

# Text Analytics Toolbox™

User's Guide



# MATLAB®

R2022a



# How to Contact MathWorks



Latest news: [www.mathworks.com](http://www.mathworks.com)  
Sales and services: [www.mathworks.com/sales\\_and\\_services](http://www.mathworks.com/sales_and_services)  
User community: [www.mathworks.com/matlabcentral](http://www.mathworks.com/matlabcentral)  
Technical support: [www.mathworks.com/support/contact\\_us](http://www.mathworks.com/support/contact_us)



Phone: 508-647-7000



The MathWorks, Inc.  
1 Apple Hill Drive  
Natick, MA 01760-2098

*Text Analytics Toolbox™ User's Guide*

© COPYRIGHT 2017–2022 by The MathWorks, Inc.

The software described in this document is furnished under a license agreement. The software may be used or copied only under the terms of the license agreement. No part of this manual may be photocopied or reproduced in any form without prior written consent from The MathWorks, Inc.

FEDERAL ACQUISITION: This provision applies to all acquisitions of the Program and Documentation by, for, or through the federal government of the United States. By accepting delivery of the Program or Documentation, the government hereby agrees that this software or documentation qualifies as commercial computer software or commercial computer software documentation as such terms are used or defined in FAR 12.212, DFARS Part 227.72, and DFARS 252.227-7014. Accordingly, the terms and conditions of this Agreement and only those rights specified in this Agreement, shall pertain to and govern the use, modification, reproduction, release, performance, display, and disclosure of the Program and Documentation by the federal government (or other entity acquiring for or through the federal government) and shall supersede any conflicting contractual terms or conditions. If this License fails to meet the government's needs or is inconsistent in any respect with federal procurement law, the government agrees to return the Program and Documentation, unused, to The MathWorks, Inc.

## Trademarks

MATLAB and Simulink are registered trademarks of The MathWorks, Inc. See [www.mathworks.com/trademarks](http://www.mathworks.com/trademarks) for a list of additional trademarks. Other product or brand names may be trademarks or registered trademarks of their respective holders.

## Patents

MathWorks products are protected by one or more U.S. patents. Please see [www.mathworks.com/patents](http://www.mathworks.com/patents) for more information.

## Revision History

March 2018	Online Only	New for Version 1.1 (Release 2018a)
September 2018	Online Only	Revised for Version 1.2 (Release 2018b)
March 2019	Online Only	Revised for Version 1.3 (Release 2019a)
September 2019	Online Only	Revised for Version 1.4 (Release 2019b)
March 2020	Online Only	Revised for Version 1.5 (Release 2020a)
September 2020	Online Only	Revised for Version 1.6 (Release 2020b)
March 2021	Online Only	Revised for Version 1.7 (Release 2021a)
September 2021	Online Only	Revised for Version 1.8 (Release 2021b)
March 2022	Online Only	Revised for Version 1.8.1 (Release 2022a)

## Text Data Preparation

### 1

<b>Extract Text Data from Files</b> .....	<b>1-2</b>
<b>Prepare Text Data for Analysis</b> .....	<b>1-11</b>
<b>Parse HTML and Extract Text Content</b> .....	<b>1-18</b>
<b>Correct Spelling in Documents</b> .....	<b>1-22</b>
<b>Create Extension Dictionary for Spelling Correction</b> .....	<b>1-24</b>
<b>Create Custom Spelling Correction Function Using Edit Distance Searchers</b> .....	<b>1-28</b>
<b>Data Sets for Text Analytics</b> .....	<b>1-34</b>

## Modeling and Prediction

### 2

<b>Create Simple Text Model for Classification</b> .....	<b>2-2</b>
<b>Analyze Text Data Using Multiword Phrases</b> .....	<b>2-7</b>
<b>Analyze Text Data Using Topic Models</b> .....	<b>2-13</b>
<b>Choose Number of Topics for LDA Model</b> .....	<b>2-19</b>
<b>Compare LDA Solvers</b> .....	<b>2-23</b>
<b>Visualize LDA Topics Using Word Clouds</b> .....	<b>2-28</b>
<b>Visualize LDA Topic Probabilities of Documents</b> .....	<b>2-30</b>
<b>Visualize Document Clusters Using LDA Model</b> .....	<b>2-35</b>
<b>Visualize LDA Topic Correlations</b> .....	<b>2-38</b>
<b>Visualize Correlations Between LDA Topics and Document Labels</b> .....	<b>2-42</b>
<b>Create Co-occurrence Network</b> .....	<b>2-48</b>

<b>Analyze Text Data Containing Emojis</b> .....	<b>2-52</b>
<b>Analyze Sentiment in Text</b> .....	<b>2-58</b>
<b>Generate Domain Specific Sentiment Lexicon</b> .....	<b>2-61</b>
<b>Train a Sentiment Classifier</b> .....	<b>2-71</b>
<b>Extract Keywords from Text Data Using RAKE</b> .....	<b>2-79</b>
<b>Extract Keywords from Text Data Using TextRank</b> .....	<b>2-82</b>
<b>Classify Text Data Using Deep Learning</b> .....	<b>2-85</b>
<b>Classify Text Data Using Convolutional Neural Network</b> .....	<b>2-93</b>
<b>Classify Text Data Using Custom Training Loop</b> .....	<b>2-100</b>
<b>Multilabel Text Classification Using Deep Learning</b> .....	<b>2-111</b>
<b>Sequence-to-Sequence Translation Using Attention</b> .....	<b>2-130</b>
<b>Language Translation Using Deep Learning</b> .....	<b>2-146</b>
<b>Classify Out-of-Memory Text Data Using Deep Learning</b> .....	<b>2-168</b>
<b>Pride and Prejudice and MATLAB</b> .....	<b>2-174</b>
<b>Word-By-Word Text Generation Using Deep Learning</b> .....	<b>2-180</b>
<b>Generate Text Using Autoencoders</b> .....	<b>2-186</b>
<b>Define Text Encoder Model Function</b> .....	<b>2-198</b>
<b>Define Text Decoder Model Function</b> .....	<b>2-205</b>
<b>Classify Out-of-Memory Text Data Using Custom Mini-Batch Datastore</b> .....	<b>2-212</b>

## **Display and Presentation**

### **3**

<b>Visualize Text Data Using Word Clouds</b> .....	<b>3-2</b>
<b>Visualize Word Embeddings Using Text Scatter Plots</b> .....	<b>3-8</b>

<b>Language Considerations</b> .....	<b>4-2</b>
Language-Independent Features .....	<b>4-4</b>
<b>Japanese Language Support</b> .....	<b>4-6</b>
Tokenization .....	<b>4-6</b>
Part of Speech Details .....	<b>4-6</b>
Named Entity Recognition .....	<b>4-7</b>
Stop Words .....	<b>4-8</b>
Lemmatization .....	<b>4-9</b>
Language-Independent Features .....	<b>4-9</b>
<b>Analyze Japanese Text Data</b> .....	<b>4-11</b>
<b>German Language Support</b> .....	<b>4-21</b>
Tokenization .....	<b>4-21</b>
Sentence Detection .....	<b>4-21</b>
Part of Speech Details .....	<b>4-22</b>
Named Entity Recognition .....	<b>4-23</b>
Stop Words .....	<b>4-24</b>
Stemming .....	<b>4-24</b>
Language-Independent Features .....	<b>4-25</b>
<b>Analyze German Text Data</b> .....	<b>4-26</b>
<b>Korean Language Support</b> .....	<b>4-37</b>
Tokenization .....	<b>4-37</b>
Part of Speech Details .....	<b>4-37</b>
Named Entity Recognition .....	<b>4-37</b>
Stop Words .....	<b>4-37</b>
Lemmatization .....	<b>4-37</b>
Language-Independent Features .....	<b>4-37</b>
<b>Language-Independent Features</b> .....	<b>4-39</b>
Word and N-Gram Counting .....	<b>4-39</b>
Modeling and Prediction .....	<b>4-39</b>

<b>Text Analytics Glossary</b> .....	<b>5-2</b>
Documents and Tokens .....	<b>5-2</b>
Preprocessing .....	<b>5-3</b>
Modeling and Prediction .....	<b>5-3</b>
Visualization .....	<b>5-5</b>



# Text Data Preparation

---

- “Extract Text Data from Files” on page 1-2
- “Prepare Text Data for Analysis” on page 1-11
- “Parse HTML and Extract Text Content” on page 1-18
- “Correct Spelling in Documents” on page 1-22
- “Create Extension Dictionary for Spelling Correction” on page 1-24
- “Create Custom Spelling Correction Function Using Edit Distance Searchers” on page 1-28
- “Data Sets for Text Analytics” on page 1-34

## Extract Text Data from Files

This example shows how to extract the text data from text, HTML, Microsoft® Word, PDF, CSV, and Microsoft Excel® files and import it into MATLAB® for analysis.

Usually, the easiest way to import text data into MATLAB is to use the `extractFileText` function. This function extracts the text data from text, PDF, HTML, and Microsoft Word files. To import text from CSV and Microsoft Excel files, use `readtable`. To extract text from HTML code, use `extractHTMLText`. To read data from PDF forms, use `readPDFFormData`.

### Text File

Extract the text from `sonnets.txt` using `extractFileText`. The file `sonnets.txt` contains Shakespeare's sonnets in plain text.

```
filename = "sonnets.txt";  
str = extractFileText(filename);
```

View the first sonnet by extracting the text between the two titles "I" and "II".

```
start = " I" + newline;  
fin = " II";  
sonnet1 = extractBetween(str,start,fin)  
  
sonnet1 =  
"  
    From fairest creatures we desire increase,  
    That thereby beauty's rose might never die,  
    But as the ripper should by time decease,  
    His tender heir might bear his memory:  
    But thou, contracted to thine own bright eyes,  
    Feed'st thy light's flame with self-substantial fuel,  
    Making a famine where abundance lies,  
    Thy self thy foe, to thy sweet self too cruel:  
    Thou that art now the world's fresh ornament,  
    And only herald to the gaudy spring,  
    Within thine own bud buriest thy content,  
    And tender churl mak'st waste in niggarding:  
    Pity the world, or else this glutton be,  
    To eat the world's due, by the grave and thee.  
"
```

For text files containing multiple documents separated by newline characters, use the `readlines` function.

```
filename = "multilineSonnets.txt";  
str = readlines(filename)  
  
str = 3x1 string  
"From fairest creatures we desire increase, That thereby beauty's rose might never die, But a  
"When forty winters shall besiege thy brow, And dig deep trenches in thy beauty's field, Thy  
"Look in thy glass and tell the face thou viewest Now is the time that face should form anotl
```



## Microsoft Word Document

Extract the text from `sonnets.docx` using `extractFileText`. The file `exampleSonnets.docx` contains Shakespeare's sonnets in a Microsoft Word document.

```
filename = "exampleSonnets.docx";
str = extractFileText(filename);
```

View the second sonnet by extracting the text between the two titles "II" and "III".

```
start = " II" + newline;
fin = " III";
sonnet2 = extractBetween(str,start,fin)
```

```
sonnet2 =
"
```

```
    When forty winters shall besiege thy brow,
    And dig deep trenches in thy beauty's field,
    Thy youth's proud livery so gazed on now,
    Will be a tatter'd weed of small worth held:
    Then being asked, where all thy beauty lies,
    Where all the treasure of thy lusty days;
    To say, within thine own deep sunken eyes,
    Were an all-eating shame, and thriftless praise.
    How much more praise deserv'd thy beauty's use,
    If thou couldst answer 'This fair child of mine
    Shall sum my count, and make my old excuse,'
    Proving his beauty by succession thine!

    This were to be new made when thou art old,
    And see thy blood warm when thou feel'st it cold.
```

```
"
```

The example Microsoft Word document uses two newline characters between each line. To replace these characters with a single newline character, use the `replace` function.

```
sonnet2 = replace(sonnet2,[newline newline],newline)
```

```
sonnet2 =
"
```

```
When forty winters shall besiege thy brow,  
And dig deep trenches in thy beauty's field,  
Thy youth's proud livery so gazed on now,  
Will be a tatter'd weed of small worth held:  
Then being asked, where all thy beauty lies,  
Where all the treasure of thy lusty days;  
To say, within thine own deep sunken eyes,  
Were an all-eating shame, and thriftless praise.  
How much more praise deserv'd thy beauty's use,  
If thou couldst answer 'This fair child of mine  
Shall sum my count, and make my old excuse,'  
Proving his beauty by succession thine!  
    This were to be new made when thou art old,  
    And see thy blood warm when thou feel'st it cold.
```

"

## PDF Files

Extract text from PDF documents and data from PDF forms.

### PDF Document

Extract the text from `sonnets.pdf` using `extractFileText`. The file `exampleSonnets.pdf` contains Shakespeare's sonnets in a PDF.

```
filename = "exampleSonnets.pdf";  
str = extractFileText(filename);
```

View the third sonnet by extracting the text between the two titles "III" and "IV". This PDF has a space before each newline character.

```
start = " III " + newline;  
fin = "IV";  
sonnet3 = extractBetween(str,start,fin)
```

```
sonnet3 =
```

"

```
    Look in thy glass and tell the face thou viewest  
    Now is the time that face should form another;  
    Whose fresh repair if now thou not renewest,  
    Thou dost beguile the world, unbless some mother.  
    For where is she so fair whose unear'd womb  
    Disdains the tillage of thy husbandry?  
    Or who is he so fond will be the tomb,  
    Of his self-love to stop posterity?  
    Thou art thy mother's glass and she in thee  
    Calls back the lovely April of her prime;  
    So thou through windows of thine age shalt see,  
    Despite of wrinkles this thy golden time.  
    But if thou live, remember'd not to be,  
    Die single and thine image dies with thee.
```

"

## PDF Form

To read text data from PDF forms, use `readPDFFormData`. The function returns a struct containing the data from the PDF form fields.

```
filename = "weatherReportForm1.pdf";
data = readPDFFormData(filename)

data = struct with fields:
    event_type: "Thunderstorm Wind"
    event_narrative: "Large tree down between Plantersville and Nettleton."
```

## HTML

Extract text from HTML files, HTML code, and the web.

### HTML File

To extract text data from a saved HTML file, use `extractFileText`.

```
filename = "exampleSonnets.html";
str = extractFileText(filename);
```

View the fourth sonnet by extracting the text between the two titles "IV" and "V".

```
start = newline + "IV" + newline;
fin = newline + "V" + newline;
sonnet4 = extractBetween(str,start,fin)

sonnet4 =
"
    Unthrifty loveliness, why dost thou spend
    Upon thy self thy beauty's legacy?
    Nature's bequest gives nothing, but doth lend,
    And being frank she lends to those are free:
    Then, beauteous niggard, why dost thou abuse
    The bounteous largess given thee to give?
    Profitless usurer, why dost thou use
    So great a sum of sums, yet canst not live?
    For having traffic with thy self alone,
    Thou of thy self thy sweet self dost deceive:
    Then how when nature calls thee to be gone,
    What acceptable audit canst thou leave?
    Thy unused beauty must be tombed with thee,
    Which, used, lives th' executor to be.
"
```

### HTML Code

To extract text data from a string containing HTML code, use `extractHTMLText`.

```
code = "<html><body><h1>THE SONNETS</h1><p>by William Shakespeare</p></body></html>";
str = extractHTMLText(code)

str =
    "THE SONNETS"
```

by William Shakespeare"

### **From the Web**

To extract text data from a web page, first read the HTML code using `webread`, and then use `extractHTMLText`.

```
url = "https://www.mathworks.com/help/textanalytics";  
code = webread(url);  
str = extractHTMLText(code)
```

```
str =  
'Text Analytics Toolbox  
  
Analyze and model text data  
  
Release Notes  
  
PDF Documentation  
  
Release Notes  
  
PDF Documentation  
  
Text Analytics Toolbox™ provides algorithms and visualizations for preprocessing, analyzing  
Text Analytics Toolbox includes tools for processing raw text from sources such as equipment  
Using machine learning techniques such as LSA, LDA, and word embeddings, you can find cluste  
Get Started  
  
Learn the basics of Text Analytics Toolbox  
  
Text Data Preparation  
  
Import text data into MATLAB® and preprocess it for analysis  
  
Modeling and Prediction  
  
Develop predictive models using topic models and word embeddings  
  
Display and Presentation  
  
Visualize text data and models using word clouds and text scatter plots  
  
Language Support  
  
Information on language support in Text Analytics Toolbox'
```

### **Parse HTML Code**

To find particular elements of HTML code, parse the code using `htmlTree` and use `findElement`. Parse the HTML code and find all the hyperlinks. The hyperlinks are nodes with element name "A".

```
tree = htmlTree(code);
selector = "A";
subtrees = findElement(tree,selector);
```

View the first 10 subtrees and extract the text using `extractHTMLText`.

```
subtrees(1:10)
```

```
ans =
    10×1 htmlTree:

    <A class="skip_link sr-only" href="#content_container">Skip to content</A>
    <A href="https://www.mathworks.com?s_tid=gn_logo" class="svg_link navbar-brand"><IMG src="/i
    <A href="https://www.mathworks.com/products.html?s_tid=gn_ps">Products</A>
    <A href="https://www.mathworks.com/solutions.html?s_tid=gn_sol">Solutions</A>
    <A href="https://www.mathworks.com/academia.html?s_tid=gn_acad">Academia</A>
    <A href="https://www.mathworks.com/support.html?s_tid=gn_supp">Support</A>
    <A href="https://www.mathworks.com/matlabcentral/?s_tid=gn_mlc">Community</A>
    <A href="https://www.mathworks.com/company/events.html?s_tid=gn_ev">Events</A>
    <A href="https://www.mathworks.com/products/get-matlab.html?s_tid=gn_getml">Get MATLAB</A>
    <A href="https://www.mathworks.com?s_tid=gn_logo" class="svg_link pull-left"><IMG src="/imag
```

```
str = extractHTMLText(subtrees);
```

View the extracted text of the first 10 hyperlinks.

```
str(1:10)
```

```
ans = 10×1 string
    "Skip to content"
    ""
    "Products"
    "Solutions"
    "Academia"
    "Support"
    "Community"
    "Events"
    "Get MATLAB"
    ""
```

To get the link targets, use `getAttributes` and specify the attribute `"href"` (hyperlink reference). Get the link targets of the first 10 subtrees.

```
attr = "href";
str = getAttribute(subtrees(1:10),attr)

str = 10×1 string
    "#content_container"
    "https://www.mathworks.com?s_tid=gn_logo"
    "https://www.mathworks.com/products.html?s_tid=gn_ps"
    "https://www.mathworks.com/solutions.html?s_tid=gn_sol"
    "https://www.mathworks.com/academia.html?s_tid=gn_acad"
    "https://www.mathworks.com/support.html?s_tid=gn_supp"
    "https://www.mathworks.com/matlabcentral/?s_tid=gn_mlc"
    "https://www.mathworks.com/company/events.html?s_tid=gn_ev"
    "https://www.mathworks.com/products/get-matlab.html?s_tid=gn_getml"
```

```
"https://www.mathworks.com?s_tid=gn_logo"
```

### CSV and Microsoft Excel Files

To extract text data from CSV and Microsoft Excel files, use `readtable` and extract the text data from the table that it returns.

Extract the table data from `factoryReposts.csv` using the `readtable` function and view the first few rows of the table.

```
T = readtable('factoryReposts.csv', 'TextType', 'string');
head(T)
```

```
ans=8x5 table
```

Description	Category
"Items are occasionally getting stuck in the scanner spools."	"Mechanical Failure"
"Loud rattling and banging sounds are coming from assembler pistons."	"Mechanical Failure"
"There are cuts to the power when starting the plant."	"Electronic Failure"
"Fried capacitors in the assembler."	"Electronic Failure"
"Mixer tripped the fuses."	"Electronic Failure"
"Burst pipe in the constructing agent is spraying coolant."	"Leak"
"A fuse is blown in the mixer."	"Electronic Failure"
"Things continue to tumble off of the belt."	"Mechanical Failure"

Extract the text data from the `event_narrative` column and view the first few strings.

```
str = T.Description;
str(1:10)
```

```
ans = 10x1 string
```

```
"Items are occasionally getting stuck in the scanner spools."
"Loud rattling and banging sounds are coming from assembler pistons."
"There are cuts to the power when starting the plant."
"Fried capacitors in the assembler."
"Mixer tripped the fuses."
"Burst pipe in the constructing agent is spraying coolant."
"A fuse is blown in the mixer."
"Things continue to tumble off of the belt."
"Falling items from the conveyor belt."
"The scanner reel is split, it will soon begin to curve."
```

### Extract Text from Multiple Files

If your text data is contained in multiple files in a folder, then you can import the text data into MATLAB using a file datastore.

Create a file datastore for the example sonnet text files. The example files are named "exampleSonnetN.txt", where N is the number of the sonnet. Specify the file name using the wildcard "\*" to find all file names of this structure. To specify the read function to be `extractFileText`, input this function to `fileDatastore` using a function handle.

```
location = fullfile(matlabroot, "examples", "textanalytics", "data", "exampleSonnet*.txt");
fds = fileDatastore(location, 'ReadFcn', @extractFileText)
```

```

fds =
  FileDatastore with properties:

        Files: {
            '...\matlab\examples\textanalytics\data\exampleSonnet1.txt';
            '...\matlab\examples\textanalytics\data\exampleSonnet2.txt';
            '...\matlab\examples\textanalytics\data\exampleSonnet3.txt'
            ... and 2 more
        }
        Folders: {
            '...\matlab\examples\textanalytics\data'
        }
        UniformRead: 0
        ReadMode: 'file'
        BlockSize: Inf
        PreviewFcn: @extractFileText
        SupportedOutputFormats: ["txt" "csv" "xlsx" "xls" "parquet" "parq" "png"
        ReadFcn: @extractFileText
        AlternateFileSystemRoots: {}

```

Loop over the files in the datastore and read each text file.

```

str = [];
while hasdata(fds)
    textData = read(fds);
    str = [str; textData];
end

```

View the extracted text.

```
str
```

```

str = 5x1 string
    " From fairest creatures we desire increase, That thereby beauty's rose might never die,
    " When forty winters shall besiege thy brow, And dig deep trenches in thy beauty's field,
    " Look in thy glass and tell the face thou viewest Now is the time that face should form a
    " Unthrifty loveliness, why dost thou spend Upon thy self thy beauty's legacy? Nature's
    "from fairest creatures we desire increase that thereby beautys rose might never die but as t

```

## See Also

[extractFileText](#) | [readPDFFormData](#) | [extractHTMLText](#) | [tokenizedDocument](#)

## Related Examples

- “Prepare Text Data for Analysis” on page 1-11
- “Create Simple Text Model for Classification” on page 2-2
- “Visualize Text Data Using Word Clouds” on page 3-2
- “Analyze Text Data Containing Emojis” on page 2-52
- “Analyze Text Data Using Topic Models” on page 2-13
- “Analyze Text Data Using Multiword Phrases” on page 2-7
- “Classify Text Data Using Deep Learning” on page 2-85

- “Train a Sentiment Classifier” on page 2-71



## Prepare Text Data for Analysis

This example shows how to create a function which cleans and preprocesses text data for analysis.

Text data can be large and can contain lots of noise which negatively affects statistical analysis. For example, text data can contain the following:

- Variations in case, for example "new" and "New"
- Variations in word forms, for example "walk" and "walking"
- Words which add noise, for example stop words such as "the" and "of"
- Punctuation and special characters
- HTML and XML tags

These word clouds illustrate word frequency analysis applied to some raw text data from factory reports, and a preprocessed version of the same text data.



### Load and Extract Text Data

Load the example data. The file `factoryReports.csv` contains factory reports, including a text description and categorical labels for each event.

```
filename = "factoryReports.csv";
data = readtable(filename, 'TextType', 'string');
```

Extract the text data from the field `Description`, and the label data from the field `Category`.

```
textData = data.Description;
labels = data.Category;
textData(1:10)

ans = 10x1 string
    "Items are occasionally getting stuck in the scanner spools."
    "Loud rattling and banging sounds are coming from assembler pistons."
    "There are cuts to the power when starting the plant."
    "Fried capacitors in the assembler."
    "Mixer tripped the fuses."
    "Burst pipe in the constructing agent is spraying coolant."
    "A fuse is blown in the mixer."
    "Things continue to tumble off of the belt."
    "Falling items from the conveyor belt."
    "The scanner reel is split, it will soon begin to curve."
```

### Create Tokenized Documents

Create an array of tokenized documents.

```
cleanedDocuments = tokenizedDocument(textData);
cleanedDocuments(1:10)

ans =
    10x1 tokenizedDocument:

    10 tokens: Items are occasionally getting stuck in the scanner spools .
    11 tokens: Loud rattling and banging sounds are coming from assembler pistons .
    11 tokens: There are cuts to the power when starting the plant .
     6 tokens: Fried capacitors in the assembler .
     5 tokens: Mixer tripped the fuses .
    10 tokens: Burst pipe in the constructing agent is spraying coolant .
     8 tokens: A fuse is blown in the mixer .
     9 tokens: Things continue to tumble off of the belt .
     7 tokens: Falling items from the conveyor belt .
    13 tokens: The scanner reel is split , it will soon begin to curve .
```

To improve lemmatization, add part of speech details to the documents using `addPartOfSpeechDetails`. Use the `addPartOfSpeech` function before removing stop words and lemmatizing.

```
cleanedDocuments = addPartOfSpeechDetails(cleanedDocuments);
```

Words like "a", "and", "to", and "the" (known as stop words) can add noise to data. Remove a list of stop words using the `removeStopWords` function. Use the `removeStopWords` function before using the `normalizeWords` function.

```
cleanedDocuments = removeStopWords(cleanedDocuments);
cleanedDocuments(1:10)
```

```
ans =
    10x1 tokenizedDocument:

     7 tokens: Items occasionally getting stuck scanner spools .
     8 tokens: Loud rattling banging sounds coming assembler pistons .
     5 tokens: cuts power starting plant .
     4 tokens: Fried capacitors assembler .
```

```

4 tokens: Mixer tripped fuses .
7 tokens: Burst pipe constructing agent spraying coolant .
4 tokens: fuse blown mixer .
6 tokens: Things continue tumble off belt .
5 tokens: Falling items conveyor belt .
8 tokens: scanner reel split , soon begin curve .

```

Lemmatize the words using `normalizeWords`.

```

cleanedDocuments = normalizeWords(cleanedDocuments, 'Style', 'lemma');
cleanedDocuments(1:10)

```

```

ans =
  10×1 tokenizedDocument:

    7 tokens: items occasionally get stuck scanner spool .
    8 tokens: loud rattle bang sound come assembler piston .
    5 tokens: cut power start plant .
    4 tokens: fry capacitor assembler .
    4 tokens: mixer trip fuse .
    7 tokens: burst pipe constructing agent spray coolant .
    4 tokens: fuse blow mixer .
    6 tokens: thing continue tumble off belt .
    5 tokens: fall item conveyor belt .
    8 tokens: scanner reel split , soon begin curve .

```

Erase the punctuation from the documents.

```

cleanedDocuments = erasePunctuation(cleanedDocuments);
cleanedDocuments(1:10)

```

```

ans =
  10×1 tokenizedDocument:

    6 tokens: items occasionally get stuck scanner spool
    7 tokens: loud rattle bang sound come assembler piston
    4 tokens: cut power start plant
    3 tokens: fry capacitor assembler
    3 tokens: mixer trip fuse
    6 tokens: burst pipe constructing agent spray coolant
    3 tokens: fuse blow mixer
    5 tokens: thing continue tumble off belt
    4 tokens: fall item conveyor belt
    6 tokens: scanner reel split soon begin curve

```

Remove words with 2 or fewer characters, and words with 15 or greater characters.

```

cleanedDocuments = removeShortWords(cleanedDocuments, 2);
cleanedDocuments = removeLongWords(cleanedDocuments, 15);
cleanedDocuments(1:10)

```

```

ans =
  10×1 tokenizedDocument:

    6 tokens: items occasionally get stuck scanner spool
    7 tokens: loud rattle bang sound come assembler piston

```

```
4 tokens: cut power start plant
3 tokens: fry capacitor assembler
3 tokens: mixer trip fuse
6 tokens: burst pipe constructing agent spray coolant
3 tokens: fuse blow mixer
5 tokens: thing continue tumble off belt
4 tokens: fall item conveyor belt
6 tokens: scanner reel split soon begin curve
```

## Create Bag-of-Words Model

Create a bag-of-words model.

```
cleanedBag = bagOfWords(cleanedDocuments)
```

```
cleanedBag =
  bagOfWords with properties:
    Counts: [480×352 double]
    Vocabulary: [1×352 string]
    NumWords: 352
    NumDocuments: 480
```

Remove words that do not appear more than two times in the bag-of-words model.

```
cleanedBag = removeInfrequentWords(cleanedBag,2)
```

```
cleanedBag =
  bagOfWords with properties:
    Counts: [480×163 double]
    Vocabulary: [1×163 string]
    NumWords: 163
    NumDocuments: 480
```

Some preprocessing steps such as `removeInfrequentWords` leaves empty documents in the bag-of-words model. To ensure that no empty documents remain in the bag-of-words model after preprocessing, use `removeEmptyDocuments` as the last step.

Remove empty documents from the bag-of-words model and the corresponding labels from `labels`.

```
[cleanedBag,idx] = removeEmptyDocuments(cleanedBag);
labels(idx) = [];
cleanedBag
```

```
cleanedBag =
  bagOfWords with properties:
    Counts: [480×163 double]
    Vocabulary: [1×163 string]
    NumWords: 163
    NumDocuments: 480
```

## Create a Preprocessing Function

It can be useful to create a function which performs preprocessing so you can prepare different collections of text data in the same way. For example, you can use a function so that you can preprocess new data using the same steps as the training data.

Create a function which tokenizes and preprocesses the text data so it can be used for analysis. The function `preprocessText`, performs the following steps:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Remove a list of stop words (such as "and", "of", and "the") using `removeStopWords`.
- 3 Lemmatize the words using `normalizeWords`.
- 4 Erase punctuation using `erasePunctuation`.
- 5 Remove words with 2 or fewer characters using `removeShortWords`.
- 6 Remove words with 15 or more characters using `removeLongWords`.

Use the example preprocessing function `preprocessText` to prepare the text data.

```
newText = "The sorting machine is making lots of loud noises.";
newDocuments = preprocessText(newText)
```

```
newDocuments =
  tokenizedDocument:

    6 tokens: sorting machine make lot loud noise
```

## Compare with Raw Data

Compare the preprocessed data with the raw data.

```
rawDocuments = tokenizedDocument(textData);
rawBag = bagOfWords(rawDocuments)
```

```
rawBag =
  bagOfWords with properties:

    Counts: [480x555 double]
    Vocabulary: [1x555 string]
    NumWords: 555
    NumDocuments: 480
```

Calculate the reduction in data.

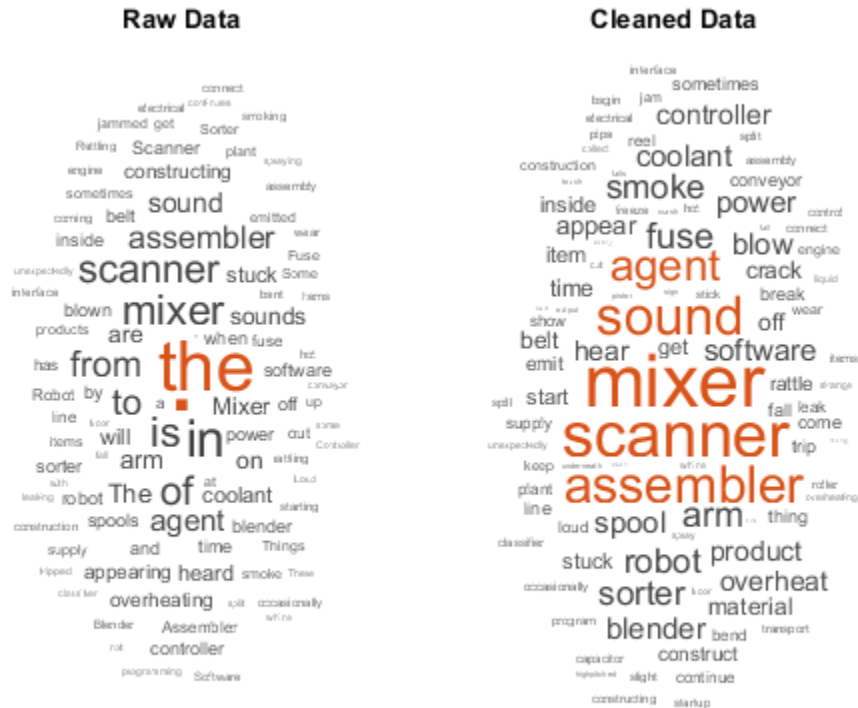
```
numWordsCleaned = cleanedBag.NumWords;
numWordsRaw = rawBag.NumWords;
reduction = 1 - numWordsCleaned/numWordsRaw

reduction = 0.7063
```

Compare the raw data and the cleaned data by visualizing the two bag-of-words models using word clouds.

```
figure
subplot(1,2,1)
```

```
wordcloud(rawBag);
title("Raw Data")
subplot(1,2,2)
wordcloud(cleanedBag);
title("Cleaned Data")
```



### Preprocessing Function

The function preprocessText, performs the following steps in order:

- 1 Tokenize the text using tokenizedDocument.
- 2 Remove a list of stop words (such as "and", "of", and "the") using removeStopWords.
- 3 Lemmatize the words using normalizeWords.
- 4 Erase punctuation using erasePunctuation.
- 5 Remove words with 2 or fewer characters using removeShortWords.
- 6 Remove words with 15 or more characters using removeLongWords.

```
function documents = preprocessText(textData)

% Tokenize the text.
documents = tokenizedDocument(textData);

% Remove a list of stop words then lemmatize the words. To improve
% lemmatization, first use addPartOfSpeechDetails.
documents = addPartOfSpeechDetails(documents);
```

```
documents = removeStopWords(documents);
documents = normalizeWords(documents, 'Style', 'lemma');

% Erase punctuation.
documents = erasePunctuation(documents);

% Remove words with 2 or fewer characters, and words with 15 or more
% characters.
documents = removeShortWords(documents,2);
documents = removeLongWords(documents,15);

end
```

## See Also

[tokenizedDocument](#) | [bagOfWords](#) | [removeStopWords](#) | [removeLongWords](#) | [removeShortWords](#) | [erasePunctuation](#) | [removeEmptyDocuments](#) | [removeInfrequentWords](#) | [normalizeWords](#) | [wordcloud](#) | [addPartOfSpeechDetails](#)

## Related Examples

- “Extract Text Data from Files” on page 1-2
- “Create Simple Text Model for Classification” on page 2-2
- “Visualize Text Data Using Word Clouds” on page 3-2
- “Analyze Text Data Containing Emojis” on page 2-52
- “Analyze Text Data Using Topic Models” on page 2-13
- “Analyze Text Data Using Multiword Phrases” on page 2-7
- “Classify Text Data Using Deep Learning” on page 2-85
- “Train a Sentiment Classifier” on page 2-71

## Parse HTML and Extract Text Content

This example shows how to parse HTML code and extract the text content from particular elements.

### Parse HTML Code

Read HTML code from the URL `https://www.mathworks.com/help/textanalytics` using `webread`.

```
url = "https://www.mathworks.com/help/textanalytics";  
code = webread(url);
```

Parse the HTML code using `htmlTree`.

```
tree = htmlTree(code);
```

View the HTML element name of the tree.

```
tree.Name
```

```
ans =  
"HTML"
```

View the child elements of the tree. The children are subtrees of `tree`.

```
tree.Children
```

```
ans =  
4x1 htmlTree:  
  
" "  
  <HEAD><TITLE>Text Analytics Toolbox Documentation</TITLE><META charset="utf-8"/><META content  
" "  
  <BODY id="responsive_offcanvas"><!-- Mobile TopNav: Start --><DIV class="header visible-xs v
```

### Extract Text from HTML Tree

To extract text directly from the HTML tree, use `extractHTMLText`.

```
str = extractHTMLText(tree)
```

```
str =  
"Text Analytics Toolbox™ provides algorithms and visualizations for preprocessing, analyzing  
Text Analytics Toolbox includes tools for processing raw text from sources such as equipment  
Using machine learning techniques such as LSA, LDA, and word embeddings, you can find cluste
```

### Find HTML Elements

To find particular elements of an HTML tree, use `findElement`. Find all the hyperlinks in the HTML tree. In HTML, hyperlinks use the "A" tag.

```
selector = "A";  
subtrees = findElement(tree,selector);
```

View the first few subtrees.



```
subtrees(1:20)
```

```
ans =
```

```
20x1 htmlTree:
```

```
<A class="svg_link navbar-brand" href="https://www.mathworks.com?s_tid=gn_logo"><IMG alt="Ma
<A class="mwa-nav_login" href="https://www.mathworks.com/login?uri=http://www.mathworks.com/
<A href="https://www.mathworks.com/products.html?s_tid=gn_ps">Products</A>
<A href="https://www.mathworks.com/solutions.html?s_tid=gn_sol">Solutions</A>
<A href="https://www.mathworks.com/academia.html?s_tid=gn_acad">Academia</A>
<A href="https://www.mathworks.com/support.html?s_tid=gn_supp">Support</A>
<A href="https://www.mathworks.com/matlabcentral/?s_tid=gn_mlc">Community</A>
<A href="https://www.mathworks.com/company/events.html?s_tid=gn_ev">Events</A>
<A href="https://www.mathworks.com/company/aboutus/contact_us.html?s_tid=gn_cntus">Contact U
<A href="https://www.mathworks.com/store?s_cid=store_top_nav&s_tid=gn_store">How to Buy<
<A href="https://www.mathworks.com/company/aboutus/contact_us.html?s_tid=gn_cntus">Contact U
<A href="https://www.mathworks.com/store?s_cid=store_top_nav&s_tid=gn_store">How to Buy<
<A class="mwa-nav_login" href="https://www.mathworks.com/login?uri=http://www.mathworks.com/
<A class="svg_link pull-left" href="https://www.mathworks.com?s_tid=gn_logo"><IMG alt="MathW
<A href="https://www.mathworks.com/products.html?s_tid=gn_ps">Products</A>
<A href="https://www.mathworks.com/solutions.html?s_tid=gn_sol">Solutions</A>
<A href="https://www.mathworks.com/academia.html?s_tid=gn_acad">Academia</A>
<A href="https://www.mathworks.com/support.html?s_tid=gn_supp">Support</A>
<A href="https://www.mathworks.com/matlabcentral/?s_tid=gn_mlc">Community</A>
<A href="https://www.mathworks.com/company/events.html?s_tid=gn_ev">Events</A>
```

Create a word cloud from the text of the hyperlinks.

```
str = extractHTMLText(subtrees);
figure
wordcloud(str);
title("Hyperlinks")
```



### Get HTML Attributes

Get the class attributes from the paragraph elements in the HTML tree.

```
subtrees = findElement(tree,'p');
attr = "class";
str = getAttribute(subtrees,attr)
```

```
str = 21x1 string array
<missing>
<missing>
"add_margin_5"
<missing>
<missing>
<missing>
<missing>
<missing>
"category_desc"
"category_desc"
"category_desc"
"category_desc"
<missing>
<missing>
<missing>
"text-center"
<missing>
<missing>
<missing>
```

```
"copyright"
<missing>
```

Create a word cloud from the text contained in paragraph elements with class "category\_desc".

```
subtrees = findElement(tree, 'p.category_desc');
str = extractHTMLText(subtrees);
figure
wordcloud(str);
```



## See Also

[htmlTree](#) | [findElement](#) | [getAttribute](#) | [extractHTMLText](#) | [tokenizedDocument](#)

## Related Examples

- "Prepare Text Data for Analysis" on page 1-11
- "Create Simple Text Model for Classification" on page 2-2
- "Visualize Text Data Using Word Clouds" on page 3-2
- "Analyze Text Data Using Topic Models" on page 2-13
- "Analyze Text Data Using Multiword Phrases" on page 2-7
- "Classify Text Data Using Deep Learning" on page 2-85
- "Train a Sentiment Classifier" on page 2-71

## Correct Spelling in Documents

This example shows how to correct spelling in documents using Hunspell.

### Load Text Data

Create an array of tokenized documents.

```
str = [  
    "Use MATLAB to correct spelling of words."  
    "Correctly spelled worrds are important for lemmatization."  
    "Text Analytics Toolbox providesfunctions for spelling correction."];  
documents = tokenizedDocument(str)  
  
documents =  
    3x1 tokenizedDocument:  
  
    8 tokens: Use MATLAB to correct spelling of words .  
    8 tokens: Correctly spelled worrds are important for lemmatization .  
    8 tokens: Text Analytics Toolbox providesfunctions for spelling correction .
```

### Correct Spelling

Correct the spelling of the documents using the `correctSpelling` function.

```
updatedDocuments = correctSpelling(documents)  
  
updatedDocuments =  
    3x1 tokenizedDocument:  
  
    9 tokens: Use MAT LAB to correct spelling of words .  
    8 tokens: Correctly spelled words are important for solemnization .  
    9 tokens: Text Analytic Toolbox provides functions for spelling correction .
```

Notice that:

- The input word "MATLAB" has been split into the two words "MAT" and "LAB".
- The input word "worrds" has been changed to "words".
- The input word "lemmatization" has been changed to "solemnization".
- The input word "Analytics" has been changed to "Analytic".
- The input word "providesfunctions" has been split into the two words "provides" and "functions".

### Specify Custom Words

To prevent the software from updating particular words, you can provide a list of known words using the 'KnownWords' option of the `correctSpelling` function.

Correct the spelling of the documents again and specify the words "MATLAB", "Analytics", and "lemmatization" as known words.

```
updatedDocuments = correctSpelling(documents, 'KnownWords', ["MATLAB" "Analytics" "lemmatization"])  
  
updatedDocuments =  
    3x1 tokenizedDocument:
```

8 tokens: Use MATLAB to correct spelling of words .  
8 tokens: Correctly spelled words are important for lemmatization .  
9 tokens: Text Analytics Toolbox provides functions for spelling correction .

Notice here that the words "MATLAB", "Analytics", and "lemmatization" remain unchanged.

## See Also

`correctSpelling` | `tokenizedDocument`

## More About

- “Create Extension Dictionary for Spelling Correction” on page 1-24
- “Create Custom Spelling Correction Function Using Edit Distance Searchers” on page 1-28
- “Prepare Text Data for Analysis” on page 1-11
- “Create Simple Text Model for Classification” on page 2-2
- “Analyze Text Data Using Topic Models” on page 2-13

## Create Extension Dictionary for Spelling Correction

This example shows how to create a Hunspell extension dictionary for spelling correction.

When using the `correctSpelling` function, the function may update some correctly spelled words. To provide a list of known words, you can use the "KnownWords" option directly with a string array of known words. Alternatively, you can specify a Hunspell extension dictionary (also known as a *personal dictionary*) that specifies lists of known words, forbidden words, and words alongside affix rules.

### Specify Known Words

Create an array of tokenized documents.

```
str = [
    "Use MATLAB to correct spelling of words."
    "Correctly spelled worrds are important for lemmatizing."
    "Text Analytics Toolbox providesfunctions for spelling correction."];
documents = tokenizedDocument(str);
```

Correct the spelling of the documents using the `correctSpelling` function.

```
updatedDocuments = correctSpelling(documents)

updatedDocuments =
    3x1 tokenizedDocument:

    9 tokens: Use MAT LAB to correct spelling of words .
    8 tokens: Correctly spelled words are important for legitimatizing .
    9 tokens: Text Analytic Toolbox provides functions for spelling correction .
```

The function has corrected the spelling of the words "worrds" and "providesfunctions", though it has also updated some correctly spelled words:

- The input word "MATLAB" has been split into the two words "MAT" and "LAB".
- The input word "lemmatizing" has been changed to "legitimizing".
- The input word "Analytics" has been changed to "Analytic".

To create a Hunspell extension dictionary containing a list of known words, create a `.dic` file containing these words with one word per line. Create an extension dictionary with name `knownWords.dic` file containing the words "MATLAB", "lemmatization", and "Analytics".

```
MATLAB
Analytics
lemmatizing
```

Correct the spelling of the documents again and specify the extension dictionary `knownWords.dic`.

```
updatedDocuments = correctSpelling(documents, 'ExtensionDictionary', 'knownWords.dic')

updatedDocuments =
    3x1 tokenizedDocument:

    8 tokens: Use MATLAB to correct spelling of words .
    8 tokens: Correctly spelled words are important for lemmatizing .
```

9 tokens: Text Analytics Toolbox provides functions for spelling correction .

### Specify Affix Rules

When specifying multiple words with the same root word (for example, specifying the words "lemmatize", "lemmatizer", "lemmatized", and so on), it can be easier to indicate a set of affix rules. Instead of specifying the same word multiple times with different affixes, you can specify particular word to inherit a set of affix rules from.

For example, create an array of tokenized documents and use the `correctSpelling` function.

```
str = [
    "A lemmatizer reduces words to their dictionary forms."
    "To lemmatize words, use the normalizeWords function."
    "Before lemmatizing, add part of speech details to the text."
    "Display lemmatized words in a word cloud."];
documents = tokenizedDocument(str);
updatedDocuments = correctSpelling(documents)

updatedDocuments =
    4x1 tokenizedDocument:

    9 tokens: A legitimatize reduces words to their dictionary forms .
    10 tokens: To legitimatize words , use the normalize Words function .
    12 tokens: Before legitimatizing , add part of speech details to the text .
    8 tokens: Display legitimatized words in a word cloud .
```

Notice that the word "normalizeWords" and variants of "lemmatize" do not get updated correctly.

Create an extension dictionary with name `knownWordsWithAffixes.dic` file containing the words "normalizeWords" and "lemmatize". For the word "lemmatize", also specify to also include valid affixes of the word "equalize" using the "/" symbol.

```
normalizeWords
lemmatize/equalize
```

Correct the spelling of the documents again and specify the extension dictionary `knownWordsWithAffixes.dic`.

```
updatedDocuments = correctSpelling(documents, 'ExtensionDictionary', 'knownWordsWithAffixes.dic')

updatedDocuments =
    4x1 tokenizedDocument:

    9 tokens: A lemmatizer reduces words to their dictionary forms .
    9 tokens: To lemmatize words , use the normalizeWords function .
    12 tokens: Before lemmatizing , add part of speech details to the text .
    8 tokens: Display lemmatized words in a word cloud .
```

Notice that the variants of "lemmatize" have not been changed. The default dictionary contains the word "equalize" and also recognizes the words "equalizer" and "equalized" via the "-r" and "-d" suffixes, respectively. By specifying the entry "lemmatize/equalize", the software recognizes the word "lemmatize" as well as other words by extension of the affixes corresponding to "equalize". For example, the words "lemmatizer" and "lemmatized".

## Specify Forbidden Words

When using the `correctSpelling` function, the function may output undesirable words, even if a more desirable word is in the dictionary. For example, for the input word "Decrease", the `correctSpelling` function may output the word "Decreases". To ensure that certain words do not appear in the output, you can specify forbidden words in the extension dictionary.

For example, create an array of tokenized documents and correct the spelling using the extension dictionary `knownWords.dic`. Note that this dictionary contains the word "MATLAB".

```
str = [  
    "Analyze text data using MATLAB."  
    "Decrease the number of typos using an extension dictionary."];  
documents = tokenizedDocument(str);  
updatedDocuments = correctSpelling(documents, 'ExtensionDictionary', 'knownWords.dic')  
  
updatedDocuments =  
    2x1 tokenizedDocument:  
  
    6 tokens: Analyze text data using MATLAB .  
    10 tokens: Decreases the number of typos using an extension dictionary .
```

Even though the word "decrease" is in the dictionary, the function may still choose other words as matches. In this case, the function chooses the word "decreases".

Create an extension dictionary with name `knownWordsWithForbiddenWords.dic` file containing the word "MATLAB" and also specify the forbidden word "decree" using the "\*" symbol. When specifying forbidden words, you must specify the root word. For example, to prevent the function from outputting the plural "decreases", specify the root word "decree".

```
MATLAB  
*decree
```

Correct the spelling of the documents using the extension dictionary `knownWordsWithForbiddenWords.dic`.

```
updatedDocuments = correctSpelling(documents, 'ExtensionDictionary', 'knownWordsWithForbiddenWords  
updatedDocuments =  
    2x1 tokenizedDocument:  
  
    6 tokens: Analyze text data using MATLAB .  
    10 tokens: Decrease the number of typos using an extension dictionary .
```

Here, the word "Decrease" is corrected to "Decrease".

## See Also

`correctSpelling` | `tokenizedDocument`

## More About

- "Correct Spelling in Documents" on page 1-22
- "Create Custom Spelling Correction Function Using Edit Distance Searchers" on page 1-28



- “Prepare Text Data for Analysis” on page 1-11
- “Create Simple Text Model for Classification” on page 2-2
- “Analyze Text Data Using Topic Models” on page 2-13

## Create Custom Spelling Correction Function Using Edit Distance Searchers

This example shows how to correct spelling using edit distance searchers and a vocabulary of known words.

Lemmatization with `normalizeWords` and `word2vec` requires correctly spelled words to work. To easily correct the spelling of words in text, use the `correctSpelling` function. To learn how to create a spelling correction function from scratch using edit distance searchers, use this example as a guide.

If you have misspelled words in a collection of text, then you can use edit distance searchers to find the nearest correctly spelled words to a given vocabulary. To correct the spelling of misspelled words in documents, replace them with the nearest neighbors in the vocabulary. Use edit distance searchers to find the nearest correctly spelled word to misspelled words according to an edit distance. For example, the number of adjacent grapheme swaps and grapheme insertions, deletions, and substitutions.

### Load Data

Create a vocabulary of known words. Download and extract the Spell Checking Oriented Word Lists (SCOWL) from <https://sourceforge.net/projects/wordlist/> into a folder in the current directory. Import the words from the downloaded data using the supporting function `scowlWordList`.

```
folderName = "scowl-2019.10.06";  
maxSize = 60;  
vocabulary = scowlWordList(folderName, 'english', maxSize);
```

View the number of words in the vocabulary.

```
numWords = numel(vocabulary)  
  
numWords = 98213
```

### Create Simple Spelling Corrector

Using the imported vocabulary, create an edit distance searcher with a maximum distance of 2. For better results, allow for adjacent grapheme swaps by setting the `'SwapCost'` option to 1. For large vocabularies, this can take a few minutes.

```
maxDist = 2;  
eds = editDistanceSearcher(vocabulary, maxDist, 'SwapCost', 1);
```

This edit distance searcher is case sensitive which means that changing the case of characters contributes to the edit distance. For example, the searcher can find the neighbor "testing" for the word "tsetting" because it has edit distance 1 (one swap), but not of the word "TSeTiNG" because it has edit distance 6.

### Correct Spelling

Correct the spelling of misspelled words in an array of tokenized documents by selecting the misspelled words and finding the nearest neighbors in the edit distance searcher.

Create a tokenized document object containing typos and spelling mistakes.

```
str = "An exmaple dccoument with typos and averyunusualword.";
document = tokenizedDocument(str)

document =
    tokenizedDocument:

        8 tokens: An exmaple dccoument with typos and averyunusualword .
```

Convert the documents to a string array of words using the string function.

```
words = string(document)

words = 1×8 string
    "An"    "exmaple"    "dccoument"    "with"    "typos"    "and"    "averyunusualword"    "."
```

Find the words that need correction. To ignore words that are correctly spelled, find the indices of the words already in the vocabulary. To ignore punctuation and complex tokens such as email addresses, find the indices of the words which do not have the token types "letters" or "other". Get the token details from the document using the tokenDetails function.

```
tetails = tokenDetails(document);
idxVocabularyWords = ismember(tetails.Token,eds.Vocabulary);

idxComplexTokens = ...
    tetails.Type ~= "letters" & ...
    tetails.Type ~= "other";

idxWordsToCheck = ...
    ~idxVocabularyWords & ...
    ~idxComplexTokens

idxWordsToCheck = 8×1 logical array

    1
    1
    1
    0
    0
    0
    1
    0
```

Find the numeric indices of the words and view the corresponding words.

```
idxWordsToCheck = find(idxWordsToCheck)

idxWordsToCheck = 4×1

    1
    2
    3
    7

wordsToCheck = words(idxWordsToCheck)
```

```
wordsToCheck = 1×4 string
    "An"    "exmaple"    "dccountment"    "averyunusualword"
```

Notice that the word "An" is flagged as a word to check. This word is flagged because the vocabulary does not contain the word "An" with an uppercase "A". A later section in the example shows how to create a case insensitive spelling corrector.

Find the nearest words and their distances using the `knnsearch` function with the edit distance searcher.

```
[idxNearestWords,d] = knnsearch(eds,wordsToCheck)
```

```
idxNearestWords = 4×1
```

```
    165
   1353
   1152
    NaN
```

```
d = 4×1
```

```
    1
    1
    2
   Inf
```

If any of the words are not found in the searcher, then the function returns index `NaN` with distance `Inf`. The word "averyunusualword" does not have a match within edit distance 2, so the function returns the index `NaN` for that word.

Find the indices of the words with positive finite edit distances.

```
idxMatches = ~isnan(idxNearestWords)
```

```
idxMatches = 4×1 logical array
```

```
    1
    1
    1
    0
```

Get the indices of the words with matches in the searcher and view the corresponding corrected words in the vocabulary.

```
idxCorrectedWords = idxNearestWords(idxMatches)
```

```
idxCorrectedWords = 3×1
```

```
    165
   1353
   1152
```

```
correctedWords = eds.Vocabulary(idxCorrectedWords)
```

```
correctedWords = 1×3 string
    "an"    "example"    "document"
```

Replace the misspelled words that have matches with the corrected words.

```
idxToCorrect = idxWordsToCheck(idxMatches);
words(idxToCorrect) = correctedWords
```

```
words = 1×8 string
    "an"    "example"    "document"    "with"    "typos"    "and"    "averyunusualword"    "."
```

To create a tokenized document of these words, use the `tokenizedDocument` function and set `'TokenizedMethod'` to `'none'`.

```
document = tokenizedDocument(words, 'TokenizeMethod', 'none')
```

```
document =
    tokenizedDocument:

    8 tokens: an example document with typos and averyunusualword .
```

The next section shows how to correct the spelling of multiple documents at once by creating a custom spelling correction function and using `docfun`.

### Create Spelling Correction Function

To correct the spelling in multiple documents at once, create a custom function using the code from the previous section and use this function with the `docfun` function.

Create a function that takes an edit distance searcher, a string array of words, and the corresponding table of token details as inputs and outputs the corrected words. The `correctSpelling` function, listed at the end of the example, corrects the spelling in a string array of words using the corresponding token details and an edit distance searcher.

To use this function with the `docfun` function, create a function handle that takes a string array of words and the corresponding table of token details as the inputs.

```
func = @(words,tetails) correctSpelling(eds,words,tetails);
```

Correct the spelling of an array of tokenized documents using `docfun` with the function handle `func`.

```
str = [
    "Here is some really badly wrirten text."
    "Some moree mitsakes here too."];
documents = tokenizedDocument(str);
updatedDocuments = docfun(func,documents)

updatedDocuments =
    2×1 tokenizedDocument:

    8 tokens: here is some really badly written text .
    6 tokens: come more mistakes here too .
```

Note that uppercase characters can get corrected to different lowercase characters. For example, the word "Some" can get corrected to "come". If multiple words in the edit distance searcher vocabulary

have the same edit distance to the input, then the function outputs the first result it found. For example, the words "come" and "some" both have edit distance 1 from the word "Some".

The next section shows how to create an spelling corrector that is case insensitive.

### **Create Case Insensitive Spelling Corrector**

To prevent differences in case clashing with other substitutions, create an edit distance searcher with the vocabulary in lower case and convert the documents to lowercase before using the edit distance searcher.

Convert the vocabulary to lowercase. This operation can introduce duplicate words, remove them by taking the unique values only.

```
vocabularyLower = lower(vocabulary);  
vocabularyLower = unique(vocabularyLower);
```

Create an edit distance searcher using the lowercase vocabulary using the same options as before. This can take a few minutes to run.

```
maxDist = 2;  
eds = editDistanceSearcher(vocabularyLower,maxDist,'SwapCost',1);
```

Use the edit distance searcher to correct the spelling of the words in tokenized document. To use the case insensitive spelling corrector, convert the documents to lowercase.

```
documentsLower = lower(documents);
```

Correct the spelling using the new edit distance searcher using same steps as before.

```
func = @(words,tetails) correctSpelling(eds,words,tetails);  
updatedDocuments = docfun(func,documentsLower)
```

```
updatedDocuments =  
    2×1 tokenizedDocument:  
  
    8 tokens: here is some really badly written text .  
    6 tokens: some more mistakes here too .
```

Here, the word "Some" in the original text is converted to "some" before being input to the spelling corrector. The corresponding word "some" is unaffected by the searcher as the word some occurs in the vocabulary.

### **Spelling Correction Function**

The `correctSpelling` function corrects the spelling in a string array of words using the corresponding token details and an edit distance searcher. You can use this function with `docfun` to correct the spelling of multiple documents at once.

```
function words = correctSpelling(eds,words,tetails)  
  
% Get indices of misspelled words ignoring complex tokens.  
idxVocabularyWords = ismember(tetails.Token,eds.Vocabulary);  
  
idxComplexTokens = ...  
    tetails.Type ~= "letters" & ...  
    tetails.Type ~= "other";
```

```
idxWordsToCheck = ...
    ~idxVocabularyWords & ...
    ~idxComplexTokens;

% Convert to numeric indices.
idxWordsToCheck = find(idxWordsToCheck);

% Find nearest words.
wordsToCheck = words(idxWordsToCheck);
idxNearestWords = knnsearch(eds,wordsToCheck);

% Find words with matches.
idxMatches = ~isnan(idxNearestWords);

% Get corrected words.
idxCorrectedWords = idxNearestWords(idxMatches);
correctedWords = eds.Vocabulary(idxCorrectedWords);

% Correct words.
idxToCorrect = idxWordsToCheck(idxMatches);
words(idxToCorrect) = correctedWords;

end
```

## See Also

[correctSpelling](#) | [editDistanceSearcher](#) | [editDistance](#) | [tokenizedDocument](#) | [tokenDetails](#) | [knnsearch](#) | [docfun](#)

## More About

- “Correct Spelling in Documents” on page 1-22
- “Create Extension Dictionary for Spelling Correction” on page 1-24
- “Prepare Text Data for Analysis” on page 1-11
- “Create Simple Text Model for Classification” on page 2-2



## Data Sets for Text Analytics

This page provides a list of different data sets that you can use to get started with text analytics applications.

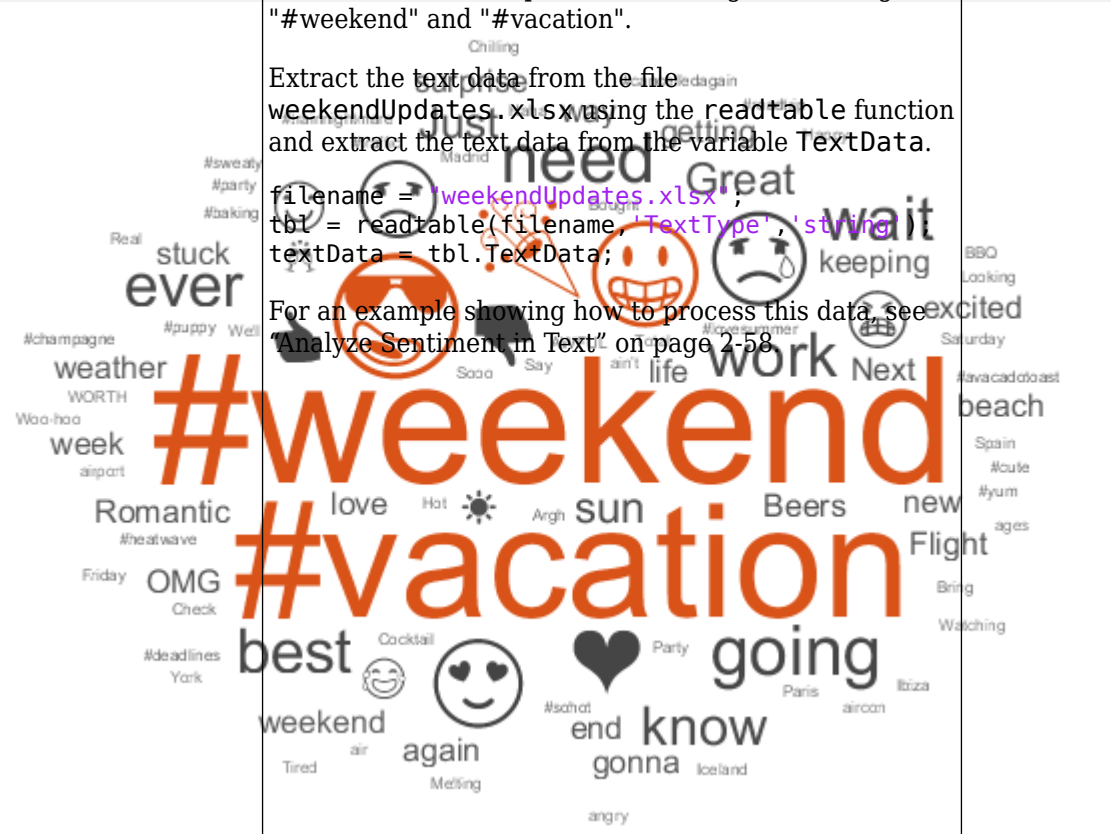
Data Set	Description	Task
<p>Factory Reports</p>	<p>The Factory Reports data set is a table containing approximately 500 reports with various attributes including a plain text description in the variable <code>Description</code> and a categorical label in the variable <code>Category</code>.</p> <p>Read the Factory Reports data from the file "factoryReports.csv". Extract the text data and the labels from the <code>Description</code> and <code>Category</code> columns, respectively.</p> <pre>filename = 'factoryReports.csv'; data = readtable(filename, 'TextType','string'); textData = data.Description; labels = data.Category;</pre> <p>For an example showing how to process this data for deep learning, see "Classify Text Data Using Deep Learning" (Deep Learning Toolbox).</p>	<p>Text classification, topic modeling</p>



Data Set	Description	Task
Shakespeare's Sonnets	<p>The file <code>sonnets.txt</code> contains all of Shakespeare's sonnets in a single text file.</p> <p>Read the Shakespeare's Sonnets data from the file <code>"sonnets.txt"</code>.</p> <pre>filename = "sonnets.txt"; textData = extractFileText(filename);</pre> <p>The sonnets are indented by two whitespace characters and are separated by two newline characters. Remove the indentations using <code>replace</code> and split the text into separate sonnets using <code>split</code>. Remove the main title from the first three elements and the sonnet titles, which appear before each sonnet.</p> <pre>textData = replace(textData, " ", ""); textData = split(textData, [newline newline]); textData = textData(5:2:end);</pre> <p>For an example showing how to process this data for deep learning, see "Generate Text Using Deep Learning" (Deep Learning Toolbox).</p>	Topic modeling, text generation

Data Set	Description	Task
	<p>The ArXiv API allows you to access the metadata of scientific e-prints submitted to <a href="https://arxiv.org">https://arxiv.org</a> including the abstract and subject areas. For more information, see <a href="https://arxiv.org/help/api">https://arxiv.org/help/api</a>.</p> <p>Import a set of abstracts and category labels from math papers using the arXiv API.</p> <pre> "https://export.arxiv.org/oai2?verb=ListRecords" + ... &amp;set=math" + &amp;metadataPrefix=arXiv"; options = weboptions('Timeout', 160); code = webread(url, options); </pre> <p>For an example showing how to parse the returned XML code and import more records, see "Multilabel Text Classification Using Deep Learning" on page 2-111.</p>	<p>Text classification, topic modeling</p> 

Data Set	Description	Task
Books from Project Gutenberg	<p>You can download many books from Project Gutenberg. For example, download the text from Alice's Adventures in Wonderland by Lewis Carroll from <a href="https://www.gutenberg.org/files/11/11-h/11-h.htm">https://www.gutenberg.org/files/11/11-h/11-h.htm</a> using the <code>webread</code> function.</p> <pre>url = "https://www.gutenberg.org/files/11/11-h/11-h.htm"; code = webread(url);</pre> <p>The HTML code contains the relevant text inside <code>&lt;p&gt;</code> (paragraph) elements. Extract the relevant text by parsing the HTML code using the <code>htmlTree</code> function and then finding all the elements with the element name <code>p</code>.</p> <pre>tree = htmlTree(code); selector = "p"; subtrees = findElement(tree, selector);</pre> <p>Extract the text data from the HTML subtrees using the <code>extractHTMLText</code> function and remove the empty elements.</p> <pre>textData = extractHTMLText(subtrees); textData(textData == "") = [];</pre> <p>For an example showing how to process this data for deep learning, see "Word-By-Word Text Generation Using Deep Learning" on page 2-180.</p>	Topic modeling, text generation

Data Set	Description	Task																					
<p>Weekend updates</p>	<p>The file <code>weekendUpdates.xlsx</code> contains example social media status updates containing the hashtags <code>#weekend</code> and <code>#vacation</code>.</p> <p>Extract the text data from the file <code>weekendUpdates.xlsx</code> using the <code>readtable</code> function and extract the text data from the variable <code>TextData</code>.</p> <pre>filename = 'weekendUpdates.xlsx'; tbl = readtable(filename, 'TextType', 'string'); textData = tbl.TextData;</pre> <p>For an example showing how to process this data, see “Analyze Sentiment in Text” on page 2-58.</p> 	<p>Sentiment analysis</p>																					
<p>Roman Numerals</p> <p>Observations</p>	<p>The CSV file <code>"romanNumerals.csv"</code> contains the decimal numbers 1-1000 in the first column and the corresponding Roman numerals in the second column.</p> <table border="1" data-bbox="435 1291 1312 1617"> <tr> <td>C</td> <td>D</td> <td>X</td> <td>X</td> <td>X</td> <td>V</td> <td>I</td> </tr> <tr> <td>C</td> <td>D</td> <td>X</td> <td>X</td> <td>X</td> <td>I</td> <td></td> </tr> <tr> <td>C</td> <td>I</td> <td>I</td> <td></td> <td></td> <td></td> <td></td> </tr> </table> <pre>filename = fullfile("romanNumerals.csv"); options = detectImportOptions(filename, ...     'TextType', 'string', ...     'ReadVariableNames', false); options.VariableNames = ["Source" "Target"]; options.VariableTypes = ["string" "string"]; data = readtable(filename, options);</pre> <p>For an example showing how to process this data for deep learning, see “Sequence-to-Sequence Translation Using Attention” on page 2-130.</p> <p>Time Steps</p>	C	D	X	X	X	V	I	C	D	X	X	X	I		C	I	I					<p>Sequence-to-sequence translation</p>
C	D	X	X	X	V	I																	
C	D	X	X	X	I																		
C	I	I																					

Data Set	Description	Task
Finance Reports	<p>The Securities and Exchange Commission (SEC) allows you to access financial reports via the Electronic Data Gathering, Analysis, and Retrieval (EDGAR) API. For more information, see <a href="https://www.sec.gov/os/accessing-edgar-data">https://www.sec.gov/os/accessing-edgar-data</a>.</p> <p>To download this data, use the function <code>financeReports</code> attached to the example “Generate Domain Specific Sentiment Lexicon” on page 2-61 as a supporting file. To access this function, open the example as a Live Script.</p> <pre> year = 2019; qtr = 4; maxLength = 2e6; textData = financeReports(year, qtr, maxLength); </pre> <p>For an example showing how to process this data, see “Generate Domain Specific Sentiment Lexicon” on page 2-61.</p>	Sentiment analysis

## See Also

### More About

- “Extract Text Data from Files” on page 1-2
- “Parse HTML and Extract Text Content” on page 1-18
- “Prepare Text Data for Analysis” on page 1-11
- “Analyze Text Data Containing Emojis” on page 2-52
- “Create Simple Text Model for Classification” on page 2-2
- “Classify Text Data Using Deep Learning” on page 2-85
- “Analyze Text Data Using Topic Models” on page 2-13
- “Analyze Sentiment in Text” on page 2-58
- “Sequence-to-Sequence Translation Using Attention” on page 2-130
- “Generate Text Using Deep Learning” (Deep Learning Toolbox)



# Modeling and Prediction

---

- “Create Simple Text Model for Classification” on page 2-2
- “Analyze Text Data Using Multiword Phrases” on page 2-7
- “Analyze Text Data Using Topic Models” on page 2-13
- “Choose Number of Topics for LDA Model” on page 2-19
- “Compare LDA Solvers” on page 2-23
- “Visualize LDA Topics Using Word Clouds” on page 2-28
- “Visualize LDA Topic Probabilities of Documents” on page 2-30
- “Visualize Document Clusters Using LDA Model” on page 2-35
- “Visualize LDA Topic Correlations” on page 2-38
- “Visualize Correlations Between LDA Topics and Document Labels” on page 2-42
- “Create Co-occurrence Network” on page 2-48
- “Analyze Text Data Containing Emojis” on page 2-52
- “Analyze Sentiment in Text” on page 2-58
- “Generate Domain Specific Sentiment Lexicon” on page 2-61
- “Train a Sentiment Classifier” on page 2-71
- “Extract Keywords from Text Data Using RAKE” on page 2-79
- “Extract Keywords from Text Data Using TextRank” on page 2-82
- “Classify Text Data Using Deep Learning” on page 2-85
- “Classify Text Data Using Convolutional Neural Network” on page 2-93
- “Classify Text Data Using Custom Training Loop” on page 2-100
- “Multilabel Text Classification Using Deep Learning” on page 2-111
- “Sequence-to-Sequence Translation Using Attention” on page 2-130
- “Language Translation Using Deep Learning” on page 2-146
- “Classify Out-of-Memory Text Data Using Deep Learning” on page 2-168
- “Pride and Prejudice and MATLAB” on page 2-174
- “Word-By-Word Text Generation Using Deep Learning” on page 2-180
- “Generate Text Using Autoencoders” on page 2-186
- “Define Text Encoder Model Function” on page 2-198
- “Define Text Decoder Model Function” on page 2-205
- “Classify Out-of-Memory Text Data Using Custom Mini-Batch Datastore” on page 2-212

## Create Simple Text Model for Classification

This example shows how to train a simple text classifier on word frequency counts using a bag-of-words model.

You can create a simple classification model which uses word frequency counts as predictors. This example trains a simple classification model to predict the category of factory reports using text descriptions.

### Load and Extract Text Data

Load the example data. The file `factoryReports.csv` contains factory reports, including a text description and categorical labels for each report.

```
filename = "factoryReports.csv";
data = readtable(filename, 'TextType', 'string');
head(data)
```

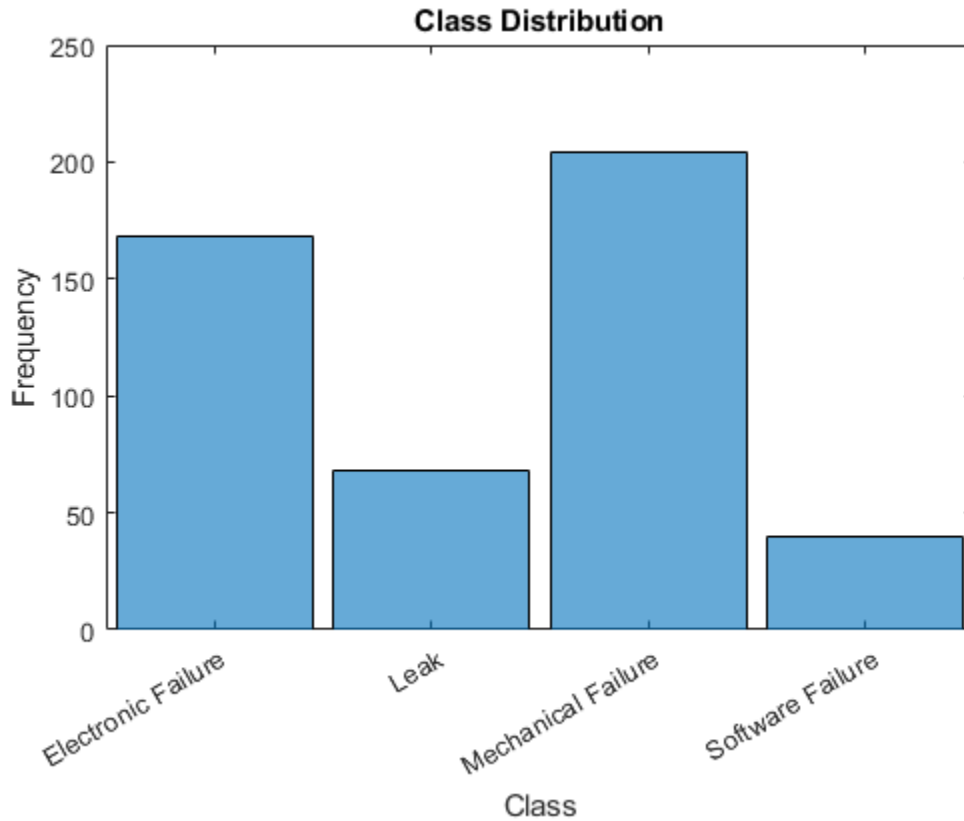
ans=8x5 table

Description	Category
"Items are occasionally getting stuck in the scanner spools."	"Mechanical Failure"
"Loud rattling and banging sounds are coming from assembler pistons."	"Mechanical Failure"
"There are cuts to the power when starting the plant."	"Electronic Failure"
"Fried capacitors in the assembler."	"Electronic Failure"
"Mixer tripped the fuses."	"Electronic Failure"
"Burst pipe in the constructing agent is spraying coolant."	"Leak"
"A fuse is blown in the mixer."	"Electronic Failure"
"Things continue to tumble off of the belt."	"Mechanical Failure"

Convert the labels in the `Category` column of the table to categorical and view the distribution of the classes in the data using a histogram.

```
data.Category = categorical(data.Category);
figure
histogram(data.Category)
xlabel("Class")
ylabel("Frequency")
title("Class Distribution")
```





Partition the data into a training partition and a held-out test set. Specify the holdout percentage to be 10%.

```
cvp = cvpartition(data.Category, 'Holdout', 0.1);
dataTrain = data(cvp.training, :);
dataTest = data(cvp.test, :);
```

Extract the text data and labels from the tables.

```
textDataTrain = dataTrain.Description;
textDataTest = dataTest.Description;
YTrain = dataTrain.Category;
YTest = dataTest.Category;
```

### Prepare Text Data for Analysis

Create a function which tokenizes and preprocesses the text data so it can be used for analysis. The function `preprocessText`, performs the following steps in order:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Remove a list of stop words (such as "and", "of", and "the") using `removeStopWords`.
- 3 Lemmatize the words using `normalizeWords`.
- 4 Erase punctuation using `erasePunctuation`.
- 5 Remove words with 2 or fewer characters using `removeShortWords`.
- 6 Remove words with 15 or more characters using `removeLongWords`.

Use the example preprocessing function `preprocessText` to prepare the text data.

```
documents = preprocessText(textDataTrain);
documents(1:5)

ans =
    5×1 tokenizedDocument:

    6 tokens: items occasionally get stuck scanner spool
    7 tokens: loud rattle bang sound come assembler piston
    4 tokens: cut power start plant
    3 tokens: fry capacitor assembler
    3 tokens: mixer trip fuse
```

Create a bag-of-words model from the tokenized documents.

```
bag = bagOfWords(documents)

bag =
    bagOfWords with properties:

        Counts: [432×336 double]
    Vocabulary: [1×336 string]
        NumWords: 336
    NumDocuments: 432
```

Remove words from the bag-of-words model that do not appear more than two times in total. Remove any documents containing no words from the bag-of-words model, and remove the corresponding entries in labels.

```
bag = removeInfrequentWords(bag,2);
[bag,idx] = removeEmptyDocuments(bag);
YTrain(idx) = [];
bag

bag =
    bagOfWords with properties:

        Counts: [432×155 double]
    Vocabulary: [1×155 string]
        NumWords: 155
    NumDocuments: 432
```

### **Train Supervised Classifier**

Train a supervised classification model using the word frequency counts from the bag-of-words model and the labels.

Train a multiclass linear classification model using `fitcecoc`. Specify the `Counts` property of the bag-of-words model to be the predictors, and the event type labels to be the response. Specify the learners to be linear. These learners support sparse data input.

```
XTrain = bag.Counts;
mdl = fitcecoc(XTrain,YTrain,'Learners','linear')

mdl =
    CompactClassificationECOC
```

```

    ResponseName: 'Y'
    ClassNames: [Electronic Failure    Leak    Mechanical Failure    Software Failure]
    ScoreTransform: 'none'
    BinaryLearners: {6×1 cell}
    CodingMatrix: [4×6 double]

```

Properties, Methods

For a better fit, you can try specifying different parameters of the linear learners. For more information on linear classification learner templates, see `templateLinear`.

### Test Classifier

Predict the labels of the test data using the trained model and calculate the classification accuracy. The classification accuracy is the proportion of the labels that the model predicts correctly.

Preprocess the test data using the same preprocessing steps as the training data. Encode the resulting test documents as a matrix of word frequency counts according to the bag-of-words model.

```

documentsTest = preprocessText(textDataTest);
XTest = encode(bag,documentsTest);

```

Predict the labels of the test data using the trained model and calculate the classification accuracy.

```

YPred = predict mdl,XTest);
acc = sum(YPred == YTest)/numel(YTest)

```

```
acc = 0.8542
```

### Predict Using New Data

Classify the event type of new factory reports. Create a string array containing the new factory reports.

```

str = [
    "Coolant is pooling underneath sorter."
    "Sorter blows fuses at start up."
    "There are some very loud rattling sounds coming from the assembler."];
documentsNew = preprocessText(str);
XNew = encode(bag,documentsNew);
labelsNew = predict mdl,XNew)

```

```

labelsNew = 3×1 categorical
    Leak
    Electronic Failure
    Mechanical Failure

```

### Example Preprocessing Function

The function `preprocessText`, performs the following steps in order:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Remove a list of stop words (such as "and", "of", and "the") using `removeStopWords`.
- 3 Lemmatize the words using `normalizeWords`.

- 4 Erase punctuation using `erasePunctuation`.
- 5 Remove words with 2 or fewer characters using `removeShortWords`.
- 6 Remove words with 15 or more characters using `removeLongWords`.

```
function documents = preprocessText(textData)

% Tokenize the text.
documents = tokenizedDocument(textData);

% Remove a list of stop words then lemmatize the words. To improve
% lemmatization, first use addPartOfSpeechDetails.
documents = addPartOfSpeechDetails(documents);
documents = removeStopWords(documents);
documents = normalizeWords(documents, 'Style', 'lemma');

% Erase punctuation.
documents = erasePunctuation(documents);

% Remove words with 2 or fewer characters, and words with 15 or more
% characters.
documents = removeShortWords(documents,2);
documents = removeLongWords(documents,15);

end
```

## See Also

`erasePunctuation` | `tokenizedDocument` | `bagOfWords` | `removeStopWords` |  
`removeLongWords` | `removeShortWords` | `normalizeWords` | `wordcloud` |  
`addPartOfSpeechDetails` | `encode`

## Related Examples

- “Analyze Text Data Using Topic Models” on page 2-13
- “Analyze Text Data Using Multiword Phrases” on page 2-7
- “Analyze Text Data Containing Emojis” on page 2-52
- “Train a Sentiment Classifier” on page 2-71
- “Classify Text Data Using Deep Learning” on page 2-85
- “Generate Text Using Deep Learning” (Deep Learning Toolbox)

## Analyze Text Data Using Multiword Phrases

This example shows how to analyze text using n-gram frequency counts.

An n-gram is a tuple of  $n$  consecutive words. For example, a bigram (the case when  $n = 2$ ) is a pair of consecutive words such as "heavy rainfall". A unigram (the case when  $n = 1$ ) is a single word. A bag-of-n-grams model records the number of times that different n-grams appear in document collections.

Using a bag-of-n-grams model, you can retain more information on word ordering in the original text data. For example, a bag-of-n-grams model is better suited for capturing short phrases which appear in the text, such as "heavy rainfall" and "thunderstorm winds".

To create a bag-of-n-grams model, use `bagOfNgrams`. You can input `bagOfNgrams` objects into other Text Analytics Toolbox functions such as `wordcloud` and `fitlda`.

### Load and Extract Text Data

Load the example data. The file `factoryReports.csv` contains factory reports, including a text description and categorical labels for each event. Remove the rows with empty reports.

```
filename = "factoryReports.csv";
data = readtable(filename, TextType="string");
```

Extract the text data from the table and view the first few reports.

```
textData = data.Description;
textData(1:5)

ans = 5x1 string
    "Items are occasionally getting stuck in the scanner spools."
    "Loud rattling and banging sounds are coming from assembler pistons."
    "There are cuts to the power when starting the plant."
    "Fried capacitors in the assembler."
    "Mixer tripped the fuses."
```

### Prepare Text Data for Analysis

Create a function which tokenizes and preprocesses the text data so it can be used for analysis. The function `preprocessText` listed at the end of the example, performs the following steps:

- 1 Convert the text data to lowercase using `lower`.
- 2 Tokenize the text using `tokenizedDocument`.
- 3 Erase punctuation using `erasePunctuation`.
- 4 Remove a list of stop words (such as "and", "of", and "the") using `removeStopWords`.
- 5 Remove words with 2 or fewer characters using `removeShortWords`.
- 6 Remove words with 15 or more characters using `removeLongWords`.
- 7 Lemmatize the words using `normalizeWords`.

Use the example preprocessing function `preprocessText` to prepare the text data.

```
documents = preprocessText(textData);
documents(1:5)
```

```
ans =
  5×1 tokenizedDocument:

    6 tokens: item occasionally get stuck scanner spool
    7 tokens: loud rattling bang sound come assembler piston
    4 tokens: cut power start plant
    3 tokens: fry capacitor assembler
    3 tokens: mixer trip fuse
```

### Create Word Cloud of Bigrams

Create a word cloud of bigrams by first creating a bag-of-n-grams model using `bagOfNgrams`, and then inputting the model to `wordcloud`.

To count the n-grams of length 2 (bigrams), use `bagOfNgrams` with the default options.

```
bag = bagOfNgrams(documents)
```

```
bag =
  bagOfNgrams with properties:

    Counts: [480×921 double]
    Vocabulary: ["item" "occasionally" "get" "stuck" "scanner" "loud" "rattl...
    Ngrams: [921×2 string]
    NgramLengths: 2
    NumNgrams: 921
    NumDocuments: 480
```

Visualize the bag-of-n-grams model using a word cloud.

```
figure
wordcloud(bag);
title("Text Data: Preprocessed Bigrams")
```





The word clouds highlight commonly co-occurring bigrams in the LDA topics. The function plots the bigrams with sizes according to their probabilities for the specified LDA topics.

### Analyze Text Using Longer Phrases

To analyze text using longer phrases, specify the `NGramLengths` option in `bagOfNgrams` to be a larger value.

When working with longer phrases, it can be useful to keep stop words in the model. For example, to detect the phrase "is not happy", keep the stop words "is" and "not" in the model.

Preprocess the text. Erase the punctuation using `erasePunctuation`, and tokenize using `tokenizedDocument`.

```
cleanTextData = erasePunctuation(textData);
documents = tokenizedDocument(cleanTextData);
```

To count the n-grams of length 3 (trigrams), use `bagOfNgrams` and specify `NGramLengths` to be 3.

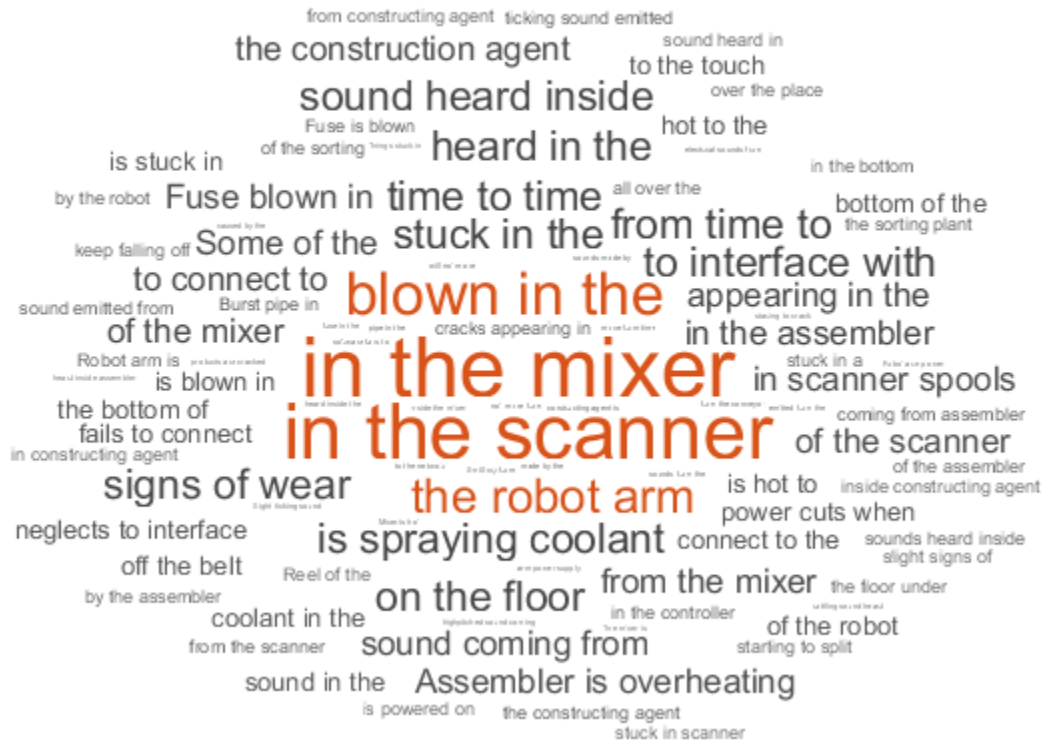
```
bag = bagOfNgrams(documents,NGramLengths=3);
```

Visualize the bag-of-n-grams model using a word cloud. The word cloud of trigrams better shows the context of the individual words.

```
figure
wordcloud(bag);
title("Text Data: Trigrams")
```



## Text Data: Trigrams



View the top 10 trigrams and their frequency counts using `topkngrams`.

```
tbl = topkngrams(bag, 10)
```

`tbl=10x3 table`

Ngram			Count	NgramLength
"in"	"the"	"mixer"	14	3
"in"	"the"	"scanner"	13	3
"blown"	"in"	"the"	9	3
"the"	"robot"	"arm"	7	3
"stuck"	"in"	"the"	6	3
"is"	"spraying"	"coolant"	6	3
"from"	"time"	"to"	6	3
"time"	"to"	"time"	6	3
"heard"	"in"	"the"	6	3
"on"	"the"	"floor"	6	3

### Example Preprocessing Function

The function `preprocessText` performs the following steps in order:

- 1 Convert the text data to lowercase using `lower`.
- 2 Tokenize the text using `tokenizedDocument`.

- 3 Erase punctuation using `erasePunctuation`.
- 4 Remove a list of stop words (such as "and", "of", and "the") using `removeStopWords`.
- 5 Remove words with 2 or fewer characters using `removeShortWords`.
- 6 Remove words with 15 or more characters using `removeLongWords`.
- 7 Lemmatize the words using `normalizeWords`.

```
function documents = preprocessText(textData)

% Convert the text data to lowercase.
cleanTextData = lower(textData);

% Tokenize the text.
documents = tokenizedDocument(cleanTextData);

% Erase punctuation.
documents = erasePunctuation(documents);

% Remove a list of stop words.
documents = removeStopWords(documents);

% Remove words with 2 or fewer characters, and words with 15 or greater
% characters.
documents = removeShortWords(documents,2);
documents = removeLongWords(documents,15);

% Lemmatize the words.
documents = addPartOfSpeechDetails(documents);
documents = normalizeWords(documents,Style="lemma");

end
```

### See Also

`tokenizedDocument` | `bagOfWords` | `removeStopWords` | `erasePunctuation` | `removeLongWords` | `removeShortWords` | `bagOfNgrams` | `normalizeWords` | `topkngrams` | `fitlda` | `ldaModel` | `wordcloud` | `addPartOfSpeechDetails`

### Related Examples

- “Create Simple Text Model for Classification” on page 2-2
- “Analyze Text Data Containing Emojis” on page 2-52
- “Analyze Text Data Using Topic Models” on page 2-13
- “Train a Sentiment Classifier” on page 2-71
- “Classify Text Data Using Deep Learning” on page 2-85
- “Generate Text Using Deep Learning” (Deep Learning Toolbox)

## Analyze Text Data Using Topic Models

This example shows how to use the Latent Dirichlet Allocation (LDA) topic model to analyze text data.

A Latent Dirichlet Allocation (LDA) model is a topic model which discovers underlying topics in a collection of documents and infers the word probabilities in topics.

### Load and Extract Text Data

Load the example data. The file `factoryReports.csv` contains factory reports, including a text description and categorical labels for each event.

```
data = readtable("factoryReports.csv", TextType="string");
head(data)
```

```
ans=8x5 table
```

Description	Category
"Items are occasionally getting stuck in the scanner spools."	"Mechanical Failure"
"Loud rattling and banging sounds are coming from assembler pistons."	"Mechanical Failure"
"There are cuts to the power when starting the plant."	"Electronic Failure"
"Fried capacitors in the assembler."	"Electronic Failure"
"Mixer tripped the fuses."	"Electronic Failure"
"Burst pipe in the constructing agent is spraying coolant."	"Leak"
"A fuse is blown in the mixer."	"Electronic Failure"
"Things continue to tumble off of the belt."	"Mechanical Failure"

Extract the text data from the field `Description`.

```
textData = data.Description;
textData(1:10)
```

```
ans = 10x1 string
```

```
"Items are occasionally getting stuck in the scanner spools."
"Loud rattling and banging sounds are coming from assembler pistons."
"There are cuts to the power when starting the plant."
"Fried capacitors in the assembler."
"Mixer tripped the fuses."
"Burst pipe in the constructing agent is spraying coolant."
"A fuse is blown in the mixer."
"Things continue to tumble off of the belt."
"Falling items from the conveyor belt."
"The scanner reel is split, it will soon begin to curve."
```

### Prepare Text Data for Analysis

Create a function which tokenizes and preprocesses the text data so it can be used for analysis. The function `preprocessText`, listed in the Preprocessing Function on page 2-0 section of the example, performs the following steps in order:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Lemmatize the words using `normalizeWords`.
- 3 Erase punctuation using `erasePunctuation`.

- 4 Remove a list of stop words (such as "and", "of", and "the") using `removeStopWords`.
- 5 Remove words with 2 or fewer characters using `removeShortWords`.
- 6 Remove words with 15 or more characters using `removeLongWords`.

Prepare the text data for analysis using the `preprocessText` function.

```
documents = preprocessText(textData);
documents(1:5)

ans =
  5×1 tokenizedDocument:

    6 tokens: item occasionally get stuck scanner spool
    7 tokens: loud rattling bang sound come assembler piston
    4 tokens: cut power start plant
    3 tokens: fry capacitor assembler
    3 tokens: mixer trip fuse
```

Create a bag-of-words model from the tokenized documents.

```
bag = bagOfWords(documents)

bag =
  bagOfWords with properties:

    Counts: [480×338 double]
    Vocabulary: [1×338 string]
    NumWords: 338
    NumDocuments: 480
```

Remove words from the bag-of-words model that have do not appear more than two times in total. Remove any documents containing no words from the bag-of-words model.

```
bag = removeInfrequentWords(bag,2);
bag = removeEmptyDocuments(bag)

bag =
  bagOfWords with properties:

    Counts: [480×158 double]
    Vocabulary: [1×158 string]
    NumWords: 158
    NumDocuments: 480
```

### Fit LDA Model

Fit an LDA model with 7 topics. For an example showing how to choose the number of topics, see "Choose Number of Topics for LDA Model" on page 2-19. To suppress verbose output, set the `Verbose` option to 0. For reproducibility, use the `rng` function with the "default" option.

```
rng("default")
numTopics = 7;
mdl = fitlda(bag,numTopics,Verbose=0);
```

If you have a large dataset, then the stochastic approximate variational Bayes solver is usually better suited as it can fit a good model in fewer passes of the data. The default solver for `fitlda` (collapsed

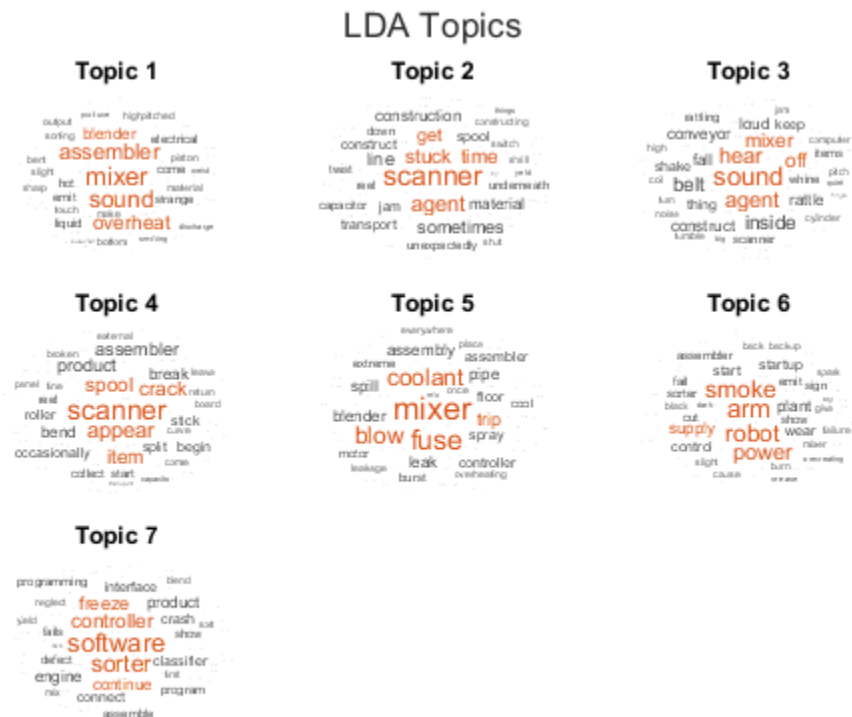
Gibbs sampling) can be more accurate at the cost of taking longer to run. To use stochastic approximate variational Bayes, set the Solver option to "savb". For an example showing how to compare LDA solvers, see "Compare LDA Solvers" on page 2-23.

### Visualize Topics Using Word Clouds

You can use word clouds to view the words with the highest probabilities in each topic. Visualize the topics using word clouds.

```
figure
t = tiledlayout("flow");
title(t,"LDA Topics")

for i = 1:numTopics
    nexttile
    wordcloud mdl,i);
    title("Topic " + i)
end
```



### View Mixtures of Topics in Documents

Create an array of tokenized documents for a set of previously unseen documents using the same preprocessing function as the training data.

```
str = [
    "Coolant is pooling underneath assembler."
    "Sorter blows fuses at start up."
```

```
"There are some very loud rattling sounds coming from the assembler."];
```

```
newDocuments = preprocessText(str);
```

Use the `transform` function to transform the documents into vectors of topic probabilities. Note that for very short documents, the topic mixtures may not be a strong representation of the document content.

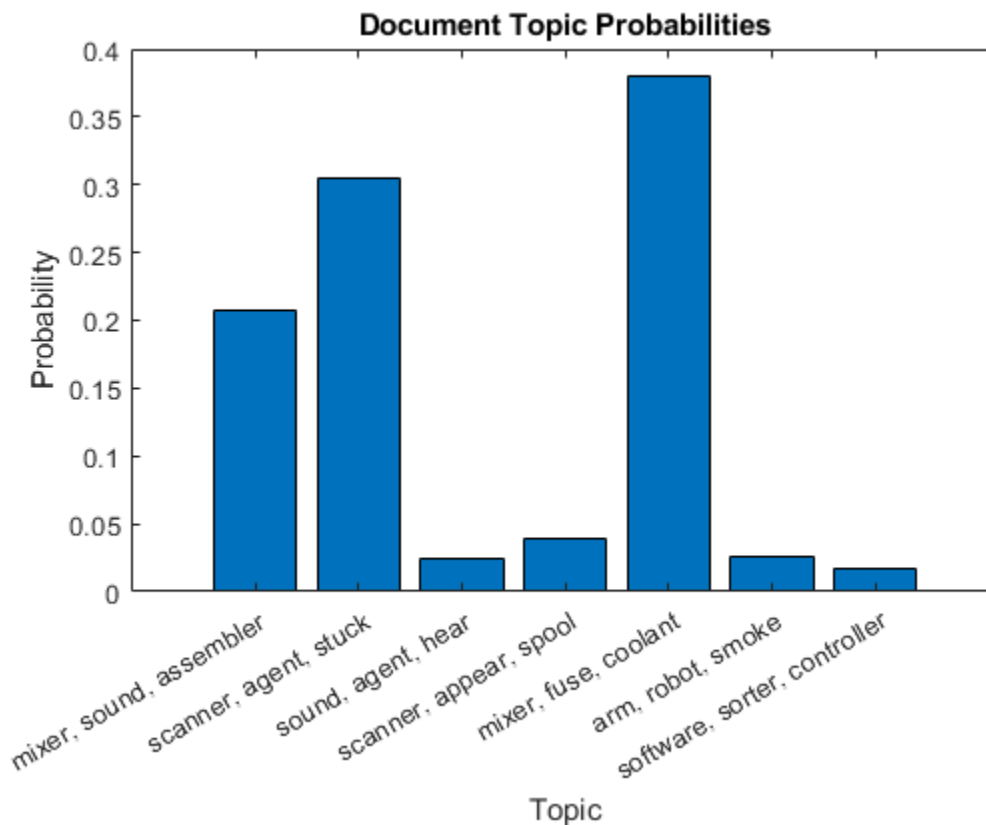
```
topicMixtures = transform mdl, newDocuments);
```

Plot the document topic probabilities of the first document in a bar chart. To label the topics, use the top three words of the corresponding topic.

```
for i = 1:numTopics
    top = topkwords(mdl,3,i);
    topWords(i) = join(top.Word, " ");
end
```

```
figure
bar(topicMixtures(1,:))
```

```
xlabel("Topic")
xticklabels(topWords);
ylabel("Probability")
title("Document Topic Probabilities")
```



Visualize multiple topic mixtures using stacked bar charts. Visualize the topic mixtures of the documents.

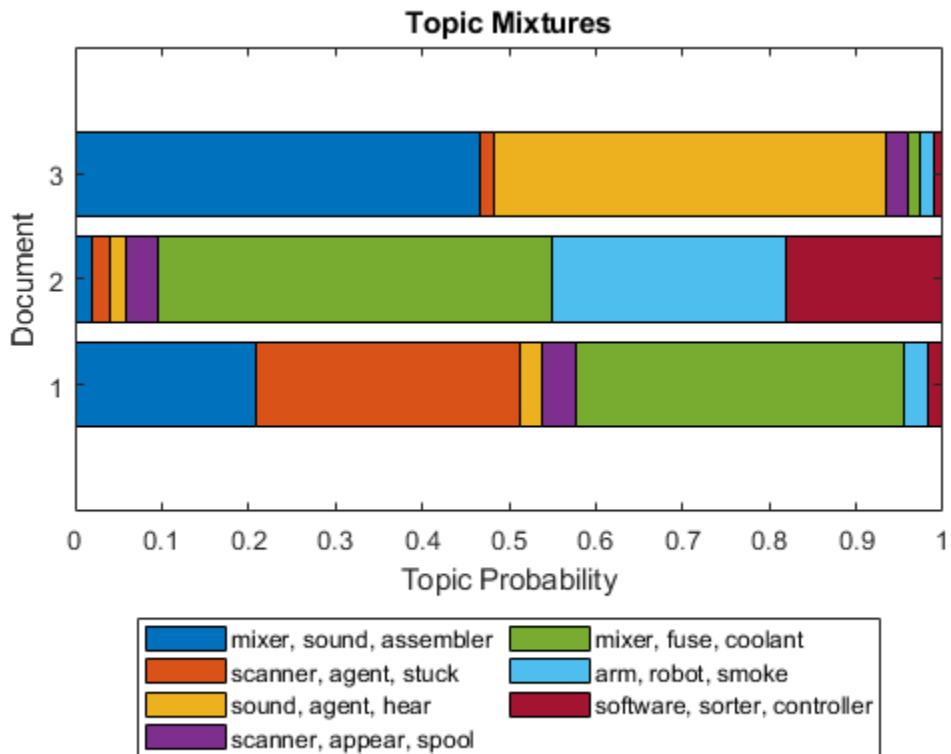
```

figure
barh(topicMixtures,"stacked")
xlim([0 1])

title("Topic Mixtures")
xlabel("Topic Probability")
ylabel("Document")

legend(topWords, ...
      Location="southoutside", ...
      NumColumns=2)

```



## Preprocessing Function

The function `preprocessText`, performs the following steps in order:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Lemmatize the words using `normalizeWords`.
- 3 Erase punctuation using `erasePunctuation`.
- 4 Remove a list of stop words (such as "and", "of", and "the") using `removeStopWords`.
- 5 Remove words with 2 or fewer characters using `removeShortWords`.
- 6 Remove words with 15 or more characters using `removeLongWords`.

```
function documents = preprocessText(textData)
```

```
% Tokenize the text.
documents = tokenizedDocument(textData);

% Lemmatize the words.
documents = addPartOfSpeechDetails(documents);
documents = normalizeWords(documents,Style="lemma");

% Erase punctuation.
documents = erasePunctuation(documents);

% Remove a list of stop words.
documents = removeStopWords(documents);

% Remove words with 2 or fewer characters, and words with 15 or greater
% characters.
documents = removeShortWords(documents,2);
documents = removeLongWords(documents,15);

end
```

### See Also

[tokenizedDocument](#) | [bagOfWords](#) | [removeStopWords](#) | [fitlda](#) | [ldaModel](#) | [wordcloud](#) | [addPartOfSpeechDetails](#) | [removeEmptyDocuments](#) | [removeInfrequentWords](#) | [transform](#)

### Related Examples

- “Visualize Document Clusters Using LDA Model” on page 2-35
- “Visualize LDA Topic Correlations” on page 2-38
- “Visualize Correlations Between LDA Topics and Document Labels” on page 2-42
- “Analyze Text Data Using Multiword Phrases” on page 2-7
- “Classify Text Data Using Deep Learning” on page 2-85



## Choose Number of Topics for LDA Model

This example shows how to decide on a suitable number of topics for a latent Dirichlet allocation (LDA) model.

To decide on a suitable number of topics, you can compare the goodness-of-fit of LDA models fit with varying numbers of topics. You can evaluate the goodness-of-fit of an LDA model by calculating the perplexity of a held-out set of documents. The perplexity indicates how well the model describes a set of documents. A lower perplexity suggests a better fit.

### Extract and Preprocess Text Data

Load the example data. The file `factoryReports.csv` contains factory reports, including a text description and categorical labels for each event. Extract the text data from the field `Description`.

```
filename = "factoryReports.csv";
data = readtable(filename, 'TextType', 'string');
textData = data.Description;
```

Tokenize and preprocess the text data using the function `preprocessText` which is listed at the end of this example.

```
documents = preprocessText(textData);
documents(1:5)
```

```
ans =
    5×1 tokenizedDocument:

    6 tokens: item occasionally get stuck scanner spool
    7 tokens: loud rattle bang sound come assembler piston
    4 tokens: cut power start plant
    3 tokens: fry capacitor assembler
    3 tokens: mixer trip fuse
```

Set aside 10% of the documents at random for validation.

```
numDocuments = numel(documents);
cvp = cvpartition(numDocuments, 'HoldOut', 0.1);
documentsTrain = documents(cvp.training);
documentsValidation = documents(cvp.test);
```

Create a bag-of-words model from the training documents. Remove the words that do not appear more than two times in total. Remove any documents containing no words.

```
bag = bagOfWords(documentsTrain);
bag = removeInfrequentWords(bag, 2);
bag = removeEmptyDocuments(bag);
```

### Choose Number of Topics

The goal is to choose a number of topics that minimize the perplexity compared to other numbers of topics. This is not the only consideration: models fit with larger numbers of topics may take longer to converge. To see the effects of the tradeoff, calculate both goodness-of-fit and the fitting time. If the optimal number of topics is high, then you might want to choose a lower value to speed up the fitting process.

Fit some LDA models for a range of values for the number of topics. Compare the fitting time and the perplexity of each model on the held-out set of test documents. The perplexity is the second output to the `logp` function. To obtain the second output without assigning the first output to anything, use the `~` symbol. The fitting time is the `TimeSinceStart` value for the last iteration. This value is in the `History` struct of the `FitInfo` property of the LDA model.

For a quicker fit, specify `'Solver'` to be `'savb'`. To suppress verbose output, set `'Verbose'` to `0`. This may take a few minutes to run.

```
numTopicsRange = [5 10 15 20 40];
for i = 1:numel(numTopicsRange)
    numTopics = numTopicsRange(i);

    mdl = fitlda(bag,numTopics, ...
        'Solver','savb', ...
        'Verbose',0);

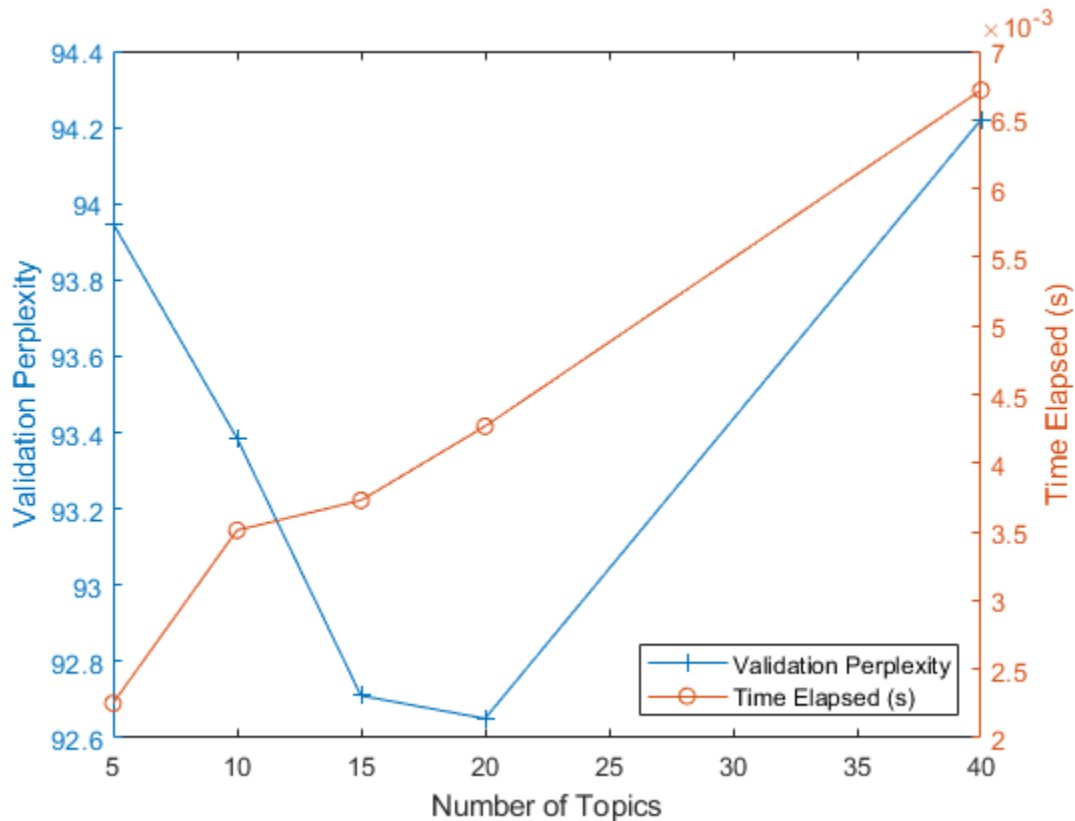
    [~,validationPerplexity(i)] = logp(mdl,documentsValidation);
    timeElapsed(i) = mdl.FitInfo.History.TimeSinceStart(end);
end
```

Show the perplexity and elapsed time for each number of topics in a plot. Plot the perplexity on the left axis and the time elapsed on the right axis.

```
figure
yyaxis left
plot(numTopicsRange,validationPerplexity,'+-')
ylabel("Validation Perplexity")

yyaxis right
plot(numTopicsRange,timeElapsed,'o-')
ylabel("Time Elapsed (s)")

legend(["Validation Perplexity" "Time Elapsed (s)"],"Location','southeast')
xlabel("Number of Topics")
```



The plot suggests that fitting a model with 10–20 topics may be a good choice. The perplexity is low compared with the models with different numbers of topics. With this solver, the elapsed time for this many topics is also reasonable. With different solvers, you may find that increasing the number of topics can lead to a better fit, but fitting the model takes longer to converge.

### Example Preprocessing Function

The function `preprocessText`, performs the following steps in order:

- 1 Convert the text data to lowercase using `lower`.
- 2 Tokenize the text using `tokenizedDocument`.
- 3 Erase punctuation using `erasePunctuation`.
- 4 Remove a list of stop words (such as "and", "of", and "the") using `removeStopWords`.
- 5 Remove words with 2 or fewer characters using `removeShortWords`.
- 6 Remove words with 15 or more characters using `removeLongWords`.
- 7 Lemmatize the words using `normalizeWords`.

```
function documents = preprocessText(textData)
% Convert the text data to lowercase.
cleanTextData = lower(textData);
% Tokenize the text.
documents = tokenizedDocument(cleanTextData);
```

```
% Erase punctuation.
documents = erasePunctuation(documents);

% Remove a list of stop words.
documents = removeStopWords(documents);

% Remove words with 2 or fewer characters, and words with 15 or greater
% characters.
documents = removeShortWords(documents,2);
documents = removeLongWords(documents,15);

% Lemmatize the words.
documents = addPartOfSpeechDetails(documents);
documents = normalizeWords(documents, 'Style', 'Lemma');
end
```

### See Also

[tokenizedDocument](#) | [bagOfWords](#) | [removeStopWords](#) | [logp](#) | [bagOfWords](#) | [fitlda](#) | [ldaModel](#) | [erasePunctuation](#) | [removeShortWords](#) | [removeLongWords](#) | [normalizeWords](#) | [addPartOfSpeechDetails](#) | [removeInfrequentWords](#) | [removeEmptyDocuments](#)

### Related Examples

- “Analyze Text Data Using Topic Models” on page 2-13
- “Compare LDA Solvers” on page 2-23
- “Visualize Document Clusters Using LDA Model” on page 2-35
- “Visualize LDA Topic Correlations” on page 2-38
- “Visualize Correlations Between LDA Topics and Document Labels” on page 2-42

## Compare LDA Solvers

This example shows how to compare latent Dirichlet allocation (LDA) solvers by comparing the goodness of fit and the time taken to fit the model.

### Import Text Data

Import a set of abstracts and category labels from math papers using the arXiv API. Specify the number of records to import using the `importSize` variable.

```
importSize = 50000;
```

Create a URL that queries records with set "math" and metadata prefix "arXiv".

```
url = "https://export.arxiv.org/oai2?verb=ListRecords" + ...
      "&set=math" + ...
      "&metadataPrefix=arXiv";
```

Extract the abstract text and the resumption token returned by the query URL using the `parseArXivRecords` function which is attached to this example as a supporting file. To access this file, open this example as a live script. Note that the arXiv API is rate limited and requires waiting between multiple requests.

```
[textData,~,resumptionToken] = parseArXivRecords(url);
```

Iteratively import more chunks of records until the required amount is reached, or there are no more records. To continue importing records from where you left off, use the resumption token from the previous result in the query URL. To adhere to the rate limits imposed by the arXiv API, add a delay of 20 seconds before each query using the `pause` function.

```
while numel(textData) < importSize
    if resumptionToken == ""
        break
    end

    url = "https://export.arxiv.org/oai2?verb=ListRecords" + ...
          "&resumptionToken=" + resumptionToken;

    pause(20)
    [textDataNew,labelsNew,resumptionToken] = parseArXivRecords(url);

    textData = [textData; textDataNew];
end
```

### Preprocess Text Data

Set aside 10% of the documents at random for validation.

```
numDocuments = numel(textData);
cvp = cvpartition(numDocuments,'HoldOut',0.1);
textDataTrain = textData(training(cvp));
textDataValidation = textData(test(cvp));
```

Tokenize and preprocess the text data using the function `preprocessText` which is listed at the end of this example.

```
documentsTrain = preprocessText(textDataTrain);
documentsValidation = preprocessText(textDataValidation);
```

Create a bag-of-words model from the training documents. Remove the words that do not appear more than two times in total. Remove any documents containing no words.

```
bag = bagOfWords(documentsTrain);
bag = removeInfrequentWords(bag,2);
bag = removeEmptyDocuments(bag);
```

For the validation data, create a bag-of-words model from the validation documents. You do not need to remove any words from the validation data because any words that do not appear in the fitted LDA models are automatically ignored.

```
validationData = bagOfWords(documentsValidation);
```

### Fit and Compare Models

For each of the LDA solvers, fit a model with 40 topics. To distinguish the solvers when plotting the results on the same axes, specify different line properties for each solver.

```
numTopics = 40;
solvers = ["cgs" "avb" "cvb0" "savb"];
lineSpecs = ["+-" "*-" "x-" "o-"];
```

Fit an LDA model using each solver. For each solver, specify the initial topic concentration 1, to validate the model once per data pass, and to not fit the topic concentration parameter. Using the data in the `FitInfo` property of the fitted LDA models, plot the validation perplexity and the time elapsed.

The stochastic solver, by default, uses a mini-batch size of 1000 and validates the model every 10 iterations. For this solver, to validate the model once per data pass, set the validation frequency to `ceil(numObservations/1000)`, where `numObservations` is the number of documents in the training data. For the other solvers, set the validation frequency to 1.

For the iterations that the stochastic solver does not evaluate the validation perplexity, the stochastic solver reports NaN in the `FitInfo` property. To plot the validation perplexity, remove the NaNs from the reported values.

```
numObservations = bag.NumDocuments;

figure
for i = 1:numel(solvers)
    solver = solvers(i);
    lineSpec = lineSpecs(i);

    if solver == "savb"
        numIterationsPerDataPass = ceil(numObservations/1000);
    else
        numIterationsPerDataPass = 1;
    end

    mdl = fitlda(bag,numTopics, ...
        'Solver',solver, ...
        'InitialTopicConcentration',1, ...
        'FitTopicConcentration',false, ...
        'ValidationData',validationData, ...
```

```

    'ValidationFrequency',numIterationsPerDataPass, ...
    'Verbose',0);

history = mdl.FitInfo.History;

timeElapsed = history.TimeSinceStart;

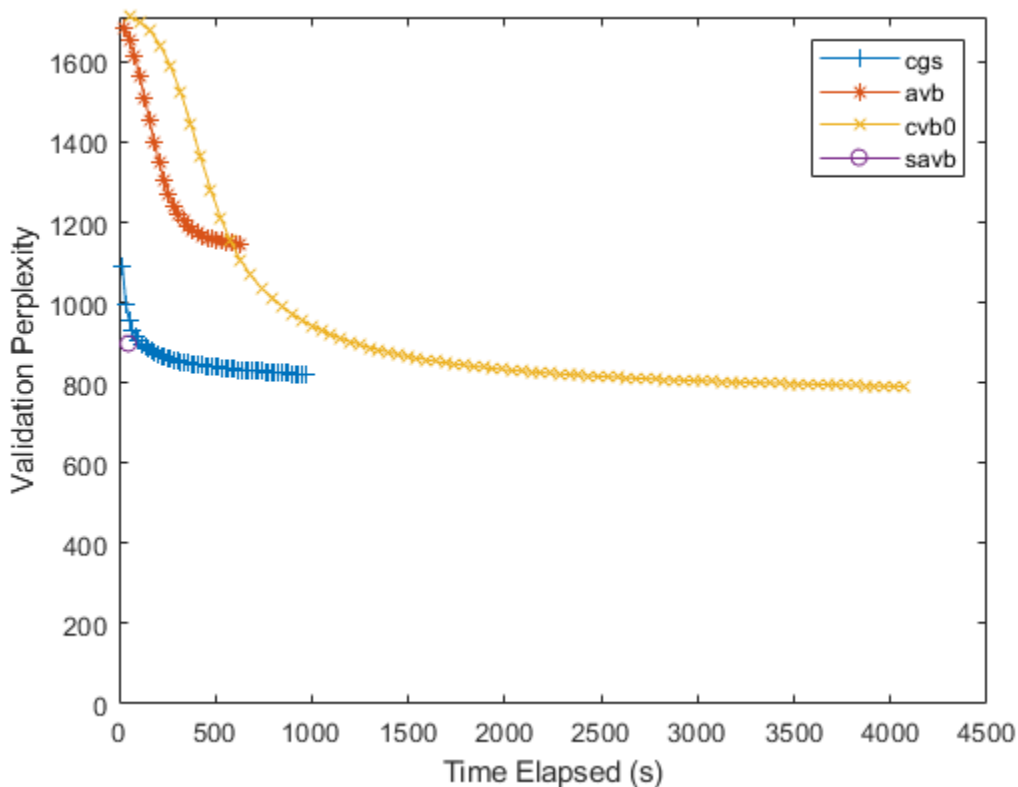
validationPerplexity = history.ValidationPerplexity;

% Remove NaNs.
idx = isnan(validationPerplexity);
timeElapsed(idx) = [];
validationPerplexity(idx) = [];

plot(timeElapsed,validationPerplexity,lineSpec)
hold on
end

hold off
xlabel("Time Elapsed (s)")
ylabel("Validation Perplexity")
ylim([0 inf])
legend(solvers)

```



For the stochastic solver, there is only one data point. This is because this solver passes through input data once. To specify more data passes, use the 'DataPassLimit' option. For the batch solvers

("cgs", "avb", and "cvb0"), to specify the number of iterations used to fit the models, use the 'IterationLimit' option.

A lower validation perplexity suggests a better fit. Usually, the solvers "savb" and "cgs" converge quickly to a good fit. The solver "cvb0" might converge to a better fit, but it can take much longer to converge.

For the `FitInfo` property, the `fitlda` function estimates the validation perplexity from the document probabilities at the maximum likelihood estimates of the per-document topic probabilities. This is usually quicker to compute, but can be less accurate than other methods. Alternatively, calculate the validation perplexity using the `logp` function. This function calculates more accurate values but can take longer to run. For an example showing how to compute the perplexity using `logp`, see “Calculate Document Log-Probabilities from Word Count Matrix”.

### Preprocessing Function

The function `preprocessText` performs the following steps:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Lemmatize the words using `normalizeWords`.
- 3 Erase punctuation using `erasePunctuation`.
- 4 Remove a list of stop words (such as "and", "of", and "the") using `removeStopWords`.
- 5 Remove words with 2 or fewer characters using `removeShortWords`.
- 6 Remove words with 15 or more characters using `removeLongWords`.

```
function documents = preprocessText(textData)

% Tokenize the text.
documents = tokenizedDocument(textData);

% Lemmatize the words.
documents = addPartOfSpeechDetails(documents);
documents = normalizeWords(documents, 'Style', 'lemma');

% Erase punctuation.
documents = erasePunctuation(documents);

% Remove a list of stop words.
documents = removeStopWords(documents);

% Remove words with 2 or fewer characters, and words with 15 or greater
% characters.
documents = removeShortWords(documents, 2);
documents = removeLongWords(documents, 15);

end
```

### See Also

`tokenizedDocument` | `bagOfWords` | `removeStopWords` | `logp` | `fitlda` | `ldaModel` | `wordcloud` | `removeInfrequentWords` | `removeEmptyDocuments` | `erasePunctuation` | `removeShortWords` | `removeLongWords` | `normalizeWords` | `addPartOfSpeechDetails`



**Related Examples**

- “Analyze Text Data Using Topic Models” on page 2-13
- “Choose Number of Topics for LDA Model” on page 2-19
- “Visualize Document Clusters Using LDA Model” on page 2-35
- “Visualize LDA Topic Correlations” on page 2-38
- “Visualize Correlations Between LDA Topics and Document Labels” on page 2-42

## Visualize LDA Topics Using Word Clouds

This example shows how to visualize the words in Latent Dirichlet Allocation (LDA) model topics.

A latent Dirichlet allocation (LDA) model is a topic model which discovers underlying topics in a collection of documents and infers word probabilities in topics. You can visualize the LDA topics using word clouds by displaying words with their corresponding topic word probabilities.

### Load LDA Model

Load the LDA model `factoryReportsLDAModel` which is trained using a data set of factory reports detailing different failure events. For an example showing how to fit an LDA model to a collection of text data, see “Analyze Text Data Using Topic Models” on page 2-13.

```
load factoryReportsLDAModel
mdl

mdl =
  ldaModel with properties:

          NumTopics: 7
      WordConcentration: 1
    TopicConcentration: 0.5755
 CorpusTopicProbabilities: [0.1587 0.1573 0.1551 0.1534 0.1340 ... ]
DocumentTopicProbabilities: [480x7 double]
  TopicWordProbabilities: [158x7 double]
          Vocabulary: ["item"      "occasionally"      "get"      ...      ]
          TopicOrder: 'initial-fit-probability'
          FitInfo: [1x1 struct]
```

### Visualize Topics Using Word Clouds

Visualize the topics using the `wordcloud` function.

```
numTopics = mdl.NumTopics;

figure
t = tiledlayout("flow");
title(t,"LDA Topics")

for i = 1:numTopics
    nexttile
    wordcloud(mdl,i);
    title("Topic " + i)
end
```



## Visualize LDA Topic Probabilities of Documents

This example shows how to visualize the topic probabilities of documents using a latent Dirichlet allocation (LDA) topic model.

A latent Dirichlet allocation (LDA) model is a topic model which discovers underlying topics in a collection of documents and infers word probabilities in topics. You can use an LDA model to transform documents into a vector of topic probabilities, also known as a topic mixture. You can visualize the LDA topics using stacked bar charts.

### Load LDA Model

Load the LDA model `factoryReportsLDAModel` which is trained using a data set of factory reports detailing different failure events. For an example showing how to fit an LDA model to a collection of text data, see “Analyze Text Data Using Topic Models” on page 2-13.

```
load factoryReportsLDAModel
mdl

mdl =
  ldaModel with properties:

          NumTopics: 7
      WordConcentration: 1
      TopicConcentration: 0.5755
  CorpusTopicProbabilities: [0.1587 0.1573 0.1551 0.1534 0.1340 ... ]
  DocumentTopicProbabilities: [480x7 double]
      TopicWordProbabilities: [158x7 double]
          Vocabulary: ["item"      "occasionally"      "get"      ...      ]
          TopicOrder: 'initial-fit-probability'
          FitInfo: [1x1 struct]
```

### Visualize Topics Using Word Clouds

Visualize the topics using the `wordcloud` function.

```
numTopics = mdl.NumTopics;
figure
t = tiledlayout("flow");
title(t,"LDA Topics")
for i = 1:numTopics
    nexttile
        wordcloud(mdl,i);
        title("Topic " + i)
end
```

## LDA Topics

Topic 1



Topic 2



Topic 3



Topic 4



Topic 5



Topic 6



Topic 7



## View Mixtures of Topics in Documents

Create an array of tokenized documents for a set of previously unseen documents using the same preprocessing function used when fitting the model.

The function `preprocessText`, listed in the Preprocessing Function on page 2-0 section of the example, performs the following steps in order:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Lemmatize the words using `normalizeWords`.
- 3 Erase punctuation using `erasePunctuation`.
- 4 Remove a list of stop words (such as "and", "of", and "the") using `removeStopWords`.
- 5 Remove words with 2 or fewer characters using `removeShortWords`.
- 6 Remove words with 15 or more characters using `removeLongWords`.

Prepare the text data for analysis using the `preprocessText` function.

```
str = [
    "Coolant is pooling underneath assembler."
    "Sorter blows fuses at start up."
    "There are some very loud rattling sounds coming from the assembler."];
documents = preprocessText(str);
```

Transform the documents into vectors of topic probabilities using the `transform` function. Note that for very short documents, the topic mixtures may not be a strong representation of the document content.

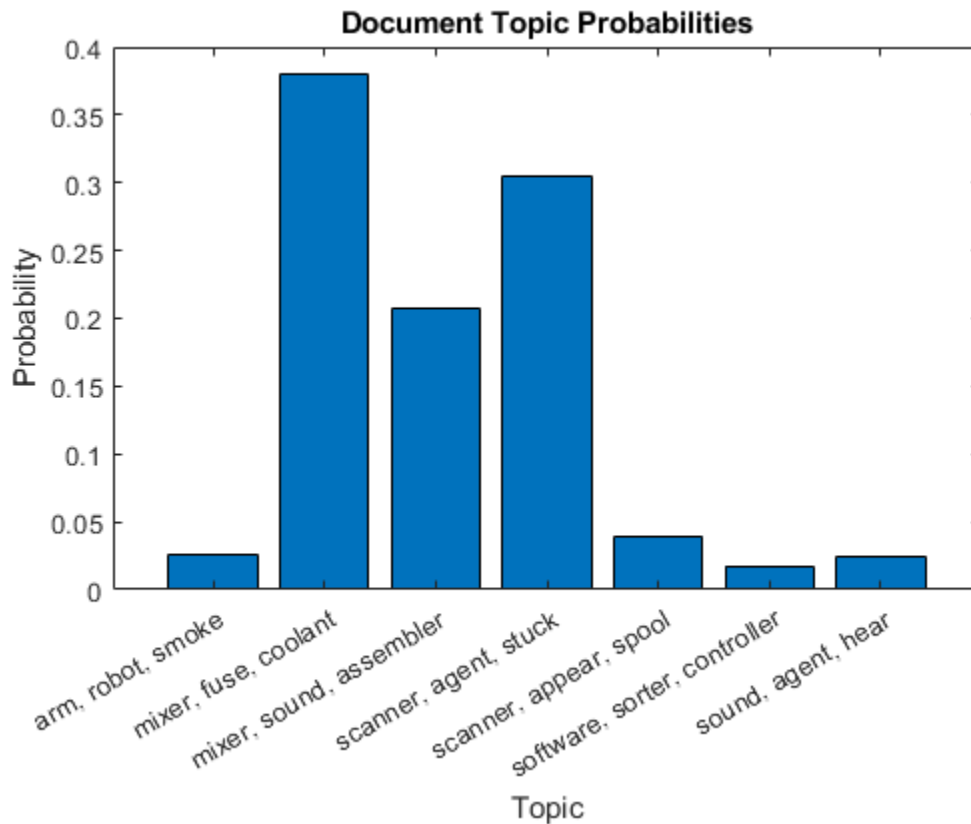
```
topicMixtures = transform mdl, documents;
```

Visualize the first topic mixture in a bar chart and label the bars using the top three words from each topic.

```
numTopics = mdl.NumTopics;
for i = 1:numTopics
    top = topkwords(mdl,3,i);
    topWords(i) = join(top.Word, ", ");
end

figure
bar(categorical(topWords),topicMixtures(1,:))

xlabel("Topic")
ylabel("Probability")
title("Document Topic Probabilities")
```



To visualize the proportions of the topics in each document, or to visualize multiple topic mixtures, use a stacked bar chart.

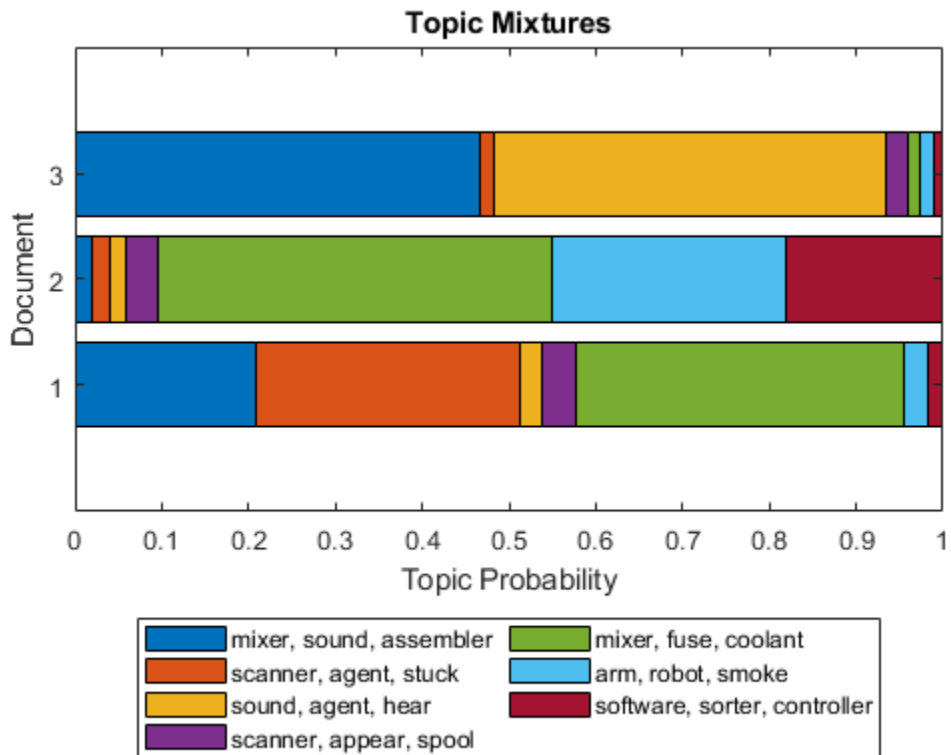
```
figure
barh(topicMixtures, "stacked")
```

```

title("Topic Mixtures")
xlabel("Topic Probability")
ylabel("Document")

legend(topWords, ...
      Location="southoutside", ...
      NumColumns=2)

```



The regions of the stacked bar chart represent the proportion of the document belonging to the corresponding topic.

### Preprocessing Function

The function `preprocessText`, performs the following steps in order:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Lemmatize the words using `normalizeWords`.
- 3 Erase punctuation using `erasePunctuation`.
- 4 Remove a list of stop words (such as "and", "of", and "the") using `removeStopWords`.
- 5 Remove words with 2 or fewer characters using `removeShortWords`.
- 6 Remove words with 15 or more characters using `removeLongWords`.

```
function documents = preprocessText(textData)
```

```
% Tokenize the text.
documents = tokenizedDocument(textData);

% Lemmatize the words.
documents = addPartOfSpeechDetails(documents);
documents = normalizeWords(documents,Style="lemma");

% Erase punctuation.
documents = erasePunctuation(documents);

% Remove a list of stop words.
documents = removeStopWords(documents);

% Remove words with 2 or fewer characters, and words with 15 or greater
% characters.
documents = removeShortWords(documents,2);
documents = removeLongWords(documents,15);

end
```

### See Also

[tokenizedDocument](#) | [fitlda](#) | [ldaModel](#) | [wordcloud](#)

### Related Examples

- “Analyze Text Data Using Topic Models” on page 2-13
- “Visualize Document Clusters Using LDA Model” on page 2-35
- “Visualize LDA Topic Correlations” on page 2-38
- “Visualize Correlations Between LDA Topics and Document Labels” on page 2-42



## Visualize Document Clusters Using LDA Model

This example shows how to visualize the clustering of documents using a Latent Dirichlet Allocation (LDA) topic model and a t-SNE plot.

A latent Dirichlet allocation (LDA) model is a topic model which discovers underlying topics in a collection of documents and infers word probabilities in topics. The vectors of per-topic word probabilities characterize the topics. You can evaluate document similarity using an LDA model by comparing the per-document topic probabilities, also known as topic mixtures.

### Load LDA Model

Load the LDA model `factoryReportsLDAModel` which is trained using a data set of factory reports detailing different failure events. For an example showing how to fit an LDA model to a collection of text data, see “Analyze Text Data Using Topic Models” on page 2-13.

```
load factoryReportsLDAModel
mdl

mdl =
  ldaModel with properties:

          NumTopics: 7
      WordConcentration: 1
    TopicConcentration: 0.5755
 CorpusTopicProbabilities: [0.1587 0.1573 0.1551 0.1534 0.1340 0.1322 0.1093]
DocumentTopicProbabilities: [480x7 double]
   TopicWordProbabilities: [158x7 double]
         Vocabulary: [1x158 string]
        TopicOrder: 'initial-fit-probability'
          FitInfo: [1x1 struct]
```

Visualize the topics using word clouds.

```
numTopics = mdl.NumTopics;

figure
tiledlayout("flow")
title("LDA Topics")

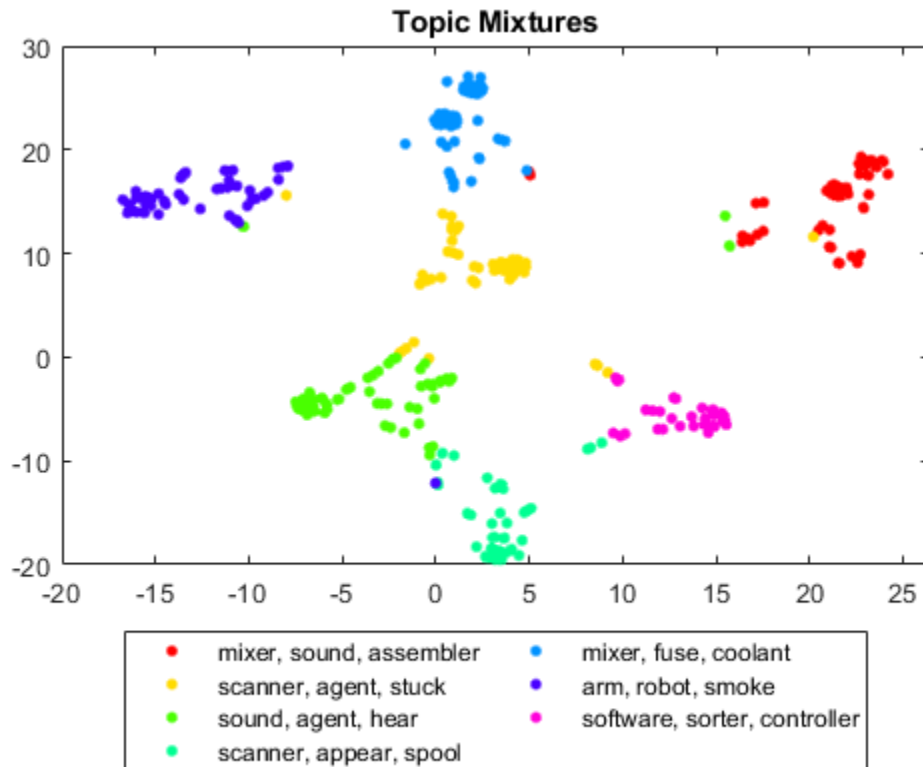
for i = 1:numTopics
    nexttile
    wordcloud(mdl,i);
    title("Topic " + i)
end
```



```

title("Topic Mixtures")
legend(topWords, ...
      Location="southoutside", ...
      NumColumns=2)

```



The t-SNE plot highlights clusters occurring in the original high-dimensional data.

## See Also

`tokenizedDocument` | `fitLda` | `LdaModel` | `wordCloud`

## Related Examples

- "Analyze Text Data Using Topic Models" on page 2-13
- "Visualize LDA Topic Correlations" on page 2-38
- "Visualize Correlations Between LDA Topics and Document Labels" on page 2-42

## Visualize LDA Topic Correlations

This example shows how to analyze correlations between topics in a Latent Dirichlet Allocation (LDA) topic model.

A latent Dirichlet allocation (LDA) model is a topic model which discovers underlying topics in a collection of documents and infers word probabilities in topics. The vectors of per-topic word probabilities characterize the topics. Using the per-topic word probabilities, you can identify correlations between the topics.

### Load LDA Model

Load the LDA model `factoryReportsLDAModel` which is trained using a data set of factory reports detailing different failure events. For an example showing how to fit an LDA model to a collection of text data, see “Analyze Text Data Using Topic Models” on page 2-13.

```
load factoryReportsLDAModel
mdl

mdl =
  ldaModel with properties:

          NumTopics: 7
      WordConcentration: 1
    TopicConcentration: 0.5755
  CorpusTopicProbabilities: [0.1587 0.1573 0.1551 0.1534 0.1340 ... ]
  DocumentTopicProbabilities: [480x7 double]
    TopicWordProbabilities: [158x7 double]
          Vocabulary: ["item"      "occasionally"      "get"      ...      ]
          TopicOrder: 'initial-fit-probability'
          FitInfo: [1x1 struct]
```

Visualize the topics using word clouds.

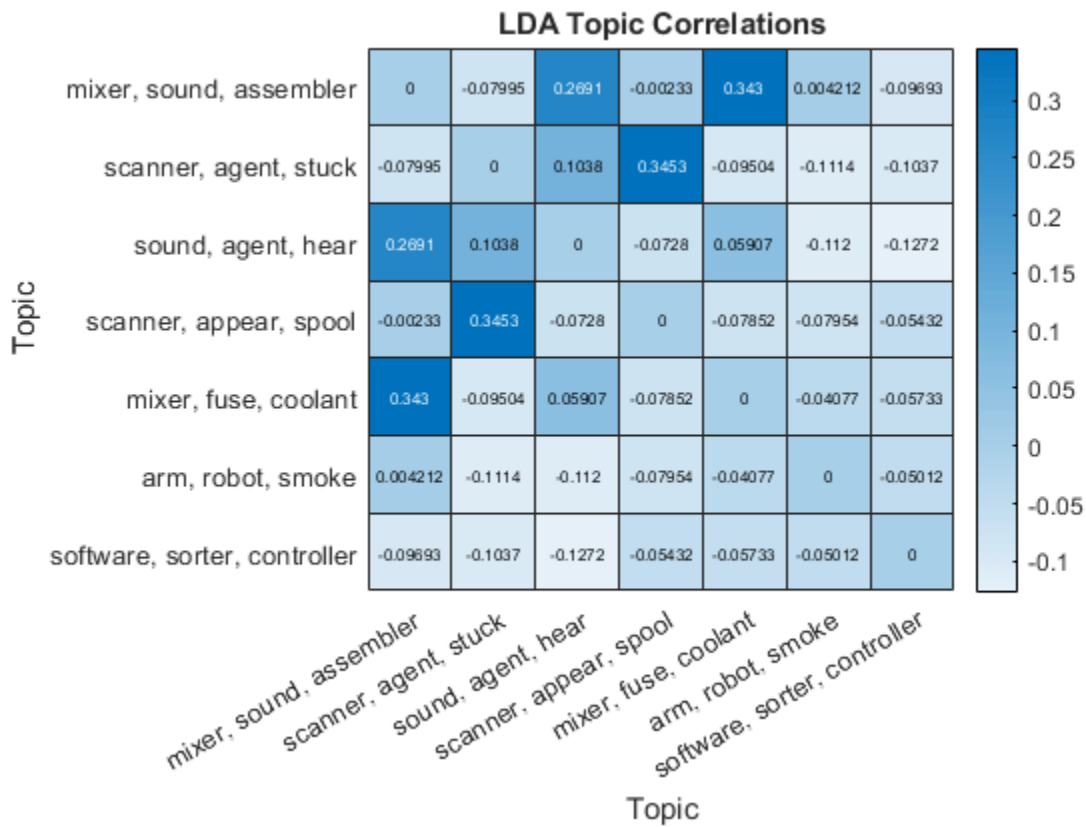
```
numTopics = mdl.NumTopics;

figure
t = tiledlayout("flow");
title(t,"LDA Topics")

for i = 1:numTopics
    nexttile
    wordcloud(mdl,i);
    title("Topic " + i)
end
```



```
xlabel("Topic")
ylabel("Topic")
```



For each topic, find the topic with the strongest correlation and display the pairs in a table with the corresponding correlation coefficient.

```
[topCorrelations,topCorrelatedTopics] = max(correlation - eye(numTopics));
```

```
tbl = table;
tbl.TopicIndex = (1:numTopics)';
tbl.Topic = topWords';
tbl.TopCorrelatedTopicIndex = topCorrelatedTopics';
tbl.TopCorrelatedTopic = topWords(topCorrelatedTopics)';
tbl.CorrelationCoefficient = topCorrelations'
```

tbl=7x5 table

TopicIndex	Topic	TopCorrelatedTopicIndex	TopCorrelatedTopic
1	"mixer, sound, assembler"	5	"mixer, fuse, coolant"
2	"scanner, agent, stuck"	4	"scanner, appear, spool"
3	"sound, agent, hear"	1	"mixer, sound, assembler"
4	"scanner, appear, spool"	2	"scanner, agent, stuck"
5	"mixer, fuse, coolant"	1	"mixer, sound, assembler"
6	"arm, robot, smoke"	1	"mixer, sound, assembler"

7 "software, sorter, controller"

7 "software, sorter"

## See Also

`tokenizedDocument` | `fitlda` | `ldaModel` | `wordcloud`

## Related Examples

- “Analyze Text Data Using Topic Models” on page 2-13
- “Visualize Document Clusters Using LDA Model” on page 2-35
- “Visualize Correlations Between LDA Topics and Document Labels” on page 2-42

## Visualize Correlations Between LDA Topics and Document Labels

This example shows how to fit a Latent Dirichlet Allocation (LDA) topic model and visualize correlations between the LDA topics and document labels.

A Latent Dirichlet Allocation (LDA) model is a topic model which discovers underlying topics in a collection of documents and infers the word probabilities in topics. Fitting an LDA model does not require labeled data. However, you can visualize correlations between the fitted LDA topics and the document labels using a parallel coordinates plot.

This example fits an LDA model to the Factory Reports data set which is a collection of factory reports detailing different failure events and identifies correlations between the LDA topics and the report category.

### Load and Extract Text Data

Load the example data. The file `factoryReports.csv` contains factory reports, including a text description and categorical labels for each event.

```
data = readtable("factoryReports.csv", TextType="string");
head(data)
```

```
ans=8x5 table
```

	Description	Category
	"Items are occasionally getting stuck in the scanner spools."	"Mechanical Failure"
	"Loud rattling and banging sounds are coming from assembler pistons."	"Mechanical Failure"
	"There are cuts to the power when starting the plant."	"Electronic Failure"
	"Fried capacitors in the assembler."	"Electronic Failure"
	"Mixer tripped the fuses."	"Electronic Failure"
	"Burst pipe in the constructing agent is spraying coolant."	"Leak"
	"A fuse is blown in the mixer."	"Electronic Failure"
	"Things continue to tumble off of the belt."	"Mechanical Failure"

Extract the text data from the field `Description`.

```
textData = data.Description;
textData(1:10)
```

```
ans = 10x1 string
```

```
"Items are occasionally getting stuck in the scanner spools."
"Loud rattling and banging sounds are coming from assembler pistons."
"There are cuts to the power when starting the plant."
"Fried capacitors in the assembler."
"Mixer tripped the fuses."
"Burst pipe in the constructing agent is spraying coolant."
"A fuse is blown in the mixer."
"Things continue to tumble off of the belt."
"Falling items from the conveyor belt."
"The scanner reel is split, it will soon begin to curve."
```

Extract the labels from the field `Category`.



```
labels = data.Category;
```

### Prepare Text Data for Analysis

Create a function which tokenizes and preprocesses the text data so it can be used for analysis. The function `preprocessText`, listed in the Preprocessing Function on page 2-0 section of the example, performs the following steps in order:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Lemmatize the words using `normalizeWords`.
- 3 Erase punctuation using `erasePunctuation`.
- 4 Remove a list of stop words (such as "and", "of", and "the") using `removeStopWords`.
- 5 Remove words with 2 or fewer characters using `removeShortWords`.
- 6 Remove words with 15 or more characters using `removeLongWords`.

Prepare the text data for analysis using the `preprocessText` function.

```
documents = preprocessText(textData);
documents(1:5)
```

```
ans = 5x1 tokenizedDocument array with properties:
    6 tokens: item occasionally get stuck scanner spool
    7 tokens: loud rattling bang sound come assembler piston
    4 tokens: cut power start plant
    3 tokens: fry capacitor assembler
    3 tokens: mixer trip fuse
```

Create a bag-of-words model from the tokenized documents.

```
bag = bagOfWords(documents)

bag = bagOfWords with properties:
    Counts: [480x338 double]
    Vocabulary: [1x338 string]
    NumWords: 338
    NumDocuments: 480
```

Remove words from the bag-of-words model that have do not appear more than two times in total. Remove any documents containing no words from the bag-of-words model.

```
bag = removeInfrequentWords(bag,2);
bag = removeEmptyDocuments(bag)

bag = bagOfWords with properties:
    Counts: [480x158 double]
    Vocabulary: [1x158 string]
    NumWords: 158
    NumDocuments: 480
```

### Fit LDA Model

Fit an LDA model with 7 topics. For an example showing how to choose the number of topics, see "Choose Number of Topics for LDA Model" on page 2-19. To suppress verbose output, set the `Verbose` option to 0. For reproducibility, set `rng` to "default".



Extract the document topic mixtures from the `DocumentTopicProbabilities` property of the LDA model.

```
topicMixtures = mdl.DocumentTopicProbabilities;
```

For the documents with each label, calculate the mean topic probabilities.

```
[groups,groupNames] = findgroups(labels);
numGroups = numel(groupNames);

for i = 1:numGroups
    idx = groups == i;
    meanTopicProbabilities(i,:) = mean(topicMixtures(idx,:));
end
```

For each topic, find the top three words.

```
numTopics = mdl.NumTopics;
for i = 1:numTopics
    top = topkwords(mdl,3,i);
    topWords(i) = join(top.Word, ", ");
end
```

Plot the per-category mean topic probabilities using a parallel coordinates plot. For readability, create a figure and increase the figure width using the `Position` property.

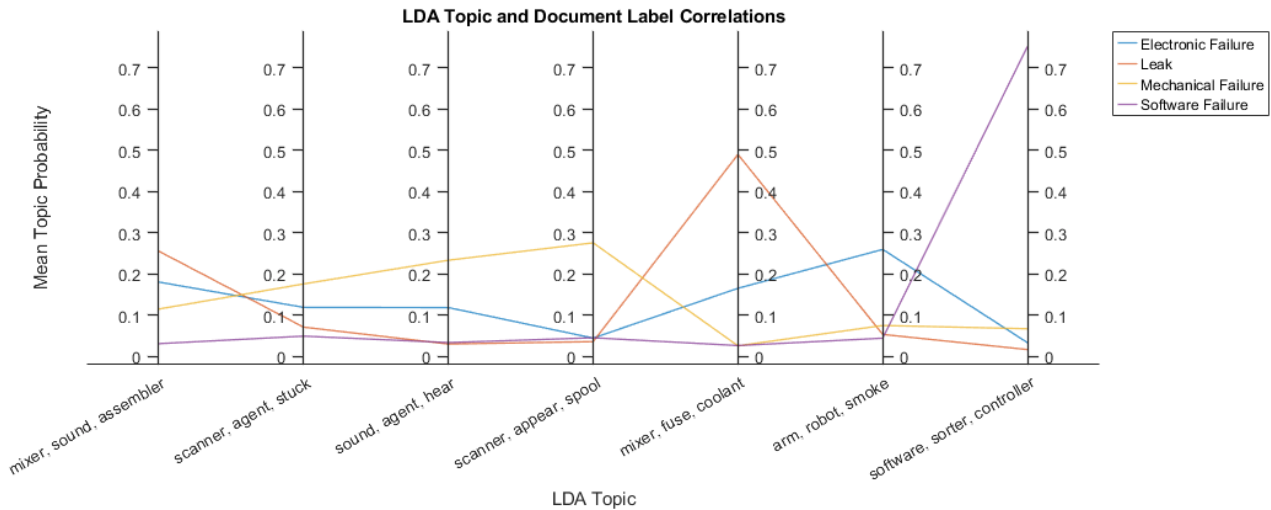
```
f = figure;
f.Position(3) = 2*f.Position(3);
```

Plot the per-category mean topic probabilities using the `parallelplot` function. Do not normalize the input data and specify the categories as the groups. Set the coordinate tick labels to the top three words of each topic.

```
p = parallelplot(meanTopicProbabilities, ...
    GroupData=groupNames, ...
    DataNormalization="none");

p.CoordinateTickLabels = topWords;

xlabel("LDA Topic")
ylabel("Mean Topic Probability")
title("LDA Topic and Document Label Correlations")
```



The parallel plot highlights the correlations between the LDA topics and the document labels. High peaks indicate a strong correlation between the corresponding topic and document label.

### Preprocessing Function

The function `preprocessText`, performs the following steps in order:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Lemmatize the words using `normalizeWords`.
- 3 Erase punctuation using `erasePunctuation`.
- 4 Remove a list of stop words (such as "and", "of", and "the") using `removeStopWords`.
- 5 Remove words with 2 or fewer characters using `removeShortWords`.
- 6 Remove words with 15 or more characters using `removeLongWords`.

```
function documents = preprocessText(textData)

% Tokenize the text.
documents = tokenizedDocument(textData);

% Lemmatize the words.
documents = addPartOfSpeechDetails(documents);
documents = normalizeWords(documents, 'Style', 'lemma');

% Erase punctuation.
documents = erasePunctuation(documents);

% Remove a list of stop words.
documents = removeStopWords(documents);

% Remove words with 2 or fewer characters, and words with 15 or greater
% characters.
documents = removeShortWords(documents, 2);
documents = removeLongWords(documents, 15);
```

end

## See Also

`tokenizedDocument` | `fitlda` | `ldaModel` | `wordcloud`

## Related Examples

- “Analyze Text Data Using Topic Models” on page 2-13
- “Visualize Document Clusters Using LDA Model” on page 2-35
- “Visualize LDA Topic Correlations” on page 2-38

## Create Co-occurrence Network

This example shows how to create a co-occurrence network using a bag-of-words model.

Given a corpus of documents, a co-occurrence network is an undirected graph, with nodes corresponding to unique words in a vocabulary and edges corresponding to the frequency of words co-occurring in a document. Use co-occurrence networks to visualize and extract information of the relationships between words in a corpus of documents. For example, you can use a co-occurrence network to discover which words commonly appear with a specified word.

### Import Text Data

Extract the text data in the file `weekendUpdates.xlsx` using `readtable`. The file `weekendUpdates.xlsx` contains status updates containing the hashtags `"#weekend"` and `"#vacation"`. Read the data using the `readtable` function and extract the text data from the `TextData` column.

```
filename = "weekendUpdates.xlsx";
tbl = readtable(filename, 'TextType', 'string');
textData = tbl.TextData;
```

View the first few observations.

```
textData(1:5)
```

```
ans = 5x1 string
    "Happy anniversary! ♥ Next stop: Paris! → #vacation"
    "Haha, BBQ on the beach, engage smug mode! ☺☺☺♥ ☺☺#vacation"
    "getting ready for Saturday night ☺☺#yum #weekend ☺☺"
    "Say it with me - I NEED A #VACATION!!! ☺"
    "☺☺Chilling ☺☺at home for the first time in ages...This is the life! ☺☺#weekend"
```

### Preprocess Text Data

Tokenize the text, convert it to lowercase, and remove the stop words.

```
documents = tokenizedDocument(textData);

documents = lower(documents);
documents = removeStopWords(documents);
```

Create a matrix of word counts using a bag-of-words model.

```
bag = bagOfWords(documents);
counts = bag.Counts;
```

To compute the word co-occurrences, multiply the word-count matrix by its transpose.

```
cooccurrence = counts.'*counts;
```

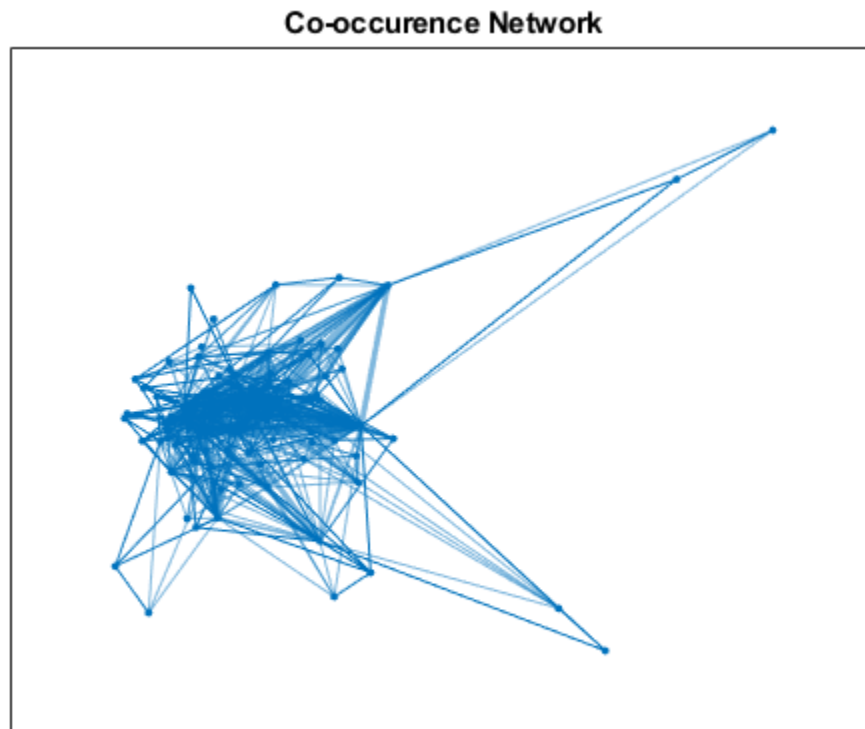
Convert the co-occurrence matrix to a network using the `graph` function.

```
G = graph(cooccurrence, bag.Vocabulary, 'omitselfloops');
```

Visualize the network using the `plot` function. Set the line thickness to a multiple of the edge weight.

```
LWidths = 5*G.Edges.Weight/max(G.Edges.Weight);
```

```
plot(G, 'LineWidth', LWidths)
title("Co-occurrence Network")
```



Find neighbors of the word "great" using the neighbors function.

```
word = "great"
```

```
word =
"great"
```

```
idx = find(bag.Vocabulary == word);
nbrs = neighbors(G,idx);
bag.Vocabulary(nbrs)'
```

```
ans = 18x1 string
"next"
"#vacation"
"[]"
"#weekend"
"@"
"excited"
"flight"
"delayed"
"stuck"
"airport"
"way"
"spend"
"[]"
```

```

"lovely"
"friends"
"_"
"mini"
"everybody"

```

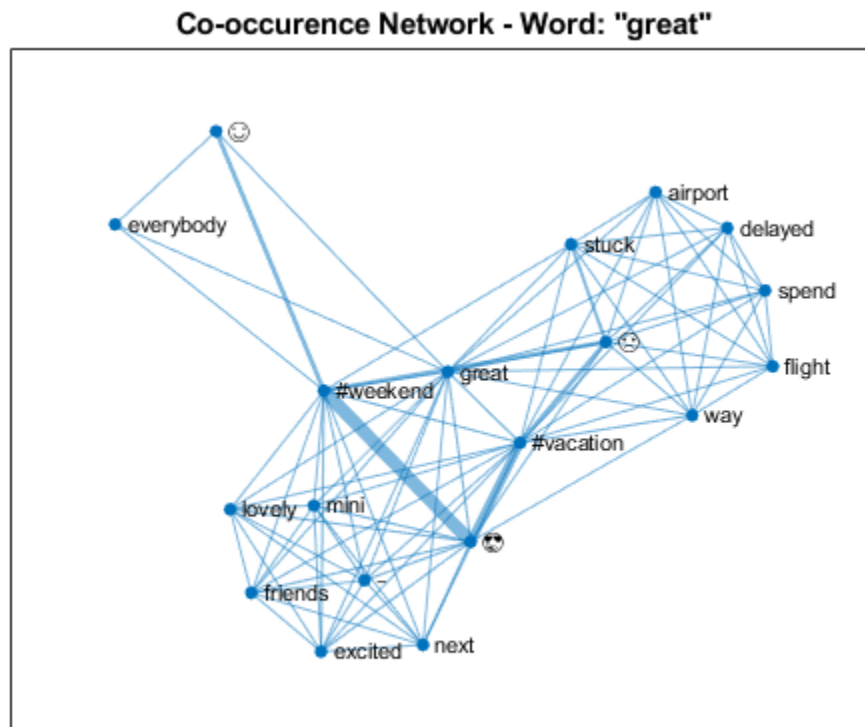
Visualize the co-occurrences of the word "great" by extracting a subgraph of this word and its neighbors.

```

H = subgraph(G,[idx; nbrs]);

LWidths = 5*H.Edges.Weight/max(H.Edges.Weight);
plot(H,'LineWidth',LWidths)
title("Co-occurrence Network - Word: " + word + "");

```



For more information about graphs and network analysis, see “Graph and Network Algorithms”.

## See Also

tokenizedDocument | bagOfWords | removeStopWords | graph

## Related Examples

- “Analyze Text Data Using Topic Models” on page 2-13
- “Analyze Text Data Using Multiword Phrases” on page 2-7



- “Analyze Text Data Containing Emojis” on page 2-52

## Analyze Text Data Containing Emojis

This example shows how to analyze text data containing emojis.

Emojis are pictorial symbols that appear inline in text. When writing text on mobile devices such as smartphones and tablets, people use emojis to keep the text short and convey emotion and feelings.

You also can use emojis to analyze text data. For example, use them to identify relevant strings of text or to visualize the sentiment or emotion of the text.

When working with text data, emojis can behave unpredictably. Depending on your system fonts, your system might not display some emojis correctly. Therefore, if an emoji is not displayed correctly, then the data is not necessarily missing. Your system might be unable to display the emoji in the current font.

### Composing Emojis

In most cases, you can read emojis from a file (for example, by using `extractFileText`, `extractHTMLText`, or `readtable`) or by copying and pasting them directly into MATLAB®. Otherwise, you must compose the emoji using Unicode UTF16 code units.

Some emojis consist of multiple Unicode UTF16 code units. For example, the "smiling face with sunglasses" emoji (☺️ with code point U+1F60E) is a single glyph but comprises two UTF16 code units "D83D" and "DE0E". Create a string containing this emoji using the `compose` function, and specify the two code units with the prefix `"\x"`.

```
emoji = compose("\xD83D\xDE0E")
```

```
emoji =  
"☺️"
```

First get the Unicode UTF16 code units of an emoji. Use `char` to get the numeric representation of the emoji, and then use `dec2hex` to get the corresponding hex value.

```
codeUnits = dec2hex(char(emoji))
```

```
codeUnits = 2×4 char array  
    'D83D'  
    'DE0E'
```

Reconstruct the composition string using the `strjoin` function with the empty delimiter `" "`.

```
formatSpec = strjoin("\x" + codeUnits, " ")
```

```
formatSpec =  
"\xD83D\xDE0E"
```

```
emoji = compose(formatSpec)
```

```
emoji =  
"☺️"
```

### Import Text Data

Extract the text data in the file `weekendUpdates.xlsx` using `readtable`. The file `weekendUpdates.xlsx` contains status updates containing the hashtags `"#weekend"` and `"#vacation"`.

```
filename = "weekendUpdates.xlsx";
tbl = readtable(filename, 'TextType', 'string');
head(tbl)
```

```
ans=8x2 table
```

ID	TextData
1	"Happy anniversary! ♥ Next stop: Paris! → #vacation"
2	"Haha, BBQ on the beach, engage smug mode! 🍷♥️ 🍷#vacation"
3	"getting ready for Saturday night 🍷#yum #weekend 🍷"
4	"Say it with me - I NEED A #VACATION!!! ☺️"
5	"🍷Chilling 🍷at home for the first time in ages...This is the life! 🍷#weekend"
6	"My last #weekend before the exam 🍷🍷"
7	"can't believe my #vacation is over 🍷so unfair"
8	"Can't wait for tennis this #weekend 🍷🍷🍷"

Extract the text data from the field TextData and view the first few status updates.

```
textData = tbl.TextData;
textData(1:5)
```

```
ans = 5x1 string
```

```
"Happy anniversary! ♥ Next stop: Paris! → #vacation"
"Haha, BBQ on the beach, engage smug mode! 🍷♥️ 🍷#vacation"
"getting ready for Saturday night 🍷#yum #weekend 🍷"
"Say it with me - I NEED A #VACATION!!! ☺️"
"🍷Chilling 🍷at home for the first time in ages...This is the life! 🍷#weekend"
```

Visualize the text data in a word cloud.

```
figure
wordcloud(textData);
```





## Extract and Visualize Emojis

Visualize all the emojis in text data using a word cloud.

Extract the emojis. First tokenize the text using `tokenizedDocument`, and then view the first few documents.

```
documents = tokenizedDocument(textData);
documents(1:5)
```

```
ans =
    5×1 tokenizedDocument:
```

```
11 tokens: Happy anniversary ! ♥ Next stop : Paris ! → #vacation
16 tokens: Haha , BBQ on the beach , engage smug mode ! 🍷🍷♥ 🍷🍷#vacation
 9 tokens: getting ready for Saturday night 🍷🍷#yum #weekend 🍷🍷
13 tokens: Say it with me - I NEED A #VACATION ! ! ! ☺
19 tokens: 🍷🍷Chilling 🍷🍷at home for the first time in ages ... This is the life ! 🍷🍷#weekend
```

The `tokenizedDocument` function automatically detects emoji and assigns the token type "emoji". View the first few token details of the documents using the `tokenDetails` function.

```
tdetails = tokenDetails(documents);
head(tdetails)
```

```
ans=8×5 table
      Token      DocumentNumber      LineNumber      Type      Language
```



## **Related Examples**

- “Analyze Text Data Using Topic Models” on page 2-13
- “Analyze Text Data Using Multiword Phrases” on page 2-7
- “Train a Sentiment Classifier” on page 2-71
- “Classify Text Data Using Deep Learning” on page 2-85
- “Generate Text Using Deep Learning” (Deep Learning Toolbox)

## Analyze Sentiment in Text

This example shows how to use the Valence Aware Dictionary and sEntiment Reasoner (VADER) algorithm for sentiment analysis.

The VADER algorithm uses a list of annotated words (the sentiment lexicon), where each word has a corresponding sentiment score. The VADER algorithm also utilizes word lists that modify the scores of proceeding words in the text:

- Boosters - words or n-grams that boost the sentiment of proceeding tokens. For example, words like "absolutely" and "amazingly".
- Dampeners - words or n-grams that dampen the sentiment of proceeding tokens. For example, words like "hardly" and "somewhat".
- Negations - words that negate the sentiment of proceeding tokens. For example, words like "not" and "isn't".

To evaluate sentiment in text, use the `vaderSentimentScores` function.

### Load Data

Extract the text data in the file `weekendUpdates.xlsx` using `readtable`. The file `weekendUpdates.xlsx` contains status updates containing the hashtags "#weekend" and "#vacation".

```
filename = "weekendUpdates.xlsx";
tbl = readtable(filename, 'TextType', 'string');
head(tbl)
```

ans=8×2 table

ID	TextData
1	"Happy anniversary! ♥ Next stop: Paris! → #vacation"
2	"Haha, BBQ on the beach, engage smug mode! ☺☺☺♥ ☺☺#vacation"
3	"getting ready for Saturday night ☺☺#yum #weekend ☺☺"
4	"Say it with me - I NEED A #VACATION!!! ☺"
5	"☺☺Chilling ☺☺at home for the first time in ages...This is the life! ☺☺#weekend"
6	"My last #weekend before the exam ☺☺☺☺"
7	"can't believe my #vacation is over ☺☺so unfair"
8	"Can't wait for tennis this #weekend ☺☺☺☺☺"

Create an array of tokenized documents from the text data and view the first few documents.

```
str = tbl.TextData;
documents = tokenizedDocument(str);
documents(1:5)
```

ans =

5x1 tokenizedDocument:

```
11 tokens: Happy anniversary ! ♥ Next stop : Paris ! → #vacation
16 tokens: Haha , BBQ on the beach , engage smug mode ! ☺☺☺♥ ☺☺#vacation
9 tokens: getting ready for Saturday night ☺☺#yum #weekend ☺☺
13 tokens: Say it with me - I NEED A #VACATION ! ! ! ☺
```



```
19 tokens: []Chilling []at home for the first time in ages ... This is the life ! []#weekend
```

## Evaluate Sentiment

Evaluate the sentiment of the tokenized documents using the `vaderSentimentLexicon` function. Scores close to 1 indicate positive sentiment, scores close to -1 indicate negative sentiment, and scores close to 0 indicate neutral sentiment.

```
compoundScores = vaderSentimentScores(documents);
```

View the scores of the first few documents.

```
compoundScores(1:5)
```

```
ans = 5×1  
  
    0.4738  
    0.9348  
    0.6705  
   -0.5067  
    0.7345
```

Visualize the text with positive and negative sentiment in word clouds.

```
idx = compoundScores > 0;  
strPositive = str(idx);  
strNegative = str(~idx);  
  
figure  
subplot(1,2,1)  
wordcloud(strPositive);  
title("Positive Sentiment")  
  
subplot(1,2,2)  
wordcloud(strNegative);  
title("Negative Sentiment")
```

**Positive Sentiment**



**Negative Sentiment**



**See Also**

`vaderSentimentScores` | `ratioSentimentScores` | `tokenizedDocument`

**More About**

- "Generate Domain Specific Sentiment Lexicon" on page 2-61
- "Train a Sentiment Classifier" on page 2-71
- "Prepare Text Data for Analysis" on page 1-11
- "Analyze Text Data Containing Emojis" on page 2-52
- "Create Simple Text Model for Classification" on page 2-2
- "Analyze Text Data Using Topic Models" on page 2-13
- "Analyze Text Data Using Multiword Phrases" on page 2-7

## Generate Domain Specific Sentiment Lexicon

This example shows how to generate a lexicon for sentiment analysis using 10-K and 10-Q financial reports.

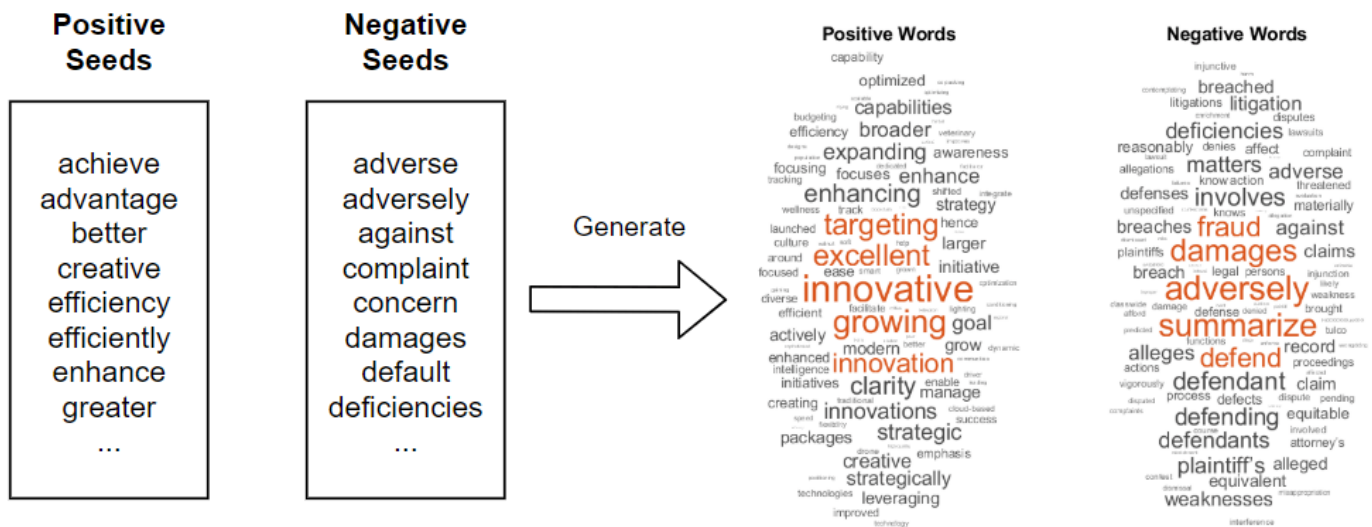
Sentiment analysis allows you to automatically summarize the sentiment in a given piece of text. For example, assign the pieces of text "This company is showing strong growth." and "This other company is accused of misleading consumers." with positive and negative sentiment, respectively. Also, for example, to assign the text "This company is showing *extremely* strong growth." a stronger sentiment score than the text "This company is showing strong growth."

Sentiment analysis algorithms such as VADER rely on annotated lists of words called sentiment lexicons. For example, VADER uses a sentiment lexicon with words annotated with a sentiment score ranging from -1 to 1, where scores close to 1 indicate strong positive sentiment, scores close to -1 indicate strong negative sentiment, and scores close to zero indicate neutral sentiment.

To analyze the sentiment of text using the VADER algorithm, use the `vaderSentimentScores` function. If the sentiment lexicon used by the `vaderSentimentScores` function does not suit the data you are analyzing, for example, if you have a domain-specific data set like medical or engineering data, then you can generate your own custom sentiment lexicon using a small set of seed words.

This example shows how to generate a sentiment lexicon given a collection of seed words using a graph-based approach based on [1 on page 2-0 ]:

- Train a word embedding that models the similarity between words using the training data.
- Create a simplified graph representing the embedding with nodes corresponding to words and edges weighted by similarity.
- To determine words with strong polarity, identify the words connected to multiple seed words through short but heavily weighted paths.



### Load Data

Download the 10-K and 10-Q financial reports data from Securities and Exchange Commission (SEC) via the Electronic Data Gathering, Analysis, and Retrieval (EDGAR) API [2 on page 2-0 ] using the

`financeReports` helper function attached to this example as a supporting file. To access this file, open this example as a Live Script. The `financeReports` function downloads 10-K and 10-Q reports for the specified year, quarter, and maximum character length.

Download a set of 20,000 reports from the fourth quarter of 2019. Depending on the sizes of the reports, this can take some time to run.

```
year = 2019;  
qtr = 4;  
textData = financeReports(year,qtr,'MaxNumReports',20000);
```

```
Downloading 10-K and 10-Q reports...  
Done.  
Elapsed time is 1799.718710 seconds.
```

Define sets of positive and negative seed words to use with this data. The seed words must appear at least once in the text data, otherwise they are ignored.

```
seedsPositive = ["achieve" "advantage" "better" "creative" "efficiency" ...  
"efficiently" "enhance" "greater" "improved" "improving" ...  
"innovation" "innovations" "innovative" "opportunities" "profitable" ...  
"profitably" "strength" "strengthen" "strong" "success"]';  
  
seedsNegative = ["adverse" "adversely" "against" "complaint" "concern" ...  
"damages" "default" "deficiencies" "disclosed" "failure" ...  
"fraud" "impairment" "litigation" "losses" "misleading" ...  
"omit" "restated" "restructuring" "termination" "weaknesses"]';
```

### Prepare Text Data

Create a function names `preprocessText` that prepares the text data for analysis. The `preprocessText` function, listed at the end of the example performs the following steps:

- Erase any URLs.
- Tokenize the text.
- Remove tokens containing digits.
- Convert the text to lower case.
- Remove any words with two or fewer characters.
- Remove any stop words.

Preprocess the text using the `preprocessText` function. Depending on the size of the text data, this can take some time to run.

```
documents = preprocessText(textData);
```

Visualize the preprocessed text data in a word cloud.

```
figure  
wordcloud(documents);
```



```
[nearestWords,dist] = vec2word(emb,wordVectors,numNeighbors);
```

To create the graph, use the `graph` function and specify pairwise source and target nodes, and specify their edge weights.

Define the source and target nodes.

```
sourceNodes = repelem(vocabulary,numNeighbors);  
targetNodes = reshape(nearestWords,1,[]);
```

Calculate the edge weights.

```
edgeWeights = reshape(dist,1,[]);
```

Create a graph connecting each word with its neighbors with edge weights corresponding to the similarity scores.

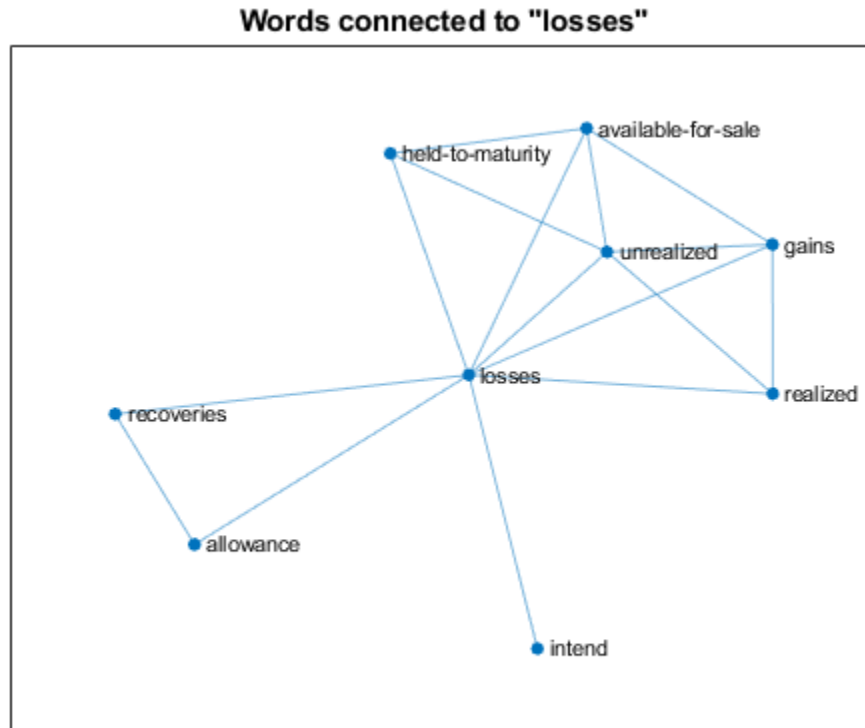
```
wordGraph = graph(sourceNodes,targetNodes,edgeWeights,vocabulary);
```

Remove the repeated edges using the `simplify` function.

```
wordGraph = simplify(wordGraph);
```

Visualize the section of the word graph connected to the word "losses".

```
word = "losses";  
idx = findnode(wordGraph,word);  
nbrs = neighbors(wordGraph,idx);  
wordSubgraph = subgraph(wordGraph,[idx; nbrs]);  
figure  
plot(wordSubgraph)  
title("Words connected to "" + word + "")")
```



### Generate Sentiment Scores

To determine words with strong polarity, identify the words connected to multiple seed words through short but heavily weighted paths.

Initialize an array of sentiment scores corresponding to each word in the vocabulary.

```
sentimentScores = zeros([1 numel(vocabulary)]);
```

Iteratively traverse the graph and update the sentiment scores.

Traverse the graph at different depths. For each depth, calculate the positive and negative polarity of the words by using the positive and negative seeds to propagate sentiment to the rest of the graph.

For each depth:

- Calculate the positive and negative polarity scores.
- Account for the difference in overall mass of positive and negative flow in the graph.
- For each node-word, normalize the difference of its two scores.

After running the algorithm, if a phrase has a higher positive than negative polarity score, then its final polarity will be positive, and negative otherwise.

Specify a maximum path length of 4.

```
maxPathLength = 4;
```

Iteratively traverse the graph and calculate the sum of the sentiment scores.

```
for depth = 1:maxPathLength

    % Calculate polarity scores.
    polarityPositive = polarityScores(seedsPositive,vocabulary,wordGraph,depth);
    polarityNegative = polarityScores(seedsNegative,vocabulary,wordGraph,depth);

    % Account for difference in overall mass of positive and negative flow
    % in the graph.
    b = sum(polarityPositive) / sum(polarityNegative);

    % Calculate new sentiment scores.
    sentimentScoresNew = polarityPositive - b * polarityNegative;
    sentimentScoresNew = normalize(sentimentScoresNew,'range',[-1,1]);

    % Add scores to sum.
    sentimentScores = sentimentScores + sentimentScoresNew;
end
```

Normalize the sentiment scores by the number of iterations.

```
sentimentScores = sentimentScores / maxPathLength;
```

Create a table containing the vocabulary and the corresponding sentiment scores.

```
tbl = table;
tbl.Token = vocabulary';
tbl.SentimentScore = sentimentScores';
```

To remove tokens with neutral sentiment from the lexicon, remove the tokens with sentiment score that have absolute value less than a threshold of 0.1.

```
thr = 0.1;
idx = abs(tbl.SentimentScore) < thr;
tbl(idx,:) = [];
```

Sort the table rows by descending sentiment score and view the first few rows.

```
tbl = sortrows(tbl, 'SentimentScore', 'descend');
head(tbl)
```

```
ans=8x2 table
      Token      SentimentScore
-----
"opportunities"    0.95633
"innovative"       0.89635
"success"          0.84362
"focused"          0.83768
"strong"           0.81042
"capabilities"     0.79174
"innovation"       0.77698
"improved"         0.77176
```

You can use this table as a custom sentiment lexicon for the `vaderSentimentScores` function.



Visualize the sentiment lexicon in word clouds. Display tokens with a positive score in one word cloud and tokens with negative scores in another. Display the words with sizes given by the absolute value their corresponding sentiment score.

```
figure
subplot(1,2,1);
idx = tbl.SentimentScore > 0;
tblPositive = tbl(idx,:);
wordcloud(tblPositive, 'Token', 'SentimentScore')
title('Positive Words')

subplot(1,2,2);
idx = tbl.SentimentScore < 0;
tblNegative = tbl(idx,:);
tblNegative.SentimentScore = abs(tblNegative.SentimentScore);
wordcloud(tblNegative, 'Token', 'SentimentScore')
title('Negative Words')
```



Export the table to a CSV file.

```
filename = "financeSentimentLexicon.csv";
writetable(tbl, filename)
```

### Analyze Sentiment in Text

To analyze the sentiment in for previously unseen text data, preprocess the text using the same preprocessing steps and use the vaderSentimentScores function.

Create a string array containing the text data and preprocess it using the `preprocessText` function.

```
textDataNew = [  
    "This innovative company is continually showing strong growth."  
    "This other company is accused of misleading consumers."];  
documentsNew = preprocessText(textDataNew);
```

Evaluate the sentiment using the `vaderSentimentScores` function. Specify the sentiment lexicon created in this example using the `'SentimentLexicon'` option.

```
compoundScores = vaderSentimentScores(documentsNew, 'SentimentLexicon', tbl)  
  
compoundScores = 2×1  
  
    0.4360  
   -0.1112
```

Positive and negative scores indicate positive and negative sentiment, respectively. The magnitude of the value corresponds to the strength of the sentiment.

## Supporting Functions

### Text Preprocessing Function

The `preprocessText` function performs the following steps:

- Erase any URLs.
- Tokenize the text.
- Remove tokens containing digits.
- Convert the text to lower case.
- Remove any words with two or fewer characters.
- Remove any stop words.

```
function documents = preprocessText(textData)  
  
% Erase URLs.  
textData = eraseURLs(textData);  
  
% Tokenize.  
documents = tokenizedDocument(textData);  
  
% Remove tokens containing digits.  
pat = textBoundary + wildcardPattern + digitsPattern + wildcardPattern + textBoundary;  
documents = replace(documents, pat, "");  
  
% Convert to lowercase.  
documents = lower(documents);  
  
% Remove short words.  
documents = removeShortWords(documents, 2);  
  
% Remove stop words.  
documents = removeStopWords(documents);  
  
end
```

## Polarity Scores Function

The `polarityScores` function returns a vector of polarity scores given a set of seed words, vocabulary, graph, and a specified depth. The function computes the sum over the maximum weighted path from every seed word to each node in the vocabulary. A high polarity score indicates phrases connected to multiple seed words via both short and strongly weighted paths.

The function performs the following steps:

- Initialize the scores of the seeds with ones and otherwise zeros.
- Loop over the seeds. For each seed, iteratively traverse the graph at different depth levels. For the first iteration, set the search space to the immediate neighbors of the seed.
- For each depth level, loop over the nodes in the search space and identify its neighbors in the graph.
- Loop over its neighbors and update the corresponding scores. The updated score is the maximum value of the current score for the seed and neighbor, and the score for the seed and search node weighted by the corresponding graph edge.
- At the end of the search for the depth level, append the neighbors to the search space. This increases the depth of the search for the next iteration.

The output polarity is the sum of the scores connected to the input seeds.

```
function polarity = polarityScores(seeds,vocabulary,wordGraph,depth)

% Remove seeds missing from vocabulary.
idx = ~ismember(seeds,vocabulary);
seeds(idx) = [];

% Initialize scores.
vocabularySize = numel(vocabulary);
scores = zeros(vocabularySize);
idx = ismember(vocabulary,seeds);
scores(idx,idx) = eye(numel(seeds));

% Loop over seeds.
for i = 1:numel(seeds)

    % Initialize search space.
    seed = seeds(i);
    idxSeed = vocabulary == seed;
    searchSpace = find(idxSeed);

    % Search at different depths.
    for d = 1:depth

        % Loop over nodes in search space.
        numNodes = numel(searchSpace);

        for k = 1:numNodes

            idxNew = searchSpace(k);

            % Find neighbors and weights.
            nbrs = neighbors(wordGraph,idxNew);
            idxWeights = findedge(wordGraph,idxNew,nbrs);
```

```
weights = wordGraph.Edges.Weight(idxWeights);

% Loop over neighbors.
for j = 1:numel(nbrs)

    % Calculate scores.
    score = scores(idxSeed,nbrs(j));
    scoreNew = scores(idxSeed,idxNew);

    % Update score.
    scores(idxSeed,nbrs(j)) = max(score,scoreNew*weights(j));
end

% Appended nodes to search space for next depth iteration.
searchSpace = [searchSpace nbrs'];
end
end
end

% Find seeds in vocabulary.
[~,idx] = ismember(seeds,vocabulary);

% Sum scores connected to seeds.
polarity = sum(scores(idx,:));

end
```

### **Bibliography**

- 1 Velikovich, Lenid. "The Viability of Web-derived Polarity Lexicons." In *Proceedings of The Annual Conference of the North American Chapter of the Association for Computational Linguistics, 2010*, pp. 777-785. 2010.
- 2 Accessing EDGAR Data. <https://www.sec.gov/os/accessing-edgar-data>

### **See Also**

## Train a Sentiment Classifier

This example shows how to train a classifier for sentiment analysis using an annotated list of positive and negative sentiment words and a pretrained word embedding.

The pretrained word embedding plays several roles in this workflow. It converts words into numeric vectors and forms the basis for a classifier. You can then use the classifier to predict the sentiment of other words using their vector representation, and use these classifications to calculate the sentiment of a piece of text. There are four steps in training and using the sentiment classifier:

- Load a pretrained word embedding.
- Load an opinion lexicon listing positive and negative words.
- Train a sentiment classifier using the word vectors of the positive and negative words.
- Calculate the mean sentiment scores of the words in a piece of text.

### Load Pretrained Word Embedding

Word embeddings map words in a vocabulary to numeric vectors. These embeddings can capture semantic details of the words so that similar words have similar vectors. They also model relationships between words through vector arithmetic. For example, the relationship *Rome is to Paris as Italy is to France* is described by the equation  $Rome - Italy + France \approx Paris$ .

Load a pretrained word embedding using the `fastTextWordEmbedding` function. This function requires Text Analytics Toolbox™ Model for *fastText English 16 Billion Token Word Embedding* support package. If this support package is not installed, then the function provides a download link.

```
emb = fastTextWordEmbedding;
```

### Load Opinion Lexicon

Load the positive and negative words from the opinion lexicon (also known as a sentiment lexicon) from <https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html>. [1] First, extract the files from the `.rar` file into a folder named `opinion-lexicon-English`, and then import the text.

Load the data using the function `readLexicon` listed at the end of this example. The output `data` is a table with variables `Word` containing the words, and `Label` containing a categorical sentiment label, `Positive` or `Negative`.

```
data = readLexicon;
```

View the first few words labeled as positive.

```
idx = data.Label == "Positive";
head(data(idx,:))
```

```
ans=8x2 table
      Word      Label
-----
"a+"      Positive
"abound"   Positive
"abounds"  Positive
"abundance" Positive
"abundant" Positive
"accessible" Positive
```

```
"accessible" Positive
"acclaim"    Positive
```

View the first few words labeled as negative.

```
idx = data.Label == "Negative";
head(data(idx,:))
```

```
ans=8x2 table
      Word      Label
-----
"2-faced"    Negative
"2-faces"    Negative
"abnormal"   Negative
"abolish"    Negative
"abominable" Negative
"abominably" Negative
"abominate"  Negative
"abomination" Negative
```

### Prepare Data for Training

To train the sentiment classifier, convert the words to word vectors using the pretrained word embedding `emb`. First remove the words that do not appear in the word embedding `emb`.

```
idx = ~isVocabularyWord(emb,data.Word);
data(idx,:) = [];
```

Set aside 10% of the words at random for testing.

```
numWords = size(data,1);
cvp = cvpartition(numWords,'HoldOut',0.1);
dataTrain = data(training(cvp),:);
dataTest = data(test(cvp),:);
```

Convert the words in the training data to word vectors using `word2vec`.

```
wordsTrain = dataTrain.Word;
XTrain = word2vec(emb,wordsTrain);
YTrain = dataTrain.Label;
```

### Train Sentiment Classifier

Train a support vector machine (SVM) classifier which classifies word vectors into positive and negative categories.

```
mdl = fitcsvm(XTrain,YTrain);
```

### Test Classifier

Convert the words in the test data to word vectors using `word2vec`.

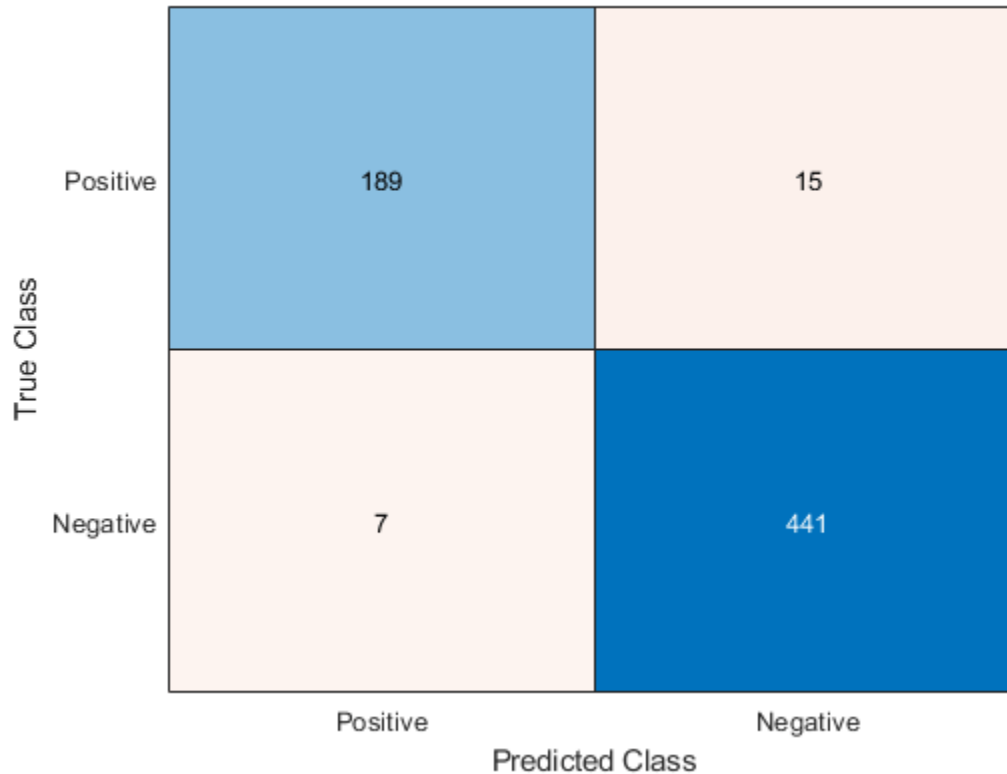
```
wordsTest = dataTest.Word;
XTest = word2vec(emb,wordsTest);
YTest = dataTest.Label;
```

Predict the sentiment labels of the test word vectors.

```
[YPred,scores] = predict mdl,XTest;
```

Visualize the classification accuracy in a confusion matrix.

```
figure
confusionchart(YTest,YPred);
```



Visualize the classifications in word clouds. Plot the words with positive and negative sentiments in word clouds with word sizes corresponding to the prediction scores.

```
figure
subplot(1,2,1)
idx = YPred == "Positive";
wordcloud(wordsTest(idx),scores(idx,1));
title("Predicted Positive Sentiment")

subplot(1,2,2)
wordcloud(wordsTest(~idx),scores(~idx,2));
title("Predicted Negative Sentiment")
```

## Predicted Positive Sentiment



## Predicted Negative Sentiment



## Calculate Sentiment of Collections of Text

To calculate the sentiment of a piece of text, for example an update on social media, predict the sentiment score of each word in the text and take the mean sentiment score.

```
filename = "weekendUpdates.xlsx";
tbl = readtable(filename, 'TextType', 'string');
textData = tbl.TextData;
textData(1:10)
```

```
ans = 10x1 string array
"Happy anniversary! ♥ Next stop: Paris! → #vacation"
"Haha, BBQ on the beach, engage smug mode! ☺☺☺♥ ☺☺#vacation"
"getting ready for Saturday night ☺☺#yum #weekend ☺☺"
"Say it with me - I NEED A #VACATION!!! ☺"
"☺☺Chilling ☺☺at home for the first time in ages...This is the life! ☺☺#weekend"
"My last #weekend before the exam ☺☺☺"
"can't believe my #vacation is over ☺☺so unfair"
"Can't wait for tennis this #weekend ☺☺☺☺☺"
"I had so much fun! ☺☺☺☺☺st trip EVER! ☺☺☺☺#vacation #weekend"
"Hot weather and air con broke in car ☺☺#sweaty #roadtrip #vacation"
```

Create a function which tokenizes and preprocesses the text data so it can be used for analysis. The function `preprocessText`, listed at the end of the example, performs the following steps in order:

- 1 Tokenize the text using `tokenizedDocument`.



- 2 Erase punctuation using `erasePunctuation`.
- 3 Remove stop words (such as "and", "of", and "the") using `removeStopWords`.
- 4 Convert to lowercase using `lower`.

Use the preprocessing function `preprocessText` to prepare the text data. This step can take a few minutes to run.

```
documents = preprocessText(textData);
```

Remove the words from the documents that do not appear in the word embedding `emb`.

```
idx = ~isVocabularyWord(emb,documents.Vocabulary);
documents = removeWords(documents,idx);
```

To visualize how well the sentiment classifier generalizes to the new text, classify the sentiments on the words that occur in the text, but not in the training data and visualize them in word clouds. Use the word clouds to manually check that the classifier behaves as expected.

```
words = documents.Vocabulary;
words(ismember(words,wordsTrain)) = [];
```

```
vec = word2vec(emb,words);
[YPred,scores] = predict mdl,vec);
```

```
figure
subplot(1,2,1)
idx = YPred == "Positive";
wordcloud(words(idx),scores(idx,1));
title("Predicted Positive Sentiment")
```

```
subplot(1,2,2)
wordcloud(words(~idx),scores(~idx,2));
title("Predicted Negative Sentiment")
```

**Predicted Positive Sentiment**



**Predicted Negative Sentiment**



To calculate the sentiment of a given piece of text, compute the sentiment score for each word in the text and calculate the mean sentiment score.

Calculate the mean sentiment score of the updates. For each document, convert the words to word vectors, predict the sentiment score on the word vectors, transform the scores using the score-to-posterior transform function and then calculate the mean sentiment score.

```
for i = 1:numel(documents)
    words = string(documents(i));
    vec = word2vec(emb,words);
    [~,scores] = predict mdl,vec;
    sentimentScore(i) = mean(scores(:,1));
end
```

View the predicted sentiment scores with the text data. Scores greater than 0 correspond to positive sentiment, scores less than 0 correspond to negative sentiment, and scores close to 0 correspond to neutral sentiment.

```
table(sentimentScore', textData)
```

```
ans=50x2 table
    Var1
```

```
textData
```

1.8382	"Happy anniversary! ♥ Next stop: Paris! → #vacation"
1.294	"Haha, BBQ on the beach, engage smug mode! ☺☺☺♥ ☺☺#vacation"
1.0922	"getting ready for Saturday night ☺☺#yum #weekend ☺☺"

```

0.094709 "Say it with me - I NEED A #VACATION!!! ☺"
1.4073 "Chilling at home for the first time in ages...This is the life! #weekend"
-0.8356 "My last #weekend before the exam "
-1.3556 "can't believe my #vacation is over so unfair"
1.4312 "Can't wait for tennis this #weekend "
3.0458 "I had so much fun! Best trip EVER! #vacation #weekend"
-0.39243 "Hot weather and air con broke in car #sweaty #roadtrip #vacation"
0.8028 "Check the out-of-office crew, we are officially ON #VACATION!! "
0.38217 "Well that wasn't how I expected this #weekend to go Total washout!! "
3.03 "So excited for my bestie to visit this #weekend! ♥ "
2.3849 "Who needs a #vacation when the weather is this good * "
-0.0006176 "I love meetings in summer that run into the weekend! Wait that was sarcasm. B
0.52992 "You know we all worked hard for this! We totes deserve this #vacation Ibi.
:

```

### Sentiment Lexicon Reading Function

This function reads the positive and negative words from the sentiment lexicon and returns a table. The table contains variables `Word` and `Label`, where `Label` contains categorical values `Positive` and `Negative` corresponding to the sentiment of each word.

```

function data = readLexicon

% Read positive words
fidPositive = fopen(fullfile('opinion-lexicon-English', 'positive-words.txt'));
C = textscan(fidPositive, '%s', 'CommentStyle', ';');
wordsPositive = string(C{1});

% Read negative words
fidNegative = fopen(fullfile('opinion-lexicon-English', 'negative-words.txt'));
C = textscan(fidNegative, '%s', 'CommentStyle', ';');
wordsNegative = string(C{1});
fclose all;

% Create table of labeled words
words = [wordsPositive; wordsNegative];
labels = categorical(nan(numel(words), 1));
labels(1: numel(wordsPositive)) = "Positive";
labels(numel(wordsPositive)+1: end) = "Negative";

data = table(words, labels, 'VariableNames', {'Word', 'Label'});

end

```

### Preprocessing Function

The function `preprocessText` performs the following steps:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Erase punctuation using `erasePunctuation`.
- 3 Remove stop words (such as "and", "of", and "the") using `removeStopWords`.
- 4 Convert to lowercase using `lower`.

```

function documents = preprocessText(textData)

% Tokenize the text.

```

```
documents = tokenizedDocument(textData);  
  
% Erase punctuation.  
documents = erasePunctuation(documents);  
  
% Remove a list of stop words.  
documents = removeStopWords(documents);  
  
% Convert to lowercase.  
documents = lower(documents);  
  
end
```

### Bibliography

- 1 Hu, Minqing, and Bing Liu. "Mining and summarizing customer reviews." In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 168-177. ACM, 2004.

### See Also

[tokenizedDocument](#) | [bagOfWords](#) | [erasePunctuation](#) | [removeWords](#) | [removeStopWords](#) | [wordcloud](#) | [word2vec](#) | [fastTextWordEmbedding](#)

### Related Examples

- "Analyze Sentiment in Text" on page 2-58
- "Generate Domain Specific Sentiment Lexicon" on page 2-61
- "Create Simple Text Model for Classification" on page 2-2
- "Analyze Text Data Containing Emojis" on page 2-52
- "Analyze Text Data Using Topic Models" on page 2-13
- "Analyze Text Data Using Multiword Phrases" on page 2-7
- "Classify Text Data Using Deep Learning" on page 2-85
- "Generate Text Using Deep Learning" (Deep Learning Toolbox)

## Extract Keywords from Text Data Using RAKE

This example shows how to extract keywords from text data using Rapid Automatic Keyword Extraction (RAKE).

The RAKE algorithm extracts keywords using a delimiter-based approach to identify candidate keywords and scores them using word co-occurrences that appear in the candidate keywords. Keywords can contain multiple tokens. Furthermore, the RAKE algorithm also merges keywords when they appear multiple times, separated by the same merging delimiter.

### Extract Keywords

Create an array of tokenized document containing the text data.

```
textData = [
    "MATLAB provides tools for scientists and engineers. MATLAB is used by scientists and engineers."
    "Analyze text and images. You can import text and images."
    "Analyze text and images. Analyze text, images, and videos in MATLAB."];
documents = tokenizedDocument(textData);
```

Extract the keywords using the `rakeKeywords` function.

```
tbl = rakeKeywords(documents)
```

`tbl=12×3 table`

	Keyword		DocumentNumber	Score
"MATLAB"	"provides"	"tools"	1	8
"MATLAB"	""	""	1	2
"scientists"	"and"	"engineers"	1	2
"engineers"	""	""	1	1
"scientists"	""	""	1	1
"Analyze"	"text"	""	2	4
"import"	"text"	""	2	4
"images"	""	""	2	1
"Analyze"	"text"	""	3	4
"MATLAB"	""	""	3	1
"images"	""	""	3	1
"videos"	""	""	3	1

If a keyword contains multiple words, then the  $i$ th element of the string array corresponds to the  $i$ th word of the keyword. If the keyword has fewer words than the longest keyword, then remaining entries of the string array are the empty string `" "`.

For readability, transform the multi-word keywords into a single string using the `join` and `strip` functions.

```
if size(tbl.Keyword,2) > 1
    tbl.Keyword = strip(join(tbl.Keyword));
end
head(tbl)
```

`ans=8×3 table`

	Keyword	DocumentNumber	Score

"MATLAB provides tools"	1	8
"MATLAB"	1	2
"scientists and engineers"	1	2
"engineers"	1	1
"scientists"	1	1
"Analyze text"	2	4
"import text"	2	4
"images"	2	1

### Specify Maximum Number of Keywords Per Document

The `rakeKeywords` function, by default, returns all identified keywords. To reduce the number of keywords, use the `'MaxNumKeywords'` option.

Extract the top three keywords for each document by setting the `'MaxNumKeywords'` option to 3.

```
tbl = rakeKeywords(documents, 'MaxNumKeywords', 3)
```

tbl=9x3 table

	Keyword		DocumentNumber	Score
"MATLAB"	"provides"	"tools"	1	8
"MATLAB"	"	"	1	2
"scientists"	"and"	"engineers"	1	2
"Analyze"	"text"	"	2	4
"import"	"text"	"	2	4
"images"	"	"	2	1
"Analyze"	"text"	"	3	4
"MATLAB"	"	"	3	1
"images"	"	"	3	1

### Specify Delimiters

Notice that in the extracted keywords above, the function extracts the multi-word keyword "scientists and engineers" from the first document, but does not extract the multi-word keyword "text and images" from the second document. This is because the RAKE algorithm uses tokens appearing between delimiters as candidate keywords, and the algorithm only merges keywords with delimiters when the merged phrase appears multiple times.

In this case, the instances of the token "text" appears within the two different multi-word keyword candidates "Analyze text" and "import text". Because, in this case, the function does not extract "text" as a separate candidate keyword, the algorithm does not consider merging candidates with the delimiter "and" and the candidate keyword "images".

You can specify the delimiters used for extracting keywords using the `'Delimiters'` and `'MergingDelimiters'` options. To specify delimiters that should not appear in extracted keywords, use the `'Delimiters'` option. To specify delimiters that can appear in extracted keywords, use the `'MergingDelimiters'` option.

Extract keywords from the same text as before and also specify the words "Analyze" and "import" as merging delimiters.

```
newDelimiters = ["Analyze" "import"];
mergingDelimiters = [stopWords newDelimiters];
```

```
tbl = rakeKeywords(documents, 'MergingDelimiters', mergingDelimiters)
```

```
tbl=12x3 table
```

	Keyword		DocumentNumber	Score
"MATLAB"	"provides"	"tools"	1	8
"MATLAB"	" "	" "	1	2
"scientists"	"and"	"engineers"	1	2
"engineers"	" "	" "	1	1
"scientists"	" "	" "	1	1
"text"	"and"	"images"	2	2
"images"	" "	" "	2	1
"text"	" "	" "	2	1
"MATLAB"	" "	" "	3	1
"images"	" "	" "	3	1
"text"	" "	" "	3	1
"videos"	" "	" "	3	1

Notice here that the function treats the tokens "text" and "images" as keywords and also extracts the merged keyword "text and images". To learn more about the RAKE algorithm, see "Rapid Automatic Keyword Extraction".

### Alternatives

You can experiment with different keyword extraction algorithms to see what works best with your data. Because the RAKE algorithm uses a delimiter-based approach to extract candidate keywords, the extracted keywords can be very long. Alternatively, you can try extracting keywords using TextRank algorithm which starts with individual tokens as candidate keywords and then merges them when appropriate. To extract keywords using TextRank, use the `textrankKeywords` function. To learn more, see "Extract Keywords from Text Data Using TextRank" on page 2-82.

### References

[1] Rose, Stuart, Dave Engel, Nick Cramer, and Wendy Cowley. "Automatic keyword extraction from individual documents." *Text mining: applications and theory* 1 (2010): 1-20.

### See Also

`tokenizedDocument` | `rakeKeywords` | `textrankKeywords` | `extractSummary`

### More About

- "Extract Keywords from Text Data Using TextRank" on page 2-82

## Extract Keywords from Text Data Using TextRank

This example shows to extract keywords from text data using TextRank.

The TextRank keyword extraction algorithm extracts keywords using a part-of-speech tag-based approach to identify candidate keywords and scores them using word co-occurrences determined by a sliding window. Keywords can contain multiple tokens. Furthermore, the TextRank keyword extraction algorithm also merges keywords when they appear consecutively in a document.

### Extract Keywords

Create an array of tokenized document containing the text data.

```
textData = [
    "MATLAB provides really useful tools for engineers. Scientists use many useful MATLAB toolboxes."
    "MATLAB and Simulink have many features. MATLAB and Simulink makes it easy to develop models."
    "You can easily import data in MATLAB. In particular, you can easily import text data."];
documents = tokenizedDocument(textData);
```

Extract the keywords using the `textrankKeywords` function.

```
tbl = textrankKeywords(documents)
```

tbl=6×3 table

	Keyword		DocumentNumber	Score
"useful"	"MATLAB"	"toolboxes"	1	4.8695
"useful"	" "	" "	1	2.3612
"MATLAB"	" "	" "	1	1.6212
"many"	"features"	" "	2	4.6152
"text"	"data"	" "	3	3.4781
"data"	" "	" "	3	1.7391

If a keyword contains multiple words, then the  $i$ th element of the string array corresponds to the  $i$ th word of the keyword. If the keyword has fewer words than the longest keyword, then remaining entries of the string array are the empty string `" "`.

For readability, transform the multi-word keywords into a single string using the `join` and `strip` functions.

```
if size(tbl.Keyword,2) > 1
    tbl.Keyword = strip(join(tbl.Keyword));
end
head(tbl)
```

ans=6×3 table

	Keyword		DocumentNumber	Score
	"useful MATLAB toolboxes"		1	4.8695
	"useful"		1	2.3612
	"MATLAB"		1	1.6212
	"many features"		2	4.6152
	"text data"		3	3.4781
	"data"		3	1.7391



### Specify Maximum Number of Keywords Per Document

The `textrankKeywords` function, by default, returns all identified keywords. To reduce the number of keywords, use the `'MaxNumKeywords'` option.

Extract the top two keywords for each document by setting the `'MaxNumKeywords'` option to 2.

```
tbl = textrankKeywords(documents, 'MaxNumKeywords', 2)
```

tbl=5x3 table

	Keyword		DocumentNumber	Score
"useful"	"MATLAB"	"toolboxes"	1	4.8695
"useful"	" "	" "	1	2.3612
"many"	"features"	" "	2	4.6152
"text"	"data"	" "	3	3.4781
"data"	" "	" "	3	1.7391

### Specify Part-of-Speech Tags

Notice that in the extracted keywords above, the function does not consider the word "import" as a keyword. This is because the TextRank keyword extraction algorithm, by default, uses tokens with the part-of-speech tags "noun", "proper-noun" and "adjective" as candidate keywords. Because the word "import" is a verb, the algorithm does not consider this as a candidate keyword. Similarly, the algorithm does not consider the adverb "easily" as a candidate keyword.

To specify which part-of-speech tags to use for identifying candidate keywords, use the `'PartOfSpeech'` option.

Extract keywords from the same text as before and also specify also specify the part-of-speech tags "adverb" and "verb".

```
newTags = ["adverb" "verb"];
tags = ["noun" "proper-noun" "adjective" newTags];
tbl = textrankKeywords(documents, 'PartOfSpeech', tags)
```

tbl=7x3 table

	Keyword		DocumentNumber	Score	
"use"	"many"	"useful"	"MATLAB"	1	5.8839
"useful"	" "	" "	" "	1	2.0169
"MATLAB"	" "	" "	" "	1	1.5478
"Simulink"	"have"	"many"	" "	2	4.5058
"Simulink"	" "	" "	" "	2	1.5161
"import"	"text"	"data"	" "	3	4.7921
"import"	"data"	" "	" "	3	3.4195

Notice here that the function treats the token "import" as a candidate keyword and merges it into the multi-word keywords "import data" and "import text data".

### Specify Windows Size

Notice that in the extracted keywords above, that the function does not extract the adverb "easily" as a keyword. This is because of the proximity of these words in the text to other candidate keywords.

The TextRank keyword extraction algorithm scores candidate keywords using the number of pairwise co-occurrences within a sliding window. To increase the window size, use the 'Window' option. Increasing the window size enables the function to find more co-occurrences between keywords which increases the keyword importance scores. This can result in finding more relevant keywords at the cost of potentially over-scoring less relevant keywords.

Extract keywords from the same text as before and also specify also specify a window size of 3.

```
tbl = textrankKeywords(documents, ...
    'PartOfSpeech', tags, ...
    'Window', 3)
```

tbl=8x3 table

Keyword				DocumentNumber	Score
"many"	"useful"	"MATLAB"	" "	1	4.2185
"really"	"useful"	" "	" "	1	2.8851
"MATLAB"	" "	" "	" "	1	1.3154
"Simulink"	" "	" "	" "	2	1.4526
"develop"	" "	" "	" "	2	1.0912
"features"	" "	" "	" "	2	1.0794
"easily"	"import"	"text"	"data"	3	5.2989
"easily"	"import"	"data"	" "	3	4.0842

Notice here that the function treats the tokens "easily" as keywords and merges it into the multi-word keywords "easily import text data" and "easily import data".

To learn more about the TextRank keyword extraction algorithm, see “TextRank Keyword Extraction”.

### Alternatives

You can experiment with different keyword extraction algorithms to see what works best with your data. Because the TextRank keywords algorithm uses a part-of-speech tag-based approach to extract candidate keywords, the extracted keywords can be short. Alternatively, you can try extracting keywords using RAKE algorithm which extracts sequences of tokens appearing between delimiters as candidate keywords. To extract keywords using RAKE, use the `rakeKeywords` function. To learn more, see “Extract Keywords from Text Data Using RAKE” on page 2-79.

### References

- [1] Mihalcea, Rada, and Paul Tarau. "TextRank: Bringing order into text." In *Proceedings of the 2004 conference on empirical methods in natural language processing*, pp. 404-411. 2004.

### See Also

`tokenizedDocument` | `rakeKeywords` | `textrankKeywords` | `extractSummary`

### More About

- “Extract Keywords from Text Data Using RAKE” on page 2-79

## Classify Text Data Using Deep Learning

This example shows how to classify text data using a deep learning long short-term memory (LSTM) network.

Text data is naturally sequential. A piece of text is a sequence of words, which might have dependencies between them. To learn and use long-term dependencies to classify sequence data, use an LSTM neural network. An LSTM network is a type of recurrent neural network (RNN) that can learn long-term dependencies between time steps of sequence data.

To input text to an LSTM network, first convert the text data into numeric sequences. You can achieve this using a word encoding which maps documents to sequences of numeric indices. For better results, also include a word embedding layer in the network. Word embeddings map words in a vocabulary to numeric vectors rather than scalar indices. These embeddings capture semantic details of the words, so that words with similar meanings have similar vectors. They also model relationships between words through vector arithmetic. For example, the relationship "Rome is to Italy as Paris is to France" is described by the equation  $Italy - Rome + Paris = France$ .

There are four steps in training and using the LSTM network in this example:

- Import and preprocess the data.
- Convert the words to numeric sequences using a word encoding.
- Create and train an LSTM network with a word embedding layer.
- Classify new text data using the trained LSTM network.

### Import Data

Import the factory reports data. This data contains labeled textual descriptions of factory events. To import the text data as strings, specify the text type to be 'string'.

```
filename = "factoryReports.csv";
data = readtable(filename, 'TextType', 'string');
head(data)
```

ans=8×5 table

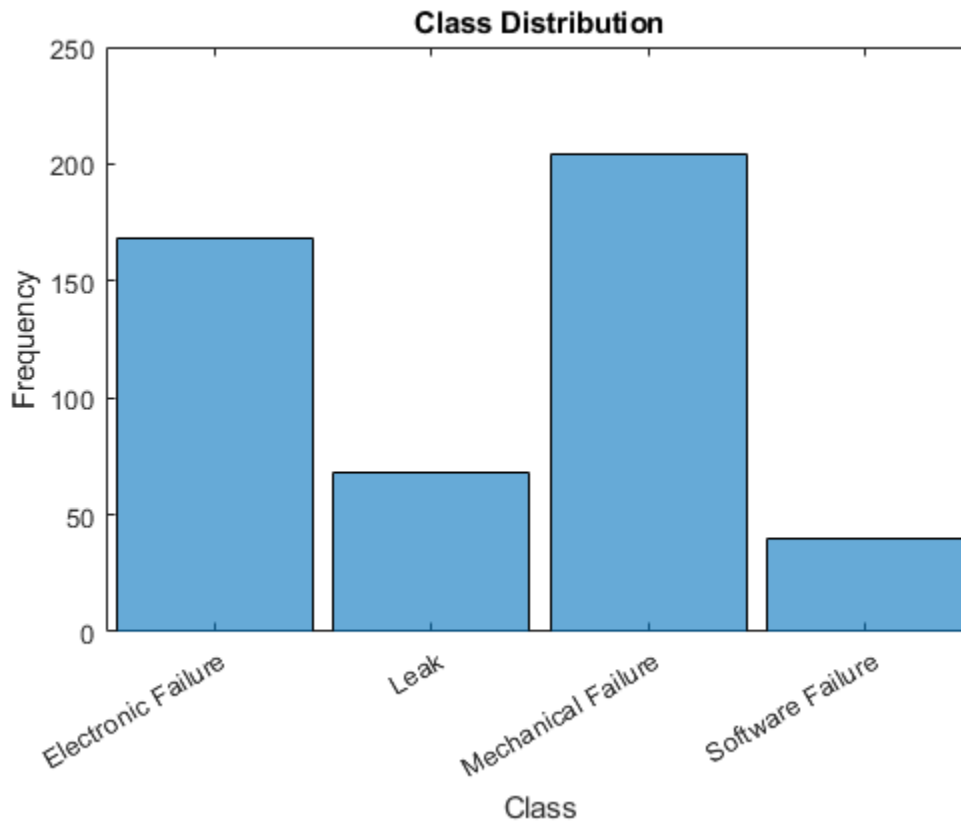
Description	Category
"Items are occasionally getting stuck in the scanner spools."	"Mechanical Failure"
"Loud rattling and banging sounds are coming from assembler pistons."	"Mechanical Failure"
"There are cuts to the power when starting the plant."	"Electronic Failure"
"Fried capacitors in the assembler."	"Electronic Failure"
"Mixer tripped the fuses."	"Electronic Failure"
"Burst pipe in the constructing agent is spraying coolant."	"Leak"
"A fuse is blown in the mixer."	"Electronic Failure"
"Things continue to tumble off of the belt."	"Mechanical Failure"

The goal of this example is to classify events by the label in the Category column. To divide the data into classes, convert these labels to categorical.

```
data.Category = categorical(data.Category);
```

View the distribution of the classes in the data using a histogram.

```
figure
histogram(data.Category);
xlabel("Class")
ylabel("Frequency")
title("Class Distribution")
```



The next step is to partition it into sets for training and validation. Partition the data into a training partition and a held-out partition for validation and testing. Specify the holdout percentage to be 20%.

```
cvp = cvpartition(data.Category, 'Holdout', 0.2);
dataTrain = data(training(cvp),:);
dataValidation = data(test(cvp),:);
```

Extract the text data and labels from the partitioned tables.

```
textDataTrain = dataTrain.Description;
textDataValidation = dataValidation.Description;
YTrain = dataTrain.Category;
YValidation = dataValidation.Category;
```

To check that you have imported the data correctly, visualize the training text data using a word cloud.

```
figure
wordcloud(textDataTrain);
title("Training Data")
```



## Convert Document to Sequences

To input the documents into an LSTM network, use a word encoding to convert the documents into sequences of numeric indices.

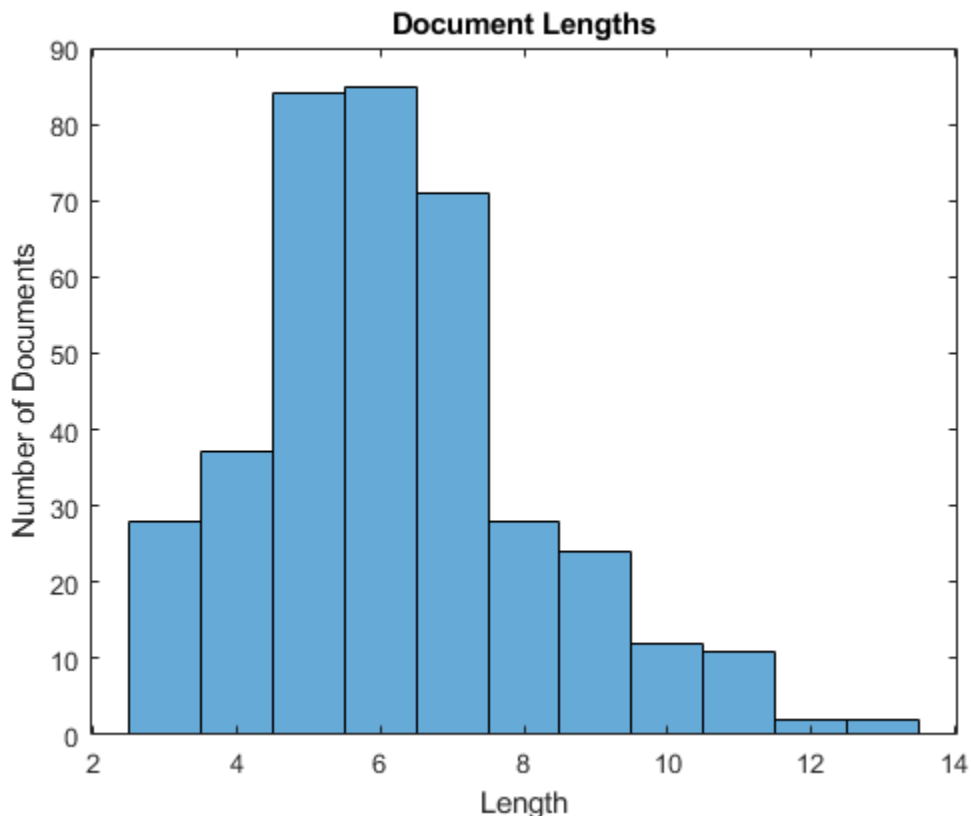
To create a word encoding, use the `wordEncoding` function.

```
enc = wordEncoding(documentsTrain);
```

The next conversion step is to pad and truncate documents so they are all the same length. The `trainingOptions` function provides options to pad and truncate input sequences automatically. However, these options are not well suited for sequences of word vectors. Instead, pad and truncate the sequences manually. If you *left-pad* and truncate the sequences of word vectors, then the training might improve.

To pad and truncate the documents, first choose a target length, and then truncate documents that are longer than it and left-pad documents that are shorter than it. For best results, the target length should be short without discarding large amounts of data. To find a suitable target length, view a histogram of the training document lengths.

```
documentLengths = doclength(documentsTrain);  
figure  
histogram(documentLengths)  
title("Document Lengths")  
xlabel("Length")  
ylabel("Number of Documents")
```



Most of the training documents have fewer than 10 tokens. Use this as your target length for truncation and padding.

Convert the documents to sequences of numeric indices using `doc2sequence`. To truncate or left-pad the sequences to have length 10, set the `'Length'` option to 10.

```
sequenceLength = 10;
XTrain = doc2sequence(enc,documentsTrain,'Length',sequenceLength);
XTrain(1:5)
```

```
ans=5x1 cell array
    {1x10 double}
    {1x10 double}
    {1x10 double}
    {1x10 double}
    {1x10 double}
```

Convert the validation documents to sequences using the same options.

```
XValidation = doc2sequence(enc,documentsValidation,'Length',sequenceLength);
```

### Create and Train LSTM Network

Define the LSTM network architecture. To input sequence data into the network, include a sequence input layer and set the input size to 1. Next, include a word embedding layer of dimension 50 and the same number of words as the word encoding. Next, include an LSTM layer and set the number of hidden units to 80. To use the LSTM layer for a sequence-to-label classification problem, set the output mode to `'last'`. Finally, add a fully connected layer with the same size as the number of classes, a softmax layer, and a classification layer.

```
inputSize = 1;
embeddingDimension = 50;
numHiddenUnits = 80;
```

```
numWords = enc.NumWords;
numClasses = numel(categories(YTrain));
```

```
layers = [ ...
    sequenceInputLayer(inputSize)
    wordEmbeddingLayer(embeddingDimension,numWords)
    lstmLayer(numHiddenUnits,'OutputMode','last')
    fullyConnectedLayer(numClasses)
    softmaxLayer
    classificationLayer]
```

```
layers =
    6x1 Layer array with layers:
```

1	''	Sequence Input	Sequence input with 1 dimensions
2	''	Word Embedding Layer	Word embedding layer with 50 dimensions and 423 unique words
3	''	LSTM	LSTM with 80 hidden units
4	''	Fully Connected	4 fully connected layer
5	''	Softmax	softmax
6	''	Classification Output	crossentropyex

### Specify Training Options

Specify the training options:

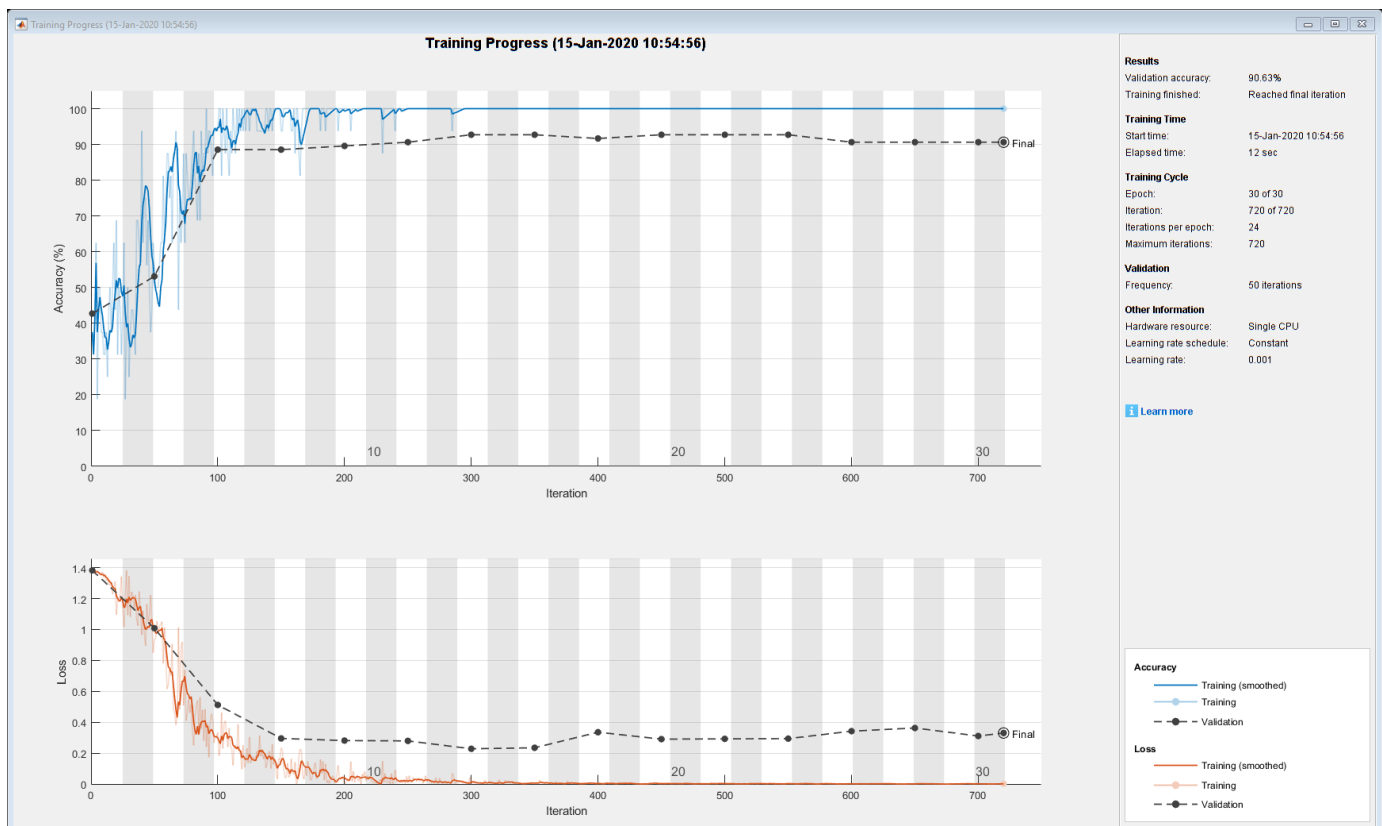
- Train using the Adam solver.
- Specify a mini-batch size of 16.
- Shuffle the data every epoch.
- Monitor the training progress by setting the 'Plots' option to 'training-progress'.
- Specify the validation data using the 'ValidationData' option.
- Suppress verbose output by setting the 'Verbose' option to false.

By default, `trainNetwork` uses a GPU if one is available. Otherwise, it uses the CPU. To specify the execution environment manually, use the 'ExecutionEnvironment' name-value pair argument of `trainingOptions`. Training on a CPU can take significantly longer than training on a GPU. Training with a GPU requires Parallel Computing Toolbox™ and a supported GPU device. For information on supported devices, see “GPU Support by Release” (Parallel Computing Toolbox).

```
options = trainingOptions('adam', ...
    'MiniBatchSize',16, ...
    'GradientThreshold',2, ...
    'Shuffle','every-epoch', ...
    'ValidationData',{XValidation,YValidation}, ...
    'Plots','training-progress', ...
    'Verbose',false);
```

Train the LSTM network using the `trainNetwork` function.

```
net = trainNetwork(XTrain,YTrain,layers,options);
```





## Predict Using New Data

Classify the event type of three new reports. Create a string array containing the new reports.

```
reportsNew = [ ...
    "Coolant is pooling underneath sorter."
    "Sorter blows fuses at start up."
    "There are some very loud rattling sounds coming from the assembler."];
```

Preprocess the text data using the preprocessing steps as the training documents.

```
documentsNew = preprocessText(reportsNew);
```

Convert the text data to sequences using `doc2sequence` with the same options as when creating the training sequences.

```
XNew = doc2sequence(enc,documentsNew,'Length',sequenceLength);
```

Classify the new sequences using the trained LSTM network.

```
labelsNew = classify(net,XNew)
```

```
labelsNew = 3×1 categorical
    Leak
    Electronic Failure
    Mechanical Failure
```

## Preprocessing Function

The function `preprocessText` performs these steps:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Convert the text to lowercase using `lower`.
- 3 Erase the punctuation using `erasePunctuation`.

```
function documents = preprocessText(textData)
```

```
% Tokenize the text.
documents = tokenizedDocument(textData);
```

```
% Convert to lowercase.
documents = lower(documents);
```

```
% Erase punctuation.
documents = erasePunctuation(documents);
```

```
end
```

## See Also

`fastTextWordEmbedding` | `wordEmbeddingLayer` | `tokenizedDocument` | `lstmLayer` | `trainNetwork` | `trainingOptions` | `doc2sequence` | `sequenceInputLayer` | `wordcloud`

## Related Examples

- “Classify Text Data Using Convolutional Neural Network” on page 2-93

- “Classify Out-of-Memory Text Data Using Deep Learning” on page 2-168
- “Generate Text Using Deep Learning” (Deep Learning Toolbox)
- “Word-By-Word Text Generation Using Deep Learning” on page 2-180
- “Create Simple Text Model for Classification” on page 2-2
- “Analyze Text Data Using Topic Models” on page 2-13
- “Analyze Text Data Using Multiword Phrases” on page 2-7
- “Train a Sentiment Classifier” on page 2-71
- “Sequence Classification Using Deep Learning” (Deep Learning Toolbox)
- “Deep Learning in MATLAB” (Deep Learning Toolbox)

## Classify Text Data Using Convolutional Neural Network

This example shows how to classify text data using a convolutional neural network.

To classify text data using convolutions, use 1-D convolutional layers that convolve over the time dimension of the input.

This example trains a network with 1-D convolutional filters of varying widths. The width of each filter corresponds to the number of words the filter can see (the n-gram length). The network has multiple branches of convolutional layers, so it can use different n-gram lengths.

### Load Data

Create a tabular text datastore from the data in `factoryReports.csv` and view the first few reports.

```
data = readtable("factoryReports.csv");
head(data)
```

ans=8x5 table

Description	Category
{'Items are occasionally getting stuck in the scanner spools.'	{'Mechanical Failure'}
{'Loud rattling and banging sounds are coming from assembler pistons.'	{'Mechanical Failure'}
{'There are cuts to the power when starting the plant.'	{'Electronic Failure'}
{'Fried capacitors in the assembler.'	{'Electronic Failure'}
{'Mixer tripped the fuses.'	{'Electronic Failure'}
{'Burst pipe in the constructing agent is spraying coolant.'	{'Leak'}
{'A fuse is blown in the mixer.'	{'Electronic Failure'}
{'Things continue to tumble off of the belt.'	{'Mechanical Failure'}

Partition the data into training and validation partitions. Use 80% of the data for training and the remaining data for validation.

```
cvp = cvpartition(data.Category, Holdout=0.2);
dataTrain = data(training(cvp),:);
dataValidation = data(test(cvp),:);
```

### Preprocess Text Data

Extract the text data from the "Description" column of the table and preprocess it using the `preprocessText` function, listed in the section [Preprocess Text Function](#) on page 2-0 of the example.

```
documentsTrain = preprocessText(dataTrain.Description);
```

Extract the labels from the "Category" column and convert them to categorical.

```
TTrain = categorical(dataTrain.Category);
```

View the class names and the number of observations.

```
classNames = unique(TTrain)
classNames = 4x1 categorical
    Electronic Failure
```

```
Leak  
Mechanical Failure  
Software Failure
```

```
numObservations = numel(TTrain)
```

```
numObservations = 384
```

Extract and preprocess the validation data using the same steps.

```
documentsValidation = preprocessText(dataValidation.Description);  
TValidation = categorical(dataValidation.Category);
```

### Convert Documents to Sequences

To input the documents into a neural network, use a word encoding to convert the documents into sequences of numeric indices.

Create a word encoding from the documents.

```
enc = wordEncoding(documentsTrain);
```

View the vocabulary size of the word encoding. The vocabulary size is the number of unique words of the word encoding.

```
numWords = enc.NumWords
```

```
numWords = 436
```

Convert the documents to sequences of integers using the `doc2sequence` function.

```
XTrain = doc2sequence(enc,documentsTrain);
```

Convert the validation documents to sequences using the word encoding created from the training data.

```
XValidation = doc2sequence(enc,documentsValidation);
```

### Define Network Architecture

Define the network architecture for the classification task.

The following steps describe the network architecture.

- Specify an input size of 1, which corresponds to the channel dimension of the integer sequence input.
- Embed the input using a word embedding of dimension 100.
- For the n-gram lengths 2, 3, 4, and 5, create blocks of layers containing a convolutional layer, a batch normalization layer, a ReLU layer, a dropout layer, and a max pooling layer.
- For each block, specify 200 convolutional filters of size 1-by- $N$  and a global max pooling layer.
- Connect the input layer to each block and concatenate the outputs of the blocks using a concatenation layer.
- To classify the outputs, include a fully connected layer with output size  $K$ , a softmax layer, and a classification layer, where  $K$  is the number of classes.

Specify the network hyperparameters.

```
embeddingDimension = 100;
ngramLengths = [2 3 4 5];
numFilters = 200;
```

First, create a layer graph containing the input layer and a word embedding layer of dimension 100. To help connect the word embedding layer to the convolution layers, set the word embedding layer name to "emb". To check that the convolution layers do not convolve the sequences to have a length of zero during training, set the `MinLength` option to the length of the shortest sequence in the training data.

```
minLength = min(doclength(documentsTrain));
layers = [
    sequenceInputLayer(1,MinLength=minLength)
    wordEmbeddingLayer(embeddingDimension,numWords,Name="emb")];
lgraph = layerGraph(layers);
```

For each of the n-gram lengths, create a block of 1-D convolution, batch normalization, ReLU, dropout, and 1-D global max pooling layers. Connect each block to the word embedding layer.

```
numBlocks = numel(ngramLengths);
for j = 1:numBlocks
    N = ngramLengths(j);

    block = [
        convolution1dLayer(N,numFilters,Name="conv"+N,Padding="same")
        batchNormalizationLayer(Name="bn"+N)
        reluLayer(Name="relu"+N)
        dropoutLayer(0.2,Name="drop"+N)
        globalMaxPooling1dLayer(Name="max"+N)];

    lgraph = addLayers(lgraph,block);
    lgraph = connectLayers(lgraph,"emb","conv"+N);
end
```

Add the concatenation layer, the fully connected layer, the softmax layer, and the classification layer.

```
numClasses = numel(classNames);

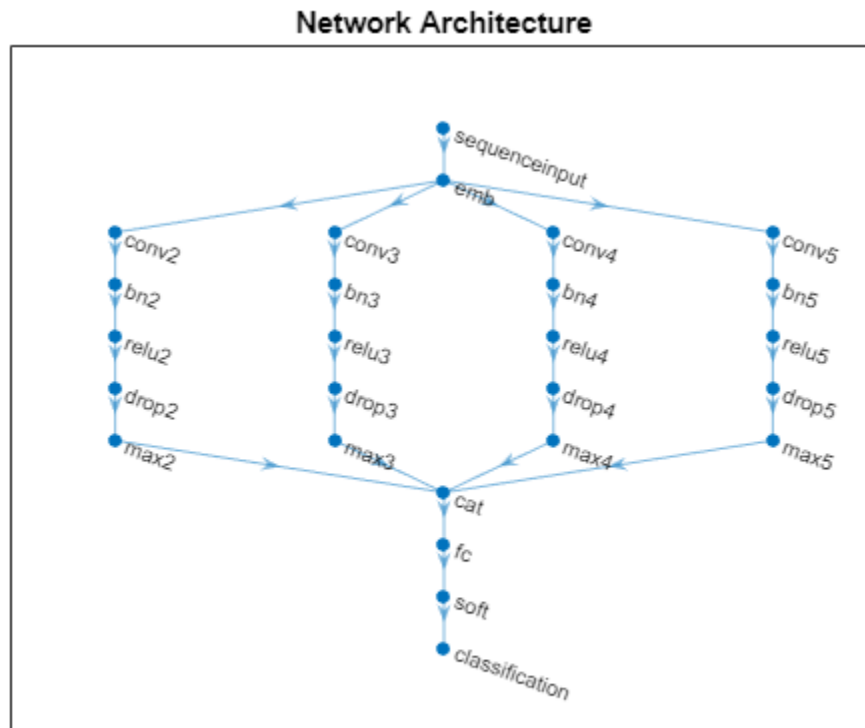
layers = [
    concatenationLayer(1,numBlocks,Name="cat")
    fullyConnectedLayer(numClasses,Name="fc")
    softmaxLayer(Name="soft")
    classificationLayer(Name="classification")];

lgraph = addLayers(lgraph,layers);
```

Connect the global max pooling layers to the concatenation layer and view the network architecture in a plot.

```
for j = 1:numBlocks
    N = ngramLengths(j);
    lgraph = connectLayers(lgraph,"max"+N,"cat/in"+j);
end

figure
plot(lgraph)
title("Network Architecture")
```



### Train Network

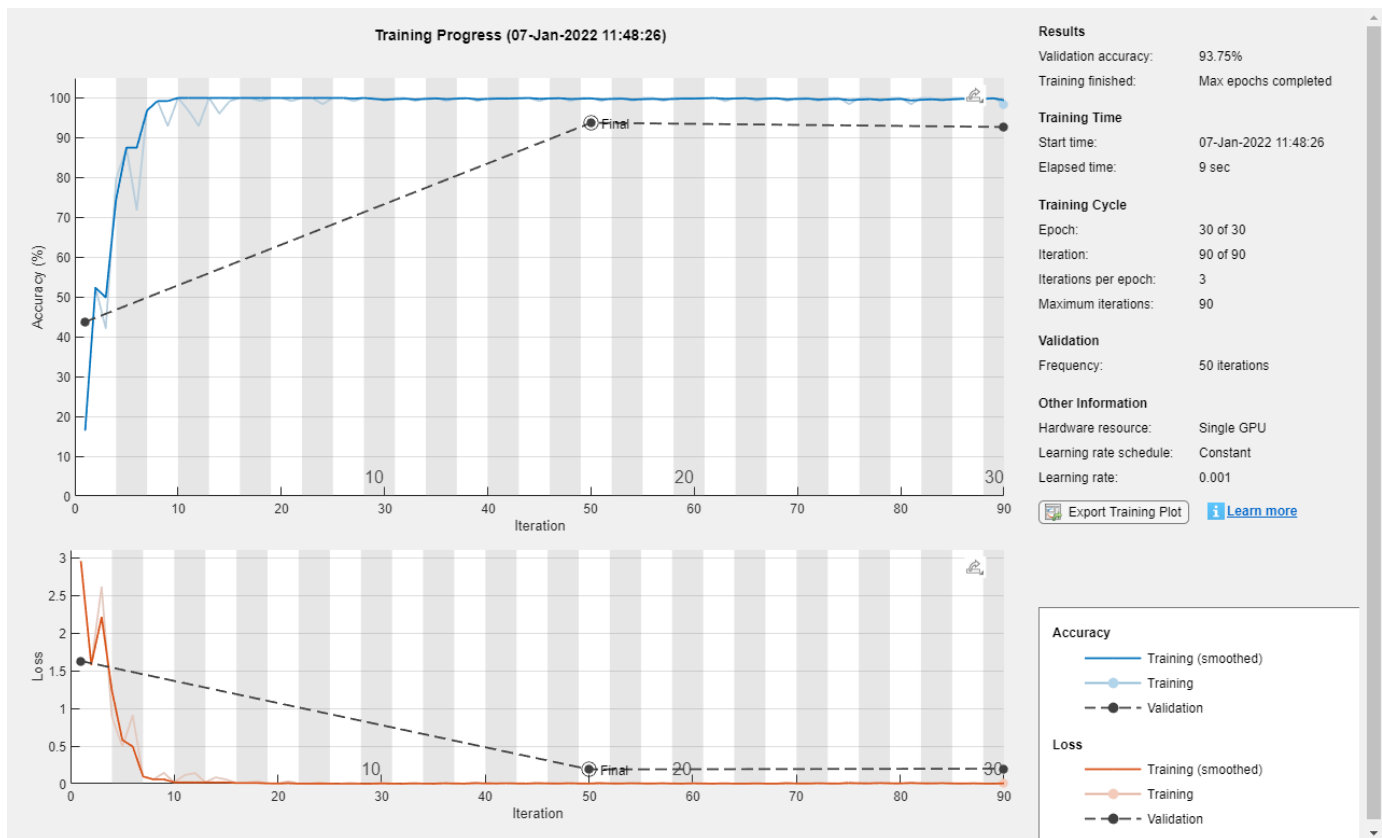
Specify the training options:

- Train with a mini-batch size of 128.
- Validate the network using the validation data.
- Return the network with the lowest validation loss.
- Display the training progress plot and suppress the verbose output.

```
options = trainingOptions("adam", ...
    MiniBatchSize=128, ...
    ValidationData={XValidation,TValidation}, ...
    OutputNetwork="best-validation-loss", ...
    Plots="training-progress", ...
    Verbose=false);
```

Train the network using the `trainNetwork` function.

```
net = trainNetwork(XTrain,TTrain,lgraph,options);
```



## Test Network

Classify the validation data using the trained network.

```
YValidation = classify(net,XValidation);
```

Visualize the predictions in a confusion chart.

```
figure
confusionchart(TValidation,YValidation)
```

True Class	Electronic Failure	30		3	
	Leak	1	13		
	Mechanical Failure	1	1	38	
	Software Failure				9
		Electronic Failure	Leak	Mechanical Failure	Software Failure
		Predicted Class			

Calculate the classification accuracy. The accuracy is the proportion of labels predicted correctly.

```
accuracy = mean(TValidation == YValidation)
```

```
accuracy = 0.9375
```

### Predict Using New Data

Classify the event type of three new reports. Create a string array containing the new reports.

```
reportsNew = [
    "Coolant is pooling underneath sorter."
    "Sorter blows fuses at start up."
    "There are some very loud rattling sounds coming from the assembler."];
```

Preprocess the text data using the preprocessing steps as the training and validation documents.

```
documentsNew = preprocessText(reportsNew);
XNew = doc2sequence(enc,documentsNew);
```

Classify the new sequences using the trained network.

```
YNew = classify(net,XNew)
```

```
YNew = 3x1 categorical
    Leak
    Electronic Failure
    Mechanical Failure
```



## Preprocess Text Function

The `preprocessTextData` function takes text data as input and performs these steps:

- 1 Tokenize the text.
- 2 Convert the text to lowercase.

```
function documents = preprocessText(textData)

documents = tokenizedDocument(textData);
documents = lower(documents);

end
```

## See Also

[fastTextWordEmbedding](#) | [wordcloud](#) | [wordEmbedding](#) | [layerGraph](#) | [convolution2dLayer](#) | [batchNormalizationLayer](#) | [trainingOptions](#) | [trainNetwork](#) | [doc2sequence](#) | [tokenizedDocument](#)

## Related Examples

- “Classify Text Data Using Deep Learning” on page 2-85
- “Classify Out-of-Memory Text Data Using Deep Learning” on page 2-168
- “Create Simple Text Model for Classification” on page 2-2
- “Analyze Text Data Using Topic Models” on page 2-13
- “Analyze Text Data Using Multiword Phrases” on page 2-7
- “Train a Sentiment Classifier” on page 2-71
- “Sequence Classification Using Deep Learning” (Deep Learning Toolbox)
- “Deep Learning in MATLAB” (Deep Learning Toolbox)

## Classify Text Data Using Custom Training Loop

This example shows how to classify text data using a deep learning bidirectional long short-term memory (BiLSTM) network with a custom training loop.

When training a deep learning network using the `trainNetwork` function, if `trainingOptions` does not provide the options you need (for example, a custom learning rate schedule), then you can define your own custom training loop using automatic differentiation. For an example showing how to classify text data using the `trainNetwork` function, see “Classify Text Data Using Deep Learning” (Deep Learning Toolbox).

This example trains a network to classify text data with the *time-based decay* learning rate schedule: for each iteration, the solver uses the learning rate given by  $\rho_t = \frac{\rho_0}{1 + kt}$ , where  $t$  is the iteration number,  $\rho_0$  is the initial learning rate, and  $k$  is the decay.

### Import Data

Import the factory reports data. This data contains labeled textual descriptions of factory events. To import the text data as strings, specify the text type to be “string”.

```
filename = "factoryReports.csv";
data = readtable(filename, TextType="string");
head(data)
```

ans=8×5 table

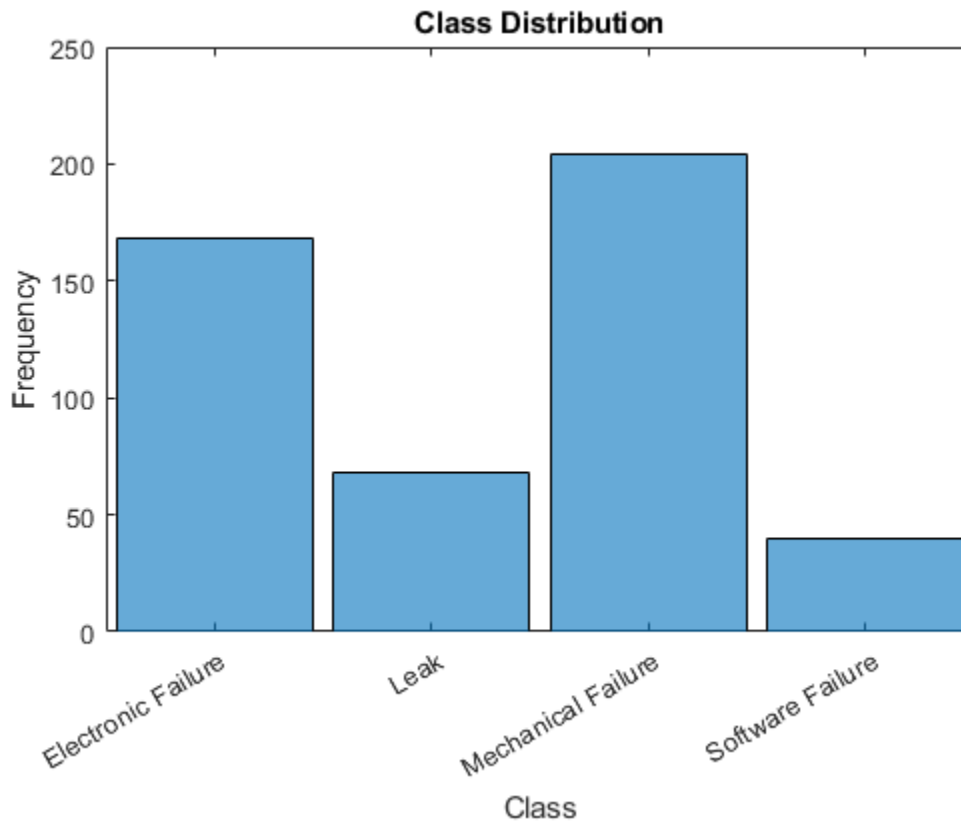
Description	Category
"Items are occasionally getting stuck in the scanner spools."	"Mechanical Failure"
"Loud rattling and banging sounds are coming from assembler pistons."	"Mechanical Failure"
"There are cuts to the power when starting the plant."	"Electronic Failure"
"Fried capacitors in the assembler."	"Electronic Failure"
"Mixer tripped the fuses."	"Electronic Failure"
"Burst pipe in the constructing agent is spraying coolant."	"Leak"
"A fuse is blown in the mixer."	"Electronic Failure"
"Things continue to tumble off of the belt."	"Mechanical Failure"

The goal of this example is to classify events by the label in the Category column. To divide the data into classes, convert these labels to categorical.

```
data.Category = categorical(data.Category);
```

View the distribution of the classes in the data using a histogram.

```
figure
histogram(data.Category);
xlabel("Class")
ylabel("Frequency")
title("Class Distribution")
```



The next step is to partition it into sets for training and validation. Partition the data into a training partition and a held-out partition for validation and testing. Specify the holdout percentage to be 20%.

```
cvp = cvpartition(data.Category, Holdout=0.2);  
dataTrain = data(training(cvp), :);  
dataValidation = data(test(cvp), :);
```

Extract the text data and labels from the partitioned tables.

```
textDataTrain = dataTrain.Description;  
textDataValidation = dataValidation.Description;  
TTrain = dataTrain.Category;  
TValidation = dataValidation.Category;
```

To check that you have imported the data correctly, visualize the training text data using a word cloud.

```
figure  
wordcloud(textDataTrain);  
title("Training Data")
```



```

9 tokens: items are occasionally getting stuck in the scanner spools
10 tokens: loud rattling and banging sounds are coming from assembler pistons
5 tokens: fried capacitors in the assembler
4 tokens: mixer tripped the fuses
9 tokens: burst pipe in the constructing agent is spraying coolant

```

Create a single datastore that contains both the documents and the labels by creating `arrayDatastore` objects, then combining them using the `combine` function.

```

dsDocumentsTrain = arrayDatastore(documentsTrain,OutputType="cell");
dsTTrain = arrayDatastore(TTrain,OutputType="cell");
dsTrain = combine(dsDocumentsTrain,dsTTrain);

```

Create an array datastore for the validation documents.

```

dsDocumentsValidation = arrayDatastore(documentsValidation,OutputType="cell");

```

### Create Word Encoding

To input the documents into a BiLSTM network, use a word encoding to convert the documents into sequences of numeric indices.

To create a word encoding, use the `wordEncoding` function.

```

enc = wordEncoding(documentsTrain)

```

```

enc =

```

```

    wordEncoding with properties:

```

```

        NumWords: 417
        Vocabulary: ["items"    "are"    "occasionally"  "getting"  "stuck"  "in"    "the"

```

### Define Network

Define the BiLSTM network architecture. To input sequence data into the network, include a sequence input layer and set the input size to 1. Next, include a word embedding layer of dimension 25 and the same number of words as the word encoding. Next, include a BiLSTM layer and set the number of hidden units to 40. To use the BiLSTM layer for a sequence-to-label classification problem, set the output mode to `"last"`. Finally, add a fully connected layer with the same size as the number of classes, and a softmax layer.

```

inputSize = 1;
embeddingDimension = 25;
numHiddenUnits = 40;

```

```

numWords = enc.NumWords;

```

```

layers = [
    sequenceInputLayer(inputSize)
    wordEmbeddingLayer(embeddingDimension,numWords)
    bilstmLayer(numHiddenUnits,OutputMode="last")
    fullyConnectedLayer(numClasses)
    softmaxLayer]

```

```

layers =
    5x1 Layer array with layers:

```

```

1 '' Sequence Input           Sequence input with 1 dimensions
2 '' Word Embedding Layer    Word embedding layer with 25 dimensions and 417 unique words
3 '' BiLSTM                  BiLSTM with 40 hidden units
4 '' Fully Connected         4 fully connected layer
5 '' Softmax                  softmax

```

Convert the layer array to a `dlnetwork` object.

```
net = dlnetwork(layers)
```

```
net =
  dlnetwork with properties:

    Layers: [5x1 nnet.cnn.layer.Layer]
  Connections: [4x2 table]
  Learnables: [6x3 table]
    State: [2x3 table]
  InputNames: {'sequenceinput'}
  OutputNames: {'softmax'}
  Initialized: 1

```

### Define Model Loss Function

Create the function `modelLoss`, listed at the end of the example, that takes a `dlnetwork` object, a mini-batch of input data with corresponding labels, and returns the loss and the gradients of the loss with respect to the learnable parameters in the network.

### Specify Training Options

Train for 30 epochs with a mini-batch size of 16.

```
numEpochs = 30;
miniBatchSize = 16;
```

Specify the options for Adam optimization. Specify an initial learn rate of 0.001 with a decay of 0.01, gradient decay factor 0.9, and squared gradient decay factor 0.999.

```
initialLearnRate = 0.001;
decay = 0.01;
gradientDecayFactor = 0.9;
squaredGradientDecayFactor = 0.999;
```

### Train Model

Create a `minibatchqueue` object that processes and manages the mini-batches of data. For each mini-batch:

- Use the custom mini-batch preprocessing function `preprocessMiniBatch` (defined at the end of this example) to convert documents to sequences and one-hot encode the labels. To pass the word encoding to the mini-batch, create an anonymous function that takes two inputs.
- Format the predictors with the dimension labels "BTC" (batch, time, channel). The `minibatchqueue` object, by default, converts the data to `dlarray` objects with underlying type `single`.
- Train on a GPU if one is available. The `minibatchqueue` object, by default, converts each output to `gpuArray` if a GPU is available. Using a GPU requires Parallel Computing Toolbox™ and a

supported GPU device. For information on supported devices, see “GPU Support by Release” (Parallel Computing Toolbox).

```
mbq = minibatchqueue(dsTrain, ...
    MiniBatchSize=miniBatchSize, ...
    MiniBatchFcn=@(X,T) preprocessMiniBatch(X,T,enc), ...
    MiniBatchFormat=["BTC" ""]);
```

Create a `minibatchqueue` object for the validation documents. For each mini-batch:

- Use the custom mini-batch preprocessing function `preprocessMiniBatchPredictors` (defined at the end of this example) to convert documents to sequences. This preprocessing function does not require label data. To pass the word encoding to the mini-batch, create an anonymous function that takes one input only.
- Format the predictors with the dimension labels "BTC" (batch, time, channel). The `minibatchqueue` object, by default, converts the data to `dlarray` objects with underlying type `single`.
- To make predictions for all observations, return any partial mini-batches.

```
mbqValidation = minibatchqueue(dsDocumentsValidation, ...
    MiniBatchSize=miniBatchSize, ...
    MiniBatchFcn=@(X) preprocessMiniBatchPredictors(X,enc), ...
    MiniBatchFormat="BTC", ...
    PartialMiniBatch="return");
```

To easily calculate the validation loss, convert the validation labels to one-hot encoded vectors and transpose the encoded labels to match the network output format.

```
TValidation = onehotencode(TValidation,2);
TValidation = TValidation';
```

Initialize the training progress plot.

```
figure
C = colororder;
lineLossTrain = animatedline(Color=C(2,:));

lineLossValidation = animatedline( ...
    LineStyle="--", ...
    Marker="o", ...
    MarkerFaceColor="black");

ylim([0 inf])
xlabel("Iteration")
ylabel("Loss")
grid on
```

Initialize the parameters for Adam.

```
trailingAvg = [];
trailingAvgSq = [];
```

Train the network. For each epoch, shuffle the data and loop over mini-batches of data. At the end of each iteration, display the training progress. At the end of each epoch, validate the network using the validation data.

For each mini-batch:

- Convert the documents to sequences of integers and one-hot encode the labels.
- Convert the data to `darray` objects with underlying type `single` and specify the dimension labels "BTC" (batch, time, channel).
- For GPU training, convert to `gpuArray` objects.
- Evaluate the model loss and gradients using `dlfeval` and the `modelLoss` function.
- Determine the learning rate for the time-based decay learning rate schedule.
- Update the network parameters using the `adamupdate` function.
- Update the training plot.

```
iteration = 0;
start = tic;

% Loop over epochs.
for epoch = 1:numEpochs

    % Shuffle data.
    shuffle(mbq);

    % Loop over mini-batches.
    while hasdata(mbq)
        iteration = iteration + 1;

        % Read mini-batch of data.
        [X,T] = next(mbq);

        % Evaluate the model loss and gradients using dlfeval and the
        % modelLoss function.
        [loss,gradients] = dlfeval(@modelLoss,net,X,T);

        % Determine learning rate for time-based decay learning rate schedule.
        learnRate = initialLearnRate/(1 + decay*iteration);

        % Update the network parameters using the Adam optimizer.
        [net,trailingAvg,trailingAvgSq] = adamupdate(net, gradients, ...
            trailingAvg, trailingAvgSq, iteration, learnRate, ...
            gradientDecayFactor, squaredGradientDecayFactor);

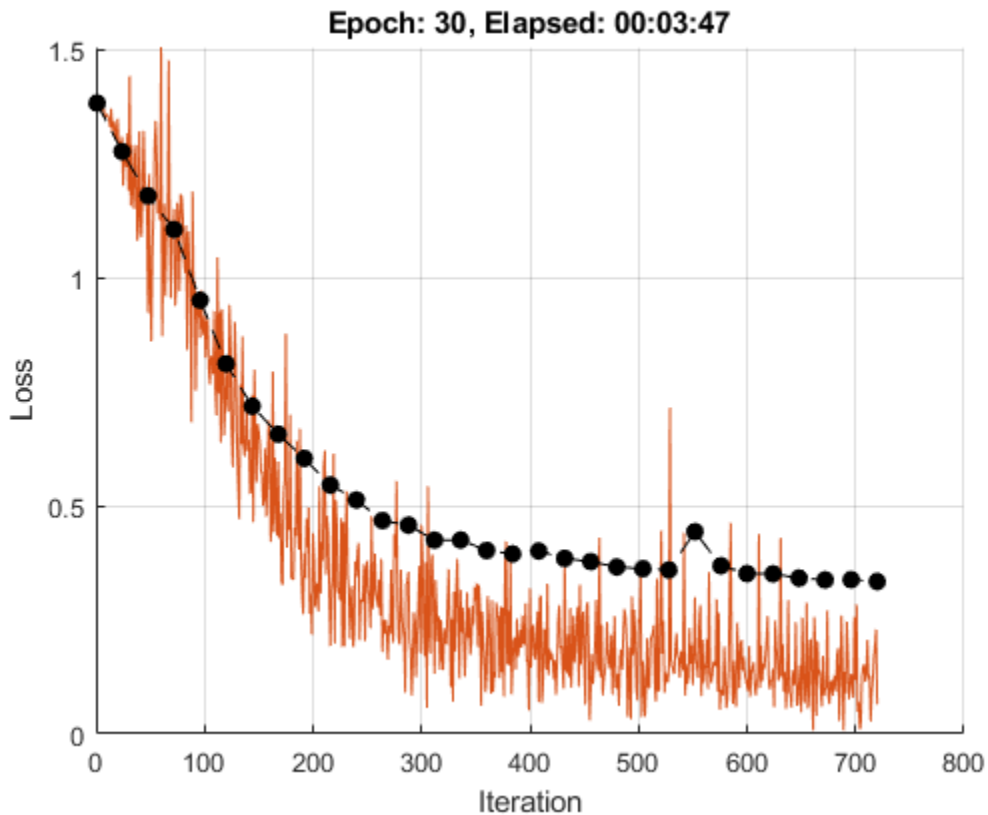
        % Display the training progress.
        D = duration(0,0,toc(start),Format="hh:mm:ss");
        loss = double(loss);
        addpoints(lineLossTrain,iteration,loss)
        title("Epoch: " + epoch + ", Elapsed: " + string(D))
        drawnow

        % Validate network.
        if iteration == 1 || ~hasdata(mbq)
            [~,scoresValidation] = modelPredictions(net,mbqValidation,classes);
            lossValidation = crossentropy(scoresValidation,TValidation);

            % Update plot.
            lossValidation = double(lossValidation);
            addpoints(lineLossValidation,iteration,lossValidation)
            drawnow
        end
    end
end
```



```
end
end
```



### Test Model

Test the classification accuracy of the model by comparing the predictions on the validation set with the true labels.

Classify the validation data using `modelPredictions` function, listed at the end of the example.

```
YNew = modelPredictions(net,mbqValidation,classes);
```

To easily calculate the validation accuracy, convert the one-hot encoded validation labels to categorical and transpose.

```
TValidation = onehotdecode(TValidation,classes,1)';
```

Evaluate the classification accuracy.

```
accuracy = mean(YNew == TValidation)
```

```
accuracy = 0.8854
```

### Predict Using New Data

Classify the event type of three new reports. Create a string array containing the new reports.

```
reportsNew = [
    "Coolant is pooling underneath sorter."
```

```
"Sorter blows fuses at start up."
"There are some very loud rattling sounds coming from the assembler."];
```

Preprocess the text data using the preprocessing steps as the training documents.

```
documentsNew = preprocessText(reportsNew);
dsNew = arrayDatastore(documentsNew,OutputType="cell");
```

Create a minibatchqueue object that processes and manages the mini-batches of data. For each mini-batch:

- Use the custom mini-batch preprocessing function `preprocessMiniBatchPredictors` (defined at the end of this example) to convert documents to sequences. This preprocessing function does not require label data. To pass the word encoding to the mini-batch, create an anonymous function that takes one input only.
- Format the predictors with the dimension labels "BTC" (batch, time, channel). The `minibatchqueue` object, by default, converts the data to `darray` objects with underlying type `single`.
- To make predictions for all observations, return any partial mini-batches.

```
mbqNew = minibatchqueue(dsNew, ...
    MiniBatchSize=miniBatchSize, ...
    MiniBatchFcn=@(X) preprocessMiniBatchPredictors(X,enc), ...
    MiniBatchFormat="BTC", ...
    PartialMiniBatch="return");
```

Classify the text data using `modelPredictions` function, listed at the end of the example and find the classes with the highest scores.

```
YNew = modelPredictions(net,mbqNew,classes)
```

```
YNew = 3×1 categorical
    Leak
    Electronic Failure
    Mechanical Failure
```

## Supporting Functions

### Text Preprocessing Function

The function `preprocessText` performs these steps:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Convert the text to lowercase using `lower`.
- 3 Erase the punctuation using `erasePunctuation`.

```
function documents = preprocessText(textData)
```

```
% Tokenize the text.
documents = tokenizedDocument(textData);
```

```
% Convert to lowercase.
documents = lower(documents);
```

```
% Erase punctuation.
```

```
documents = erasePunctuation(documents);
```

```
end
```

### Mini-Batch Preprocessing Function

The `preprocessMiniBatch` function converts a mini-batch of documents to sequences of integers and one-hot encodes label data.

```
function [X,T] = preprocessMiniBatch(dataX,dataT,enc)
```

```
% Preprocess predictors.
```

```
X = preprocessMiniBatchPredictors(dataX,enc);
```

```
% Extract labels from cell and concatenate.
```

```
T = cat(1,dataT{1:end});
```

```
% One-hot encode labels.
```

```
T = onehotencode(T,2);
```

```
% Transpose the encoded labels to match the network output.
```

```
T = T';
```

```
end
```

### Mini-Batch Predictors Preprocessing Function

The `preprocessMiniBatchPredictors` function converts a mini-batch of documents to sequences of integers.

```
function X = preprocessMiniBatchPredictors(dataX,enc)
```

```
% Extract documents from cell and concatenate.
```

```
documents = cat(4,dataX{1:end});
```

```
% Convert documents to sequences of integers.
```

```
X = doc2sequence(enc,documents);
```

```
X = cat(1,X{:});
```

```
end
```

### Model Loss Function

The `modelLoss` function takes a `dlnetwork` object `net`, a mini-batch of input data `X` with corresponding target labels `T` and returns the gradients of the loss with respect to the learnable parameters in `net`, and the loss. To compute the gradients automatically, use the `dlgradient` function.

```
function [loss,gradients] = modelLoss(net,X,T)
```

```
Y = forward(net,X);
```

```
loss = crossentropy(Y,T);
```

```
gradients = dlgradient(loss,net.Learnables);
```

```
end
```

### Model Predictions Function

The `modelPredictions` function takes a `dlnetwork` object `net`, a mini-batch queue, and outputs the model predictions and scores by iterating over mini-batches in the queue.

```
function [predictions,scores] = modelPredictions(net,mbq,classes)

% Initialize predictions.
predictions = [];
scores = [];

% Reset mini-batch queue.
reset(mbq);

% Loop over mini-batches.
while hasdata(mbq)

    % Make predictions.
    X = next(mbq);
    Y = predict(net,X);

    scores = [scores Y];

    Y = onehotdecode(Y,classes,1)';
    predictions = [predictions; Y];
end

end
```

### See Also

[wordEmbeddingLayer](#) | [tokenizedDocument](#) | [lstmLayer](#) | [doc2sequence](#) | [sequenceInputLayer](#) | [wordcloud](#) | [dlfeval](#) | [dlgradient](#) | [dlarray](#)

### Related Examples

- “Define Custom Training Loops, Loss Functions, and Networks” (Deep Learning Toolbox)
- “Classify Text Data Using Deep Learning” on page 2-85
- “Create Simple Text Model for Classification” on page 2-2
- “Analyze Text Data Using Topic Models” on page 2-13
- “Analyze Text Data Using Multiword Phrases” on page 2-7
- “Train a Sentiment Classifier” on page 2-71
- “Sequence Classification Using Deep Learning” (Deep Learning Toolbox)
- “Deep Learning in MATLAB” (Deep Learning Toolbox)

## Multilabel Text Classification Using Deep Learning

This example shows how to classify text data that has multiple independent labels.

For classification tasks where there can be multiple independent labels for each observation—for example, tags on an scientific article—you can train a deep learning model to predict probabilities for each independent class. To enable a network to learn multilabel classification targets, you can optimize the loss of each class independently using binary cross-entropy loss.

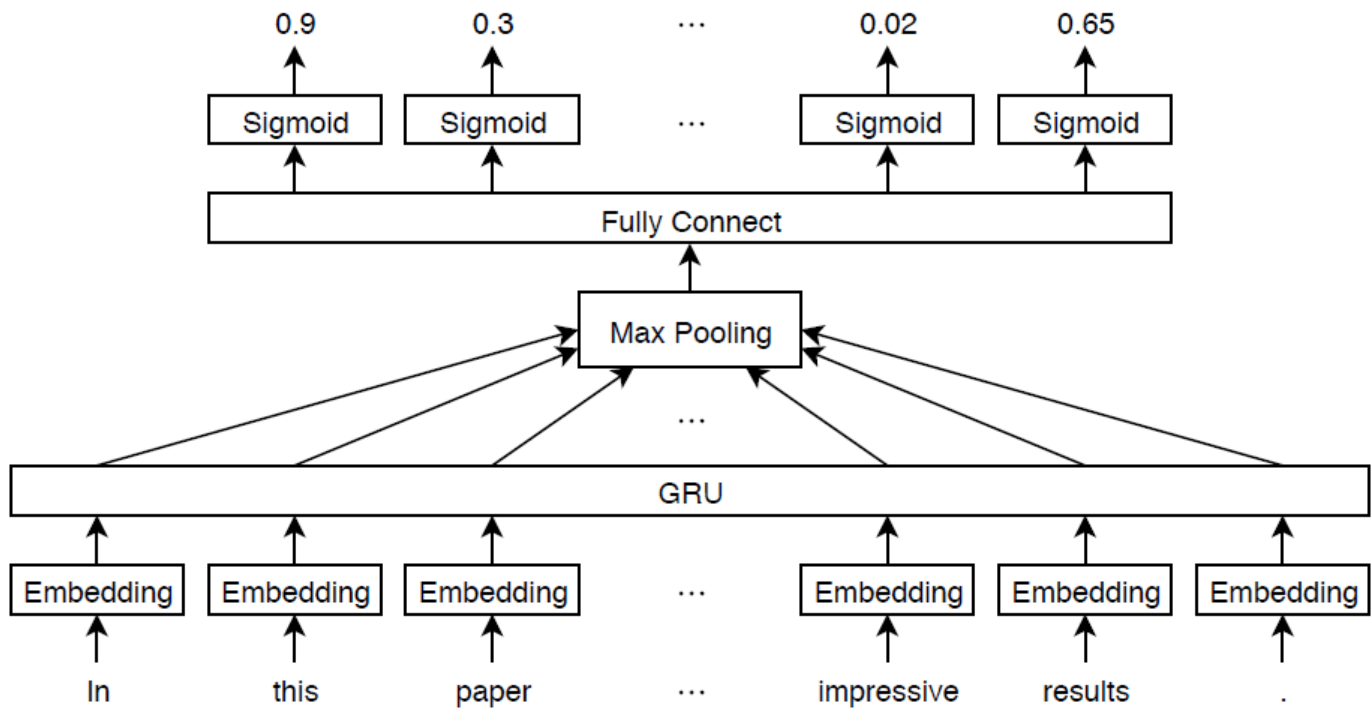
This example defines a deep learning model that classifies subject areas given the abstracts of mathematical papers collected using the arXiv API [1]. The model consists of a word embedding and GRU, max pooling operation, fully connected, and sigmoid operations.

To measure the performance of multilabel classification, you can use the labeling F-score [2]. The labeling F-score evaluates multilabel classification by focusing on per-text classification with partial matches. The measure is the normalized proportion of matching labels against the total number of true and predicted labels.

This example defines the following model:

- A word embedding that maps a sequence of words to a sequence of numeric vectors.
- A GRU operation that learns dependencies between the embedding vectors.
- A max pooling operation that reduces a sequence of feature vectors to a single feature vector.
- A fully connected layer that maps the features to the binary outputs.
- A sigmoid operation for learning the binary cross entropy loss between the outputs and the target labels.

This diagram shows a piece of text propagating through the model architecture and outputting a vector of probabilities. The probabilities are independent, so they need not sum to one.



### Import Text Data

Import a set of abstracts and category labels from math papers using the arXiv API. Specify the number of records to import using the `importSize` variable.

```
importSize = 50000;
```

Create a URL that queries records with set "math" and metadata prefix "arXiv".

```
url = "https://export.arxiv.org/oai2?verb=ListRecords" + ...
      "&set=math" + ...
      "&metadataPrefix=arXiv";
```

Extract the abstract text, category labels, and the resumption token returned by the query URL using the `parseArXivRecords` function which is attached to this example as a supporting file. To access this file, open this example as a live script. Note that the arXiv API is rate limited and requires waiting between multiple requests.

```
[textData, labelsAll, resumptionToken] = parseArXivRecords(url);
```

Iteratively import more chunks of records until the required amount is reached, or there are no more records. To continue importing records from where you left off, use the resumption token from the previous result in the query URL. To adhere to the rate limits imposed by the arXiv API, add a delay of 20 seconds before each query using the `pause` function.

```
while numel(textData) < importSize
    if resumptionToken == ""
        break
    end
```

```

url = "https://export.arxiv.org/oai2?verb=ListRecords" + ...
    "&resumptionToken=" + resumptionToken;

pause(20)
[textDataNew, labelsNew, resumptionToken] = parseArXivRecords(url);

textData = [textData; textDataNew];
labelsAll = [labelsAll; labelsNew];
end

```

### Preprocess Text Data

Tokenize and preprocess the text data using the `preprocessText` function, listed at the end of the example.

```

documentsAll = preprocessText(textData);
documentsAll(1:5)

ans =
    5×1 tokenizedDocument:

    72 tokens: describe new algorithm  $(k, \ell)$  pebble game color obtain characterization family
    22 tokens: show determinant stirling cycle number count unlabeled acyclic singlesource autom
    18 tokens: paper show compute  $\lambda_{\alpha}$  norm alpha dyadic grid result consequence d
    62 tokens: partial cube isometric subgraphs hypercubes structure graph define mean semicubes
    29 tokens: paper present algorithm compute hecke eigensystems hilbertsiegel cusp form real q

```

Remove labels that do not belong to the "math" set.

```

for i = 1:numel(labelsAll)
    labelsAll{i} = labelsAll{i}(startsWith(labelsAll{i}, "math."));
end

```

Visualize some of the classes in a word cloud. Find the documents corresponding to the following:

- Abstracts tagged with "Combinatorics" and not tagged with "Statistics Theory"
- Abstracts tagged with "Statistics Theory" and not tagged with "Combinatorics"
- Abstracts tagged with both "Combinatorics" and "Statistics Theory"

Find the document indices for each of the groups using the `ismember` function.

```

idxC0 = cellfun(@(lbls) ismember("math.CO", lbls) && ~ismember("math.ST", lbls), labelsAll);
idxST = cellfun(@(lbls) ismember("math.ST", lbls) && ~ismember("math.CO", lbls), labelsAll);
idxCOST = cellfun(@(lbls) ismember("math.CO", lbls) && ismember("math.ST", lbls), labelsAll);

```

Visualize the documents for each group in a word cloud.

```

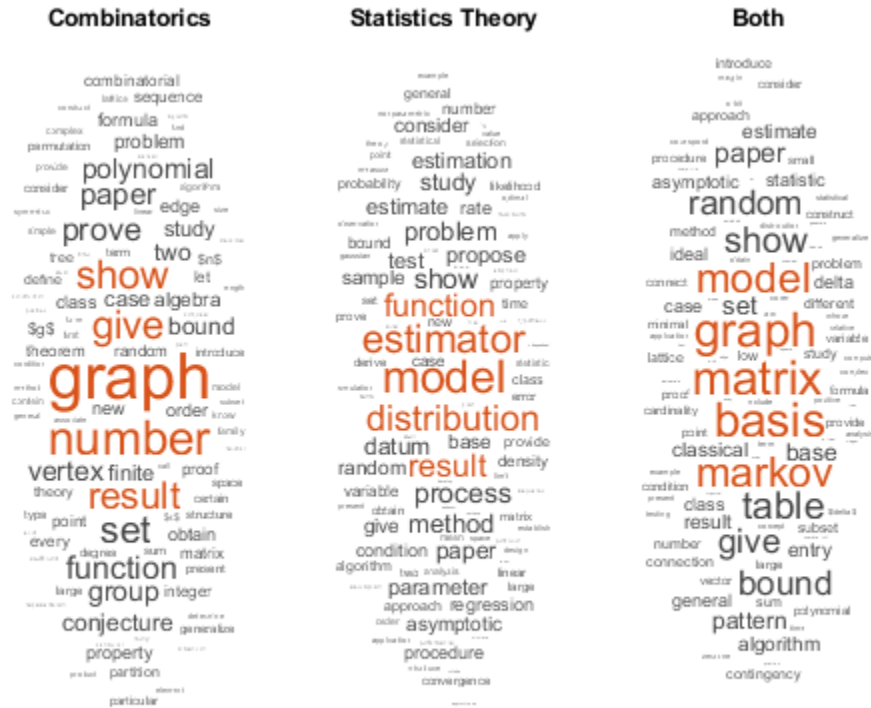
figure
subplot(1,3,1)
wordcloud(documentsAll(idxC0));
title("Combinatorics")

subplot(1,3,2)
wordcloud(documentsAll(idxST));
title("Statistics Theory")

subplot(1,3,3)

```

```
wordcloud(documentsAll(idxCOST));
title("Both")
```



View the number of classes.

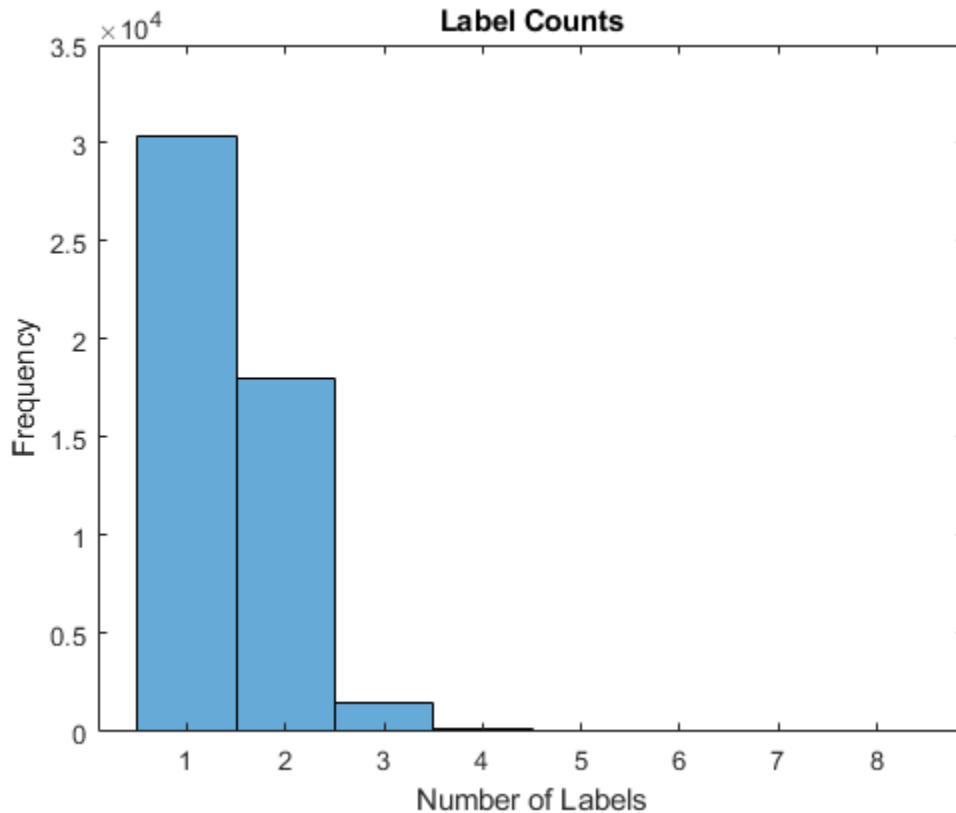
```
classNames = unique(cat(1, labelsAll{:}));
numClasses = numel(classNames)
```

```
numClasses = 32
```

Visualize the number of per-document labels using a histogram.

```
labelCounts = cellfun(@numel, labelsAll);
figure
histogram(labelCounts)
xlabel("Number of Labels")
ylabel("Frequency")
title("Label Counts")
```





### Prepare Text Data for Deep Learning

Partition the data into training and validation partitions using the `cvpartition` function. Hold out 10% of the data for validation by setting the `HoldOut` option to 0.1.

```
cvp = cvpartition(numel(documentsAll),HoldOut=0.1);
documentsTrain = documentsAll(training(cvp));
documentsValidation = documentsAll(test(cvp));
```

```
labelsTrain = labelsAll(training(cvp));
labelsValidation = labelsAll(test(cvp));
```

Create a word encoding object that encodes the training documents as sequences of word indices. Specify a vocabulary of the 5000 words by setting the `Order` option to "frequency", and the `MaxNumWords` option to 5000.

```
enc = wordEncoding(documentsTrain,Order="frequency",MaxNumWords=5000)
```

```
enc =
  wordEncoding with properties:
    NumWords: 5000
    Vocabulary: [1x5000 string]
```

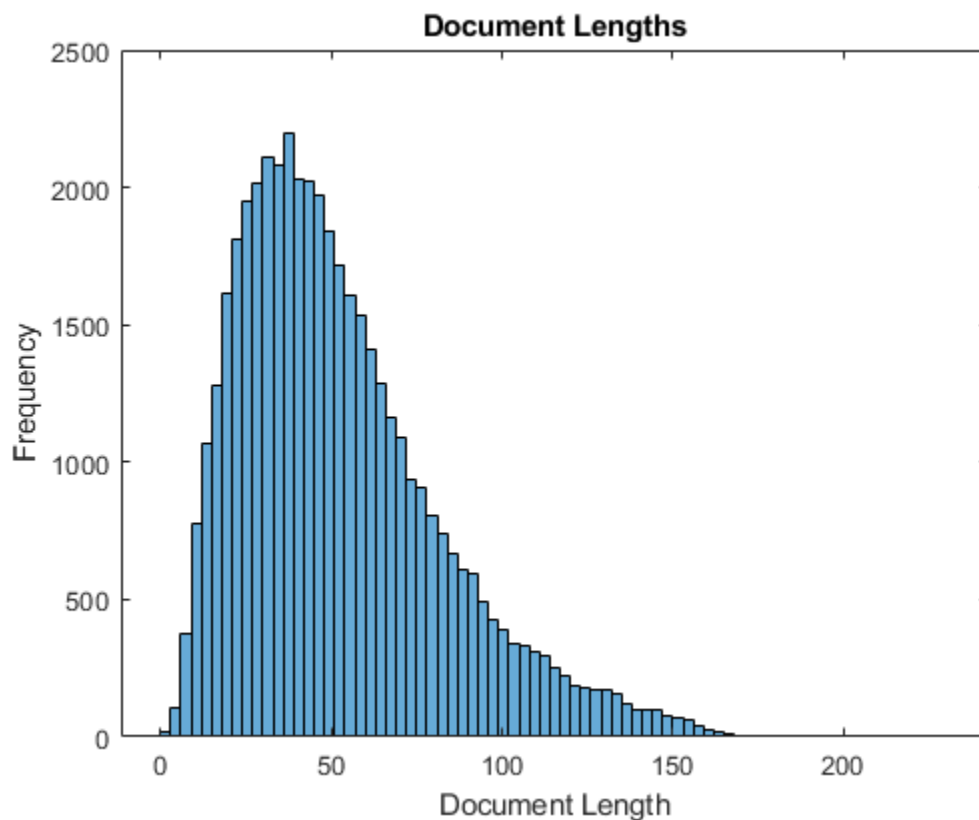
To improve training, use the following techniques:

- 1 When training, truncate the documents to a length that reduces the amount of padding used and does not discard too much data.
- 2 Train for one epoch with the documents sorted by length in ascending order, then shuffle the data each epoch. This technique is known as *sortagrad*.

To choose a sequence length for truncation, visualize the document lengths in a histogram and choose a value that captures most of the data.

```
documentLengths = doclength(documentsTrain);
```

```
figure
histogram(documentLengths)
xlabel("Document Length")
ylabel("Frequency")
title("Document Lengths")
```



Most of the training documents have fewer than 175 tokens. Use 175 tokens as the target length for truncation and padding.

```
maxSequenceLength = 175;
```

To use the sortagrad technique, sort the documents by length in ascending order.

```
[~,idx] = sort(documentLengths);
documentsTrain = documentsTrain(idx);
labelsTrain = labelsTrain(idx);
```

## Define and Initialize Model Parameters

Define the parameters for each of the operations and include them in a struct. Use the format `parameters.OperationName.ParameterName`, where `parameters` is the struct, `OperationName` is the name of the operation (for example "fc"), and `ParameterName` is the name of the parameter (for example, "Weights").

Create a struct `parameters` containing the model parameters. Initialize the bias with zeros. Use the following weight initializers for the operations:

- For the embedding, initialize the weights using the `initializeGaussian` function.
- For the GRU operation, initialize the weights and bias using the `initializeGlorot` and `initializeZeros` functions, respectively.
- For the fully connect operation, initialize the weights and bias using the `initializeGaussian` and `initializeZeros` functions, respectively.

The initialization functions `initializeGlorot`, `initializeGaussian`, and `initializeZeros` are attached to the example as supporting files. To access these functions, open the example as a live script.

Initialize the learnable parameters for the embedding.

```
embeddingDimension = 300;
numHiddenUnits = 250;
inputSize = enc.NumWords + 1;

parameters = struct;

sz = [embeddingDimension inputSize];
mu = 0;
sigma = 0.01;
parameters.emb.Weights = initializeGaussian(sz,mu,sigma);
```

Initialize the learnable parameters for the GRU operation using.

```
sz = [3*numHiddenUnits embeddingDimension];
numOut = 3*numHiddenUnits;
numIn = embeddingDimension;
parameters.gru.InputWeights = initializeGlorot(sz,numOut,numIn);

sz = [3*numHiddenUnits numHiddenUnits];
numOut = 3*numHiddenUnits;
numIn = numHiddenUnits;
parameters.gru.RecurrentWeights = initializeGlorot(sz,numOut,numIn);

sz = [3*numHiddenUnits 1];
parameters.gru.Bias = initializeZeros(sz);
```

Initialize the learnable parameters for the fully connect operation.

```
sz = [numClasses numHiddenUnits];
mu = 0;
sigma = 0.01;
parameters.fc.Weights = initializeGaussian(sz,mu,sigma);

sz = [numClasses 1];
parameters.fc.Bias = initializeZeros(sz);
```

View the parameters struct.

```
parameters

parameters = struct with fields:
    emb: [1×1 struct]
    gru: [1×1 struct]
    fc: [1×1 struct]
```

View the parameters for the GRU operation.

```
parameters.gru

ans = struct with fields:
    InputWeights: [750×300 dlarray]
    RecurrentWeights: [750×250 dlarray]
    Bias: [750×1 dlarray]
```

### **Define Model Function**

Create the function `model`, listed at the end of the example, which computes the outputs of the deep learning model described earlier. The function `model` takes as input the input data and the model parameters. The network outputs the predictions for the labels.

### **Define Model Loss Function**

Create the function `modelLoss`, listed at the end of the example, which takes as input a mini-batch of input data and the corresponding targets, and returns the loss, the gradients of the loss with respect to the learnable parameters, and the network outputs.

### **Specify Training Options**

Train for 5 epochs with a mini-batch size of 256.

```
numEpochs = 5;
miniBatchSize = 256;
```

Train using the Adam optimizer, with a learning rate of 0.01, and specify gradient decay and squared gradient decay factors of 0.5 and 0.999, respectively.

```
learnRate = 0.01;
gradientDecayFactor = 0.5;
squaredGradientDecayFactor = 0.999;
```

Clip the gradients with a threshold of 1 using  $L_2$  norm gradient clipping.

```
gradientThreshold = 1;
```

To convert a vector of probabilities to labels, use the labels with probabilities higher than a specified threshold. Specify a label threshold of 0.5.

```
labelThreshold = 0.5;
```

Validate the network every epoch.

```

numObservationsTrain = numel(documentsTrain);
numIterationsPerEpoch = floor(numObservationsTrain/miniBatchSize);
validationFrequency = numIterationsPerEpoch;

```

### Train Model

Initialize the training progress plot. Create animated lines for the F-score and the loss.

```

figure
C = colororder;

subplot(2,1,1)
lineFScoreTrain = animatedline(Color=C(1,:));
lineFScoreValidation = animatedline( ...
    LineStyle="--", ...
    Marker="o", ...
    MarkerFaceColor="black");
ylim([0 1])
xlabel("Iteration")
ylabel("Labeling F-Score")
grid on

subplot(2,1,2)
lineLossTrain = animatedline(Color=C(2,:));
lineLossValidation = animatedline( ...
    LineStyle="--", ...
    Marker="o", ...
    MarkerFaceColor="black");
ylim([0 inf])
xlabel("Iteration")
ylabel("Loss")
grid on

```

Initialize parameters for the Adam optimizer.

```

trailingAvg = [];
trailingAvgSq = [];

```

Prepare the validation data. Create a one-hot encoded matrix where non-zero entries correspond to the labels of each observation.

```

numObservationsValidation = numel(documentsValidation);
TValidation = zeros(numClasses, numObservationsValidation,"single");
for i = 1:numObservationsValidation
    [~,idx] = ismember(labelsValidation{i},classNames);
    TValidation(idx,i) = 1;
end

```

Train the model using a custom training loop.

For each epoch, loop over mini-batches of data. At the end of each epoch, shuffle the data. At the end of each iteration, update the training progress plot.

For each mini-batch:

- Convert the documents to sequences of word indices and convert the labels to dummy variables.
- Convert the sequences to `darray` objects with underlying type `single` and specify the dimension labels "BTC" (batch, time, channel).

- Train on a GPU if one is available. This requires Parallel Computing Toolbox™. Using a GPU requires Parallel Computing Toolbox™ and a supported GPU device. For information on supported devices, see “GPU Support by Release” (Parallel Computing Toolbox).
- For GPU training, convert to `gpuArray` objects.
- Evaluate the model loss and gradients using `dlfeval` and the `modelLoss` function.
- Clip the gradients.
- Update the network parameters using the `adamupdate` function.
- If necessary, validate the network using the `modelPredictions` function, listed at the end of the example.
- Update the training plot.

```

iteration = 0;
start = tic;

% Loop over epochs.
for epoch = 1:numEpochs

    % Loop over mini-batches.
    for i = 1:numIterationsPerEpoch
        iteration = iteration + 1;
        idx = (i-1)*miniBatchSize+1:i*miniBatchSize;

        % Read mini-batch of data and convert the labels to dummy
        % variables.
        documents = documentsTrain(idx);
        labels = labelsTrain(idx);

        % Convert documents to sequences.
        len = min(maxSequenceLength,max(doclength(documents)));
        X = doc2sequence(enc,documents, ...
            PaddingValue=inputSize, ...
            Length=len);
        X = cat(1,X{:});

        % Dummify labels.
        T = zeros(numClasses,miniBatchSize,"single");
        for j = 1:miniBatchSize
            [~,idx2] = ismember(labels{j},classNames);
            T(idx2,j) = 1;
        end

        % Convert mini-batch of data to dlarray.
        X = dlarray(X,"BTC");

        % If training on a GPU, then convert data to gpuArray.
        if canUseGPU
            X = gpuArray(X);
        end

        % Evaluate the model loss, gradients, and predictions using dlfeval and the
        % modelLoss function.
        [loss,gradients,Y] = dlfeval(@modelLoss,X,T,parameters);

        % Gradient clipping.
        gradients = dlupdate(@(g) thresholdL2Norm(g,gradientThreshold),gradients);

```

```

% Update the network parameters using the Adam optimizer.
[parameters,trailingAvg,trailingAvgSq] = adamupdate(parameters,gradients, ...
    trailingAvg,trailingAvgSq,iteration,learnRate, ...
    gradientDecayFactor,squaredGradientDecayFactor);

% Display the training progress.
subplot(2,1,1)
D = duration(0,0,toc(start),Format="hh:mm:ss");
title("Epoch: " + epoch + ", Elapsed: " + string(D))

% Loss.
loss = double(loss);
addpoints(lineLossTrain,iteration,loss)

% Labeling F-score.
Y = Y > labelThreshold;
score = labelingFScore(Y,T);
addpoints(lineFScoreTrain,iteration,double(gather(score)))

drawnow

% Display validation metrics.
if iteration == 1 || mod(iteration,validationFrequency) == 0
    YValidation = modelPredictions(parameters,enc,documentsValidation,miniBatchSize,maxS

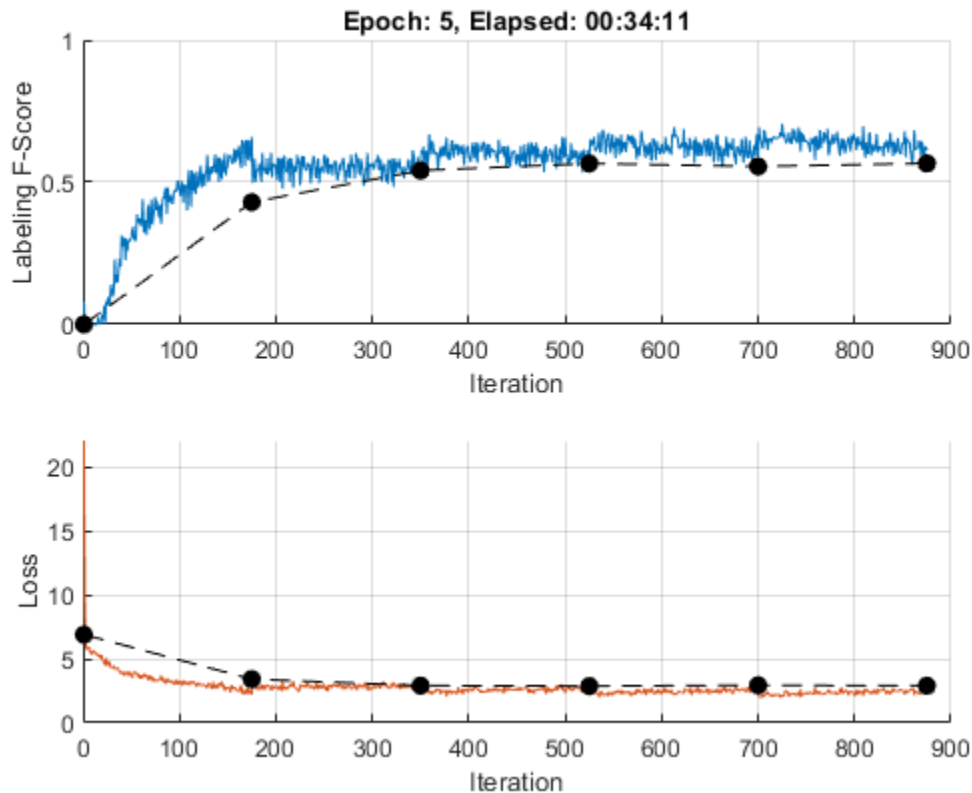
    % Loss.
    lossValidation = crossentropy(YValidation,TValidation, ...
        TargetCategories="independent", ...
        DataFormat="CB");
    lossValidation = double(lossValidation);
    addpoints(lineLossValidation,iteration,lossValidation)

    % Labeling F-score.
    YValidation = YValidation > labelThreshold;
    score = labelingFScore(YValidation,TValidation);
    score = double(score);
    addpoints(lineFScoreValidation,iteration,score)

    drawnow
end
end

% Shuffle data.
idx = randperm(numObservationsTrain);
documentsTrain = documentsTrain(idx);
labelsTrain = labelsTrain(idx);
end

```



### Test Model

To make predictions on a new set of data, use the `modelPredictions` function, listed at the end of the example. The `modelPredictions` function takes as input the model parameters, a word encoding, and an array of tokenized documents, and outputs the model predictions corresponding to the specified mini-batch size and the maximum sequence length.

```
YValidation = modelPredictions(parameters,enc,documentsValidation,miniBatchSize,maxSequenceLength)
```

To evaluate the performance, calculate the labeling F-score using the `labelingFScore` function, listed at the end of the example. The labeling F-score evaluates multilabel classification by focusing on per-text classification with partial matches. To convert the network outputs to an array of labels, find the labels with scores higher than the specified label threshold.

```
score = labelingFScore(YValidation > labelThreshold,TValidation)
```

```
score = single
       0.5663
```

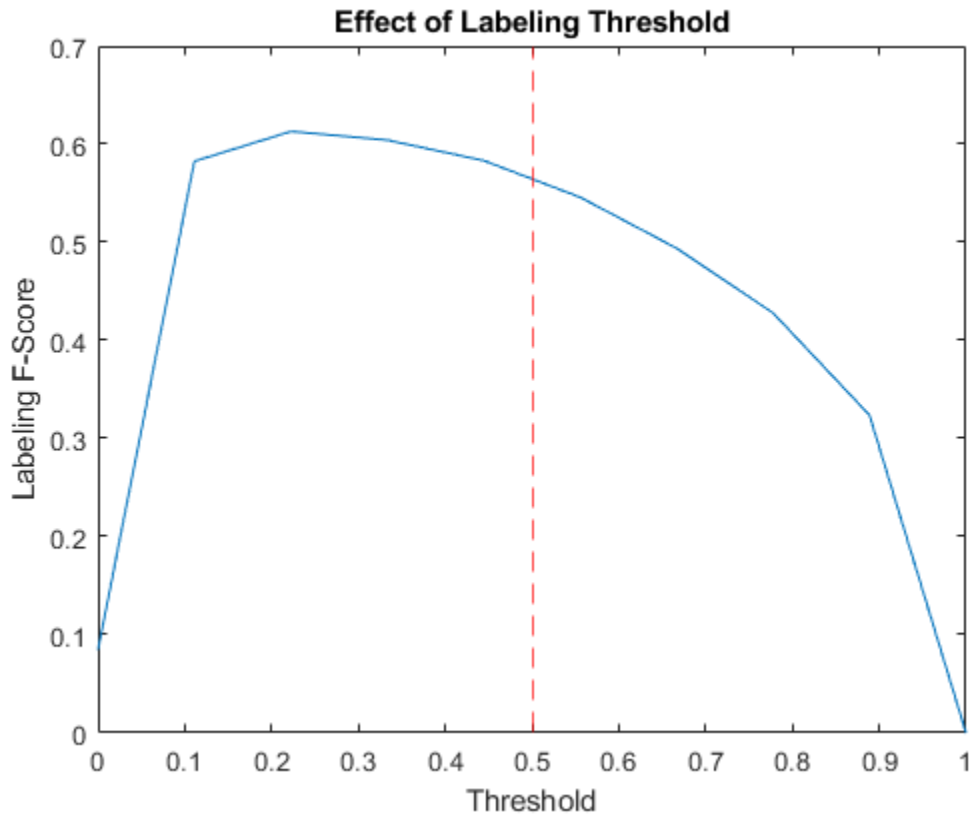
View the effect of the labeling threshold on the labeling F-score by trying a range of values for the threshold and comparing the results.

```
thr = linspace(0,1,10);
score = zeros(size(thr));
for i = 1:numel(thr)
    YPredValidationThr = YValidation >= thr(i);
    score(i) = labelingFScore(YPredValidationThr,TValidation);
end
```



```
end
```

```
figure
plot(thr,score)
xline(labelThreshold,"r--");
xlabel("Threshold")
ylabel("Labeling F-Score")
title("Effect of Labeling Threshold")
```



### Visualize Predictions

To visualize the correct predictions of the classifier, calculate the numbers of true positives. A true positive is an instance of a classifier correctly predicting a particular class for an observation.

```
Y = YValidation > labelThreshold;
T = TValidation;
```

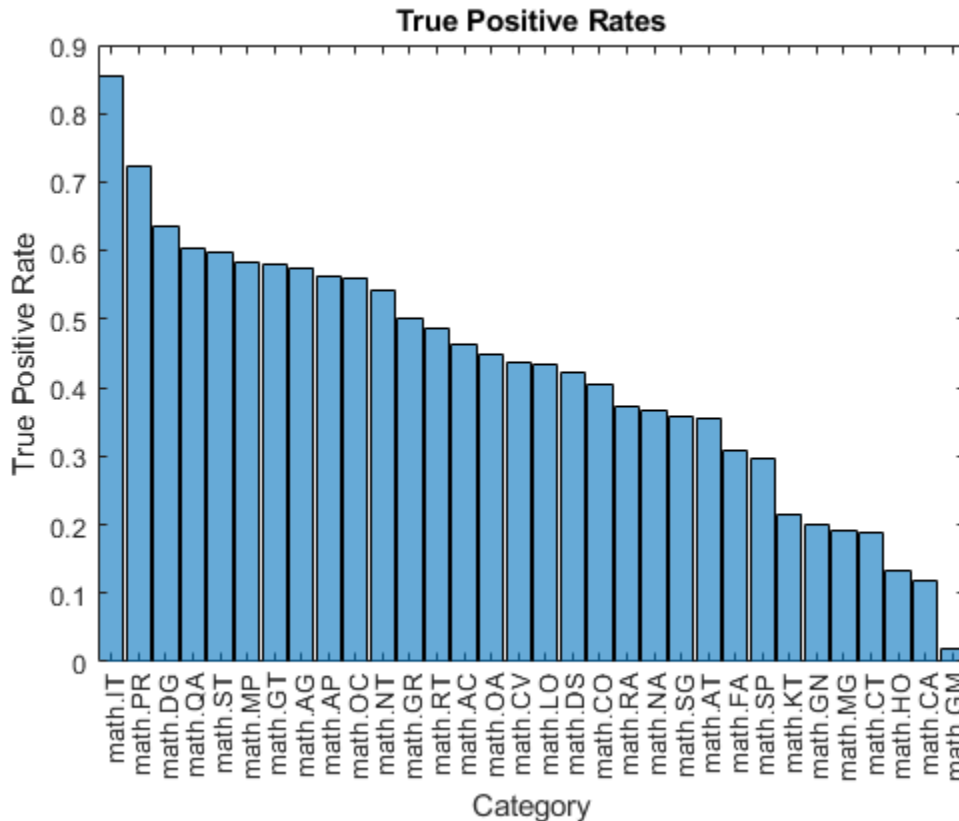
```
numTruePositives = sum(T & Y,2);
```

```
numObservationsPerClass = sum(T,2);
truePositiveRates = numTruePositives ./ numObservationsPerClass;
```

Visualize the numbers of true positives for each class in a histogram.

```
figure
truePositiveRates = extractdata(truePositiveRates);
[~,idx] = sort(truePositiveRates,"descend");
histogram(Categories=classNames(idx),BinCounts=truePositiveRates(idx))
```

```
xlabel("Category")
ylabel("True Positive Rate")
title("True Positive Rates")
```



Visualize the instances where the classifier predicts incorrectly by showing the distribution of true positives, false positives, and false negatives. A false positive is an instance of a classifier assigning a particular incorrect class to an observation. A false negative is an instance of a classifier failing to assign a particular correct class to an observation.

Create a confusion matrix showing the true positive, false positive, and false negative counts:

- For each class, display the true positive counts on the diagonal.
- For each pair of classes  $(i, j)$ , display the number of instances of a false positive for  $j$  when the instance is also a false negative for  $i$ .

That is, the confusion matrix with elements given by:

$$\text{TPFN}_{ij} = \begin{cases} \text{numTruePositives}(i), & \text{if } i = j \\ \text{numFalsePositives}(j | i \text{ is a false negative}), & \text{if } i \neq j \end{cases}$$

Calculate the false negatives and false positives.

```
falseNegatives = T & ~Y;
falsePositives = ~T & Y;
```

Calculate the off-diagonal elements.

```
falseNegatives = permute(falseNegatives,[3 2 1]);
numConditionalFalsePositives = sum(falseNegatives & falsePositives, 2);
numConditionalFalsePositives = squeeze(numConditionalFalsePositives);
```

```
tpfnMatrix = numConditionalFalsePositives;
```

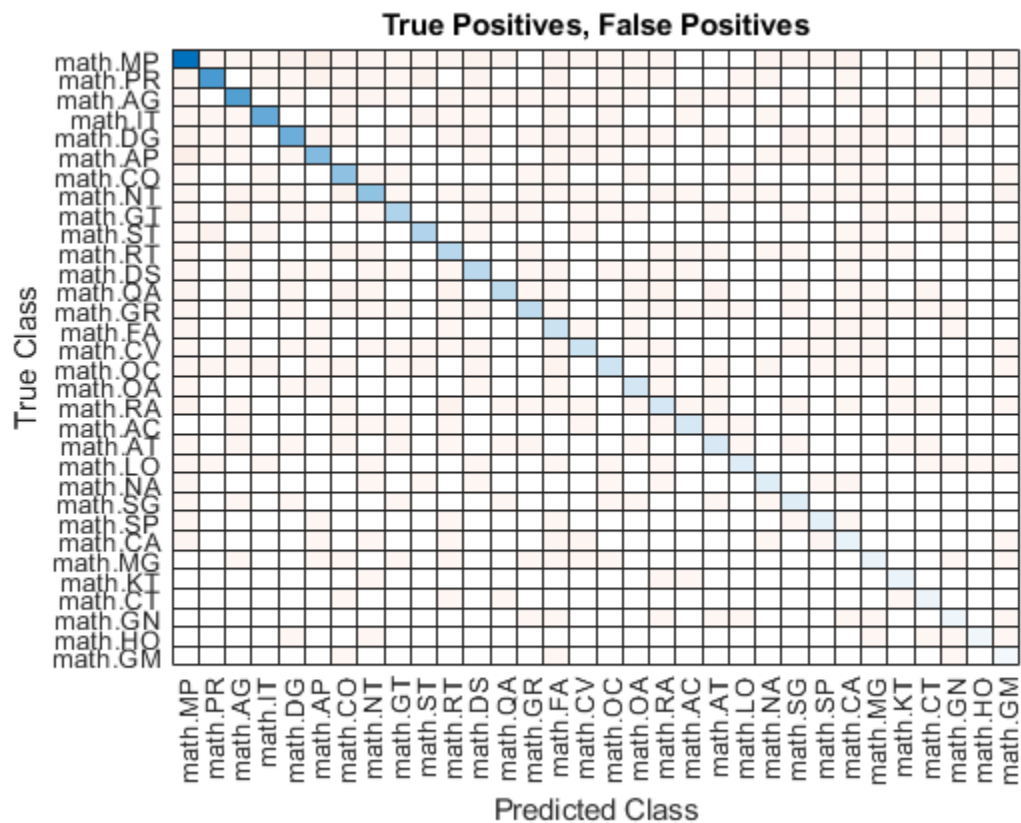
Set the diagonal elements to the true positive counts.

```
idxDiagonal = 1:numClasses+1:numClasses^2;
tpfnMatrix(idxDiagonal) = numTruePositives;
```

Visualize the true positive and false positive counts in a confusion matrix using the `confusionchart` function and sort the matrix such that the elements on the diagonal are in descending order.

```
figure
```

```
tpfnMatrix = extractdata(tpfnMatrix);
cm = confusionchart(tpfnMatrix,classNames);
sortClasses(cm,"descending-diagonal");
title("True Positives, False Positives")
```



To view the matrix in more detail, open [this example](#) as a live script and open the figure in a new window.

### Preprocess Text Function

The `preprocessText` function tokenizes and preprocesses the input text data using the following steps:

- 1 Tokenize the text using the `tokenizedDocument` function. Extract mathematical equations as a single token using the `RegularExpressions` option by specifying the regular expression `"\$.*?\$"`, which captures text appearing between two "\$" symbols.
- 2 Erase the punctuation using the `erasePunctuation` function.
- 3 Convert the text to lowercase using the `lower` function.
- 4 Remove the stop words using the `removeStopWords` function.
- 5 Lemmatize the text using the `normalizeWords` function with the `Style` option set to "lemma".

```
function documents = preprocessText(textData)

% Tokenize the text.
regularExpressions = table;
regularExpressions.Pattern = "\$.*?\$";
regularExpressions.Type = "equation";

documents = tokenizedDocument(textData,RegularExpressions=regularExpressions);

% Erase punctuation.
documents = erasePunctuation(documents);

% Convert to lowercase.
documents = lower(documents);

% Lemmatize.
documents = addPartOfSpeechDetails(documents);
documents = normalizeWords(documents,Style="lemma");

% Remove stop words.
documents = removeStopWords(documents);

% Remove short words.
documents = removeShortWords(documents,2);

end
```

### **Model Function**

The function `model` takes as input the input data `X` and the model parameters `parameters`, and returns the predictions for the labels.

```
function Y = model(X,parameters)

% Embedding
weights = parameters.emb.Weights;
X = embed(X,weights);

% GRU
inputWeights = parameters.gru.InputWeights;
recurrentWeights = parameters.gru.RecurrentWeights;
bias = parameters.gru.Bias;

numHiddenUnits = size(inputWeights,1)/3;
hiddenState = dlarray(zeros([numHiddenUnits 1]));

Y = gru(X,hiddenState,inputWeights,recurrentWeights,bias);
```

```

% Max pooling along time dimension
Y = max(Y,[],3);

% Fully connect
weights = parameters.fc.Weights;
bias = parameters.fc.Bias;
Y = fullyconnect(Y,weights,bias);

% Sigmoid
Y = sigmoid(Y);

end

```

### Model Loss Function

The `modelLoss` function takes as input a mini-batch of input data  $X$  with corresponding targets  $T$  containing the labels and returns the loss, the gradients of the loss with respect to the learnable parameters, and the network outputs.

```

function [loss,gradients,Y] = modelLoss(X,T,parameters)

Y = model(X,parameters);

loss = crossentropy(Y,T,TargetCategories="independent");

gradients = dlgradient(loss,parameters);

end

```

### Model Predictions Function

The `modelPredictions` function takes as input the model parameters, a word encoding, an array of tokenized documents, a mini-batch size, and a maximum sequence length, and returns the model predictions by iterating over mini-batches of the specified size.

```

function Y = modelPredictions(parameters,enc,documents,miniBatchSize,maxSequenceLength)

inputSize = enc.NumWords + 1;

numObservations = numel(documents);
numIterations = ceil(numObservations / miniBatchSize);

numFeatures = size(parameters.fc.Weights,1);
Y = zeros(numFeatures,numObservations,"like",parameters.fc.Weights);

for i = 1:numIterations

    idx = (i-1)*miniBatchSize+1:min(i*miniBatchSize,numObservations);

    len = min(maxSequenceLength,max(doclength(documents(idx))));
    X = doc2sequence(enc,documents(idx), ...
        PaddingValue=inputSize, ...
        Length=len);
    X = cat(1,X{:});

    X = dlarray(X,"BTC");

    Y(:,idx) = model(X,parameters);

```

```
end
```

```
end
```

### Labeling F-Score Function

The labeling F-score function [2] evaluates multilabel classification by focusing on per-text classification with partial matches. The measure is the normalized proportion of matching labels against the total number of true and predicted labels given by

$$\frac{1}{N} \sum_{n=1}^N \left( \frac{2 \sum_{c=1}^C Y_{nc} T_{nc}}{\sum_{c=1}^C (Y_{nc} + T_{nc})} \right),$$

where  $N$  and  $C$  correspond to the number of observations and classes, respectively, and  $Y$  and  $T$  correspond to the predictions and targets, respectively.

```
function score = labelingFScore(Y,T)
numObservations = size(T,2);
scores = (2 * sum(Y .* T)) ./ sum(Y + T);
score = sum(scores) / numObservations;
end
```

### Gradient Clipping Function

The `thresholdL2Norm` function scales the input gradients so that their  $L_2$  norm values equal the specified gradient threshold when the  $L_2$  norm value of the gradient of a learnable parameter is larger than the specified threshold.

```
function gradients = thresholdL2Norm(gradients,gradientThreshold)
gradientNorm = sqrt(sum(gradients(:).^2));
if gradientNorm > gradientThreshold
    gradients = gradients * (gradientThreshold / gradientNorm);
end
end
```

### References

- 1 arXiv. "arXiv API." Accessed January 15, 2020. <https://arxiv.org/help/api>
- 2 Sokolova, Marina, and Guy Lapalme. "A Sytematic Analysis of Performance Measures for Classification Tasks." *Information Processing & Management* 45, no. 4 (2009): 427-437.

### See Also

`tokenizedDocument` | `fullyconnect` | `d\lupdate` | `adamupdate` | `d\larray` | `d\lfeval` | `d\lgradient` | `wordEncoding` | `doc2sequence`

### Related Examples

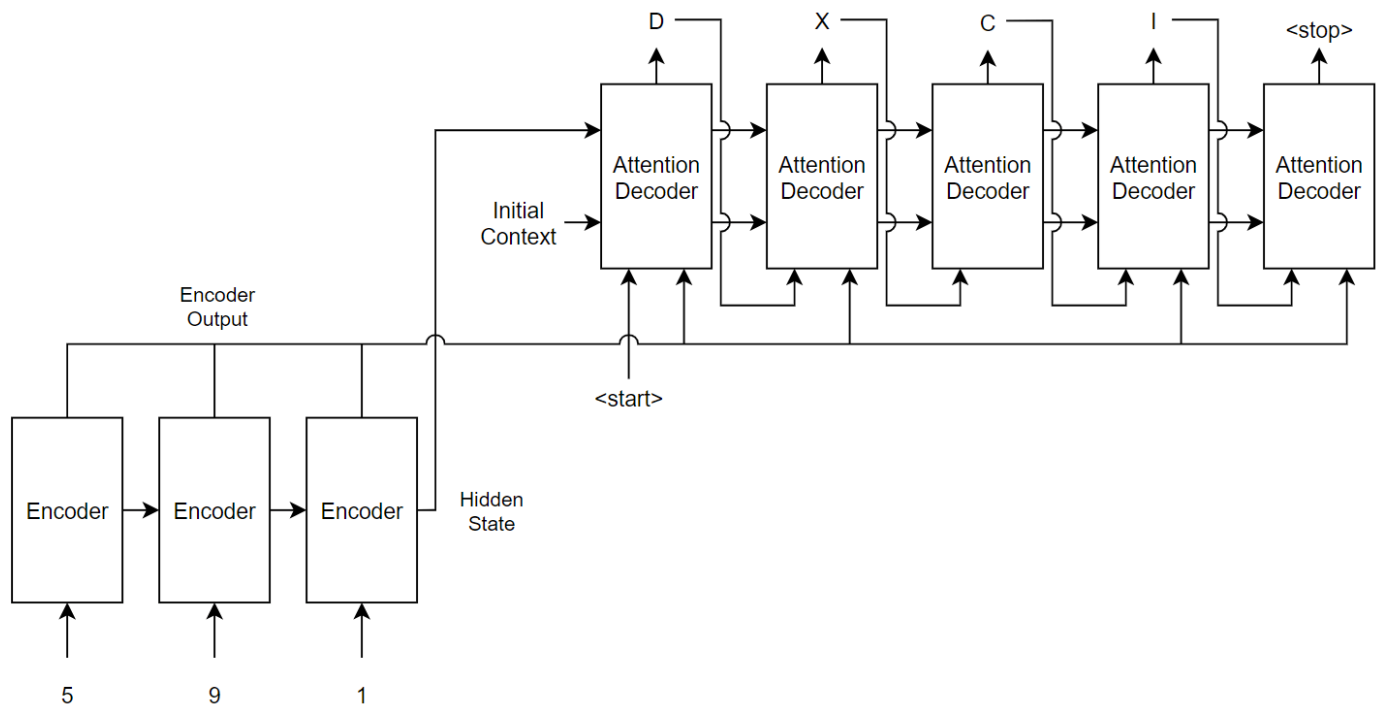
- "Classify Text Data Using Deep Learning" on page 2-85

- “Create Simple Text Model for Classification” on page 2-2
- “Deep Learning in MATLAB” (Deep Learning Toolbox)

## Sequence-to-Sequence Translation Using Attention

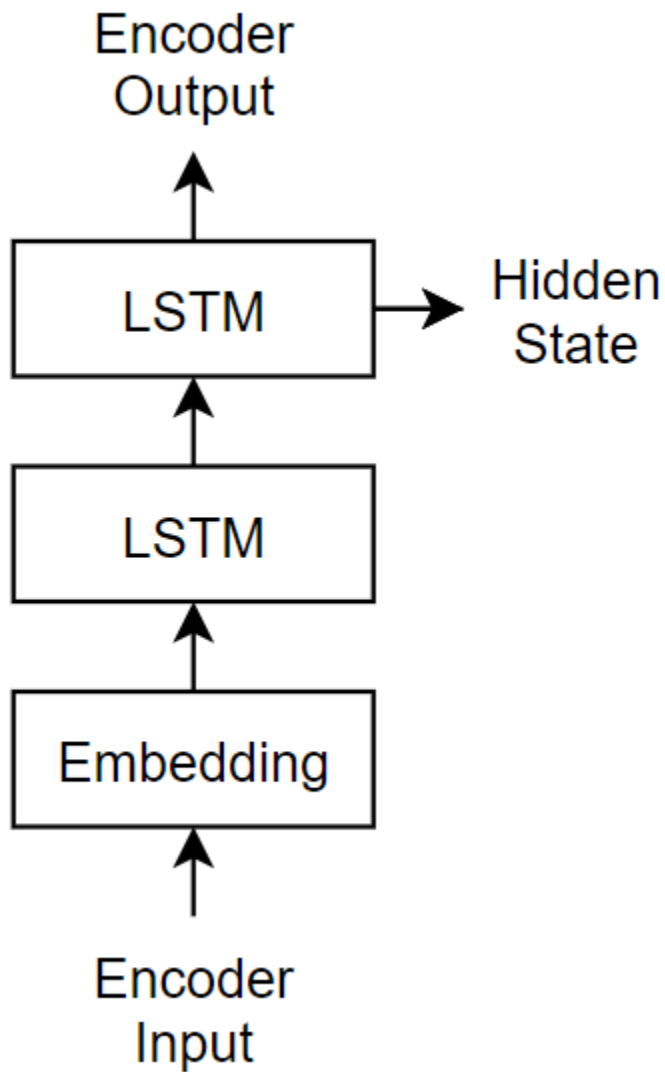
This example shows how to convert decimal strings to Roman numerals using a recurrent sequence-to-sequence encoder-decoder model with attention.

Recurrent encoder-decoder models have proven successful at tasks like abstractive text summarization and neural machine translation. The model consists of an *encoder* which typically processes input data with a recurrent layer such as LSTM, and a *decoder* which maps the encoded input into the desired output, typically with a second recurrent layer. Models that incorporate *attention mechanisms* into the models allows the decoder to focus on parts of the encoded input while generating the translation.

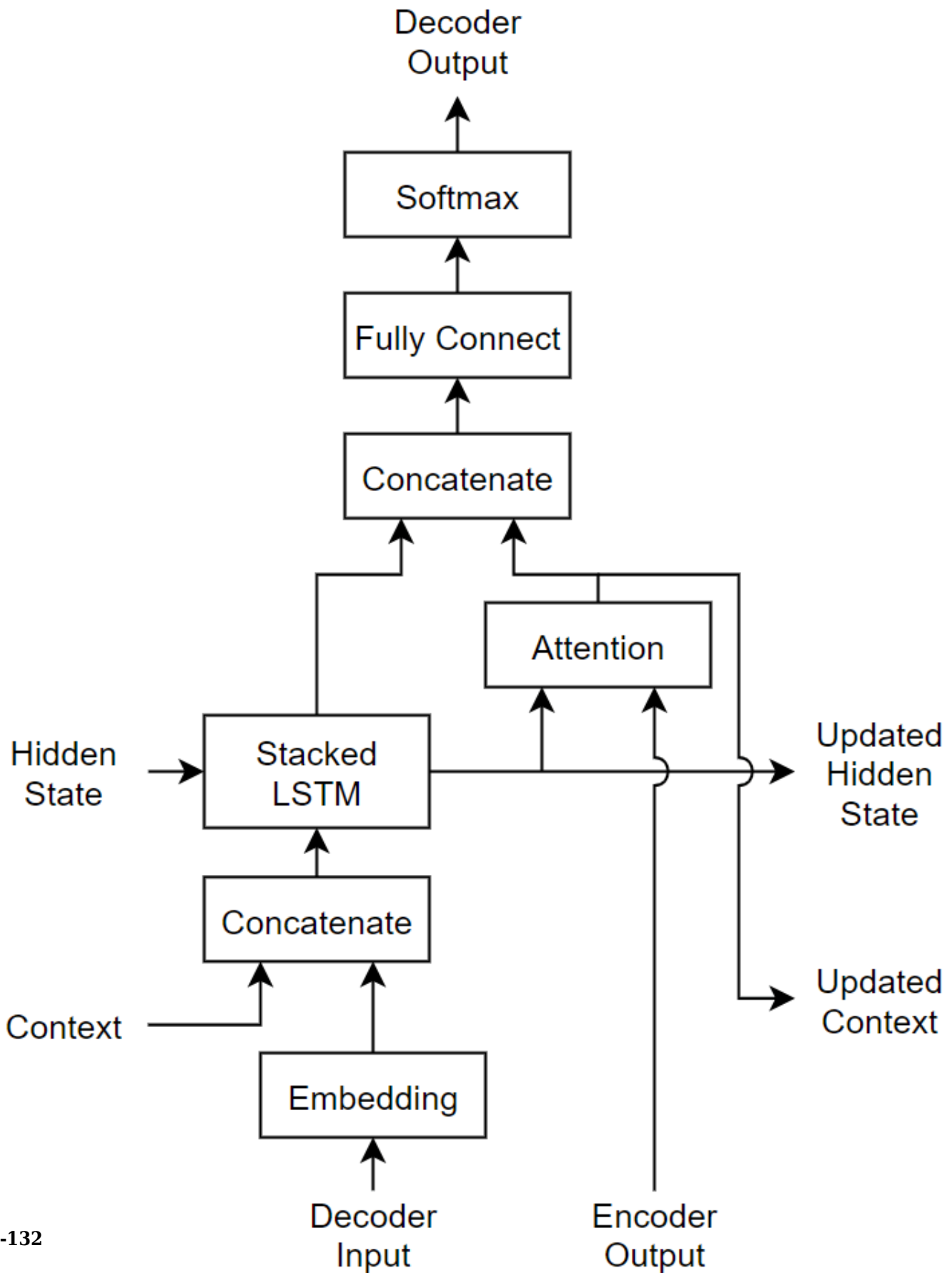


For the encoder model, this example uses a simple network consisting of an embedding followed by two LSTM operations. Embedding is a method of converting categorical tokens into numeric vectors.





For the decoder model, this example uses a network very similar to the encoder that contains two LSTMs. However, an important difference is that the decoder contains an attention mechanism. The attention mechanism allows the decoder to *attend* to specific parts of the encoder output.



## Load Training Data

Download the decimal-Roman numeral pairs from "romanNumerals.csv"

```
filename = fullfile("romanNumerals.csv");

options = detectImportOptions(filename, ...
    TextType="string", ...
    ReadVariableNames=false);
options.VariableNames = ["Source" "Target"];
options.VariableTypes = ["string" "string"];

data = readtable(filename,options);
```

Split the data into training and test partitions containing 50% of the data each.

```
idx = randperm(size(data,1),500);
dataTrain = data(idx,:);
dataTest = data;
dataTest(idx,:) = [];
```

View some of the decimal-Roman numeral pairs.

```
head(dataTrain)
```

```
ans=8x2 table
    Source    Target
    -----
    "713"    "DCCXIII"
    "752"    "DCCLII"
    "434"    "CDXXXIV"
    "641"    "DCXLI"
    "68"     "LXVIII"
    "87"     "LXXXVII"
    "469"    "CDLXIX"
    "715"    "DCCXV"
```

## Preprocess Data

Preprocess the text data using the `transformText` function, listed at the end of the example. The `transformText` function preprocesses and tokenizes the input text for translation by splitting the text into characters and adding start and stop tokens. To translate text by splitting the text into words instead of characters, skip the first step.

```
startToken = "<start>";
stopToken = "<stop>";

strSource = dataTrain.Source;
documentsSource = transformText(strSource,startToken,stopToken);
```

Create a `wordEncoding` object that maps tokens to a numeric index and vice-versa using a vocabulary.

```
encSource = wordEncoding(documentsSource);
```

Using the word encoding, convert the source text data to numeric sequences.

```
sequencesSource = doc2sequence(encSource,documentsSource,PaddingDirection="none");
```

Convert the target data to sequences using the same steps.

```
strTarget = dataTrain.Target;  
documentsTarget = transformText(strTarget,startToken,stopToken);  
encTarget = wordEncoding(documentsTarget);  
sequencesTarget = doc2sequence(encTarget,documentsTarget,PaddingDirection="none");
```

Sort the sequences by length. Training with the sequences sorted by increasing sequence length results in batches with sequences of approximately the same sequence length and ensures smaller sequence batches are used to update the model before longer sequence batches.

```
sequenceLengths = cellfun(@(sequence) size(sequence,2),sequencesSource);  
[~,idx] = sort(sequenceLengths);  
sequencesSource = sequencesSource(idx);  
sequencesTarget = sequencesTarget(idx);
```

Create `arrayDatastore` objects containing the source and target data and combine them using the `combine` function.

```
sequencesSourceDs = arrayDatastore(sequencesSource,OutputType="same");  
sequencesTargetDs = arrayDatastore(sequencesTarget,OutputType="same");
```

```
sequencesDs = combine(sequencesSourceDs,sequencesTargetDs);
```

### **Initialize Model Parameters**

Initialize the model parameters. For both the encoder and decoder, specify an embedding dimension of 128, two LSTM layers with 200 hidden units, and dropout layers with random dropout with probability 0.05.

```
embeddingDimension = 128;  
numHiddenUnits = 200;  
dropout = 0.05;
```

### **Initialize Encoder Model Parameters**

Initialize the weights of the encoding embedding using the Gaussian using the `initializeGaussian` function which is attached to this example as a supporting file. Specify a mean of 0 and a standard deviation of 0.01. To learn more, see “Gaussian Initialization” (Deep Learning Toolbox).

```
inputSize = encSource.NumWords + 1;  
sz = [embeddingDimension inputSize];  
mu = 0;  
sigma = 0.01;  
parameters.encoder.emb.Weights = initializeGaussian(sz,mu,sigma);
```

Initialize the learnable parameters for the encoder LSTM operations:

- Initialize the input weights with the Glorot initializer using the `initializeGlorot` function which is attached to this example as a supporting file. To learn more, see “Glorot Initialization” (Deep Learning Toolbox).
- Initialize the recurrent weights with the orthogonal initializer using the `initializeOrthogonal` function which is attached to this example as a supporting file. To learn more, see “Orthogonal Initialization” (Deep Learning Toolbox).

- Initialize the bias with the unit forget gate initializer using the `initializeUnitForgetGate` function which is attached to this example as a supporting file. To learn more, see “Unit Forget Gate Initialization” (Deep Learning Toolbox).

Initialize the learnable parameters for the first encoder LSTM operation.

```
sz = [4*numHiddenUnits embeddingDimension];
numOut = 4*numHiddenUnits;
numIn = embeddingDimension;

parameters.encoder.lstm1.InputWeights = initializeGlorot(sz,numOut,numIn);
parameters.encoder.lstm1.RecurrentWeights = initializeOrthogonal([4*numHiddenUnits numHiddenUnits]);
parameters.encoder.lstm1.Bias = initializeUnitForgetGate(numHiddenUnits);
```

Initialize the learnable parameters for the second encoder LSTM operation.

```
sz = [4*numHiddenUnits numHiddenUnits];
numOut = 4*numHiddenUnits;
numIn = numHiddenUnits;

parameters.encoder.lstm2.InputWeights = initializeGlorot(sz,numOut,numIn);
parameters.encoder.lstm2.RecurrentWeights = initializeOrthogonal([4*numHiddenUnits numHiddenUnits]);
parameters.encoder.lstm2.Bias = initializeUnitForgetGate(numHiddenUnits);
```

### Initialize Decoder Model Parameters

Initialize the weights of the encoding embedding using the Gaussian using the `initializeGaussian` function. Specify a mean of 0 and a standard deviation of 0.01.

```
outputSize = encTarget.NumWords + 1;
sz = [embeddingDimension outputSize];
mu = 0;
sigma = 0.01;
parameters.decoder.emb.Weights = initializeGaussian(sz,mu,sigma);
```

Initialize the weights of the attention mechanism using the Glorot initializer using the `initializeGlorot` function.

```
sz = [numHiddenUnits numHiddenUnits];
numOut = numHiddenUnits;
numIn = numHiddenUnits;
parameters.decoder.attn.Weights = initializeGlorot(sz,numOut,numIn);
```

Initialize the learnable parameters for the decoder LSTM operations:

- Initialize the input weights with the Glorot initializer using the `initializeGlorot` function.
- Initialize the recurrent weights with the orthogonal initializer using the `initializeOrthogonal` function.
- Initialize the bias with the unit forget gate initializer using the `initializeUnitForgetGate` function.

Initialize the learnable parameters for the first decoder LSTM operation.

```
sz = [4*numHiddenUnits embeddingDimension+numHiddenUnits];
numOut = 4*numHiddenUnits;
numIn = embeddingDimension + numHiddenUnits;
```

```
parameters.decoder.lstm1.InputWeights = initializeGlorot(sz,numOut,numIn);
parameters.decoder.lstm1.RecurrentWeights = initializeOrthogonal([4*numHiddenUnits numHiddenUnits]);
parameters.decoder.lstm1.Bias = initializeUnitForgetGate(numHiddenUnits);
```

Initialize the learnable parameters for the second decoder LSTM operation.

```
sz = [4*numHiddenUnits numHiddenUnits];
numOut = 4*numHiddenUnits;
numIn = numHiddenUnits;
```

```
parameters.decoder.lstm2.InputWeights = initializeGlorot(sz,numOut,numIn);
parameters.decoder.lstm2.RecurrentWeights = initializeOrthogonal([4*numHiddenUnits numHiddenUnits]);
parameters.decoder.lstm2.Bias = initializeUnitForgetGate(numHiddenUnits);
```

Initialize the learnable parameters for the decoder fully connected operation:

- Initialize the weights with the Glorot initializer.
- Initialize the bias with zeros using the `initializeZeros` function which is attached to this example as a supporting file. To learn more, see “Zeros Initialization” (Deep Learning Toolbox).

```
sz = [outputSize 2*numHiddenUnits];
numOut = outputSize;
numIn = 2*numHiddenUnits;
```

```
parameters.decoder.fc.Weights = initializeGlorot(sz,numOut,numIn);
parameters.decoder.fc.Bias = initializeZeros([outputSize 1]);
```

### **Define Model Functions**

Create the functions `modelEncoder` and `modelDecoder`, listed at the end of the example, that compute the outputs of the encoder and decoder models, respectively.

The `modelEncoder` function, listed in the Encoder Model Function on page 2-0 section of the example, takes the input data, the model parameters, the optional mask that is used to determine the correct outputs for training and returns the model outputs and the LSTM hidden state.

The `modelDecoder` function, listed in the Decoder Model Function on page 2-0 section of the example, takes the input data, the model parameters, the context vector, the LSTM initial hidden state, the outputs of the encoder, and the dropout probability and outputs the decoder output, the updated context vector, the updated LSTM state, and the attention scores.

### **Define Model Loss Function**

Create the function `modelLoss`, listed in the Model LossFunction on page 2-0 section of the example, that takes the encoder and decoder model parameters, a mini-batch of input data and the padding masks corresponding to the input data, and the dropout probability and returns the loss and the gradients of the loss with respect to the learnable parameters in the models.

### **Specify Training Options**

Train with a mini-batch size of 32 for 75 epochs with a learning rate of 0.002.

```
miniBatchSize = 32;
numEpochs = 75;
learnRate = 0.002;
```

Initialize the options from Adam.

```
gradientDecayFactor = 0.9;
squaredGradientDecayFactor = 0.999;
```

### Train Model

Train the model using a custom training loop. Use `minibatchqueue` to process and manage mini-batches of images during training. For each mini-batch:

- Use the custom mini-batch preprocessing function `preprocessMiniBatch` (defined at the end of this example) to find the lengths of all sequence in the mini-batch and pad the sequences to the same length as the longest sequence, for the source and target sequences, respectively.
- Permute the second and third dimensions of the padded sequences.
- Return the mini-batch variables unformatted `dlarray` objects with underlying data type `single`. All other outputs are arrays of data type `single`.
- Train on a GPU if one is available. Return all mini-batch variables on the GPU if one is available. Using a GPU requires Parallel Computing Toolbox™ and a supported GPU device. For information on supported devices, see GPU Support by Release.

The `minibatchqueue` object returns four output arguments for each mini-batch: the source sequences, the target sequences, the lengths of all source sequences in the mini-batch, and the sequence mask of the target sequences.

```
numMiniBatchOutputs = 4;

mbq = minibatchqueue(sequencesDs,numMiniBatchOutputs,...
    MiniBatchSize=miniBatchSize,...
    MiniBatchFcn=@(x,t) preprocessMiniBatch(x,t,inputSize,outputSize));
```

Initialize the training progress plot.

```
figure
C = colororder;
lineLossTrain = animatedline(Color=C(2,:));
ylim([0 inf])
xlabel("Iteration")
ylabel("Loss")
grid on
```

Initialize the values for the `adamupdate` function.

```
trailingAvg = [];
trailingAvgSq = [];
```

Train the model. For each mini-batch:

- Read a mini-batch of padded sequences.
- Compute loss and gradients.
- Update the encoder and decoder model parameters using the `adamupdate` function.
- Update the training progress plot.

```
iteration = 0;
start = tic;

% Loop over epochs.
for epoch = 1:numEpochs
```

```

reset(mbq);

% Loop over mini-batches.
while hasdata(mbq)

    iteration = iteration + 1;

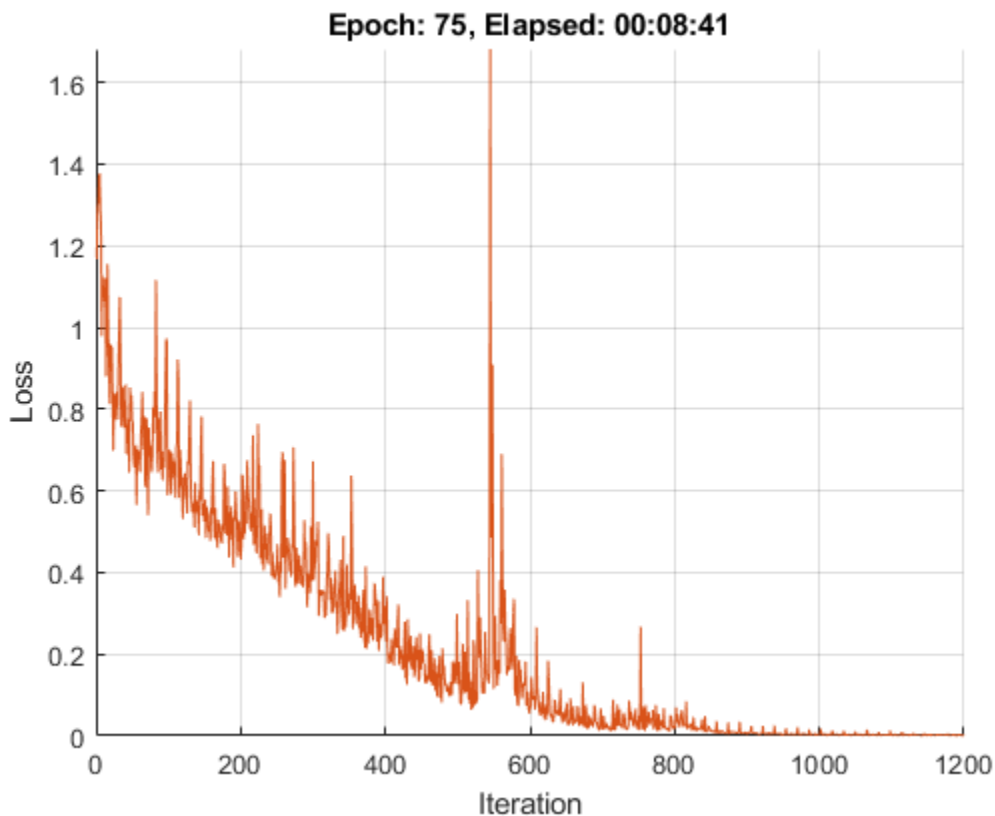
    [X,T,sequenceLengthsSource,maskSequenceTarget] = next(mbq);

    % Compute loss and gradients.
    [loss,gradients] = dlfeval(@modelLoss,parameters,X,T,sequenceLengthsSource,...
        maskSequenceTarget,dropout);

    % Update parameters using adamupdate.
    [parameters,trailingAvg,trailingAvgSq] = adamupdate(parameters,gradients,trailingAvg,tra
        iteration,learnRate,gradientDecayFactor,squaredGradientDecayFactor);

    % Display the training progress. Normalize loss by sequence length.
    D = duration(0,0,toc(start),Format="hh:mm:ss");
    loss = double(loss);
    loss = loss ./ size(T,3);
    addpoints(lineLossTrain,iteration,loss)
    title("Epoch: " + epoch + ", Elapsed: " + string(D))
    drawnow
end
end

```





## Generate Translations

To generate translations for new data using the trained model, convert the text data to numeric sequences using the same steps as when training and input the sequences into the encoder-decoder model and convert the resulting sequences back into text using the token indices.

Preprocess the text data using the same steps as when training. Use the `transformText` function, listed at the end of the example, to split the text into characters and add the start and stop tokens.

```
strSource = dataTest.Source;
strTarget = dataTest.Target;
```

Translate the text using the `modelPredictions` function.

```
maxSequenceLength = 10;
delimiter = "";
```

```
strTranslated = translateText(parameters, strSource, maxSequenceLength, miniBatchSize, ...
    encSource, encTarget, startToken, stopToken, delimiter);
```

Create a table containing the test source text, target text, and translations.

```
tbl = table;
tbl.Source = strSource;
tbl.Target = strTarget;
tbl.Translated = strTranslated;
```

View a random selection of the translations.

```
idx = randperm(size(dataTest,1), miniBatchSize);
tbl(idx, :)
```

*ans=32×3 table*

Source	Target	Translated
"318"	"CCCXVIII"	"CCCXVIII"
"461"	"CDLXI"	"DCLII"
"71"	"LXXI"	"DCXXI"
"520"	"DXX"	"DCX"
"490"	"CDXC"	"CDX"
"526"	"DXXVI"	"DCXVI"
"898"	"DCCCXCVIII"	"DCCCXCVIII"
"67"	"LXVII"	"DCLXVII"
"4"	"IV"	"CLII"
"53"	"LIII"	"DCXIII"
"116"	"CXVI"	"CXVI"
"362"	"CCCLXII"	"CCCLIII"
"709"	"DCCIX"	"DCCII"
"291"	"CCXCI"	"CCXCI"
"390"	"CCCXC"	"CCCX"
"336"	"CCCXXXVI"	"CCCXXVI"
⋮		

### Text Transformation Function

The `transformText` function preprocesses and tokenizes the input text for translation by splitting the text into characters and adding start and stop tokens. To translate text by splitting the text into words instead of characters, skip the first step.

```
function documents = transformText(str,startToken,stopToken)

str = strip(replace(str," "," "));
str = startToken + str + stopToken;
documents = tokenizedDocument(str,CustomTokens=[startToken stopToken]);

end
```

### Mini-Batch Preprocessing Function

The `preprocessMiniBatch` function, described in the Train Model section of the example, preprocesses the data for training. The function preprocesses the data using the following steps:

- 1 Determine the lengths of all source and target sequences in the mini-batch
- 2 Pad the sequences to the same length as the longest sequence in the mini-batch using the `padsequences` function.
- 3 Permute the last two dimensions of the sequences

```
function [X,T,sequenceLengthsSource,maskTarget] = preprocessMiniBatch(sequencesSource,sequencesTarget,
sequenceLengthsSource = cellfun(@(x) size(x,2),sequencesSource);

X = padsequences(sequencesSource,2,PaddingValue=inputSize);
X = permute(X,[1 3 2]);

[T,maskTarget] = padsequences(sequencesTarget,2,PaddingValue=outputSize);
T = permute(T,[1 3 2]);
maskTarget = permute(maskTarget,[1 3 2]);

end
```

### Model Loss Function

The `modelLoss` function takes the encoder and decoder model parameters, a mini-batch of input data and the padding masks corresponding to the input data, and the dropout probability and returns the loss and the gradients of the loss with respect to the learnable parameters in the models.

```
function [loss,gradients] = modelLoss(parameters,X,T,...
sequenceLengthsSource,maskTarget,dropout)

% Forward through encoder.
[Z,hiddenState] = modelEncoder(parameters.encoder,X,sequenceLengthsSource);

% Decoder Output.
doTeacherForcing = rand < 0.5;
sequenceLength = size(T,3);
Y = decoderPredictions(parameters.decoder,Z,T,hiddenState,dropout,...
doTeacherForcing,sequenceLength);

% Masked loss.
Y = Y(:,:,1:end-1);
```

```

T = extractdata(gather(T(:,:,2:end)));
T = onehotencode(T,1,ClassNames=1:size(Y,1));

maskTarget = maskTarget(:,:,2:end);
maskTarget = repmat(maskTarget,[size(Y,1),1,1]);

loss = crossentropy(Y,T,Mask=maskTarget,Dataformat="CBT");

% Update gradients.
gradients = dlgradient(loss,parameters);

end

```

### Encoder Model Function

The function `modelEncoder` takes the input data, the model parameters, the optional mask that is used to determine the correct outputs for training and returns the model output and the LSTM hidden state.

If `sequenceLengths` is empty, then the function does not mask the output. Specify an empty value for `sequenceLengths` when using the `modelEncoder` function for prediction.

```

function [Z,hiddenState] = modelEncoder(parametersEncoder,X,sequenceLengths)

% Embedding.
weights = parametersEncoder.emb.Weights;
Z = embed(X,weights,DataFormat="CBT");

% LSTM 1.
inputWeights = parametersEncoder.lstm1.InputWeights;
recurrentWeights = parametersEncoder.lstm1.RecurrentWeights;
bias = parametersEncoder.lstm1.Bias;

numHiddenUnits = size(recurrentWeights, 2);
initialHiddenState = darray(zeros([numHiddenUnits 1]));
initialCellState = darray(zeros([numHiddenUnits 1]));

Z = lstm(Z,initialHiddenState,initialCellState,inputWeights, ...
    recurrentWeights,bias,DataFormat="CBT");

% LSTM 2.
inputWeights = parametersEncoder.lstm2.InputWeights;
recurrentWeights = parametersEncoder.lstm2.RecurrentWeights;
bias = parametersEncoder.lstm2.Bias;

[Z, hiddenState] = lstm(Z,initialHiddenState,initialCellState, ...
    inputWeights,recurrentWeights,bias,DataFormat="CBT");

% Masking for training.
if ~isempty(sequenceLengths)
    miniBatchSize = size(Z,2);
    for n = 1:miniBatchSize
        hiddenState(:,n) = Z(:,n,sequenceLengths(n));
    end
end

end

```

### Decoder Model Function

The function `modelDecoder` takes the input data, the model parameters, the context vector, the LSTM initial hidden state, the outputs of the encoder, and the dropout probability and outputs the decoder output, the updated context vector, the updated LSTM state, and the attention scores.

```
function [Y, context, hiddenState, attentionScores] = modelDecoder(parametersDecoder,X,context,
    hiddenState,Z,dropout)

% Embedding.
weights = parametersDecoder.emb.Weights;
X = embed(X,weights,DataFormat="CBT");

% RNN input.
sequenceLength = size(X,3);
Y = cat(1, X, repmat(context, [1 1 sequenceLength]));

% LSTM 1.
inputWeights = parametersDecoder.lstm1.InputWeights;
recurrentWeights = parametersDecoder.lstm1.RecurrentWeights;
bias = parametersDecoder.lstm1.Bias;

initialCellState = dlarray(zeros(size(hiddenState)));

Y = lstm(Y,hiddenState,initialCellState,inputWeights,recurrentWeights,bias,DataFormat="CBT");

% Dropout.
mask = rand(size(Y),"like",Y) > dropout;
Y = Y.*mask;

% LSTM 2.
inputWeights = parametersDecoder.lstm2.InputWeights;
recurrentWeights = parametersDecoder.lstm2.RecurrentWeights;
bias = parametersDecoder.lstm2.Bias;

[Y, hiddenState] = lstm(Y,hiddenState,initialCellState,inputWeights,recurrentWeights,bias,DataFo

% Attention.
weights = parametersDecoder.attn.Weights;
[attentionScores, context] = attention(hiddenState, Z, weights);

% Concatenate.
Y = cat(1, Y, repmat(context,[1 1 sequenceLength]));

% Fully connect.
weights = parametersDecoder.fc.Weights;
bias = parametersDecoder.fc.Bias;
Y = fullyconnect(Y,weights,bias,DataFormat="CBT");

% Softmax.
Y = softmax(Y,DataFormat="CBT");

end
```

## Attention Function

The `attention` function returns the attention scores according to Luong "general" scoring and the updated context vector. The energy at each time step is the dot product of the hidden state and the learnable attention weights times the encoder output.

```
function [attentionScores, context] = attention(hiddenState, encoderOutputs, weights)

% Initialize attention energies.
[miniBatchSize, sequenceLength] = size(encoderOutputs, 2:3);
attentionEnergies = zeros([sequenceLength miniBatchSize], "like", hiddenState);

% Attention energies.
hWX = hiddenState .* pagetimes(weights, encoderOutputs);
for tt = 1:sequenceLength
    attentionEnergies(tt, :) = sum(hWX(:, :, tt), 1);
end

% Attention scores.
attentionScores = softmax(attentionEnergies, DataFormat="CB");

% Context.
encoderOutputs = permute(encoderOutputs, [1 3 2]);
attentionScores = permute(attentionScores, [1 3 2]);
context = pagetimes(encoderOutputs, attentionScores);
context = squeeze(context);

end
```

## Decoder Model Predictions Function

The `decoderModelPredictions` function returns the predicted sequence  $Y$  given the input sequence, target sequence, hidden state, dropout probability, flag to enable teacher forcing, and the sequence length.

```
function Y = decoderPredictions(parametersDecoder, Z, T, hiddenState, dropout, ...
    doTeacherForcing, sequenceLength)

% Convert to dlarray.
T = dlarray(T);

% Initialize context.
miniBatchSize = size(T, 2);
numHiddenUnits = size(Z, 1);
context = zeros([numHiddenUnits miniBatchSize], "like", Z);

if doTeacherForcing
    % Forward through decoder.
    Y = modelDecoder(parametersDecoder, T, context, hiddenState, Z, dropout);
else
    % Get first time step for decoder.
    decoderInput = T(:, :, 1);

    % Initialize output.
    numClasses = numel(parametersDecoder.fc.Bias);
    Y = zeros([numClasses miniBatchSize sequenceLength], "like", decoderInput);

    % Loop over time steps.
```

```

    for t = 1:sequenceLength
        % Forward through decoder.
        [Y(:, :, t), context, hiddenState] = modelDecoder(parametersDecoder, decoderInput, context,
            hiddenState, Z, dropout);

        % Update decoder input.
        [~, decoderInput] = max(Y(:, :, t), [], 1);
    end
end
end

```

### Text Translation Function

The `translateText` function translates an array of text by iterating over mini-batches. The function takes as input the model parameters, the input string array, a maximum sequence length, the mini-batch size, the source and target word encoding objects, the start and stop tokens, and the delimiter for assembling the output.

```

function strTranslated = translateText(parameters, strSource, maxSequenceLength, miniBatchSize, ...
    encSource, encTarget, startToken, stopToken, delimiter)

% Transform text.
documentsSource = transformText(strSource, startToken, stopToken);
sequencesSource = doc2sequence(encSource, documentsSource, ...
    PaddingDirection="right", ...
    PaddingValue=encSource.NumWords + 1);

% Convert to dlarray.
X = cat(3, sequencesSource{:});
X = permute(X, [1 3 2]);
X = dlarray(X);

% Initialize output.
numObservations = numel(strSource);
strTranslated = strings(numObservations, 1);

% Loop over mini-batches.
numIterations = ceil(numObservations / miniBatchSize);
for i = 1:numIterations
    idxMiniBatch = (i-1)*miniBatchSize+1:min(i*miniBatchSize, numObservations);
    miniBatchSize = numel(idxMiniBatch);

    % Encode using model encoder.
    sequenceLengths = [];
    [Z, hiddenState] = modelEncoder(parameters.encoder, X(:, idxMiniBatch, :), sequenceLengths);

    % Decoder predictions.
    doTeacherForcing = false;
    dropout = 0;
    decoderInput = repmat(word2ind(encTarget, startToken), [1 miniBatchSize]);
    decoderInput = dlarray(decoderInput);
    Y = decoderPredictions(parameters.decoder, Z, decoderInput, hiddenState, dropout, ...
        doTeacherForcing, maxSequenceLength);
    [~, idxPred] = max(extractdata(Y), [], 1);

    % Keep translating flag.
    idxStop = word2ind(encTarget, stopToken);

```

```

keepTranslating = idxPred ~= idxStop;

% Loop over time steps.
t = 1;
while t <= maxSequenceLength && any(keepTranslating(:, :, t))

    % Update output.
    newWords = ind2word(encTarget, idxPred(:, :, t))';
    idxUpdate = idxMiniBatch(keepTranslating(:, :, t));
    strTranslated(idxUpdate) = strTranslated(idxUpdate) + delimiter + newWords(keepTranslating(:, :, t));

    t = t + 1;
end
end
end

```

## See Also

word2ind | tokenizedDocument | wordEncoding | dlarray | adamupdate | dlupdate | dlfeval  
 | dlgradient | crossentropy | softmax | lstm | doc2sequence

## More About

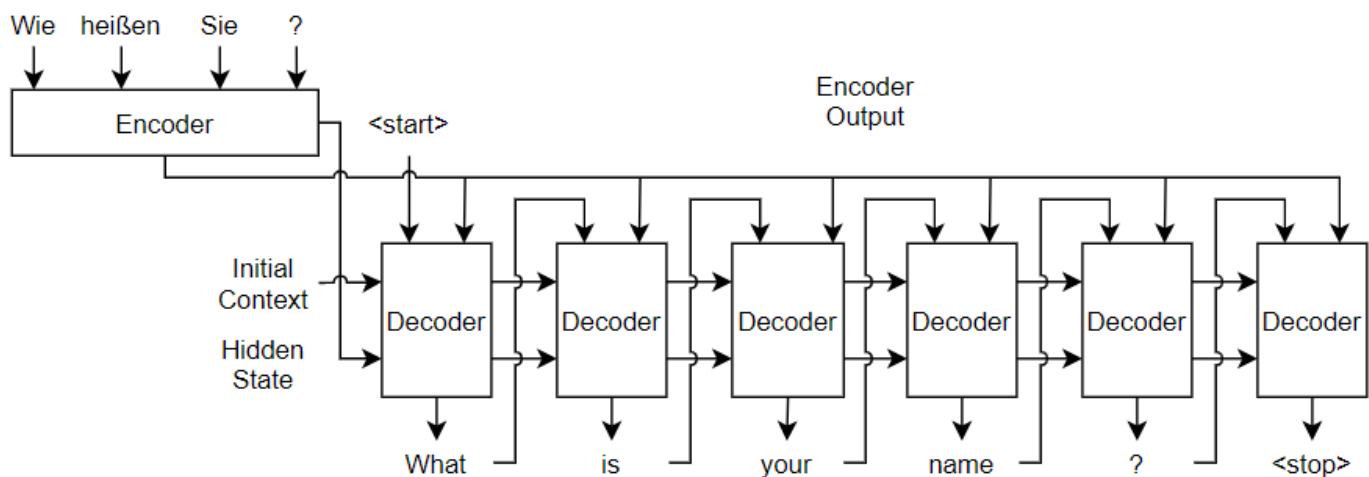
- “Define Custom Training Loops, Loss Functions, and Networks” (Deep Learning Toolbox)
- “Prepare Text Data for Analysis” on page 1-11
- “Analyze Text Data Using Topic Models” on page 2-13
- “Classify Text Data Using Deep Learning” on page 2-85
- “Classify Text Data Using Convolutional Neural Network” on page 2-93
- “Train a Sentiment Classifier” on page 2-71
- “Visualize Word Embeddings Using Text Scatter Plots” on page 3-8

## Language Translation Using Deep Learning

This example shows how to train a German to English language translator using a recurrent sequence-to-sequence encoder-decoder model with attention.

Recurrent encoder-decoder models have proven successful at tasks such as abstractive text summarization and neural machine translation. These models consist of an *encoder*, which typically processes input data with a recurrent layer such as an LSTM layer, and a *decoder* which maps the encoded input into the desired output, typically also with a recurrent layer. Models that incorporate *attention mechanisms* into the models allow the decoder to focus on parts of the encoded input while generating the translation one time step at a time. This example implements *Bahdanau attention* [1] on page 2-0 using the custom layer `attentionLayer`, attached to this example as a supporting file. To access this layer, open this example as a live script.

This diagram shows the structure of a language translation model. The input text, specified as a sequence of words, is passed through the encoder, which outputs an encoded version of the input sequence and a hidden state used to initialize the decoder state. The decoder makes predictions one word at a time using previous prediction as input and also outputs update state and context values.



For more information and the details about the encoder and decoder networks used in this example, see the Define Encoder and Decoder Networks on page 2-0 section of the example.

Predicting the most likely word for each step in the sequence can lead to suboptimal results. Any incorrect predictions can cause even more incorrect predictions in later time steps. For example, for the target text "An eagle flew by.", if the decoder predicts the first word of a translation as "A", then the probability of predicting "eagle" for the next word becomes much more unlikely because of the low probability of the phrase "a eagle" appearing in English text. The translation generation process differs for training and prediction. This example uses different approaches to stabilize training and the predictions:

- To stabilize training, you can randomly use the target values as inputs to the decoder. In particular, you can adjust the probability used to inject the target values as training progresses. For example, you can train using the target values at a much higher rate at the start of training, then decay the probability such that towards the end of training the model uses only the previous predictions. This technique is known as *scheduled sampling* [2] on page 2-0. For more information, see the Decoder Predictions Function on page 2-0 section of the example.



- To improve the predictions at translation time, for each time step, you can consider the top  $K$  predictions for some positive integer  $K$  and explore different sequences of predictions to identify the best combination. This technique is known as *beam search*. For more information, see the Beam Search Function on page 2-0 section of the example.

This example shows how to load and preprocess text data to train a German to English language translator, define the encoder and decoder networks, train the model using a custom training loop, and generate translations using beam search.

**Note:** Language translation is a computationally intensive task. Training on the full data set used in this example can take many hours to run. To make the example run quicker, you can reduce training time at the cost of accuracy of predictions with previously unseen data by discarding a portion of the training data. Removing observations can speed up training because it reduces the amount of data to process in an epoch and reduces the vocabulary size of the training data.

To shorten the time it takes to run the example, discard 70% of the data. Note that discarding large amounts of data negatively affects the accuracy of the learned model. For more accurate results, reduce the amount of discarded data. To speed up the example, increase the amount of discarded data.

```
discardProp = 0.70;
```

### Load Training Data

Download and extract the English-German Tab-delimited Bilingual Sentence Pairs data set. The data comes from <http://www.manythings.org/anki/> and <https://tatoeba.org>, and is provided under the Tatoeba Terms of Use and the CC-BY license.

```
downloadFolder = tempdir;
url = "http://www.manythings.org/anki/deu-eng.zip";
filename = fullfile(downloadFolder, "deu-eng.zip");
dataFolder = fullfile(downloadFolder, "deu-eng");

if ~exist(dataFolder, "dir")
    fprintf("Downloading English-German Tab-delimited Bilingual Sentence Pairs data set (7.6 MB)
    websave(filename, url);
    unzip(filename, dataFolder);
    fprintf("Done.\n")
end
```

Create a table that contains the sentence pairs specified as strings. Read the tab-delimited sentences pairs using `readtable`. Specify the German text as the source and the English text as the target.

```
filename = fullfile(dataFolder, "deu.txt");

opts = delimitedTextImportOptions(...
    Delimiter="\t", ...
    VariableNames=["Target" "Source" "License"], ...
    SelectedVariableNames=["Source" "Target"], ...
    VariableTypes=["string" "string" "string"], ...
    Encoding="UTF-8");
```

View the first few sentence pairs in the data.

```
data = readtable(filename, opts);
head(data)
```

```
ans=8x2 table
      Source      Target
-----
"Geh."          "Go."
"Hallo!"        "Hi."
"Grüß Gott!"    "Hi."
"Lauf!"         "Run!"
"Lauf!"         "Run."
"Potzdonner!"   "Wow!"
"Donnerwetter!" "Wow!"
"Feuer!"        "Fire!"
```

Training on the full dataset can take a long time to run. To reduce training time at the cost of accuracy, you can discard a portion of the training data. Removing observations can speed up training because it reduces the amount of data to process in an epoch as well as reducing the vocabulary size of the training data.

Discard a portion of the data according to the `discardProp` variable defined at the start of the example. Note that discarding large amounts of data negatively affects the accuracy of the learned model. For more accurate results, reduce the amount of discarded data by setting `discardProp` to a lower value.

```
idx = size(data,1) - floor(discardProp*size(data,1)) + 1;
data(idx:end,:) = [];
```

View the number of remaining observations.

```
size(data,1)
```

```
ans = 68124
```

Split the data into training and test partitions containing 90% and 10% of the data, respectively.

```
trainingProp = 0.9;
idx = randperm(size(data,1), floor(trainingProp*size(data,1)));
dataTrain = data(idx,:);
dataTest = data;
dataTest(idx,:) = [];
```

View the first few rows of the training data.

```
head(dataTrain)
```

```
ans=8x2 table
      Source      Target
-----
"Tom erschoss Mary."      "Tom shot Mary."
"Ruf mich bitte an."      "Call me, please."
"Kann das einer nachprüfen?"  "Can someone check this?"
"Das lasse ich mir nicht gefallen!"  "I won't stand for it."
"Ich mag Englisch nicht."      "I don't like English."
"Er ist auf dem Laufenden."      "He is up to date."
"Sie sieht glücklich aus."      "She seems happy."
"Wo wurden sie geboren?"      "Where were they born?"
```

View the number of training observations.

```
numObservationsTrain = size(dataTrain,1)
numObservationsTrain = 61311
```

### Preprocess Data

Preprocess the text data using the `preprocessText` function, listed at the end of the example. The `preprocessText` function preprocesses and tokenizes the input text for translation by splitting the text into words and adding start and stop tokens.

```
documentsGerman = preprocessText(dataTrain.Source);
```

Create a `wordEncoding` object that maps tokens to a numeric index and vice versa using a vocabulary.

```
encGerman = wordEncoding(documentsGerman);
```

Convert the target data to sequences using the same steps.

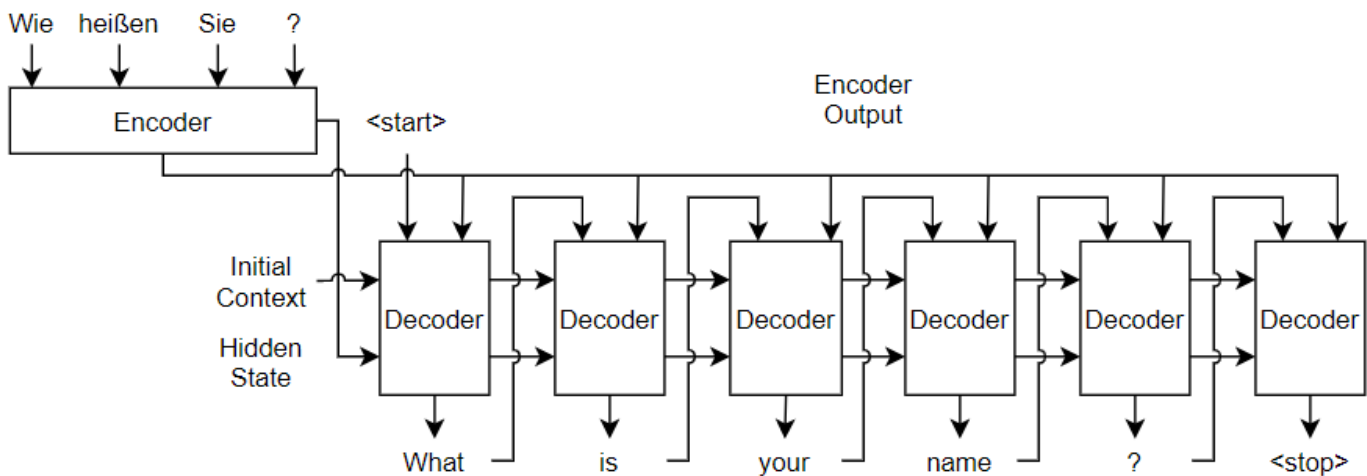
```
documentsEnglish = preprocessText(dataTrain.Target);
encEnglish = wordEncoding(documentsEnglish);
```

View the vocabulary sizes of the source and target encodings.

```
numWordsGerman = encGerman.NumWords
numWordsGerman = 12117
numWordsEnglish = encEnglish.NumWords
numWordsEnglish = 7226
```

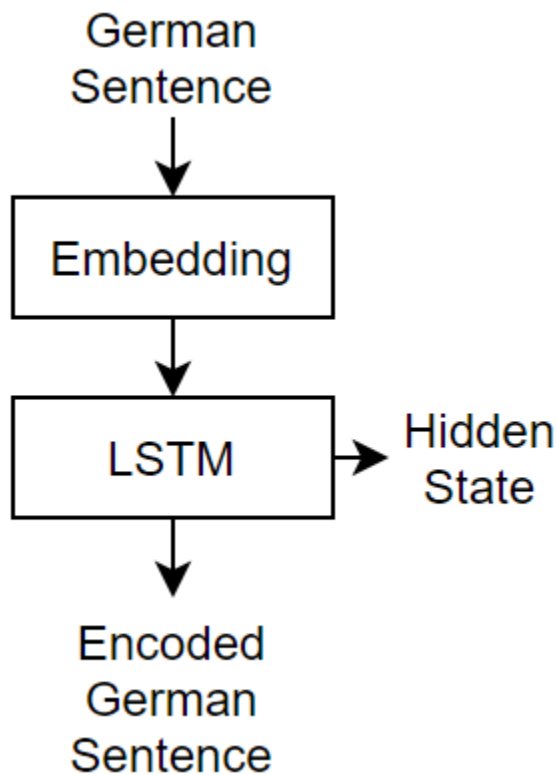
### Define Encoder and Decoder Networks

This diagram shows the structure of a language translation model. The input text, specified as a sequence of words, is passed through the encoder, which outputs an encoded version of the input sequence and a hidden state used to initialize the decoder state. The decoder makes predictions one word at a time using previous the prediction as input and also outputs updated state and context values.

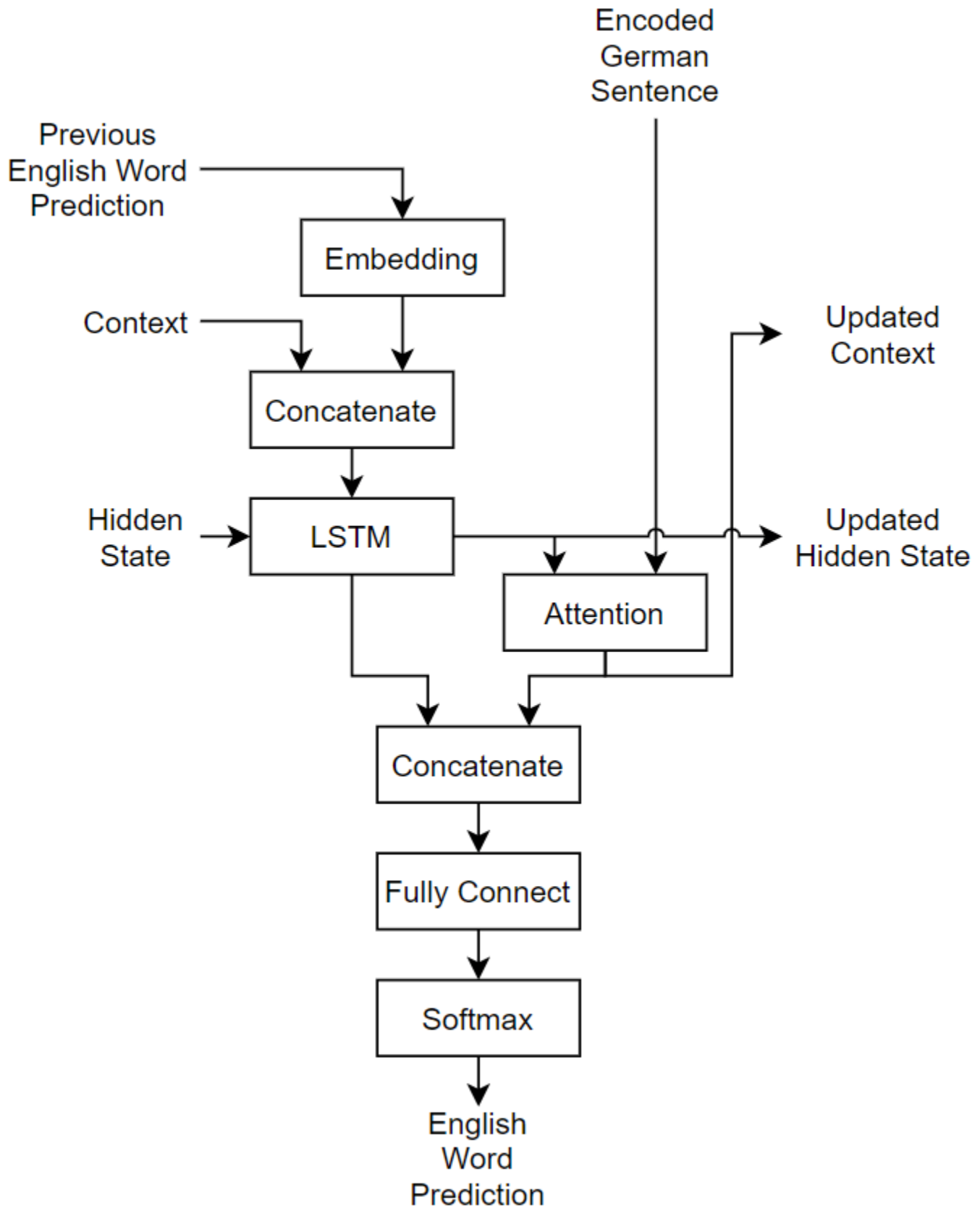


Create the encoder and decoder networks using the `languageTranslationLayers` function, attached to this example as a supporting file. To access this function, open the example as a live script.

For the encoder network, the `languageTranslationLayers` function defines a simple network consisting of an embedding layer followed by an LSTM layer. An embedding operation converts categorical tokens into numeric vectors, where the numeric vectors are learned by the network.



For the decoder network, the `languageTranslationLayers` function defines a network that passes the input data concatenated with the input context through an LSTM layer, and takes the updated hidden state and the encoder output and passes it through an attention mechanism to determine the context vector. The LSTM output and the context vector are then concatenated and passed through a fully connected and a softmax layer for classification.



Create the encoder and decoder networks using the `languageTranslationLayers` function, attached to this example as a supporting file. To access this function, open the example as a live script. Specify an embedding dimension of 128, and 128 hidden units in the LSTM layers.

```
embeddingDimension = 128;  
numHiddenUnits = 128;
```

```
[lgraphEncoder,lgraphDecoder] = languageTranslationLayers(embeddingDimension,numHiddenUnits,numW
```

To train the network in a custom training loop, convert the encoder and decoder networks to `dlnetwork` objects.

```
netEncoder = dlnetwork(lgraphEncoder);  
netDecoder = dlnetwork(lgraphDecoder);
```

The decoder has multiple outputs including the context output of the attention layer, which is also passed to another layer. Specify the network outputs using the `OutputNames` property of the decoder `dlnetwork` object.

```
netDecoder.OutputNames = ["softmax" "context" "lstm2/hidden" "lstm2/cell"];
```

### Define Model Loss Function

Create the function `modelLoss`, listed in the Model Loss Function on page 2-0 section of the example, which takes as input the encoder and decoder model parameters, a mini-batch of input data and the padding masks corresponding to the input data, and the dropout probability and returns the loss, the gradients of the loss with respect to the learnable parameters in the models, and the model predictions.

### Specify Training Options

Train with a mini-batch size of 64 for 15 epochs and a learning rate of 0.005.

```
miniBatchSize = 64;  
numEpochs = 15;  
learnRate = 0.005;
```

Initialize the options for Adam optimization.

```
gradientDecayFactor = 0.9;  
squaredGradientDecayFactor = 0.999;
```

Train using gradually decaying values of  $\epsilon$  for scheduled sampling. Start with a value of  $\epsilon = 0.5$  and linearly decay to end with a value of  $\epsilon = 0$ . For more information about scheduled sampling, see the Decoder Predictions Function on page 2-0 section of the example.

```
epsilonStart = 0.5;  
epsilonEnd = 0;
```

Train using SortaGrad [3] on page 2-0 , which is a strategy to improve training of ragged sequences by training for one epoch with the sequences sorted by sequence then shuffling once per epoch thereafter.

Sort the training sequences by sequence length.

```
sequenceLengths = doclength(documentsGerman);  
[~,idx] = sort(sequenceLengths);
```

```
documentsGerman = documentsGerman(idx);
documentsEnglish = documentsEnglish(idx);
```

### Train Model

Train the model using a custom training loop.

Create array datastores for the source and target data using the `arrayDatastore` function. Combine the datastores using the `combine` function.

```
adsSource = arrayDatastore(documentsGerman);
adsTarget = arrayDatastore(documentsEnglish);
cbs = combine(adsSource,adsTarget);
```

Create a mini-batch queue to automatically prepare mini-batches for training.

- Preprocess the training data using the `preprocessMiniBatch` function, which returns a mini-batch of source sequences, target sequences, the corresponding mask, and the initial start token.
- Output `darray` objects with the format "CTB" (channel, time, batch).
- Discard any partial mini-batches.

```
mbq = minibatchqueue(cbs,4, ...
    MiniBatchSize=miniBatchSize, ...
    MiniBatchFcn=@(X,Y) preprocessMiniBatch(X,Y,encGerman,encEnglish), ...
    MiniBatchFormat=["CTB" "CTB" "CTB" "CTB"], ...
    PartialMiniBatch="discard");
```

Initialize the training progress plot.

```
figure
C = colororder;
lineLossTrain = animatedline(Color=C(2,:));

xlabel("Iteration")
ylabel("Loss")
ylim([0 inf])
grid on
```

For the encoder and decoder networks, initialize the values for Adam optimization.

```
trailingAvgEncoder = [];
trailingAvgSqEncoder = [];
trailingAvgDecoder = [];
trailingAvgSqDecoder = [];
```

Create an array of  $\epsilon$  values for scheduled sampling.

```
numIterationsPerEpoch = floor(numObservationsTrain/miniBatchSize);
numIterations = numIterationsPerEpoch * numEpochs;
epsilon = linspace(epsilonStart,epsilonEnd,numIterations);
```

Train the model. For each iteration:

- Read a mini-batch of data from the mini-batch queue.
- Compute the model loss and gradients.
- Update the encoder and decoder networks using the `adamupdate` function.

- Update the training progress plot and display an example translation using the `ind2str` function, attached to this example as a supporting file. To access this function, open this example as a live script.
- If the iteration yields the lowest training loss, then save the network.

At the end of each epoch, shuffle the mini-batch queue.

For large data sets, training can take many hours to run.

```

iteration = 0;
start = tic;
lossMin = inf;
reset(mbq)

% Loop over epochs.
for epoch = 1:numEpochs

    % Loop over mini-batches.
    while hasdata(mbq)
        iteration = iteration + 1;

        % Read mini-batch of data.
        [X,T,maskT,decoderInput] = next(mbq);

        % Compute loss and gradients.
        [loss,gradientsEncoder,gradientsDecoder,YPred] = dlfeval(@modelLoss,netEncoder,netDecoder,decoderInput,X,T,maskT);

        % Update network learnable parameters using adamupdate.
        [netEncoder, trailingAvgEncoder, trailingAvgSqEncoder] = adamupdate(netEncoder,gradientsEncoder,gradientsDecoder,iteration,learnRate,gradientDecayFactor,squaredGradientDecayFactor);

        [netDecoder, trailingAvgDecoder, trailingAvgSqDecoder] = adamupdate(netDecoder,gradientsEncoder,gradientsDecoder,iteration,learnRate,gradientDecayFactor,squaredGradientDecayFactor);

        % Generate translation for plot.
        if iteration == 1 || mod(iteration,10) == 0
            strGerman = ind2str(X(:,1,:),encGerman);
            strEnglish = ind2str(T(:,1,:),encEnglish,Mask=maskT);
            strTranslated = ind2str(YPred(:,1,:),encEnglish);
        end

        % Display training progress.
        D = duration(0,0,toc(start),Format="hh:mm:ss");
        loss = double(gather(extractdata(loss)));
        addpoints(lineLossTrain,iteration,loss)
        title( ...
            "Epoch: " + epoch + ", Elapsed: " + string(D) + newline + ...
            "Source: " + strGerman + newline + ...
            "Target: " + strEnglish + newline + ...
            "Training Translation: " + strTranslated)

        drawnow

    % Save best network.
    if loss < lossMin
        lossMin = loss;
        netBest.netEncoder = netEncoder;
    end
end

```

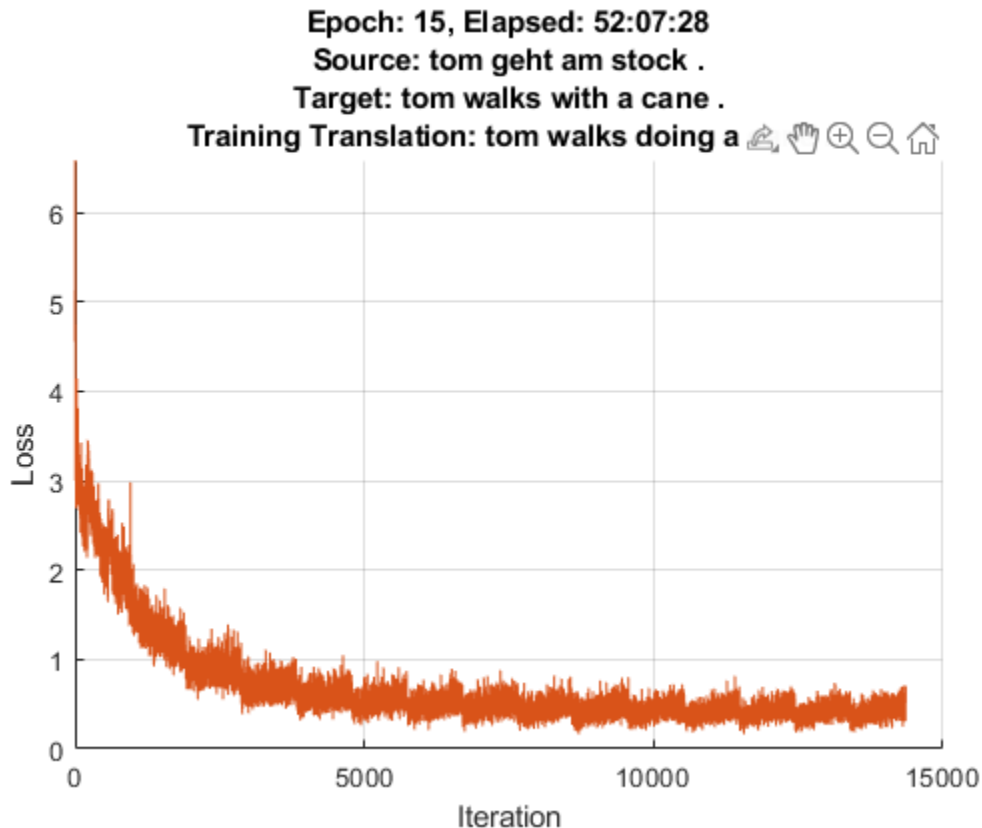


```

        netBest.netDecoder = netDecoder;
        netBest.loss = loss;
        netBest.iteration = iteration;
        netBest.D = D;
    end
end

% Shuffle.
shuffle(mbq);
end

```



The plot shows two translations of the source text. The target is the target translation provided by the training data that the network attempts to reproduce. The training translation is the predicted translation, which uses information from the target text via the scheduled sampling mechanism.

Add the word encodings to the `netBest` structure and save the structure in a MAT file.

```

netBest.encGerman = encGerman;
netBest.encEnglish = encEnglish;

D = datetime("now",Format="yyyy_MM_dd__HH_mm_ss");
filename = "net_best__" + string(D) + ".mat";
save(filename,"netBest");

```

Extract the best network from `netBest`.

```

netEncoder = netBest.netEncoder;
netDecoder = netBest.netDecoder;

```

## Test Model

To evaluate the quality of the translations, use the BiLingual Evaluation Understudy (BLEU) scoring algorithm [4] on page 2-0 .

Translate the test data using the `translateText` function listed at the end of the example.

```
strTranslatedTest = translateText(netEncoder,netDecoder,encGerman,encEnglish,dataTest.Source);
```

View a random selection of the test source text, target text, and predicted translations in a table.

```
numObservationsTest = size(dataTest,1);
idx = randperm(numObservationsTest,8);
tbl = table;
tbl.Source = dataTest.Source(idx);
tbl.Target = dataTest.Target(idx);
tbl.Translated = strTranslatedTest(idx)
```

`tbl=8x3 table`

Source	Target	Translated
"Er sieht krank aus."	"He seems ill."	"he looks sick ."
"Ich werde das Buch holen."	"I'll get the book."	"i'll get the book ."
"Ruhst du dich jemals aus?"	"Do you ever rest?"	"do you look out of ?"
"Was willst du?"	"What are you after?"	"what do you want want"
"Du hast keinen Beweis."	"You have no proof."	"you have no proof ."
"Macht es, wann immer ihr wollt."	"Do it whenever you want."	"do it you like it ."
"Tom ist gerade nach Hause gekommen."	"Tom has just come home."	"tom just came home h"
"Er lügt nie."	"He never tells a lie."	"he never lie lies ."

To evaluate the quality of the translations using the BLEU similarity score, first preprocess the text data using the same steps as for training. Specify empty start and stop tokens, as these are not used in the translation.

```
candidates = preprocessText(strTranslatedTest,StartToken="",StopToken="");
references = preprocessText(dataTest.Target,StartToken="",StopToken="");
```

The `bleuEvaluationScore` function, by default, evaluates the similarity scores by comparing n-grams of length one through four (multiword phrases with four or fewer words or single words). If the candidate or reference documents have fewer than four tokens, then the resulting BLEU evaluation score is zero. To ensure that `bleuEvaluationScore` returns nonzero scores for these short candidate documents, set the n-gram weights to a vector with fewer elements than the number of words in candidate.

Determine the length of the shortest candidate document.

```
minLength = min([doclength(candidates); doclength(references)])
```

```
minLength = 2
```

If the shortest document has fewer than four tokens, then set the n-gram weights to a vector with a length matching the shortest document with equal weights that sum to one. Otherwise, specify n-gram weights of `[0.25 0.25 0.25 0.25]`. Note that if `minLength` is 1 (and consequently the n-gram weights is also 1), then the `bleuEvaluationScore` function can return less meaningful results as it only compares individual words (unigrams) and does not compare any n-grams (multiword phrases).

```

if minLength < 4
    ngramWeights = ones(1,minLength) / minLength;
else
    ngramWeights = [0.25 0.25 0.25 0.25];
end

```

Calculate the BLEU evaluation scores by iterating over the translations and using the `bleuEvaluationScore` function.

```

for i = 1:numObservationsTest
    score(i) = bleuEvaluationScore(candidates(i),references(i),NgramWeights=ngramWeights);
end

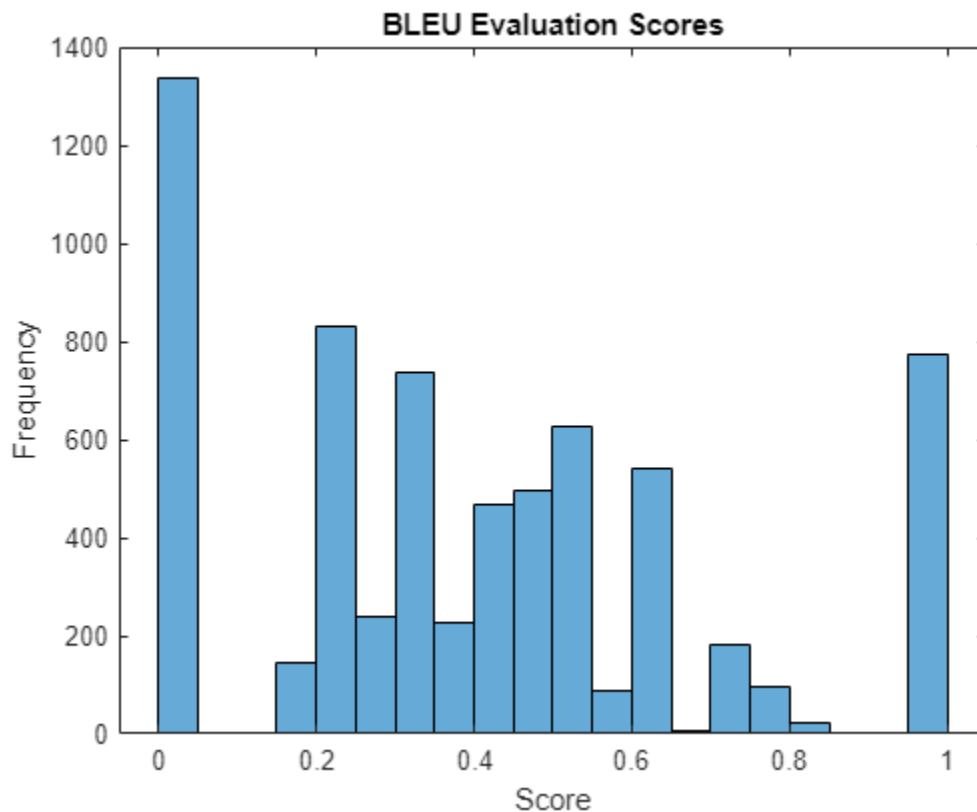
```

Visualize the BLEU evaluation scores in a histogram.

```

figure
histogram(score);
title("BLEU Evaluation Scores")
xlabel("Score")
ylabel("Frequency")

```



View a table of some of the best translations.

```

[~,idxSorted] = sort(score,"descend");
idx = idxSorted(1:8);
tbl = table;
tbl.Source = dataTest.Source(idx);

```

```
tbl.Target = dataTest.Target(idx);
tbl.Translated = strTranslatedTest(idx)
```

tbl=8x3 table

Source	Target	Translated
"Legen Sie sich hin!"	"Lie low."	"lie low ."
"Ich gähnte."	"I yawned."	"i yawned ."
"Küsse Tom!"	"Kiss Tom."	"kiss tom ."
"Küssen Sie Tom!"	"Kiss Tom."	"kiss tom ."
"Nimm Tom."	"Take Tom."	"take tom ."
"Komm bald."	"Come soon."	"come soon ."
"Ich habe es geschafft."	"I made it."	"i made it ."
"Ich sehe Tom."	"I see Tom."	"i see tom ."

View a table of some of the worst translations.

```
idx = idxSorted(end-7:end);
tbl = table;
tbl.Source = dataTest.Source(idx);
tbl.Target = dataTest.Target(idx);
tbl.Translated = strTranslatedTest(idx)
```

tbl=8x3 table

Source	Target
"Diese Schnecken kann man essen."	"These snails are edible."
"Sie stehen noch zu Verfügung."	"They're still available."
"Diese Schraube passt zu dieser Mutter."	"This bolt fits this nut."
"Diese Puppe gehört mir."	"This doll belongs to me."
"Das ist eine japanische Puppe."	"This is a Japanese doll."
"Das ist eine Kreuzung, an der alle Fahrzeuge anhalten müssen."	"This is a four-way stop."
"Diese Sendung ist eine Wiederholung."	"This program is a rerun."
"Die heutige Folge ist eine Wiederholung."	"Today's show is a rerun."

### Generate Translations

Generate translations for new data using the `translateText` function.

```
strGermanNew = [
    "Wie geht es Dir heute?"
    "Wie heißen Sie?"
    "Das Wetter ist heute gut."];
```

Translate the text using the `translateText`, function listed at the end of the example.

```
strTranslatedNew = translateText(netEncoder,netDecoder,encGerman,encEnglish,strGermanNew)
```

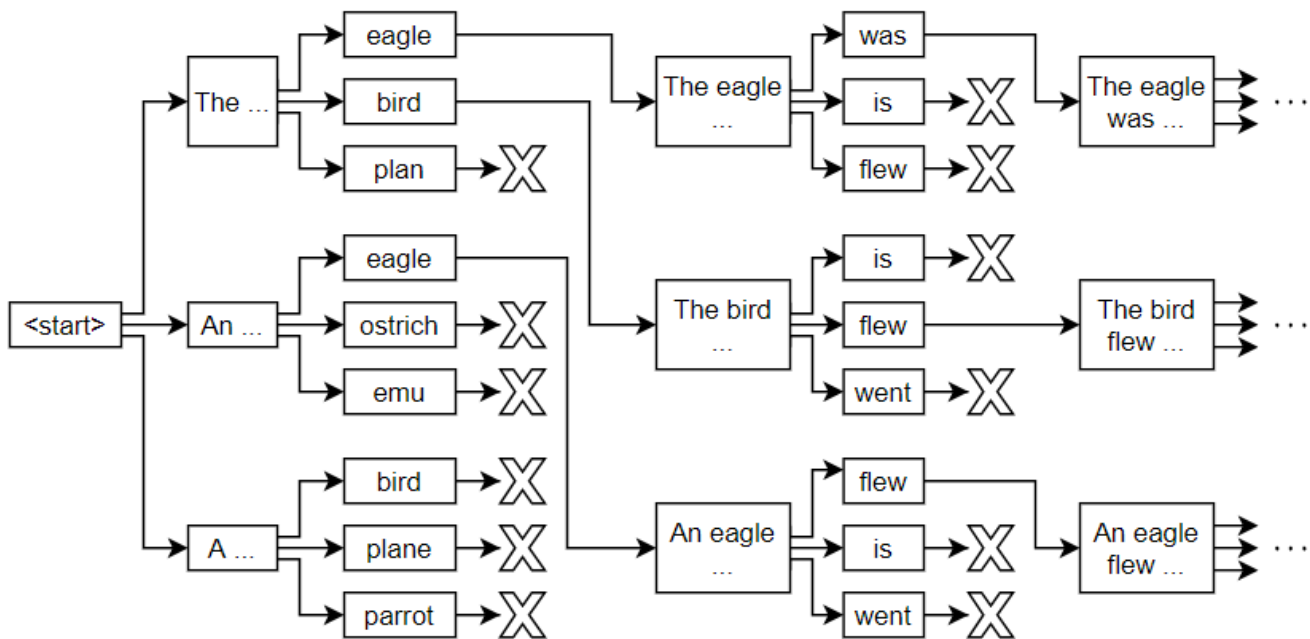
```
strTranslatedNew = 3x1 string
    "how do you feel today ?"
    "what's your your name ? ? ? ? ?"
    "the is is today . . today . ."
```

## Prediction Functions

### Beam Search Function

Beam search is a technique for exploring different combinations of time-step predictions to help find the best prediction sequence. The premise of beam search is for each time-step prediction, identify the top  $K$  predictions for some positive integer  $K$  (also known as the *beam index* or the *beam width*), and maintain the top  $K$  predicted sequences so far at each time step.

This diagram shows the structure of an example beam search with beam index  $K = 3$ . For each prediction, the top three sequences are maintained.



The `beamSearch` function takes as input the input data  $X$ , the encoder and decoder networks, and the target word encoding, and returns the predicted translated words using the beam search algorithm with a beam index of 3 and a maximum sequence length of 10. You can also specify optional arguments using name-value arguments:

- `BeamIndex` — Beam index. The default is 3.
- `MaxNumWords` — Maximum sequence length. The default is 10.

```
function str = beamSearch(X,netEncoder,netDecoder,encEnglish,args)
```

```
% Parse input arguments.
arguments
    X
    netEncoder
    netDecoder
    encEnglish

    args.BeamIndex = 3;
    args.MaxNumWords = 10;
end
```

```
beamIndex = args.BeamIndex;
maxNumWords = args.MaxNumWords;
startToken = "<start>";
stopToken = "<stop>";

% Encoder predictions.
[Z, hiddenState, cellState] = predict(netEncoder,X);

% Initialize context.
miniBatchSize = size(X,2);
numHiddenUnits = size(Z,1);
context = zeros([numHiddenUnits miniBatchSize],"like",Z);
context = darray(context,"CB");

% Initialize candidates.
candidates = struct;
candidates.Words = startToken;
candidates.Score = 0;
candidates.StopFlag = false;
candidates.HiddenState = hiddenState;
candidates.CellState = cellState;

% Loop over words.
t = 0;
while t < maxNumWords
    t = t + 1;

    candidatesNew = [];

    % Loop over candidates.
    for i = 1:numel(candidates)

        % Stop generating when stop token is predicted.
        if candidates(i).StopFlag
            continue
        end

        % Candidate details.
        words = candidates(i).Words;
        score = candidates(i).Score;
        hiddenState = candidates(i).HiddenState;
        cellState = candidates(i).CellState;

        % Predict next token.
        decoderInput = word2ind(encEnglish,words(end));
        decoderInput = darray(decoderInput,"CBT");

        [YPred,context,hiddenState,cellState] = predict(netDecoder,decoderInput,hiddenState,cellState);
        Outputs=["softmax" "context" "lstm2/hidden" "lstm2/cell"];

        % Find top predictions.
        [scoresTop,idxTop] = maxk(extractdata(YPred),beamIndex);
        idxTop = gather(idxTop);

        % Loop over top predictions.
        for j = 1:beamIndex
            candidate = struct;
```

```

        % Determine candidate word and score.
        candidateWord = ind2word(encEnglish,idxTop(j));
        candidateScore = scoresTop(j);

        % Set stop translating flag.
        if candidateWord == stopToken
            candidate.StopFlag = true;
        else
            candidate.StopFlag = false;
        end

        % Update candidate details.
        candidate.Words = [words candidateWord];
        candidate.Score = score + log(candidateScore);
        candidate.HiddenState = hiddenState;
        candidate.CellState = cellState;

        % Add to new candidates.
        candidatesNew = [candidatesNew candidate];
    end
end

% Get top candidates.
[~,idx] = maxk([candidatesNew.Score],beamIndex);
candidates = candidatesNew(idx);

% Stop predicting when all candidates have stop token.
if all([candidates.StopFlag])
    break
end
end

% Get top candidate.
words = candidates(1).Words;

% Convert to string scalar.
words(ismember(words,[startToken stopToken])) = [];
str = join(words);

end

```

### Translate Text Function

The `translateText` function takes as input the encoder and decoder networks, an input string, and source and target word encodings and returns the translated text.

```

function strTranslated = translateText(netEncoder,netDecoder,encGerman,encEnglish,strGerman,args)

% Parse input arguments.
arguments
    netEncoder
    netDecoder
    encGerman
    encEnglish
    strGerman

    args.BeamIndex = 3;

```

```

end

beamIndex = args.BeamIndex;

% Preprocess text.
documentsGerman = preprocessText(strGerman);
X = preprocessPredictors(documentsGerman,encGerman);
X = dlarray(X,"CTB");

% Loop over observations.
numObservations = numel(strGerman);
strTranslated = strings(numObservations,1);
for n = 1:numObservations

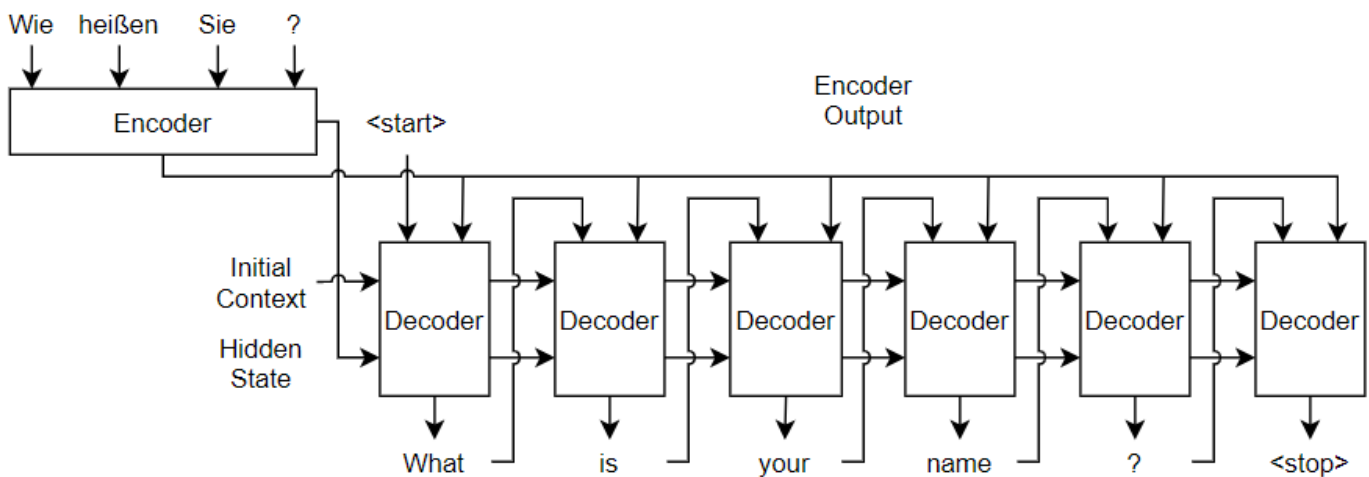
    % Translate text.
    strTranslated(n) = beamSearch(X(:,n,:),netEncoder,netDecoder,encEnglish,BeamIndex=beamIndex)
end
end

```

## Model Functions

### Model Loss Function

The `modelLoss` function takes as input the encoder network, decoder network, mini-batches of predictors `X`, targets `T`, padding mask corresponding to the targets `maskT`, and  $\epsilon$  value for scheduled sampling. The function returns the loss, the gradients of the loss with respect to the learnable parameters in the networks `gradientsE` and `gradientsD`, and the decoder predictions `YPred` encoded as sequences of one-hot vectors.



```

function [loss,gradientsE,gradientsD,YPred] = modelLoss(netEncoder,netDecoder,X,T,maskT,decoderInput,epsilon);

% Forward through encoder.
[Z, hiddenState, cellState] = forward(netEncoder,X);

% Decoder output.
Y = decoderPredictions(netDecoder,Z,T,hiddenState,cellState,decoderInput,epsilon);

% Sparse cross-entropy loss.
loss = sparseCrossEntropy(Y,T,maskT);

```



```
% Update gradients.
[gradientsE,gradientsD] = dlgradient(loss,netEncoder.Learnables,netDecoder.Learnables);

% For plotting, return loss normalized by sequence length.
sequenceLength = size(T,3);
loss = loss ./ sequenceLength;

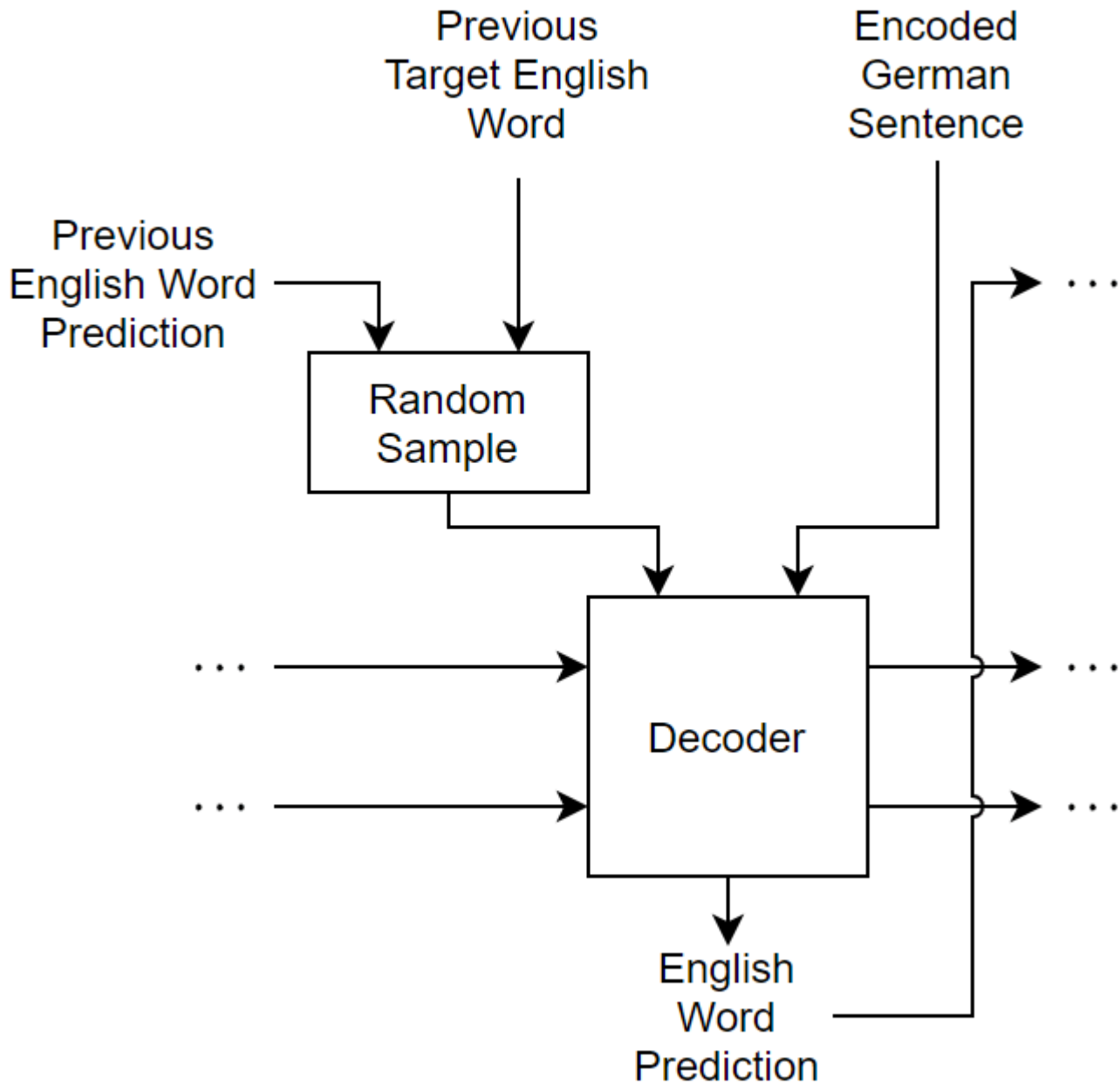
% For plotting example translations, return the decoder output.
YPred = onehotdecode(Y,1:size(Y,1),1,"single");

end
```

### Decoder Predictions Function

The `decoderPredictions` function takes as input, the decoder network, the encoder output  $Z$ , the targets  $T$ , the decoder input hidden and cell state values, and the  $\epsilon$  value for scheduled sampling.

To stabilize training, you can randomly use the target values as inputs to the decoder. In particular, you can adjust the probability used to inject the target values as training progresses. For example, you can train using the target values at a much higher rate at the start of training, then decay the probability such that towards the end of training the model uses only the previous predictions. This technique is known as *scheduled sampling* [2] on page 2-0 . This diagram shows the sampling mechanism incorporated into one time step of a decoder prediction.



The decoder makes predictions one time step at a time. At each time step, the input is randomly selected according to the  $\epsilon$  value for scheduled sampling. In particular, the function uses the target value as input with probability  $\epsilon$  and uses the previous prediction otherwise.

```
function Y = decoderPredictions(netDecoder,Z,T,hiddenState,cellState,decoderInput,epsilon)

% Initialize context.
numHiddenUnits = size(Z,1);
miniBatchSize = size(Z,2);
context = zeros([numHiddenUnits miniBatchSize],"like",Z);
context = dlarray(context,"CB");
```

```

% Initialize output.
idx = (netDecoder.Learnables.Layer == "fc" & netDecoder.Learnables.Parameter=="Bias");
numClasses = numel(netDecoder.Learnables.Value{idx});
sequenceLength = size(T,3);
Y = zeros([numClasses miniBatchSize sequenceLength],"like",Z);
Y = dlarray(Y,"CBT");

% Forward start token through decoder.
[Y(:, :, 1), context, hiddenState, cellState] = forward(netDecoder, decoderInput, hiddenState, cellState);

% Loop over remaining time steps.
for t = 2:sequenceLength

    % Scheduled sampling. Randomly select previous target or previous
    % prediction.
    if rand < epsilon
        % Use target value.
        decoderInput = T(:, :, t-1);
    else
        % Use previous prediction.
        [~, Yhat] = max(Y(:, :, t-1), [], 1);
        decoderInput = Yhat;
    end

    % Forward through decoder.
    [Y(:, :, t), context, hiddenState, cellState] = forward(netDecoder, decoderInput, hiddenState, cellState);
end
end

```

### Sparse Cross-Entropy Loss

The `sparseCrossEntropy` function calculates the cross-entropy loss between the predictions `Y` and targets `T` with the target mask `maskT`, where `Y` is an array of probabilities and `T` is encoded as a sequence of integer values.

```

function loss = sparseCrossEntropy(Y,T,maskT)

% Initialize loss.
[~,miniBatchSize,sequenceLength] = size(Y);
loss = zeros([miniBatchSize sequenceLength],"like",Y);

% To prevent calculating log of 0, bound away from zero.
precision = underlyingType(Y);
Y(Y < eps(precision)) = eps(precision);

% Loop over time steps.
for n = 1:miniBatchSize
    for t = 1:sequenceLength
        idx = T(1,n,t);
        loss(n,t) = -log(Y(idx,n,t));
    end
end

% Apply masking.
maskT = squeeze(maskT);
loss = loss .* maskT;

```

```
% Calculate sum and normalize.  
loss = sum(loss, "all");  
loss = loss / miniBatchSize;
```

```
end
```

## Preprocessing Functions

### Text Preprocessing Function

The `preprocessText` function preprocesses the input text for translation by converting the text to lowercase, adding start and stop tokens, and tokenizing.

```
function documents = preprocessText(str, args)
```

```
arguments
```

```
    str  
    args.StartToken = "<start>";  
    args.StopToken = "<stop>";
```

```
end
```

```
startToken = args.StartToken;  
stopToken = args.StopToken;
```

```
str = lower(str);  
str = startToken + str + stopToken;  
documents = tokenizedDocument(str, CustomTokens=[startToken stopToken]);
```

```
end
```

### Mini-Batch Preprocessing Function

The `preprocessMiniBatch` function preprocesses tokenized documents for training. The function encodes mini-batches of documents as sequences of numeric indices and pads the sequences to have the same length.

```
function [XSource, XTarget, mask, decoderInput] = preprocessMiniBatch(dataSource, dataTarget, encGerman)
```

```
documentsGerman = cat(1, dataSource{:});  
XSource = preprocessPredictors(documentsGerman, encGerman);
```

```
documentsEnglish = cat(1, dataTarget{:});  
sequencesTarget = doc2sequence(encEnglish, documentsEnglish, PaddingDirection="none");
```

```
[XTarget, mask] = padsequences(sequencesTarget, 2, PaddingValue=1);
```

```
decoderInput = XTarget(:, 1, :);  
XTarget(:, 1, :) = [];  
mask(:, 1, :) = [];
```

```
end
```

### Predictors Preprocessing Function

The `preprocessPredictors` function preprocesses source documents for training or prediction. The function encodes an array of tokenized documents as sequences of numeric indices.

```
function XSource = preprocessPredictors(documentsGerman, encGerman)
```

```
sequencesSource = doc2sequence(encGerman,documentsGerman,PaddingDirection="none");  
XSource = padsequences(sequencesSource,2);
```

```
end
```

## Bibliography

- 1 Chorowski, Jan, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and Yoshua Bengio. "Attention-Based Models for Speech Recognition." Preprint, submitted June 24, 2015. <https://arxiv.org/abs/1506.07503>.
- 2 Bengio, Samy, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. "Scheduled Sampling for Sequence Prediction with Recurrent Neural Networks." Preprint, submitted September 23, 2015. <https://arxiv.org/abs/1506.03099>.
- 3 Amodei, Dario, Sundaram Ananthanarayanan, Rishita Anubhai, Jingliang Bai, Eric Battenberg, Carl Case, Jared Casper et al. "Deep Speech 2: End-to-End Speech Recognition in English and Mandarin." In *Proceedings of Machine Learning Research* 48 (2016): 173-182.
- 4 Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. "BLEU: A Method for Automatic Evaluation of Machine Translation." In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics* (2002): 311-318.

## See Also

tokenizedDocument | wordEncoding | dlarray | adamupdate | dlfeval | dlgradient | doc2sequence

## More About

- "Define Custom Training Loops, Loss Functions, and Networks" (Deep Learning Toolbox)
- "Prepare Text Data for Analysis" on page 1-11
- "Analyze Text Data Using Topic Models" on page 2-13
- "Classify Text Data Using Deep Learning" on page 2-85
- "Classify Text Data Using Convolutional Neural Network" on page 2-93
- "Train a Sentiment Classifier" on page 2-71
- "Visualize Word Embeddings Using Text Scatter Plots" on page 3-8

## Classify Out-of-Memory Text Data Using Deep Learning

This example shows how to classify out-of-memory text data with a deep learning network using a transformed datastore.

A transformed datastore transforms or processes data read from an underlying datastore. You can use a transformed datastore as a source of training, validation, test, and prediction data sets for deep learning applications. Use transformed datastores to read out-of-memory data or to perform specific preprocessing operations when reading batches of data.

When training the network, the software creates mini-batches of sequences of the same length by padding, truncating, or splitting the input data. The `trainingOptions` function provides options to pad and truncate input sequences, however, these options are not well suited for sequences of word vectors. Furthermore, this function does not support padding data in a custom datastore. Instead, you must pad and truncate the sequences manually. If you *left-pad* and truncate the sequences of word vectors, then the training might improve.

The “Classify Text Data Using Deep Learning” on page 2-85 example manually truncates and pads all the documents to the same length. This process adds lots of padding to very short documents and discards lots of data from very long documents.

Alternatively, to prevent adding too much padding or discarding too much data, create a transformed datastore that inputs mini-batches into the network. The datastore created in this example converts mini-batches of documents to sequences or word indices and left-pads each mini-batch to the length of the longest document in the mini-batch.

### Load Pretrained Word Embedding

The datastore requires a word embedding to convert documents to sequences of vectors. Load a pretrained word embedding using `fastTextWordEmbedding`. This function requires Text Analytics Toolbox™ Model for *fastText English 16 Billion Token Word Embedding* support package. If this support package is not installed, then the function provides a download link.

```
emb = fastTextWordEmbedding;
```

### Load Data

Create a tabular text datastore from the data in `factoryReports.csv`. Specify to read the data from the “Description” and “Category” columns only.

```
filenameTrain = "factoryReports.csv";
textName = "Description";
labelName = "Category";
ttdsTrain = tabularTextDatastore(filenameTrain, 'SelectedVariableNames', [textName labelName]);
```

View a preview of the datastore.

```
preview(ttdsTrain)
```

```
ans=8×2 table
```

	Description	Category
	'Items are occasionally getting stuck in the scanner spools.'	{'Mechanical Failure'}
	'Loud rattling and banging sounds are coming from assembler pistons.'	{'Mechanical Failure'}
	'There are cuts to the power when starting the plant.'	{'Electronic Failure'}

```

{'Fried capacitors in the assembler.' } {'Electronic Failure'}
{'Mixer tripped the fuses.' } {'Electronic Failure'}
{'Burst pipe in the constructing agent is spraying coolant.' } {'Leak'}
{'A fuse is blown in the mixer.' } {'Electronic Failure'}
{'Things continue to tumble off of the belt.' } {'Mechanical Failure'}

```

## Transform Datastore

Create a custom transform function that converts data read from the datastore to a table containing the predictors and the responses. The `transformText` function takes the data read from a `tabularTextDatastore` object and returns a table of predictors and responses. The predictors are  $C$ -by- $S$  arrays of word vectors given by the word embedding `emb`, where  $C$  is the embedding dimension and  $S$  is the sequence length. The responses are categorical labels over the classes.

To get the class names, read the labels from the training data using the `readLabels` function, listed at the end of the example, and find the unique class names.

```

labels = readLabels(tdsTrain, labelName);
classNames = unique(labels);
numObservations = numel(labels);

```

Because tabular text datastores can read multiple rows of data in a single read, you can process a full mini-batch of data in the transform function. To ensure that the transform function processes a full mini-batch of data, set the read size of the tabular text datastore to the mini-batch size that will be used for training.

```

miniBatchSize = 64;
tdsTrain.ReadSize = miniBatchSize;

```

To convert the output of the tabular text data to sequences for training, transform the datastore using the `transform` function.

```

tdsTrain = transform(tdsTrain, @(data) transformText(data, emb, classNames))

```

```

tdsTrain =
    TransformedDatastore with properties:

        UnderlyingDatastore: [1x1 matlab.io.datastore.TabularTextDatastore]
    SupportedOutputFormats: ["txt" "csv" "xlsx" "xls" "parquet" "parq" "png"]
        Transforms: {@(data)transformText(data,emb,classNames)}
    IncludeInfo: 0

```

Preview of the transformed datastore. The predictors are  $C$ -by- $S$  arrays, where  $S$  is the sequence length and  $C$  is the number of features (the embedding dimension). The responses are the categorical labels.

```

preview(tdsTrain)

```

```

ans=8x2 table
    predictors      responses
    _____  _____
    {300x11 single} Mechanical Failure
    {300x11 single} Mechanical Failure
    {300x11 single} Electronic Failure
    {300x11 single} Electronic Failure

```

```
{300×11 single} Electronic Failure
{300×11 single} Leak
{300×11 single} Electronic Failure
{300×11 single} Mechanical Failure
```

### Create and Train LSTM Network

Define the LSTM network architecture. To input sequence data into the network, include a sequence input layer and set the input size to the embedding dimension. Next, include an LSTM layer with 180 hidden units. To use the LSTM layer for a sequence-to-label classification problem, set the output mode to 'last'. Finally, add a fully connected layer with output size equal to the number of classes, a softmax layer, and a classification layer.

```
numFeatures = emb.Dimension;
numHiddenUnits = 180;
numClasses = numel(classNames);
layers = [ ...
    sequenceInputLayer(numFeatures)
    lstmLayer(numHiddenUnits, 'OutputMode', 'last')
    fullyConnectedLayer(numClasses)
    softmaxLayer
    classificationLayer];
```

Specify the training options. Specify the solver to be 'adam' and the gradient threshold to be 2. The datastore does not support shuffling, so set 'Shuffle', to 'never'. Validate the network once per epoch. To monitor the training progress, set the 'Plots' option to 'training-progress'. To suppress verbose output, set 'Verbose' to false.

By default, `trainNetwork` uses a GPU if one is available. To specify the execution environment manually, use the 'ExecutionEnvironment' name-value pair argument of `trainingOptions`. Training on a CPU can take significantly longer than training on a GPU. Training using a GPU requires Parallel Computing Toolbox™ and a supported GPU device. For information on supported devices, see “GPU Support by Release” (Parallel Computing Toolbox).

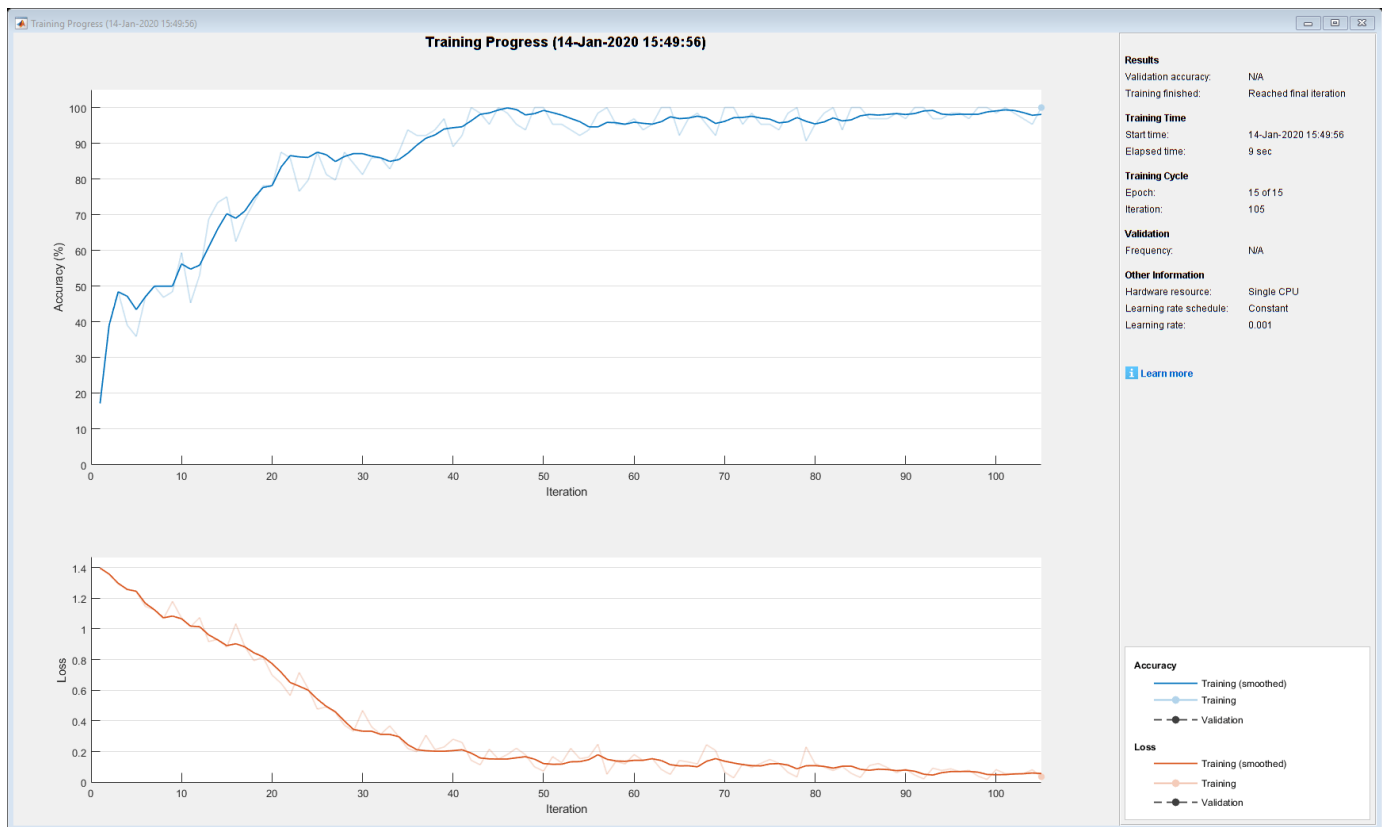
```
numIterationsPerEpoch = floor(numObservations / miniBatchSize);

options = trainingOptions('adam', ...
    'MaxEpochs',15, ...
    'MiniBatchSize',miniBatchSize, ...
    'GradientThreshold',2, ...
    'Shuffle','never', ...
    'Plots','training-progress', ...
    'Verbose',false);
```

Train the LSTM network using the `trainNetwork` function.

```
net = trainNetwork(tdsTrain,layers,options);
```





## Predict Using New Data

Classify the event type of three new reports. Create a string array containing the new reports.

```
reportsNew = [ ...
    "Coolant is pooling underneath sorter."
    "Sorter blows fuses at start up."
    "There are some very loud rattling sounds coming from the assembler."];
```

Preprocess the text data using the preprocessing steps as the training documents.

```
documentsNew = preprocessText(reportsNew);
```

Convert the text data to sequences of embedding vectors using doc2sequence.

```
XNew = doc2sequence(emb,documentsNew);
```

Classify the new sequences using the trained LSTM network.

```
labelsNew = classify(net,XNew)
```

```
labelsNew = 3x1 categorical
    Leak
    Electronic Failure
    Mechanical Failure
```

### Transform Text Function

The `transformText` function takes the data read from a `tabularTextDatastore` object and returns a table of predictors and responses. The predictors are  $C$ -by- $S$  arrays of word vectors given by the word embedding `emb`, where  $C$  is the embedding dimension and  $S$  is the sequence length. The responses are categorical labels over the classes in `classNames`.

```
function dataTransformed = transformText(data,emb,classNames)

% Preprocess documents.
textData = data(:,1);
documents = preprocessText(textData);

% Convert to sequences.
predictors = doc2sequence(emb,documents);

% Read labels.
labels = data(:,2);
responses = categorical(labels,classNames);

% Convert data to table.
dataTransformed = table(predictors,responses);

end
```

### Preprocessing Function

The function `preprocessText` performs these steps:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Convert the text to lowercase using `lower`.
- 3 Erase the punctuation using `erasePunctuation`.

```
function documents = preprocessText(textData)

documents = tokenizedDocument(textData);
documents = lower(documents);
documents = erasePunctuation(documents);

end
```

### Read Labels Function

The `readLabels` function creates a copy of the `tabularTextDatastore` object `ttds` and reads the labels from the `labelName` column.

```
function labels = readLabels(ttds,labelName)

ttdsNew = copy(ttds);
ttdsNew.SelectedVariableNames = labelName;
tbl = readall(ttdsNew);
labels = tbl.(labelName);
```

end

## See Also

`fastTextWordEmbedding` | `wordEmbeddingLayer` | `doc2sequence` | `tokenizedDocument` | `lstmLayer` | `trainNetwork` | `trainingOptions` | `sequenceInputLayer` | `transform`

## Related Examples

- “Create Simple Text Model for Classification” on page 2-2
- “Analyze Text Data Using Topic Models” on page 2-13
- “Analyze Text Data Using Multiword Phrases” on page 2-7
- “Train a Sentiment Classifier” on page 2-71
- “Sequence Classification Using Deep Learning” (Deep Learning Toolbox)
- “Deep Learning in MATLAB” (Deep Learning Toolbox)

## Pride and Prejudice and MATLAB

This example shows how to train a deep learning LSTM network to generate text using character embeddings.

To train a deep learning network for text generation, train a sequence-to-sequence LSTM network to predict the next character in a sequence of characters. To train the network to predict the next character, specify the responses to be the input sequences shifted by one time step.

To use character embeddings, convert each training observation to a sequence of integers, where the integers index into a vocabulary of characters. Include a word embedding layer in the network which learns an embedding of the characters and maps the integers to vectors.

### Load Training Data

Read the HTML code from The Project Gutenberg EBook of Pride and Prejudice, by Jane Austen and parse it using `webread` and `htmlTree`.

```
url = "https://www.gutenberg.org/files/1342/1342-h/1342-h.htm";  
code = webread(url);  
tree = htmlTree(code);
```

Extract the paragraphs by finding the `p` elements. Specify to ignore paragraph elements with class `"toc"` using the CSS selector `' :not(.toc) '`.

```
paragraphs = findElement(tree, 'p:not(.toc)');
```

Extract the text data from the paragraphs using `extractHTMLText`. and remove the empty strings.

```
textData = extractHTMLText(paragraphs);  
textData(textData == "") = [];
```

Remove strings shorter than 20 characters.

```
idx = strlength(textData) < 20;  
textData(idx) = [];
```

Visualize the text data in a word cloud.

```
figure  
wordcloud(textData);  
title("Pride and Prejudice")
```



```
Y = categorical(charactersShifted);  
  
XTrain{i} = X;  
YTrain{i} = Y;  
end
```

During training, by default, the software splits the training data into mini-batches and pads the sequences so that they have the same length. Too much padding can have a negative impact on the network performance.

To prevent the training process from adding too much padding, you can sort the training data by sequence length, and choose a mini-batch size so that sequences in a mini-batch have a similar length.

Get the sequence lengths for each observation.

```
numObservations = numel(XTrain);  
for i=1:numObservations  
    sequence = XTrain{i};  
    sequenceLengths(i) = size(sequence,2);  
end
```

Sort the data by sequence length.

```
[~,idx] = sort(sequenceLengths);  
XTrain = XTrain(idx);  
YTrain = YTrain(idx);
```

### Create and Train LSTM Network

Define the LSTM architecture. Specify a sequence-to-sequence LSTM classification network with 400 hidden units. Set the input size to be the feature dimension of the training data. For sequences of character indices, the feature dimension is 1. Specify a word embedding layer with dimension 200 and specify the number of words (which correspond to characters) to be the highest character value in the input data. Set the output size of the fully connected layer to be the number of categories in the responses. To help prevent overfitting, include a dropout layer after the LSTM layer.

The word embedding layer learns an embedding of characters and maps each character to a 200-dimension vector.

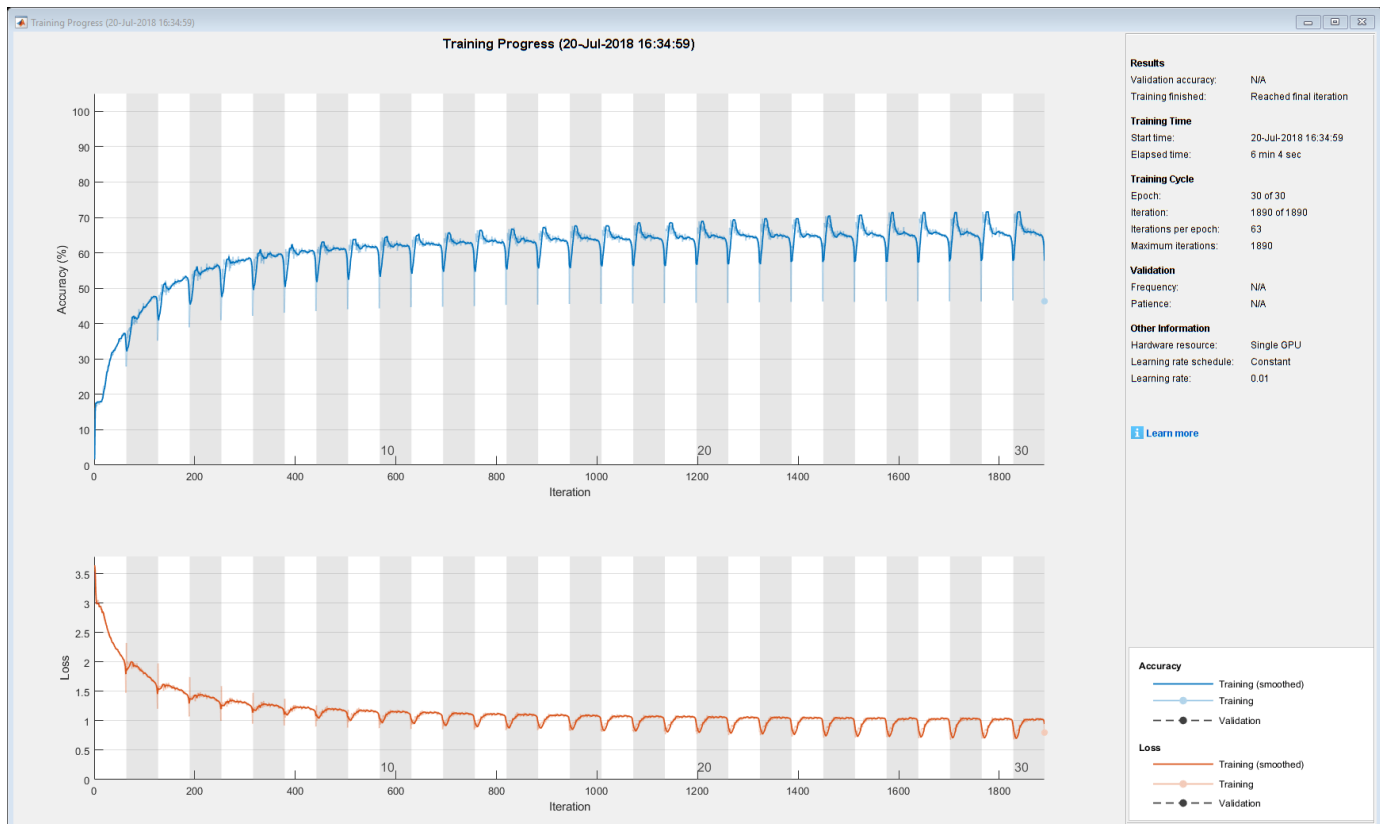
```
inputSize = size(XTrain{1},1);  
numClasses = numel(categories([YTrain{:}]));  
numCharacters = max([textData{:}]);  
  
layers = [  
    sequenceInputLayer(inputSize)  
    wordEmbeddingLayer(200,numCharacters)  
    lstmLayer(400,'OutputMode','sequence')  
    dropoutLayer(0.2);  
    fullyConnectedLayer(numClasses)  
    softmaxLayer  
    classificationLayer];
```

Specify the training options. Specify to train with a mini-batch size of 32 and initial learn rate 0.01. To prevent the gradients from exploding, set the gradient threshold to 1. To ensure the data remains sorted, set 'Shuffle' to 'never'. To monitor the training progress, set the 'Plots' option to 'training-progress'. To suppress verbose output, set 'Verbose' to false.

```
options = trainingOptions('adam', ...
    'MiniBatchSize',32,...
    'InitialLearnRate',0.01, ...
    'GradientThreshold',1, ...
    'Shuffle','never', ...
    'Plots','training-progress', ...
    'Verbose',false);
```

Train the network.

```
net = trainNetwork(XTrain,YTrain,layers,options);
```



## Generate New Text

Generate the first character of the text by sampling a character from a probability distribution according to the first characters of the text in the training data. Generate the remaining characters by using the trained LSTM network to predict the next sequence using the current sequence of generated text. Keep generating characters one-by-one until the network predicts the "end of text" character.

Sample the first character according to the distribution of the first characters in the training data.

```
initialCharacters = extractBefore(textData,2);
firstCharacter = datasample(initialCharacters,1);
generatedText = firstCharacter;
```

Convert the first character to a numeric index.

```
X = double(char(firstCharacter));
```

For the remaining predictions, sample the next character according to the prediction scores of the network. The prediction scores represent the probability distribution of the next character. Sample the characters from the vocabulary of characters given by the class names of the output layer of the network. Get the vocabulary from the classification layer of the network.

```
vocabulary = string(net.Layers(end).ClassNames);
```

Make predictions character by character using `predictAndUpdateState`. For each prediction, input the index of the previous character. Stop predicting when the network predicts the end of text character or when the generated text is 500 characters long. For large collections of data, long sequences, or large networks, predictions on the GPU are usually faster to compute than predictions on the CPU. Otherwise, predictions on the CPU are usually faster to compute. For single time step predictions, use the CPU. To use the CPU for prediction, set the `'ExecutionEnvironment'` option of `predictAndUpdateState` to `'cpu'`.

```
maxLength = 500;
while strlen(generatedText) < maxLength
    % Predict the next character scores.
    [net,characterScores] = predictAndUpdateState(net,X,'ExecutionEnvironment','cpu');

    % Sample the next character.
    newCharacter = datasample(vocabulary,1,'Weights',characterScores);

    % Stop predicting at the end of text.
    if newCharacter == endOfTextCharacter
        break
    end

    % Add the character to the generated text.
    generatedText = generatedText + newCharacter;

    % Get the numeric index of the character.
    X = double(char(newCharacter));
end
```

Reconstruct the generated text by replacing the special characters with their corresponding whitespace and new line characters.

```
generatedText = replace(generatedText,[newlineCharacter whitespaceCharacter],[newline " "])
```

```
generatedText =
```

```
"I wish Mr. Darcy, upon latter of my sort sincerely fixed in the regard to relanth. We were to ;
```

To generate multiple pieces of text, reset the network state between generations using `resetState`.

```
net = resetState(net);
```

## See Also

`wordEmbeddingLayer` | `doc2sequence` | `tokenizedDocument` | `lstmLayer` | `trainNetwork` | `trainingOptions` | `sequenceInputLayer` | `wordcloud` | `extractHTMLText` | `findElement` | `htmlTree`

## Related Examples

- “Generate Text Using Deep Learning” (Deep Learning Toolbox)



- “Word-By-Word Text Generation Using Deep Learning” on page 2-180
- “Create Simple Text Model for Classification” on page 2-2
- “Analyze Text Data Using Topic Models” on page 2-13
- “Analyze Text Data Using Multiword Phrases” on page 2-7
- “Train a Sentiment Classifier” on page 2-71
- “Sequence Classification Using Deep Learning” (Deep Learning Toolbox)
- “Deep Learning in MATLAB” (Deep Learning Toolbox)

## Word-By-Word Text Generation Using Deep Learning

This example shows how to train a deep learning LSTM network to generate text word-by-word.

To train a deep learning network for word-by-word text generation, train a sequence-to-sequence LSTM network to predict the next word in a sequence of words. To train the network to predict the next word, specify the responses to be the input sequences shifted by one time step.

This example reads text from a website. It reads and parses the HTML code to extract the relevant text, then uses a custom mini-batch datastore `documentGenerationDatastore` to input the documents to the network as mini-batches of sequence data. The datastore converts documents to sequences of numeric word indices. The deep learning network is an LSTM network that contains a word embedding layer.

A mini-batch datastore is an implementation of a datastore with support for reading data in batches. You can use a mini-batch datastore as a source of training, validation, test, and prediction data sets for deep learning applications. Use mini-batch datastores to read out-of-memory data or to perform specific preprocessing operations when reading batches of data.

You can adapt the custom mini-batch datastore specified by `documentGenerationDatastore.m` to your data by customizing the functions. This file is eattached to this example as a supporting file. To access this file, open the example as a live script. For an example showing how to create your own custom mini-batch datastore, see “Develop Custom Mini-Batch Datastore” (Deep Learning Toolbox).

### Load Training Data

Load the training data. Read the HTML code from Alice's Adventures in Wonderland by Lewis Carroll from Project Gutenberg.

```
url = "https://www.gutenberg.org/files/11/11-h/11-h.htm";
code = webread(url);
```

### Parse HTML Code

The HTML code contains the relevant text inside `<p>` (paragraph) elements. Extract the relevant text by parsing the HTML code using `htmlTree` and then finding all the elements with element name `"p"`.

```
tree = htmlTree(code);
selector = "p";
subtrees = findElement(tree,selector);
```

Extract the text data from the HTML subtrees using `extractHTMLText` and view the first 10 paragraphs.

```
textData = extractHTMLText(subtrees);
textData(1:10)
```

```
ans = 10x1 string
    "Alice was beginning to get very tired of sitting by her sister on the bank, and of having no
    "So she was considering in her own mind (as well as she could, for the hot day made her feel
    "There was nothing so very remarkable in that; nor did Alice think it so very much out of the
    "In another moment down went Alice after it, never once considering how in the world she was
    "The rabbit-hole went straight on like a tunnel for some way, and then dipped suddenly down,
    "Either the well was very deep, or she fell very slowly, for she had plenty of time as she wa
    "“Well!” thought Alice to herself, “after such a fall as this, I shall think nothing of tumb
```



### Prepare Data for Training

Create a datastore that contains the data for training using `documentGenerationDatastore`. For the predictors, this datastore converts the documents into sequences of word indices using a word encoding. The first word index for each document corresponds to a "start of text" token. The "start of text" token is given by the string `"startOfText"`. For the responses, the datastore returns categorical sequences of the words shifted by one.

Tokenize the text data using `tokenizedDocument`.

```
documents = tokenizedDocument(textData);
```

Create a document generation datastore using the tokenized documents.

```
ds = documentGenerationDatastore(documents);
```

To reduce the amount of padding added to the sequences, sort the documents in the datastore by sequence length.

```
ds = sort(ds);
```

### Create and Train LSTM Network

Define the LSTM network architecture. To input sequence data into the network, include a sequence input layer and set the input size to 1. Next, include a word embedding layer of dimension 100 and the same number of words as the word encoding. Next, include an LSTM layer and specify the hidden size to be 100. Finally, add a fully connected layer with the same size as the number of classes, a softmax layer, and a classification layer. The number of classes is the number of words in the vocabulary plus an extra class for the "end of text" class.

```
inputSize = 1;
embeddingDimension = 100;
numWords = numel(ds.Encoding.Vocabulary);
numClasses = numWords + 1;

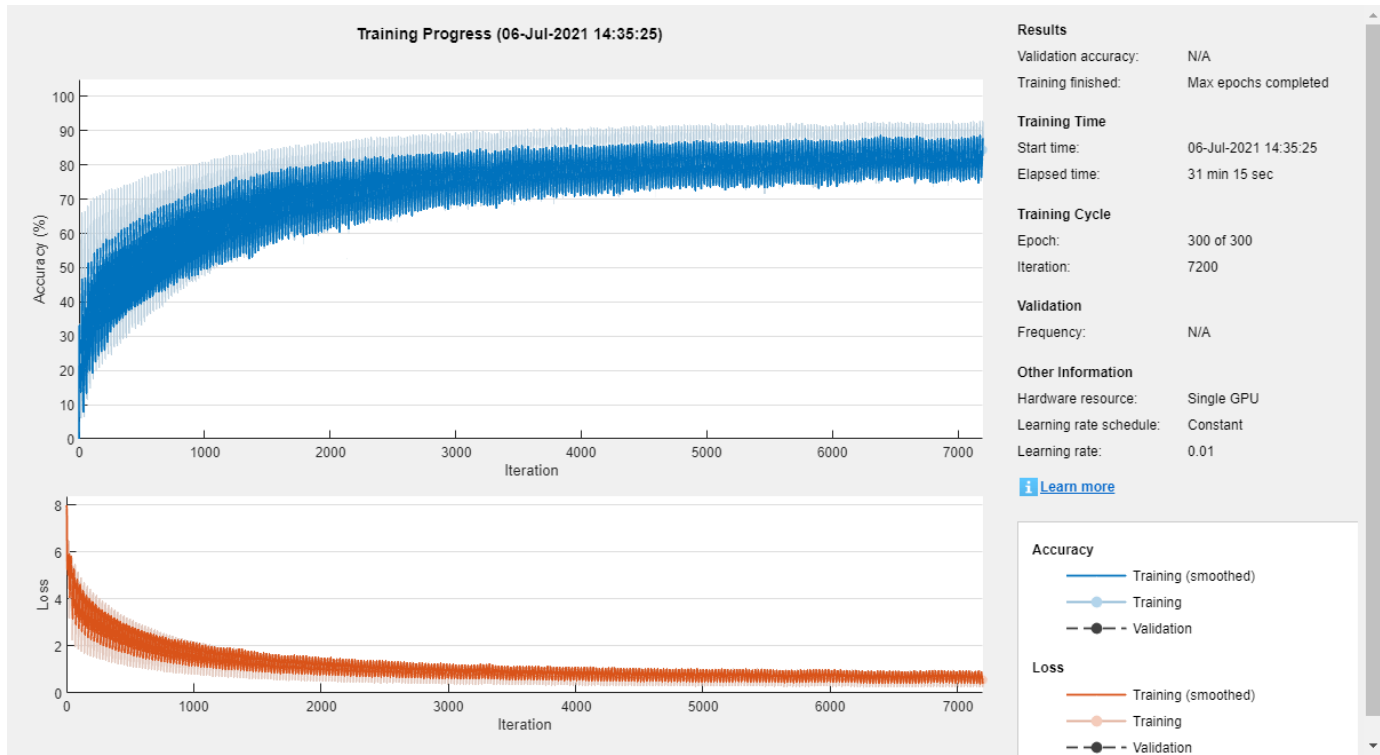
layers = [
    sequenceInputLayer(inputSize)
    wordEmbeddingLayer(embeddingDimension,numWords)
    lstmLayer(100)
    dropoutLayer(0.2)
    fullyConnectedLayer(numClasses)
    softmaxLayer
    classificationLayer];
```

Specify the training options. Specify the solver to be `'adam'`. Train for 300 epochs with learn rate 0.01. Set the mini-batch size to 32. To keep the data sorted by sequence length, set the `'Shuffle'` option to `'never'`. To monitor the training progress, set the `'Plots'` option to `'training-progress'`. To suppress verbose output, set `'Verbose'` to `false`.

```
options = trainingOptions('adam', ...
    'MaxEpochs',300, ...
    'InitialLearnRate',0.01, ...
    'MiniBatchSize',32, ...
    'Shuffle','never', ...
    'Plots','training-progress', ...
    'Verbose',false);
```

Train the network using `trainNetwork`.

```
net = trainNetwork(ds, layers, options);
```



## Generate New Text

Generate the first word of the text by sampling a word from a probability distribution according to the first words of the text in the training data. Generate the remaining words by using the trained LSTM network to predict the next time step using the current sequence of generated text. Keep generating words one-by-one until the network predicts the "end of text" word.

To make the first prediction using the network, input the index that represents the "start of text" token. Find the index by using the `word2ind` function with the word encoding used by the document datastore.

```
enc = ds.Encoding;
wordIndex = word2ind(enc, "startOfText")
```

```
wordIndex = 1
```

For the remaining predictions, sample the next word according to the prediction scores of the network. The prediction scores represent the probability distribution of the next word. Sample the words from the vocabulary given by the class names of the output layer of the network.

```
vocabulary = string(net.Layers(end).Classes);
```

Make predictions word by word using `predictAndUpdateState`. For each prediction, input the index of the previous word. Stop predicting when the network predicts the end of text word or when the generated text is 500 characters long. For large collections of data, long sequences, or large networks, predictions on the GPU are usually faster to compute than predictions on the CPU. Otherwise, predictions on the CPU are usually faster to compute. For single time step predictions, use

the CPU. To use the CPU for prediction, set the 'ExecutionEnvironment' option of predictAndUpdateState to 'cpu'.

```
generatedText = "";
maxLength = 500;
while strlength(generatedText) < maxLength
    % Predict the next word scores.
    [net,wordScores] = predictAndUpdateState(net,wordIndex,'ExecutionEnvironment','cpu');

    % Sample the next word.
    newWord = datasample(vocabulary,1,'Weights',wordScores);

    % Stop predicting at the end of text.
    if newWord == "EndOfText"
        break
    end

    % Add the word to the generated text.
    generatedText = generatedText + " " + newWord;

    % Find the word index for the next input.
    wordIndex = word2ind(enc,newWord);
end
```

The generation process introduces whitespace characters between each prediction, which means that some punctuation characters appear with unnecessary spaces before and after. Reconstruct the generated text by removing the spaces before and after the appropriate punctuation characters.

Remove the spaces that appear before the specified punctuation characters.

```
punctuationCharacters = [". " ", " "' " ) " : " ? " ! "];
generatedText = replace(generatedText," " + punctuationCharacters,punctuationCharacters);
```

Remove the spaces that appear after the specified punctuation characters.

```
punctuationCharacters = [{" " ""}];
generatedText = replace(generatedText,punctuationCharacters + " ",punctuationCharacters)

generatedText =
" " " Just about as much right, " said the Duchess, " and that's all the least, " said the Hatter.
```

To generate multiple pieces of text, reset the network state between generations using resetState.

```
net = resetState(net);
```

## See Also

wordEmbeddingLayer | doc2sequence | tokenizedDocument | lstmLayer | trainNetwork | trainingOptions | sequenceInputLayer | wordcloud | extractHTMLText | findElement | htmlTree

## Related Examples

- "Generate Text Using Deep Learning" (Deep Learning Toolbox)
- "Create Simple Text Model for Classification" on page 2-2
- "Analyze Text Data Using Topic Models" on page 2-13

- “Analyze Text Data Using Multiword Phrases” on page 2-7
- “Train a Sentiment Classifier” on page 2-71
- “Sequence Classification Using Deep Learning” (Deep Learning Toolbox)
- “Deep Learning in MATLAB” (Deep Learning Toolbox)

## Generate Text Using Autoencoders

This example shows how to generate text data using autoencoders.

An autoencoder is a type of deep learning network that is trained to replicate its input. An autoencoder consists of two smaller networks: an encoder and a decoder. The encoder maps the input data to a feature vector in some latent space. The decoder reconstructs data using vectors in this latent space.

The training process is unsupervised. In other words, the model does not require labeled data. To generate text, you can use the decoder to reconstruct text from arbitrary input.

This example trains an autoencoder to generate text. The encoder uses a word embedding and an LSTM operation to map the input text into latent vectors. The decoder uses an LSTM operation and the same embedding to reconstruct the text from the latent vectors.

### Load Data

The file `sonnets.txt` contains all of Shakespeare's sonnets in a single text file.

Read the Shakespeare's Sonnets data from the file `"sonnets.txt"`.

```
filename = "sonnets.txt";
textData = fileread(filename);
```

The sonnets are indented by two whitespace characters. Remove the indentations using `replace` and split the text into separate lines using the `split` function. Remove the header from the first nine elements and the short sonnet titles.

```
textData = replace(textData, "  ", "");
textData = split(textData, newline);
textData(1:9) = [];
textData(strlen(textData)<5) = [];
```

### Prepare Data

Create a function that tokenizes and preprocesses the text data. The function `preprocessText`, listed at the end of the example, performs these steps:

- 1 Prepends and appends each input string with the specified start and stop tokens, respectively.
- 2 Tokenize the text using `tokenizedDocument`.

Preprocess the text data and specify the start and stop tokens `"<start>"` and `"<stop>"`, respectively.

```
startToken = "<start>";
stopToken = "<stop>";
documents = preprocessText(textData, startToken, stopToken);
```

Create a word encoding object from the tokenized documents.

```
enc = wordEncoding(documents);
```

When training a deep learning model, the input data must be a numeric array containing sequences of a fixed length. Because the documents have different lengths, you must pad the shorter sequences with a padding value.



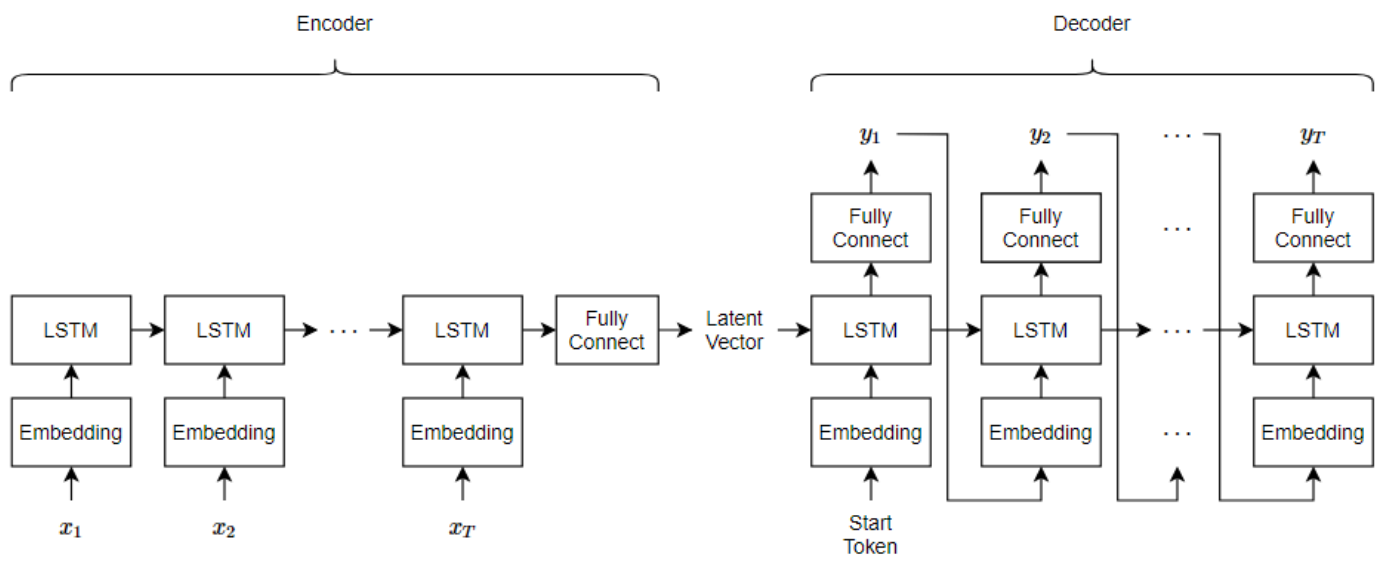
Recreate the word encoding to also include a padding token and determine the index of that token.

```
paddingToken = "<pad>";
newVocabulary = [enc.Vocabulary paddingToken];
enc = wordEncoding(newVocabulary);
paddingIdx = word2ind(enc, paddingToken)

paddingIdx = 3595
```

### Initialize Model Parameters

Initialize the parameters for the following model.



Here,  $T$  is the sequence length,  $x_1, \dots, x_T$  is the input sequence of word indices, and  $y_1, \dots, y_T$  is the reconstructed sequence.

The encoder maps sequences of word indices to a latent vector by converting the input to sequences of word vectors using an embedding, inputting the word vector sequences into an LSTM operation, and applying a fully connected operation to the last time step of the LSTM output. The decoder reconstructs the input using an LSTM initialized the encoder output. For each time step, the decoder predicts the next time step and uses the output for the next time-step predictions. Both the encoder and the decoder use the same embedding.

Specify the dimensions of the parameters.

```
embeddingDimension = 100;
numHiddenUnits = 150;
latentDimension = 75;
vocabularySize = enc.NumWords;
```

Create a struct for the parameters.

```
parameters = struct;
```

Initialize the weights of the embedding using the Gaussian using the `initializeGaussian` function which is attached to this example as a supporting file. Specify a mean of 0 and a standard deviation of 0.01. To learn more, see "Gaussian Initialization" (Deep Learning Toolbox).

```
sz = [embeddingDimension vocabularySize];
mu = 0;
sigma = 0.01;
parameters.emb.Weights = initializeGaussian(sz,mu,sigma);
```

Initialize the learnable parameters for the encoder LSTM operation:

- Initialize the input weights with the Glorot initializer using the `initializeGlorot` function which is attached to this example as a supporting file. To learn more, see “Glorot Initialization” (Deep Learning Toolbox).
- Initialize the recurrent weights with the orthogonal initializer using the `initializeOrthogonal` function which is attached to this example as a supporting file. To learn more, see “Orthogonal Initialization” (Deep Learning Toolbox).
- Initialize the bias with the unit forget gate initializer using the `initializeUnitForgetGate` function which is attached to this example as a supporting file. To learn more, see “Unit Forget Gate Initialization” (Deep Learning Toolbox).

```
sz = [4*numHiddenUnits embeddingDimension];
numOut = 4*numHiddenUnits;
numIn = embeddingDimension;
```

```
parameters.lstmEncoder.InputWeights = initializeGlorot(sz,numOut,numIn);
parameters.lstmEncoder.RecurrentWeights = initializeOrthogonal([4*numHiddenUnits numHiddenUnits]
parameters.lstmEncoder.Bias = initializeUnitForgetGate(numHiddenUnits);
```

Initialize the learnable parameters for the encoder fully connected operation:

- Initialize the weights with the Glorot initializer.
- Initialize the bias with zeros using the `initializeZeros` function which is attached to this example as a supporting file. To learn more, see “Zeros Initialization” (Deep Learning Toolbox).

```
sz = [latentDimension numHiddenUnits];
numOut = latentDimension;
numIn = numHiddenUnits;
```

```
parameters.fcEncoder.Weights = initializeGlorot(sz,numOut,numIn);
parameters.fcEncoder.Bias = initializeZeros([latentDimension 1]);
```

Initialize the learnable parameters for the decoder LSTM operation:

- Initialize the input weights with the Glorot initializer.
- Initialize the recurrent weights with the orthogonal initializer.
- Initialize the bias with the unit forget gate initializer.

```
sz = [4*latentDimension embeddingDimension];
numOut = 4*latentDimension;
numIn = embeddingDimension;
```

```
parameters.lstmDecoder.InputWeights = initializeGlorot(sz,numOut,numIn);
parameters.lstmDecoder.RecurrentWeights = initializeOrthogonal([4*latentDimension latentDimension]
parameters.lstmDecoder.Bias = initializeZeros([4*latentDimension 1]);
```

Initialize the learnable parameters for the decoder fully connected operation:

- Initialize the weights with the Glorot initializer.

- Initialize the bias with zeros.

```
sz = [vocabularySize latentDimension];
numOut = vocabularySize;
numIn = latentDimension;

parameters.fcDecoder.Weights = initializeGlorot(sz,numOut,numIn);
parameters.fcDecoder.Bias = initializeZeros([vocabularySize 1]);
```

To learn more about weight initialization, see “Initialize Learnable Parameters for Model Function” (Deep Learning Toolbox).

### Define Model Encoder Function

Create the function `modelEncoder`, listed in the Encoder Model Function on page 2-0 section of the example, that computes the output of the encoder model. The `modelEncoder` function, takes as input sequences of word indices, the model parameters, and the sequence lengths, and returns the corresponding latent feature vector. To learn more about defining a model encoder function, see “Define Text Encoder Model Function” (Deep Learning Toolbox).

### Define Model Decoder Function

Create the function `modelDecoder`, listed in the Decoder Model Function on page 2-0 section of the example, that computes the output of the decoder model. The `modelDecoder` function, takes as input sequences of word indices, the model parameters, and the sequence lengths, and returns the corresponding latent feature vector. To learn more about defining a model decoder function, see “Define Text Decoder Model Function” (Deep Learning Toolbox).

### Define Model Loss Function

The `modelLoss` function, listed in the Model Loss Function on page 2-0 section of the example, takes as input the model learnable parameters, the input data and a vector of sequence lengths for masking, and returns the loss, and the gradients of the loss with respect to the learnable parameters. To learn more about defining a model loss function, see “Define Model Loss Function for Custom Training Loop” (Deep Learning Toolbox).

### Specify Training Options

Specify the options for training.

Train for 100 epochs with a mini-batch size of 128.

```
miniBatchSize = 128;
numEpochs = 100;
```

Train with a learning rate of 0.01.

```
learnRate = 0.01;
```

### Train Network

Train the network using a custom training loop.

Initialize the parameters for the Adam optimizer.

```
trailingAvg = [];
trailingAvgSq = [];
```

Initialize the training progress plot. Create an animated line that plots the loss against the corresponding iteration.

```
figure
C = colororder;
lineLossTrain = animatedline(Color=C(2,:));
xlabel("Iteration")
ylabel("Loss")
ylim([0 inf])
grid on
```

Train the model. For the first epoch, shuffle the data and loop over mini-batches of data.

For each mini-batch:

- Convert the text data to sequences of word indices.
- Convert the data to `darray`.
- For GPU training, convert the data to `gpuArray` objects.
- Compute loss and gradients.
- Update the learnable parameters using the `adamupdate` function.
- Update the training progress plot.

Train on a GPU if one is available. Using a GPU requires Parallel Computing Toolbox™ and a supported GPU device. For information on supported devices, see “GPU Support by Release” (Parallel Computing Toolbox).

Training can take some time to run.

```
numObservations = numel(documents);
numIterationsPerEpoch = floor(numObservations / miniBatchSize);

iteration = 0;
start = tic;

for epoch = 1:numEpochs

    % Shuffle.
    idx = randperm(numObservations);
    documents = documents(idx);

    for i = 1:numIterationsPerEpoch
        iteration = iteration + 1;

        % Read mini-batch.
        idx = (i-1)*miniBatchSize+1:i*miniBatchSize;
        documentsBatch = documents(idx);

        % Convert to sequences.
        X = doc2sequence(enc,documentsBatch, ...
            PaddingDirection="right", ...
            PaddingValue=paddingIdx);

        X = cat(1,X{:});

        % Convert to darray.
```

```

X = dldarray(X, "BTC");

% If training on a GPU, then convert data to gpuArray.
if canUseGPU
    X = gpuArray(X);
end

% Calculate sequence lengths.
sequenceLengths = doclength(documentsBatch);

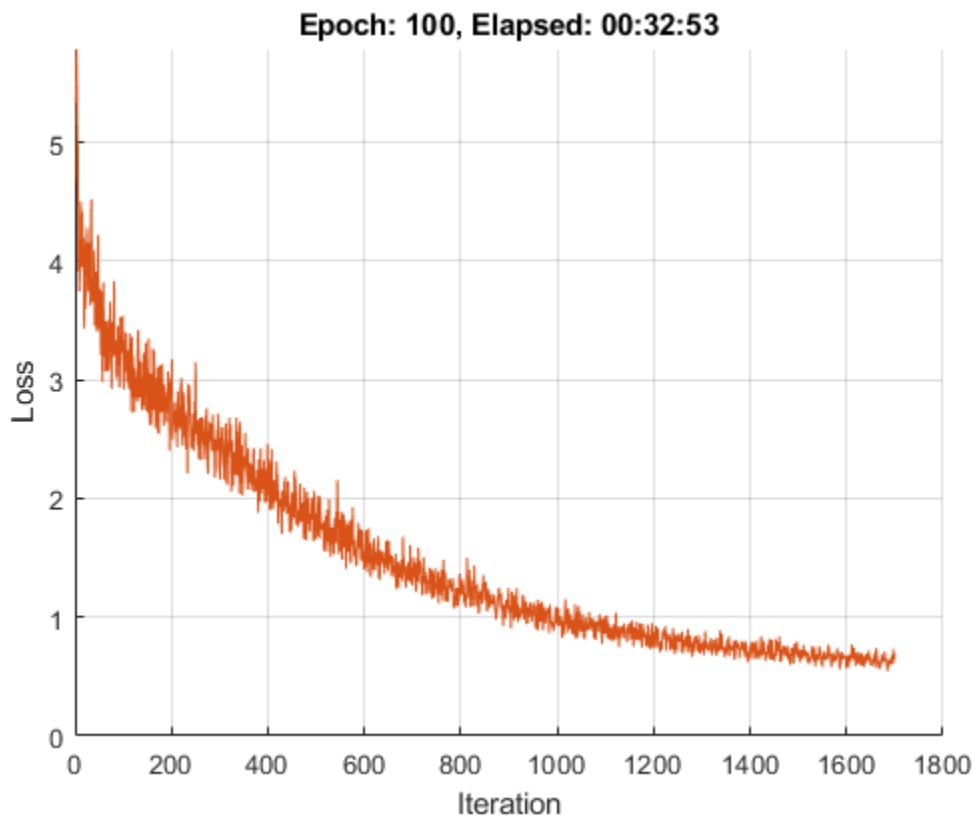
% Evaluate model loss and gradients.
[loss,gradients] = dlfeval(@modelLoss, parameters, X, sequenceLengths);

% Update learnable parameters.
[parameters,trailingAvg,trailingAvgSq] = adamupdate(parameters,gradients, ...
    trailingAvg,trailingAvgSq,iteration,learnRate);

% Display the training progress.
D = duration(0,0,toc(start),Format="hh:mm:ss");
loss = double(loss);
addpoints(lineLossTrain,iteration,loss)
title("Epoch: " + epoch + ", Elapsed: " + string(D))

drawnow
end
end

```



### Generate Text

Generate text using closed loop generation by initializing the decoder with different random states. Closed loop generation is when the model generates data one time-step at a time and uses the previous prediction as input for the next prediction.

Specify to generate 3 sequences of length 16.

```
numGenerations = 3;  
sequenceLength = 16;
```

Create an array of random values to initialize the decoder state.

```
Z = darray(randn(latentDimension,numGenerations),"CB");
```

If predicting on a GPU, then convert data to `gpuArray`.

```
if canUseGPU  
    Z = gpuArray(Z);  
end
```

Make predictions using the `modelPredictions` function, listed at the end of the example. The `modelPredictions` function returns the output scores of the decoder given the model parameters, decoder initial state, maximum sequence length, word encoding, start token, and mini-batch size.

```
Y = modelDecoderPredictions(parameters,Z,sequenceLength,enc,startToken,miniBatchSize);
```

Find the word indices with the highest scores.

```
[~,idx] = max(Y,[],1);  
idx = squeeze(idx);
```

Convert the numeric indices to words and join them using the `join` function.

```
strGenerated = join(enc.Vocabulary(idx));
```

Extract the text before the first stop token using the `extractBefore` function. To prevent the function from returning missing when there are no stop tokens, append a stop token to the end of each sequence.

```
strGenerated = extractBefore(strGenerated+stopToken,stopToken);
```

Remove padding tokens.

```
strGenerated = erase(strGenerated,paddingToken);
```

The generation process introduces whitespace characters between each prediction, which means that some punctuation characters appear with unnecessary spaces before and after. Reconstruct the generated text by removing the spaces before and after the appropriate punctuation characters.

Remove the spaces that appear before the specified punctuation characters.

```
punctuationCharacters = [". " ", " "'" " ) " : " ; " ? " ! "];  
strGenerated = replace(strGenerated," " + punctuationCharacters,punctuationCharacters);
```

Remove the spaces that appear after the specified punctuation characters.

```
punctuationCharacters = [" (" "'"];  
strGenerated = replace(strGenerated,punctuationCharacters + " ",punctuationCharacters);
```

Remove leading and trailing white space using the `strip` function and view the generated text.

```
strGenerated = strip(strGenerated)

strGenerated = 3×1 string
    "me whose fool black grounded less waning travels less pine pine sing cool thrive kindness t
    "perjur'd outward a looks black, here might."
    "birds him antique side his hours age,"
```

### Encoder Model Function

The `modelEncoder` function, takes as input the model parameters, sequences of word indices, and the sequence lengths, and returns the corresponding latent feature vector.

Because the input data contains padded sequences of different lengths, the padding can have adverse effects on loss calculations. For the LSTM operation, instead of returning the output of the last time step of the sequence (which likely corresponds to the LSTM state after processing lots of padding values), determine the actual last time step given by the `sequenceLengths` input.

```
function Z = modelEncoder(parameters,X,sequenceLengths)

% Embedding.
weights = parameters.emb.Weights;
Z = embed(X,weights);

% LSTM.
inputWeights = parameters.lstmEncoder.InputWeights;
recurrentWeights = parameters.lstmEncoder.RecurrentWeights;
bias = parameters.lstmEncoder.Bias;

numHiddenUnits = size(recurrentWeights,2);
hiddenState = zeros(numHiddenUnits,1,"like",X);
cellState = zeros(numHiddenUnits,1,"like",X);

Z1 = lstm(Z,hiddenState,cellState,inputWeights,recurrentWeights,bias);

% Output mode 'last' with masking.
miniBatchSize = size(Z1,2);
Z = zeros(numHiddenUnits,miniBatchSize,"like",Z1);

for n = 1:miniBatchSize
    t = sequenceLengths(n);
    Z(:,n) = Z1(:,n,t);
end

% Fully connect.
weights = parameters.fcEncoder.Weights;
bias = parameters.fcEncoder.Bias;
Z = fullyconnect(Z,weights,bias,DataFormat="CB");

end
```

### Decoder Model Function

The `modelDecoder` function, takes as input the model parameters, sequences of word indices, and the network state, and returns the decoded sequences.

Because the `lstm` function is *stateful* (when given a time series as input, the function propagates and updates the state between each time step) and that the `embed` and `fullyconnect` functions are time-distributed by default (when given a time series as input, the functions operate on each time step independently), the `modelDecoder` function supports both sequence and single time-step inputs.

```
function [Y,state] = modelDecoder(parameters,X,state)

% Embedding.
weights = parameters.emb.Weights;
X = embed(X,weights);

% LSTM.
inputWeights = parameters.lstmDecoder.InputWeights;
recurrentWeights = parameters.lstmDecoder.RecurrentWeights;
bias = parameters.lstmDecoder.Bias;

hiddenState = state.HiddenState;
cellState = state.CellState;

[Y,hiddenState,cellState] = lstm(X,hiddenState,cellState, ...
    inputWeights,recurrentWeights,bias);

state.HiddenState = hiddenState;
state.CellState = cellState;

% Fully connect.
weights = parameters.fcDecoder.Weights;
bias = parameters.fcDecoder.Bias;
Y = fullyconnect(Y,weights,bias);

% Softmax.
Y = softmax(Y);

end
```

### **Model Loss Function**

The `modelLoss` function takes as input the model learnable parameters, the input data `X`, and a vector of sequence lengths for masking, and returns the loss and the gradients of the loss with respect to the learnable parameters.

To calculate the masked loss, the model loss function uses the `maskedCrossEntropy` function, listed at the end of the example. To train the decoder to predict the next time-step of the sequence, specify the targets to be the input sequences shifted by one time-step.

To learn more about defining a model loss function, see “Define Model Loss Function for Custom Training Loop” (Deep Learning Toolbox).

```
function [loss,gradients] = modelLoss(parameters,X,sequenceLengths)

% Model encoder.
Z = modelEncoder(parameters,X,sequenceLengths);

% Initialize LSTM state.
state = struct;
state.HiddenState = Z;
```



```

state.CellState = zeros(size(Z), "like", Z);

% Teacher forcing.
Y = modelDecoder(parameters, X, state);

% Loss.
Y = Y(:, :, 1:end-1);
T = X(:, :, 2:end);
loss = mean(maskedCrossEntropy(Y, T, sequenceLengths));

% Gradients.
gradients = dlgradient(loss, parameters);

% Normalize loss for plotting.
sequenceLength = size(X, 3);
loss = loss / sequenceLength;

end

```

### Model Predictions Function

The `modelPredictions` function returns the output scores of the decoder given the model parameters, decoder initial state, maximum sequence length, word encoding, start token, and mini-batch size.

```

function Y = modelDecoderPredictions(parameters, Z, maxLength, enc, startToken, miniBatchSize)

numObservations = size(Z, 2);
numIterations = ceil(numObservations / miniBatchSize);

startTokenIdx = word2ind(enc, startToken);
vocabularySize = enc.NumWords;

Y = zeros(vocabularySize, numObservations, maxLength, "like", Z);

% Loop over mini-batches.
for i = 1:numIterations
    idxMiniBatch = (i-1)*miniBatchSize+1:min(i*miniBatchSize, numObservations);
    miniBatchSize = numel(idxMiniBatch);

    % Initialize state.
    state = struct;
    state.HiddenState = Z(:, idxMiniBatch);
    state.CellState = zeros(size(Z(:, idxMiniBatch)), "like", Z);

    % Initialize decoder input.
    decoderInput = dlarray repmat(startTokenIdx, [1 miniBatchSize]), "CBT");

    % Loop over time steps.
    for t = 1:maxLength
        % Predict next time step.
        [Y(:, idxMiniBatch, t), state] = modelDecoder(parameters, decoderInput, state);

        % Closed loop generation.
        [~, idx] = max(Y(:, idxMiniBatch, t));
        decoderInput = dlarray(idx, "CB");
    end
end
end

```

end

### Masked Cross Entropy Loss Function

The `maskedCrossEntropy` function calculates the loss between the specified input sequences and target sequences ignoring any time steps containing padding using the specified vector of sequence lengths.

```
function maskedLoss = maskedCrossEntropy(Y,T,sequenceLengths)

numClasses = size(Y,1);
miniBatchSize = size(Y,2);
sequenceLength = size(Y,3);

maskedLoss = zeros(sequenceLength,miniBatchSize,"like",Y);

for t = 1:sequenceLength
    T1 = single(oneHot(T(:, :, t), numClasses));

    mask = (t <= sequenceLengths)';

    maskedLoss(t, :) = mask .* crossentropy(Y(:, :, t), T1);
end

maskedLoss = sum(maskedLoss, 1);

end
```

### Text Preprocessing Function

The function `preprocessText` performs these steps:

- 1 Prepends and appends each input string with the specified start and stop tokens, respectively.
- 2 Tokenize the text using `tokenizedDocument`.

```
function documents = preprocessText(textData, startToken, stopToken)

% Add start and stop tokens.
textData = startToken + textData + stopToken;

% Tokenize the text.
documents = tokenizedDocument(textData, 'CustomTokens', [startToken stopToken]);

end
```

### One-Hot Encoding Function

The `oneHot` function converts an array of numeric indices to one-hot encoded vectors.

```
function oh = oneHot(idx, outputSize)

miniBatchSize = numel(idx);
oh = zeros(outputSize, miniBatchSize);

for n = 1:miniBatchSize
    c = idx(n);
    oh(c, n) = 1;
end
```

end

end

## See Also

[wordEncoding](#) | [word2ind](#) | [doc2sequence](#) | [tokenizedDocument](#)

## More About

- “Sequence-to-Sequence Translation Using Attention” on page 2-130
- “Generate Text Using Deep Learning” (Deep Learning Toolbox)
- “Define Text Encoder Model Function” on page 2-198
- “Define Text Decoder Model Function” on page 2-205
- “Define Custom Training Loops, Loss Functions, and Networks” (Deep Learning Toolbox)

## Define Text Encoder Model Function

This example shows how to define a text encoder model function.

In the context of deep learning, an encoder is the part of a deep learning network that maps the input to some latent space. You can use these vectors for various tasks. For example,

- Classification by applying a softmax operation to the encoded data and using cross entropy loss.
- Sequence-to-sequence translation by using the encoded vector as a context vector.

### Load Data

The file `sonnets.txt` contains all of Shakespeare's sonnets in a single text file.

Read the Shakespeare's Sonnets data from the file `"sonnets.txt"`.

```
filename = "sonnets.txt";
textData = fileread(filename);
```

The sonnets are indented by two whitespace characters. Remove the indentations using `replace` and `split` the text into separate lines using the `split` function. Remove the header from the first nine elements and the short sonnet titles.

```
textData = replace(textData, "  ", "");
textData = split(textData, newline);
textData(1:9) = [];
textData(strlen(textData)<5) = [];
```

### Prepare Data

Create a function that tokenizes and preprocesses the text data. The function `preprocessText`, listed at the end of the example, performs these steps:

- 1 Prepends and appends each input string with the specified start and stop tokens, respectively.
- 2 Tokenize the text using `tokenizedDocument`.

Preprocess the text data and specify the start and stop tokens `"<start>"` and `"<stop>"`, respectively.

```
startToken = "<start>";
stopToken = "<stop>";
documents = preprocessText(textData, startToken, stopToken);
```

Create a word encoding object from the tokenized documents.

```
enc = wordEncoding(documents);
```

When training a deep learning model, the input data must be a numeric array containing sequences of a fixed length. Because the documents have different lengths, you must pad the shorter sequences with a padding value.

Recreate the word encoding to also include a padding token and determine the index of that token.

```
paddingToken = "<pad>";
newVocabulary = [enc.Vocabulary paddingToken];
enc = wordEncoding(newVocabulary);
paddingIdx = word2ind(enc, paddingToken)
```

```
paddingIdx = 3595
```

### Initialize Model Parameters

The goal of the encoder is to map sequences of word indices to vectors in some latent space.

Initialize the parameters for the following model.



This model uses three operations:

- The embedding maps word indices in the range 1 though `vocabularySize` to vectors of dimension `embeddingDimension`, where `vocabularySize` is the number of words in the encoding vocabulary and `embeddingDimension` is the number of components learned by the embedding.
- The LSTM operation takes as input sequences of word vectors and outputs 1-by-`numHiddenUnits` vectors, where `numHiddenUnits` is the number of hidden units in the LSTM operation.
- The fully connected operation multiplies the input by a weight matrix adding bias and outputs vectors of size `latentDimension`, where `latentDimension` is the dimension of the latent space.

Specify the dimensions of the parameters.

```
embeddingDimension = 100;
numHiddenUnits = 150;
latentDimension = 50;
vocabularySize = enc.NumWords;
```

Create a struct for the parameters.

```
parameters = struct;
```

Initialize the weights of the embedding using the Gaussian using the `initializeGaussian` function which is attached to this example as a supporting file. Specify a mean of 0 and a standard deviation of 0.01. To learn more, see “Gaussian Initialization” (Deep Learning Toolbox).

```
mu = 0;
sigma = 0.01;
parameters.emb.Weights = initializeGaussian([embeddingDimension vocabularySize],mu,sigma);
```

Initialize the learnable parameters for the encoder LSTM operation:

- Initialize the input weights with the Glorot initializer using the `initializeGlorot` function which is attached to this example as a supporting file. To learn more, see “Glorot Initialization” (Deep Learning Toolbox).
- Initialize the recurrent weights with the orthogonal initializer using the `initializeOrthogonal` function which is attached to this example as a supporting file. To learn more, see “Orthogonal Initialization” (Deep Learning Toolbox).

- Initialize the bias with the unit forget gate initializer using the `initializeUnitForgetGate` function which is attached to this example as a supporting file. To learn more, see “Unit Forget Gate Initialization” (Deep Learning Toolbox).

The sizes of the learnable parameters depend on the size of the input. Because the inputs to the LSTM operation are sequences of word vectors from the embedding operation, the number of input channels is `embeddingDimension`.

- The input weight matrix has size `4*numHiddenUnits-by-inputSize`, where `inputSize` is the dimension of the input data.
- The recurrent weight matrix has size `4*numHiddenUnits-by-numHiddenUnits`.
- The bias vector has size `4*numHiddenUnits-by-1`.

```
sz = [4*numHiddenUnits embeddingDimension];
numOut = 4*numHiddenUnits;
numIn = embeddingDimension;
```

```
parameters.lstmEncoder.InputWeights = initializeGlorot(sz,numOut,numIn);
parameters.lstmEncoder.RecurrentWeights = initializeOrthogonal([4*numHiddenUnits numHiddenUnits]);
parameters.lstmEncoder.Bias = initializeUnitForgetGate(numHiddenUnits);
```

Initialize the learnable parameters for the encoder fully connected operation:

- Initialize the weights with the Glorot initializer.
- Initialize the bias with zeros using the `initializeZeros` function which is attached to this example as a supporting file. To learn more, see “Zeros Initialization” (Deep Learning Toolbox).

The sizes of the learnable parameters depend on the size of the input. Because the inputs to the fully connected operation are the outputs of the LSTM operation, the number of input channels is `numHiddenUnits`. To make the fully connected operation output vectors with size `latentDimension`, specify an output size of `latentDimension`.

- The weights matrix has size `outputSize-by-inputSize`, where `outputSize` and `inputSize` correspond to the output and input dimensions, respectively.
- The bias vector has size `outputSize-by-1`.

```
sz = [latentDimension numHiddenUnits];
numOut = latentDimension;
numIn = numHiddenUnits;
```

```
parameters.fcEncoder.Weights = initializeGlorot(sz,numOut,numIn);
parameters.fcEncoder.Bias = initializeZeros([latentDimension 1]);
```

### Define Model Encoder Function

Create the function `modelEncoder`, listed in the Encoder Model Function on page 2-0 section of the example, that computes the output of the encoder model. The `modelEncoder` function, takes as input sequences of word indices, the model parameters, and the sequence lengths, and returns the corresponding latent feature vector.

### Prepare Mini-Batch of Data

To train the model using a custom training loop, you must iterate over mini-batches of data and convert it into the format required for the encoder model and the model gradients functions. This

section of the example illustrates the steps needed for preparing a mini-batch of data inside the custom training loop.

Prepare an example mini-batch of data. Select a mini-batch of 32 documents from `documents`. This represents the mini-batch of data used in an iteration of a custom training loop.

```
miniBatchSize = 32;
idx = 1:miniBatchSize;
documentsBatch = documents(idx);
```

Convert the documents to sequences using the `doc2sequence` function and specify to right-pad the sequences with the word index corresponding to the padding token.

```
X = doc2sequence(enc,documentsBatch, ...
    PaddingDirection="right", ...
    PaddingValue=paddingIdx);
```

The output of the `doc2sequence` function is a cell array, where each element is a row vector of word indices. Because the encoder model function requires numeric input, concatenate the rows of the data using the `cat` function and specify to concatenate along the first dimension. The output has size `miniBatchSize-by-sequenceLength`, where `sequenceLength` is the length of the longest sequence in the mini-batch.

```
X = cat(1,X{:});
size(X)
```

```
ans = 1×2
    32    14
```

Convert the data to a `darray` with format "BTC" (batch, time, channel). The software automatically rearranges the output to have format "CTB" so the output has size `1-by-miniBatchSize-by-sequenceLength`.

```
X = darray(X, 'BTC');
size(X)
```

```
ans = 1×3
    1    32    14
```

For masking, calculate the unpadded sequence lengths of the input data using the `doclength` function with the mini-batch of documents as input.

```
sequenceLengths = doclength(documentsBatch);
```

This code snippet shows an example of preparing a mini-batch in a custom training loop.

```
iteration = 0;

% Loop over epochs.
for epoch = 1:numEpochs

    % Loop over mini-batches.
    for i = 1:numIterationsPerEpoch
```

```
iteration = iteration + 1;

% Read mini-batch.
idx = (i-1)*miniBatchSize+1:i*miniBatchSize;
documentsBatch = documents(idx);

% Convert to sequences.
X = doc2sequence(enc,documentsBatch, ...
    PaddingDirection="right", ...
    PaddingValue=paddingIdx);
X = cat(1,X{:});

% Convert to dlarray.
X = dlarray(X,"BTC");

% Calculate sequence lengths.
sequenceLengths = doclength(documentsBatch);

% Evaluate model gradients.
% ...

% Update learnable parameters.
% ...
end
end
```

### Use Model Function in Model Loss Function

When training a deep learning model with a custom training loop, you must calculate the loss and the gradients of the loss with respect to the learnable parameters. This calculation depends on the output of a forward pass of the model function.

To perform a forward pass of the encoder, use the `modelEncoder` function directly with the parameters, data, and sequence lengths as input. The output is a `latentDimension-by-miniBatchSize` matrix.

```
Z = modelEncoder(parameters,X,sequenceLengths);
size(Z)
```

```
ans = 1×2
    50    32
```

This code snippet shows an example of using a model encoder function inside the model gradients function.

```
function [loss,gradients] = modelLoss(parameters,X,sequenceLengths)

    Z = modelEncoder(parameters,X,sequenceLengths);

    % Calculate loss.
    % ...

    % Calculate gradients.
    % ...

end
```



This code snippet shows an example of evaluating the model gradients in a custom training loop.

```

iteration = 0;

% Loop over epochs.
for epoch = 1:numEpochs

    % Loop over mini-batches.
    for i = 1:numIterationsPerEpoch
        iteration = iteration + 1;

        % Prepare mini-batch.
        % ...

        % Evaluate model gradients.
        [loss,gradients] = dlfeval(@modelLoss, parameters, X, sequenceLengths);

        % Update learnable parameters.
        [parameters,trailingAvg,trailingAvgSq] = adamupdate(parameters,gradients, ...
            trailingAvg,trailingAvgSq,iteration);
    end
end

```

### Encoder Model Function

The `modelEncoder` function, takes as input the model parameters, sequences of word indices, and the sequence lengths, and returns the corresponding latent feature vector.

Because the input data contains padded sequences of different lengths, the padding can have adverse effects on loss calculations. For the LSTM operation, instead of returning the output of the last time step of the sequence (which likely corresponds to the LSTM state after processing lots of padding values), determine the actual last time step given by the `sequenceLengths` input.

```

function Z = modelEncoder(parameters,X,sequenceLengths)

% Embedding.
weights = parameters.emb.Weights;
Z = embed(X,weights);

% LSTM.
inputWeights = parameters.lstmEncoder.InputWeights;
recurrentWeights = parameters.lstmEncoder.RecurrentWeights;
bias = parameters.lstmEncoder.Bias;

numHiddenUnits = size(recurrentWeights,2);
hiddenState = zeros(numHiddenUnits,1,"like",X);
cellState = zeros(numHiddenUnits,1,"like",X);

Z1 = lstm(Z,hiddenState,cellState,inputWeights,recurrentWeights,bias);

% Output mode "last" with masking.
miniBatchSize = size(Z1,2);
Z = zeros(numHiddenUnits,miniBatchSize,"like",Z1);
Z = dlarray(Z,"CB");

for n = 1:miniBatchSize
    t = sequenceLengths(n);
    Z(:,n) = Z1(:,n,t);
end

```

```
end

% Fully connect.
weights = parameters.fcEncoder.Weights;
bias = parameters.fcEncoder.Bias;
Z = fullyconnect(Z,weights,bias);

end
```

### **Preprocessing Function**

The function `preprocessText` performs these steps:

- 1 Prepends and appends each input string with the specified start and stop tokens, respectively.
- 2 Tokenize the text using `tokenizedDocument`.

```
function documents = preprocessText(textData,startToken,stopToken)

% Add start and stop tokens.
textData = startToken + textData + stopToken;

% Tokenize the text.
documents = tokenizedDocument(textData,'CustomTokens',[startToken stopToken]);

end
```

### **See Also**

`wordEncoding` | `word2ind` | `doc2sequence` | `tokenizedDocument`

### **More About**

- “Classify Text Data Using Deep Learning” on page 2-85
- “Classify Text Data Using Convolutional Neural Network” on page 2-93
- “Sequence-to-Sequence Translation Using Attention” on page 2-130
- “Generate Text Using Autoencoders” on page 2-186
- “Define Text Decoder Model Function” on page 2-205
- “Define Custom Training Loops, Loss Functions, and Networks” (Deep Learning Toolbox)

## Define Text Decoder Model Function

This example shows how to define a text decoder model function.

In the context of deep learning, a decoder is the part of a deep learning network that maps a latent vector to some sample space. You can use decode the vectors for various tasks. For example,

- Text generation by initializing a recurrent network with the encoded vector.
- Sequence-to-sequence translation by using the encoded vector as a context vector.
- Image captioning by using the encoded vector as a context vector.

### Load Data

Load the encoded data from `sonnetsEncoded.mat`. This MAT file contains the word encoding, a mini-batch of sequences  $X$ , and the corresponding encoded data  $Z$  output by the encoder used in the example "Define Text Encoder Model Function" (Deep Learning Toolbox).

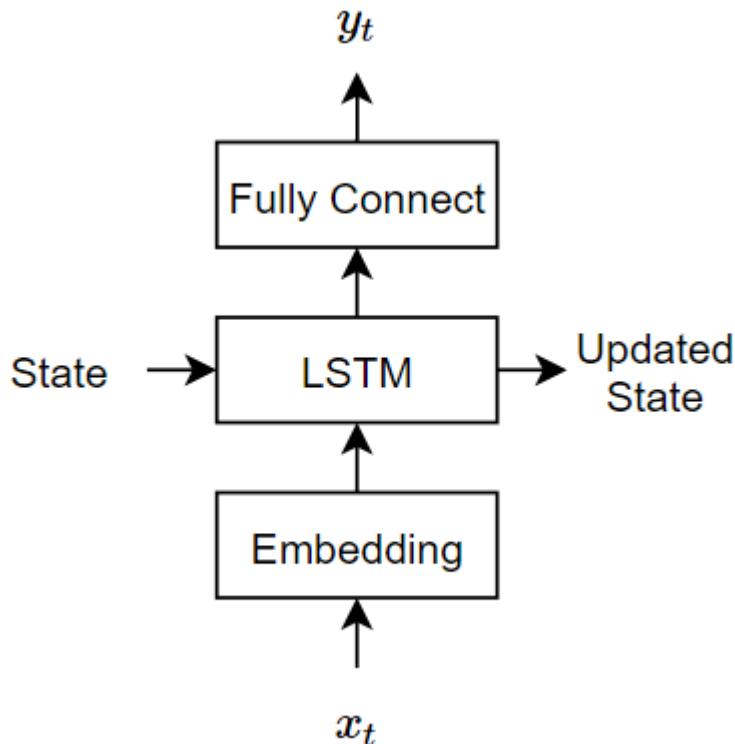
```
s = load("sonnetsEncoded.mat");
enc = s.enc;
X = s.X;
Z = s.Z;
```

```
[latentDimension,miniBatchSize] = size(Z,1:2);
```

### Initialize Model Parameters

The goal of the decoder is to generate sequences given some initial input data and network state.

Initialize the parameters for the following model.



The decoder reconstructs the input using an LSTM initialized the encoder output. For each time step, the decoder predicts the next time step and uses the output for the next time-step predictions. Both the encoder and the decoder use the same embedding.

This model uses three operations:

- The embedding maps word indices in the range 1 though `vocabularySize` to vectors of dimension `embeddingDimension`, where `vocabularySize` is the number of words in the encoding vocabulary and `embeddingDimension` is the number of components learned by the embedding.
- The LSTM operation takes as input a single word vector and outputs 1-by-`numHiddenUnits` vector, where `numHiddenUnits` is the number of hidden units in the LSTM operation. The initial state of the LSTM network (the state at the first time-step) is the encoded vector, so the number of hidden units must match the latent dimension of the encoder.
- The fully connected operation multiplies the input by a weight matrix adding bias and outputs vectors of size `vocabularySize`.

Specify the dimensions of the parameters. The embedding sizes must match the encoder.

```
embeddingDimension = 100;  
vocabularySize = enc.NumWords;  
numHiddenUnits = latentDimension;
```

Create a struct for the parameters.

```
parameters = struct;
```

Initialize the weights of the embedding using the Gaussian using the `initializeGaussian` function which is attached to this example as a supporting file. Specify a mean of 0 and a standard deviation of 0.01. To learn more, see “Gaussian Initialization” (Deep Learning Toolbox).

```
sz = [embeddingDimension vocabularySize];  
mu = 0;  
sigma = 0.01;  
parameters.emb.Weights = initializeGaussian(sz,mu,sigma);
```

Initialize the learnable parameters for the decoder LSTM operation:

- Initialize the input weights with the Glorot initializer using the `initializeGlorot` function which is attached to this example as a supporting file. To learn more, see “Glorot Initialization” (Deep Learning Toolbox).
- Initialize the recurrent weights with the orthogonal initializer using the `initializeOrthogonal` function which is attached to this example as a supporting file. To learn more, see “Orthogonal Initialization” (Deep Learning Toolbox).
- Initialize the bias with the unit forget gate initializer using the `initializeUnitForgetGate` function which is attached to this example as a supporting file. To learn more, see “Unit Forget Gate Initialization” (Deep Learning Toolbox).

The sizes of the learnable parameters depend on the size of the input. Because the inputs to the LSTM operation are sequences of word vectors from the embedding operation, the number of input channels is `embeddingDimension`.

- The input weight matrix has size `4*numHiddenUnits-by-inputSize`, where `inputSize` is the dimension of the input data.

- The recurrent weight matrix has size  $4*\text{numHiddenUnits}$ -by- $\text{numHiddenUnits}$ .
- The bias vector has size  $4*\text{numHiddenUnits}$ -by-1.

```
sz = [4*numHiddenUnits embeddingDimension];
numOut = 4*numHiddenUnits;
numIn = embeddingDimension;
```

```
parameters.lstmDecoder.InputWeights = initializeGlorot(sz,numOut,numIn);
parameters.lstmDecoder.RecurrentWeights = initializeOrthogonal([4*numHiddenUnits numHiddenUnits]);
parameters.lstmDecoder.Bias = initializeUnitForgetGate(numHiddenUnits);
```

Initialize the learnable parameters for the encoder fully connected operation:

- Initialize the weights with the Glorot initializer.
- Initialize the bias with zeros using the `initializeZeros` function which is attached to this example as a supporting file. To learn more, see “Zeros Initialization” (Deep Learning Toolbox).

The sizes of the learnable parameters depend on the size of the input. Because the inputs to the fully connected operation are the outputs of the LSTM operation, the number of input channels is `numHiddenUnits`. To make the fully connected operation output vectors with size `latentDimension`, specify an output size of `latentDimension`.

- The weights matrix has size `outputSize`-by-`inputSize`, where `outputSize` and `inputSize` correspond to the output and input dimensions, respectively.
- The bias vector has size `outputSize`-by-1.

To make the fully connected operation output vectors with size `vocabularySize`, specify an output size of `vocabularySize`.

```
sz = [vocabularySize numHiddenUnits];
mu = 0;
sigma = 1;
parameters.fcDecoder.Weights = initializeGaussian(sz,mu,sigma);
parameters.fcDecoder.Bias = initializeZeros([vocabularySize 1]);
```

### Define Model Decoder Function

Create the function `modelDecoder`, listed in the Decoder Model Function on page 2-0 section of the example, that computes the output of the decoder model. The `modelDecoder` function, takes as input sequences of word indices, the model parameters, and the sequence lengths, and returns the corresponding latent feature vector.

### Use Model Function in Model Loss Function

When training a deep learning model with a custom training loop, you must calculate the loss and gradients of the loss with respect to the learnable parameters. This calculation depends on the output of a forward pass of the model function.

There are two common approaches to generating text data with a decoder:

- 1 Closed loop — For each time step, make predictions using the previous prediction as input.
- 2 Open loop — For each time step, make predictions using inputs from an external source (for example, training targets).

### Closed Loop Generation

Closed loop generation is when the model generates data one time-step at a time and uses the previous prediction as input for the next prediction. Unlike open loop generation, this process does not require any input between predictions and is best suited for scenarios without supervision. For example, a language translation model that generates output text in one go.

Initialize the hidden state of the LSTM network with the encoder output  $Z$ .

```
state = struct;
state.HiddenState = Z;
state.CellState = zeros(size(Z), 'like', Z);
```

For the first time step, use an array of start tokens as input for the decoder. For simplicity, extract an array of start tokens from the first time-step of the training data.

```
decoderInput = X(:, :, 1);
```

Preallocate the decoder output to have size `numClasses-by-miniBatchSize-by-sequenceLength` with the same datatype as `d1X`, where `sequenceLength` is the desired length of the generation, for example, the length of the training targets. For this example, specify a sequence length of 16.

```
sequenceLength = 16;
Y = zeros(vocabularySize, miniBatchSize, sequenceLength, "like", X);
Y = d1array(Y, "CBT");
```

For each time step, predict the next time step of the sequence using the `modelDecoder` function. After each prediction, find the indices corresponding to the maximum values of the decoder output and use these indices as the decoder input for the next time step.

```
for t = 1:sequenceLength
    [Y(:, :, t), state] = modelDecoder(parameters, decoderInput, state);

    [~, idx] = max(Y(:, :, t));
    decoderInput = idx;
end
```

The output is a `vocabularySize-by-miniBatchSize-by-sequenceLength` array.

```
size(Y)
ans = 1×3
      3595      32      16
```

This code snippet shows an example of performing closed loop generation in a model gradients function.

```
function [loss, gradients] = modelLoss(parameters, X, sequenceLengths)

    % Encode input.
    Z = modelEncoder(parameters, X, sequenceLengths);

    % Initialize LSTM state.
    state = struct;
    state.HiddenState = Z;
    state.CellState = zeros(size(Z), "like", Z);
```

```

% Initialize decoder input.
decoderInput = X(:,:,1);

% Closed loop prediction.
sequenceLength = size(X,3);
Y = zeros(numClasses,miniBatchSize,sequenceLength,"like",X);
for t = 1:sequenceLength
    [Y(:,:,t), state] = modelDecoder(parameters,decoderInput,state);

    [~,idx] = max(Y(:,:,t));
    decoderInput = idx;
end

% Calculate loss.
% ...

% Calculate gradients.
% ...

end

```

### Open Loop Generation: Teacher Forcing

When training with closed loop generation, predicting the most likely word for each step in the sequence can lead to suboptimal results. For example, in an image captioning workflow, if the decoder predicts the first word of a caption is "a" when given an image of an elephant, then the probability of predicting "elephant" for the next word becomes much more unlikely because of the extremely low probability of the phrase "a elephant" appearing in English text.

To help the network converge faster, you can use *teacher forcing*: use the target values as input to the decoder instead of the previous predictions. Using teacher forcing helps the network to learn characteristics from the later time steps of the sequences without having to wait for the network to correctly generate the earlier time steps of the sequences.

To perform teacher forcing, use the `modelEncoder` function directly with the target sequence as input.

Initialize the hidden state of the LSTM network with the encoder output `Z`.

```

state = struct;
state.HiddenState = Z;
state.CellState = zeros(size(Z),"like",Z);

```

Make predictions using the target sequence as input.

```

Y = modelDecoder(parameters,X,state);

```

The output is a `vocabularySize-by-miniBatchSize-by-sequenceLength` array, where `sequenceLength` is the length of the input sequences.

```

size(Y)

```

```

ans = 1×3

```

```

    3595

```

```

     32

```

```

     14

```

This code snippet shows an example of performing teacher forcing in a model gradients function.

```
function [loss,gradients] = modelLoss(parameters,X,sequenceLengths)

    % Encode input.
    Z = modelEncoder(parameters,X);

    % Initialize LSTM state.
    state = struct;
    state.HiddenState = Z;
    state.CellState = zeros(size(Z),"like",Z);

    % Teacher forcing.
    Y = modelDecoder(parameters,X,state);

    % Calculate loss.
    % ...

    % Calculate gradients.
    % ...

end
```

### Decoder Model Function

The `modelDecoder` function, takes as input the model parameters, sequences of word indices, and the network state, and returns the decoded sequences.

Because the `lstm` function is *stateful* (when given a time series as input, the function propagates and updates the state between each time step) and that the `embed` and `fullyconnect` functions are time-distributed by default (when given a time series as input, the functions operate on each time step independently), the `modelDecoder` function supports both sequence and single time-step inputs.

```
function [Y,state] = modelDecoder(parameters,X,state)

% Embedding.
weights = parameters.emb.Weights;
X = embed(X,weights);

% LSTM.
inputWeights = parameters.lstmDecoder.InputWeights;
recurrentWeights = parameters.lstmDecoder.RecurrentWeights;
bias = parameters.lstmDecoder.Bias;

hiddenState = state.HiddenState;
cellState = state.CellState;

[Y,hiddenState,cellState] = lstm(X,hiddenState,cellState, ...
    inputWeights,recurrentWeights,bias);

state.HiddenState = hiddenState;
state.CellState = cellState;

% Fully connect.
weights = parameters.fcDecoder.Weights;
bias = parameters.fcDecoder.Bias;
Y = fullyconnect(Y,weights,bias);
```



end

## See Also

[wordEncoding](#) | [word2ind](#) | [doc2sequence](#) | [tokenizedDocument](#)

## More About

- “Classify Text Data Using Deep Learning” on page 2-85
- “Classify Text Data Using Convolutional Neural Network” on page 2-93
- “Sequence-to-Sequence Translation Using Attention” on page 2-130
- “Generate Text Using Autoencoders” on page 2-186
- “Define Text Encoder Model Function” on page 2-198
- “Define Custom Training Loops, Loss Functions, and Networks” (Deep Learning Toolbox)

## Classify Out-of-Memory Text Data Using Custom Mini-Batch Datastore

This example shows how to classify out-of-memory text data with a deep learning network using a custom mini-batch datastore.

A mini-batch datastore is an implementation of a datastore with support for reading data in batches. You can use a mini-batch datastore as a source of training, validation, test, and prediction data sets for deep learning applications. Use mini-batch datastores to read out-of-memory data or to perform specific preprocessing operations when reading batches of data.

When training the network, the software creates mini-batches of sequences of the same length by padding, truncating, or splitting the input data. The `trainingOptions` function provides options to pad and truncate input sequences, however, these options are not well suited for sequences of word vectors. Furthermore, this function does not support padding data in a custom datastore. Instead, you must pad and truncate the sequences manually. If you *left-pad* and truncate the sequences of word vectors, then the training might improve.

The “Classify Text Data Using Deep Learning” on page 2-85 example manually truncates and pads all the documents to the same length. This process adds lots of padding to very short documents and discards lots of data from very long documents.

Alternatively, to prevent adding too much padding or discarding too much data, create a custom mini-batch datastore that inputs mini-batches into the network. The custom mini-batch datastore `textDatastore.m` converts mini-batches of documents to sequences or word indices and left-pads each mini-batch to the length of the longest document in the mini-batch. For sorted data, this datastore can help reduce the amount of padding added to the data since documents are not padded to a fixed length. Similarly, the datastore does not discard any data from the documents.

This example uses the custom mini-batch datastore `textDatastore`, attached to this example as a supporting file. To access this file, open the example as a live script. You can adapt this datastore to your data by customizing the functions. For an example showing how to create your own custom mini-batch datastore, see “Develop Custom Mini-Batch Datastore” (Deep Learning Toolbox).

### Load Pretrained Word Embedding

The datastore `textDatastore` requires a word embedding to convert documents to sequences of vectors. Load a pretrained word embedding using `fastTextWordEmbedding`. This function requires Text Analytics Toolbox™ Model for *fastText English 16 Billion Token Word Embedding* support package. If this support package is not installed, then the function provides a download link.

```
emb = fastTextWordEmbedding;
```

### Create Mini-Batch Datastore of Documents

Create a datastore that contains the data for training. The custom mini-batch datastore `textDatastore` reads predictors and labels from a CSV file. For the predictors, the datastore converts the documents into sequences of word indices and for the responses, the datastore returns a categorical label for each document. For more information about creating custom mini-batch datastores, see “Develop Custom Mini-Batch Datastore” (Deep Learning Toolbox).

For the training data, specify the CSV file “`factoryReports.csv`” and that the text and labels are in the columns “`Description`” and “`Category`” respectively.

```
filenameTrain = "factoryReports.csv";
textName = "Description";
labelName = "Category";
dsTrain = textDatastore(filenameTrain, textName, labelName, emb)
```

```
dsTrain =
```

```
textDatastore with properties:
```

```
ClassNames: ["Electronic Failure" "Leak" "Mechanical Failure" "Software Fai
Datastore: [1x1 matlab.io.datastore.TransformedDatastore]
EmbeddingDimension: 300
LabelName: "Category"
MiniBatchSize: 128
NumClasses: 4
NumObservations: 480
```

### Create and Train LSTM Network

Define the LSTM network architecture. To input sequence data into the network, include a sequence input layer and set the input size to the embedding dimension. Next, include an LSTM layer with 180 hidden units. To use the LSTM layer for a sequence-to-label classification problem, set the output mode to 'last'. Finally, add a fully connected layer with output size equal to the number of classes, a softmax layer, and a classification layer.

```
numFeatures = dsTrain.EmbeddingDimension;
numHiddenUnits = 180;
numClasses = dsTrain.NumClasses;
```

```
layers = [ ...
    sequenceInputLayer(numFeatures)
    lstmLayer(numHiddenUnits, 'OutputMode', 'last')
    fullyConnectedLayer(numClasses)
    softmaxLayer
    classificationLayer];
```

Specify the training options. Specify the solver to be 'adam' and the gradient threshold to be 2. The datastore `textDatastore.m` does not support shuffling, so set 'Shuffle', to 'never'. For an example showing how to implement a datastore with support for shuffling, see “Develop Custom Mini-Batch Datastore” (Deep Learning Toolbox). To monitor the training progress, set the 'Plots' option to 'training-progress'. To suppress verbose output, set 'Verbose' to false.

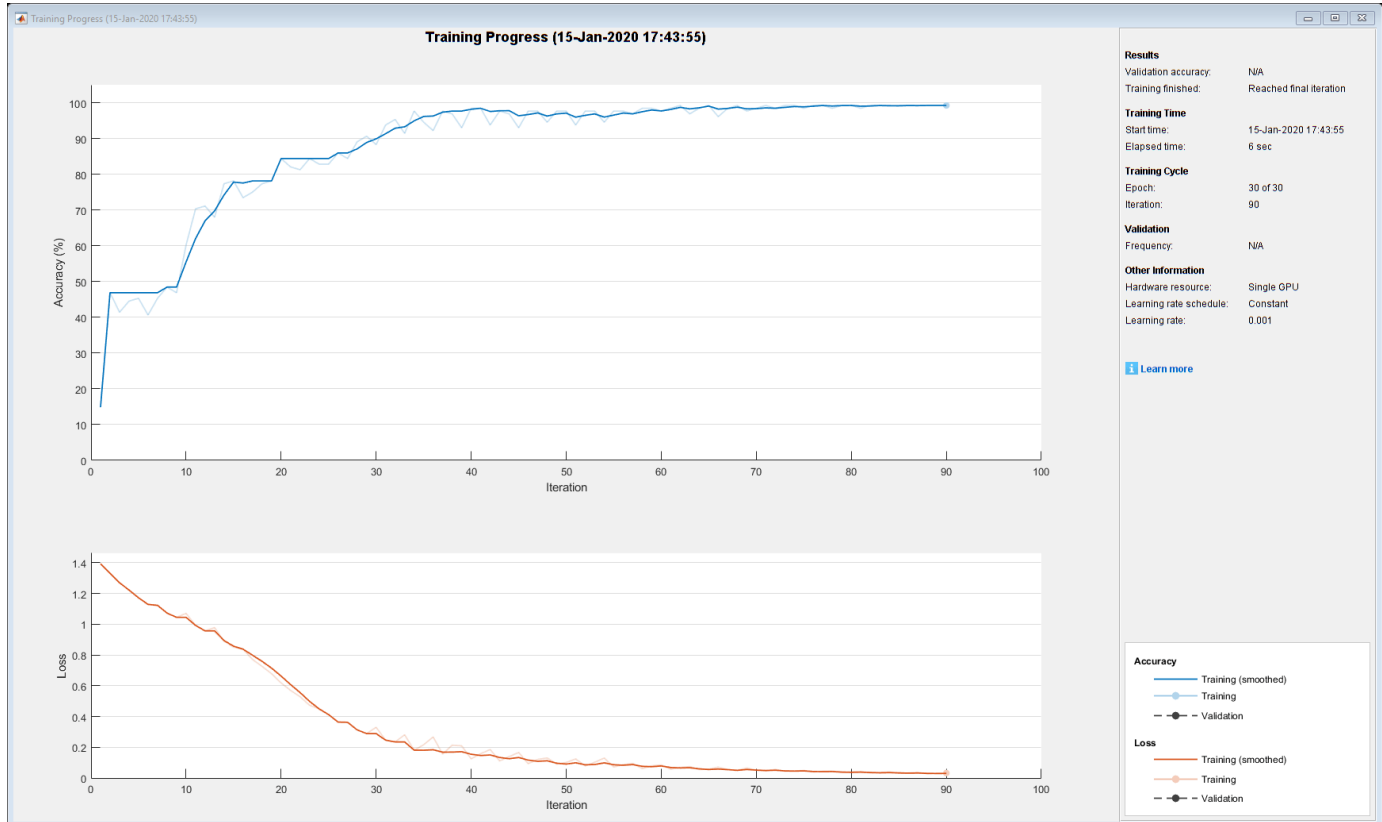
By default, `trainNetwork` uses a GPU if one is available. To specify the execution environment manually, use the 'ExecutionEnvironment' name-value pair argument of `trainingOptions`. Training on a CPU can take significantly longer than training on a GPU. Training using a GPU requires Parallel Computing Toolbox™ and a supported GPU device. For information on supported devices, see “GPU Support by Release” (Parallel Computing Toolbox).

```
miniBatchSize = 128;
numObservations = dsTrain.NumObservations;
numIterationsPerEpoch = floor(numObservations / miniBatchSize);
```

```
options = trainingOptions('adam', ...
    'MiniBatchSize', miniBatchSize, ...
    'GradientThreshold', 2, ...
    'Shuffle', 'never', ...
    'Plots', 'training-progress', ...
    'Verbose', false);
```

Train the LSTM network using the `trainNetwork` function.

```
net = trainNetwork(dsTrain, layers, options);
```



### Predict Using New Data

Classify the event type of three new reports. Create a string array containing the new reports.

```
reportsNew = [  
    "Coolant is pooling underneath sorter."  
    "Sorter blows fuses at start up."  
    "There are some very loud rattling sounds coming from the assembler."];
```

Preprocess the text data using the preprocessing steps as the datastore `textDatastore`.

```
documents = tokenizedDocument(reportsNew);  
documents = lower(documents);  
documents = erasePunctuation(documents);  
predictors = doc2sequence(emb, documents);
```

Classify the new sequences using the trained LSTM network.

```
labelsNew = classify(net, predictors)
```

```
labelsNew = 3x1 categorical  
    Leak  
    Electronic Failure
```

Mechanical Failure

## See Also

[wordEmbeddingLayer](#) | [doc2sequence](#) | [tokenizedDocument](#) | [lstmLayer](#) | [trainNetwork](#) | [trainingOptions](#) | [sequenceInputLayer](#) | [wordcloud](#) | [extractHTMLText](#) | [findElement](#) | [htmlTree](#)

## Related Examples

- “Generate Text Using Deep Learning” (Deep Learning Toolbox)
- “Create Simple Text Model for Classification” on page 2-2
- “Analyze Text Data Using Topic Models” on page 2-13
- “Analyze Text Data Using Multiword Phrases” on page 2-7
- “Train a Sentiment Classifier” on page 2-71
- “Sequence Classification Using Deep Learning” (Deep Learning Toolbox)
- “Deep Learning in MATLAB” (Deep Learning Toolbox)



# Display and Presentation

---

- “Visualize Text Data Using Word Clouds” on page 3-2
- “Visualize Word Embeddings Using Text Scatter Plots” on page 3-8

## Visualize Text Data Using Word Clouds

This example shows how to visualize text data using word clouds.

Text Analytics Toolbox extends the functionality of the `wordcloud` (MATLAB) function. It adds support for creating word clouds directly from string arrays and creating word clouds from bag-of-words models and LDA topics.

Load the example data. The file `factoryReports.csv` contains factory reports, including a text description and categorical labels for each event.

```
filename = "factoryReports.csv";  
tbl = readtable(filename, 'TextType', 'string');
```

Extract the text data from the `Description` column.

```
textData = tbl.Description;  
textData(1:10)
```

```
ans = 10x1 string  
"Items are occasionally getting stuck in the scanner spools."  
"Loud rattling and banging sounds are coming from assembler pistons."  
"There are cuts to the power when starting the plant."  
"Fried capacitors in the assembler."  
"Mixer tripped the fuses."  
"Burst pipe in the constructing agent is spraying coolant."  
"A fuse is blown in the mixer."  
"Things continue to tumble off of the belt."  
"Falling items from the conveyor belt."  
"The scanner reel is split, it will soon begin to curve."
```

Create a word cloud from the reports.

```
figure  
wordcloud(textData);  
title("Factory Reports")
```



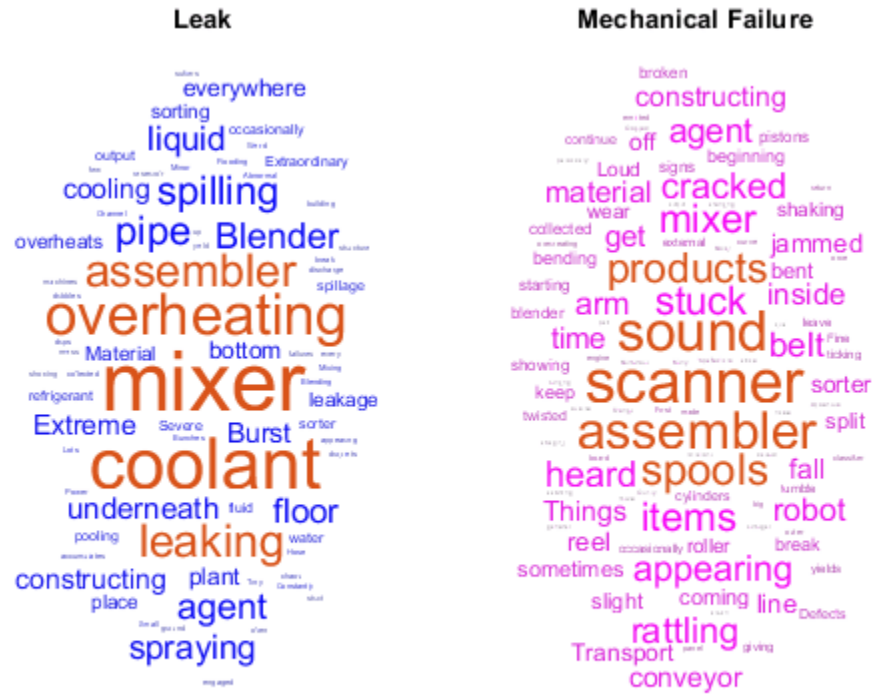


Compare the words in the reports with labels "Leak" and "Mechanical Failure". Create word clouds of the reports for each of these labels. Specify the word colors to be blue and magenta for each word cloud respectively.

```
figure
labels = tbl.Category;

subplot(1,2,1)
idx = labels == "Leak";
wordcloud(textData(idx), 'Color', 'blue');
title("Leak")

subplot(1,2,2)
idx = labels == "Mechanical Failure";
wordcloud(textData(idx), 'Color', 'magenta');
title("Mechanical Failure")
```



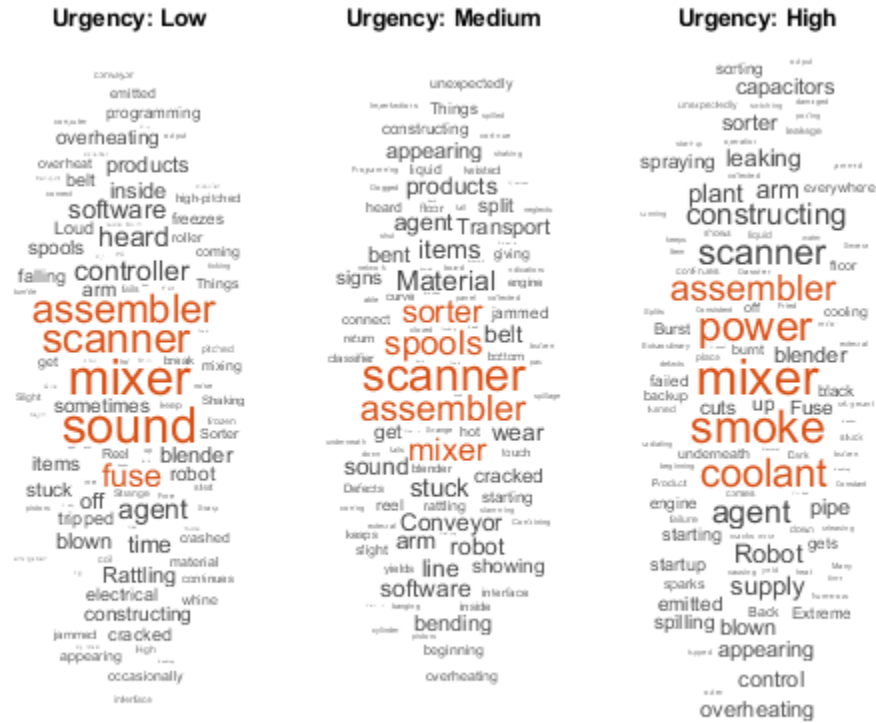
Compare the words in the reports with urgency "Low", "Medium", and "High".

```
figure
urgency = tbl.Urgency;

subplot(1,3,1)
idx = urgency == "Low";
wordcloud(textData(idx));
title("Urgency: Low")

subplot(1,3,2)
idx = urgency == "Medium";
wordcloud(textData(idx));
title("Urgency: Medium")

subplot(1,3,3)
idx = urgency == "High";
wordcloud(textData(idx));
title("Urgency: High")
```



Compare the words in the reports with cost reported in hundreds of dollars to the reports with costs reported in thousands of dollars. Create word clouds of the reports for each of these amounts with highlight color blue and red respectively.

```
cost = tbl.Cost;
idx = cost > 100;
figure
wordcloud(textData(idx), 'HighlightColor', 'blue');
title("Cost > $100")
```





## Visualize Word Embeddings Using Text Scatter Plots

This example shows how to visualize word embeddings using 2-D and 3-D t-SNE and text scatter plots.

Word embeddings map words in a vocabulary to real vectors. The vectors attempt to capture the semantics of the words, so that similar words have similar vectors. Some embeddings also capture relationships between words like "Italy is to France as Rome is to Paris". In vector form, this relationship is  $Italy - Rome + Paris = France$ .

### Load Pretrained Word Embedding

Load a pretrained word embedding using `fastTextWordEmbedding`. This function requires Text Analytics Toolbox™ Model for fastText English 16 Billion Token Word Embedding support package. If this support package is not installed, then the function provides a download link.

```
emb = fastTextWordEmbedding

emb =
  wordEmbedding with properties:

    Dimension: 300
    Vocabulary: [1×999994 string]
```

Explore the word embedding using `word2vec` and `vec2word`. Convert the words *Italy*, *Rome*, and *Paris* to vectors using `word2vec`.

```
italy = word2vec(emb, "Italy");
rome = word2vec(emb, "Rome");
paris = word2vec(emb, "Paris");
```

Compute the vector given by  $italy - rome + paris$ . This vector encapsulates the semantic meaning of the word *Italy*, without the semantics of the word *Rome*, and also includes the semantics of the word *Paris*.

```
vec = italy - rome + paris

vec = 1×300 single row vector

    0.1606    -0.0690     0.1183    -0.0349     0.0672     0.0907    -0.1820    -0.0080     0.0320    -0.0
```

Find the closest words in the embedding to `vec` using `vec2word`.

```
word = vec2word(emb, vec)

word =
  "France"
```

### Create 2-D Text Scatter Plot

Visualize the word embedding by creating a 2-D text scatter plot using `tsne` and `textscatter`.

Convert the first 5000 words to vectors using `word2vec`. `V` is a matrix of word vectors of length 300.

```
words = emb.Vocabulary(1:5000);
V = word2vec(emb,words);
size(V)
```

```
ans = 1×2
```

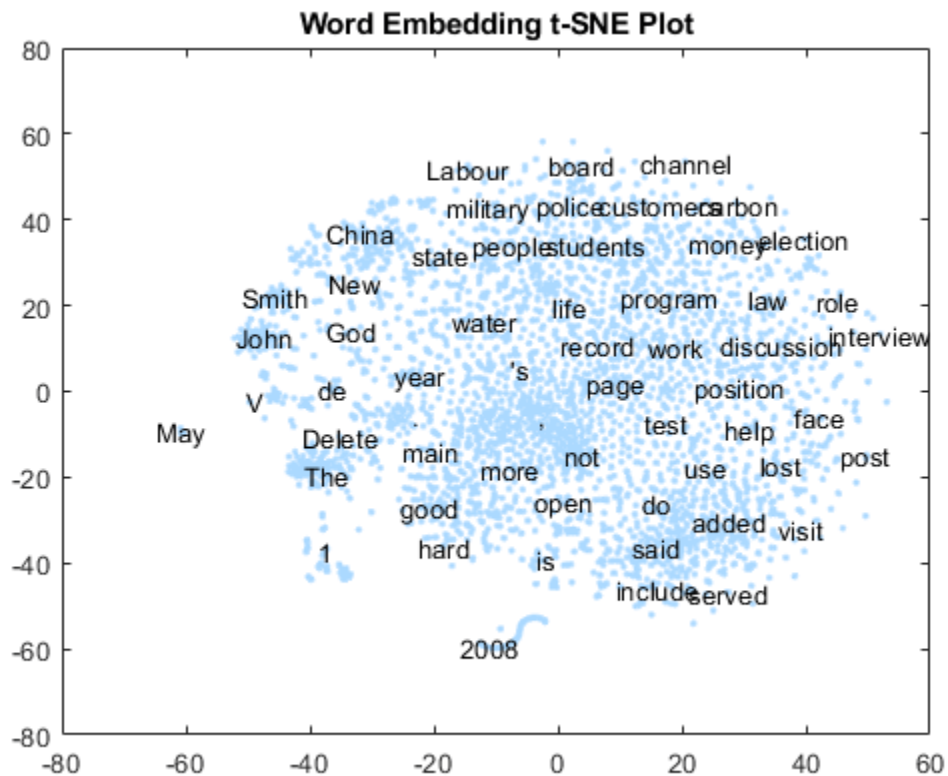
```
5000    300
```

Embed the word vectors in two-dimensional space using `tsne`. This function may take a few minutes to run. If you want to display the convergence information, then set the `'Verbose'` name-value pair to 1.

```
XY = tsne(V);
```

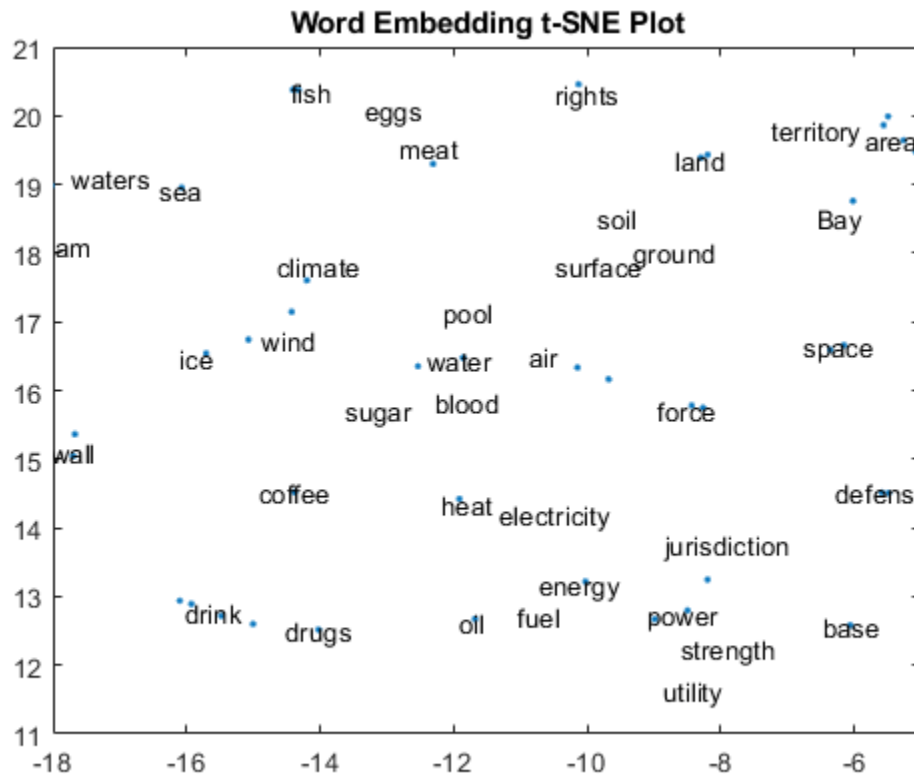
Plot the words at the coordinates specified by `XY` in a 2-D text scatter plot. For readability, `textscatter`, by default, does not display all of the input words and displays markers instead.

```
figure
textscatter(XY,words)
title("Word Embedding t-SNE Plot")
```



Zoom in on a section of the plot.

```
xlim([-18 -5])
ylim([11 21])
```



### Create 3-D Text Scatter Plot

Visualize the word embedding by creating a 3-D text scatter plot using `tsne` and `textscatter`.

Convert the first 5000 words to vectors using `word2vec`. `V` is a matrix of word vectors of length 300.

```
words = emb.Vocabulary(1:5000);
V = word2vec(emb,words);
size(V)
```

```
ans = 1x2
```

```
5000    300
```

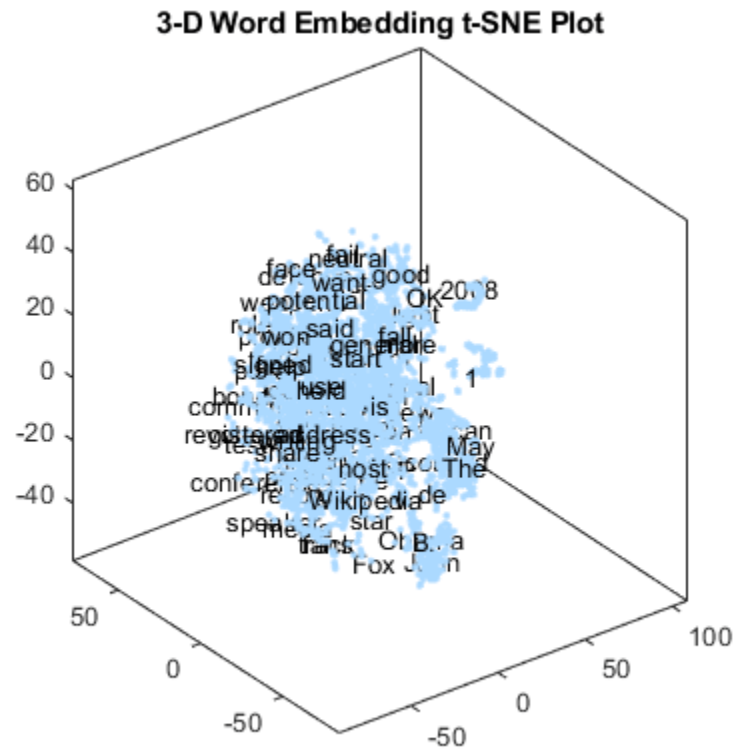
Embed the word vectors in a three-dimensional space using `tsne` by specifying the number of dimensions to be three. This function may take a few minutes to run. If you want to display the convergence information, then you can set the `'Verbose'` name-value pair to 1.

```
XYZ = tsne(V,'NumDimensions',3);
```

Plot the words at the coordinates specified by `XYZ` in a 3-D text scatter plot.

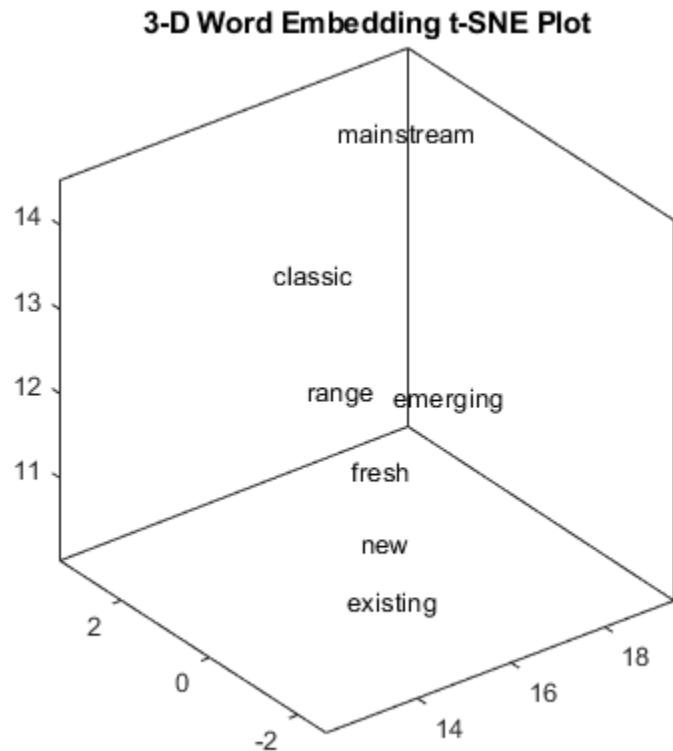
```
figure
ts = textscatter3(XYZ,words);
title("3-D Word Embedding t-SNE Plot")
```





Zoom in on a section of the plot.

```
xlim([12.04 19.48])  
ylim([-2.66 3.40])  
zlim([10.03 14.53])
```



### Perform Cluster Analysis

Convert the first 5000 words to vectors using `word2vec`. `V` is a matrix of word vectors of length 300.

```
words = emb.Vocabulary(1:5000);
V = word2vec(emb,words);
size(V)
```

```
ans = 1×2
```

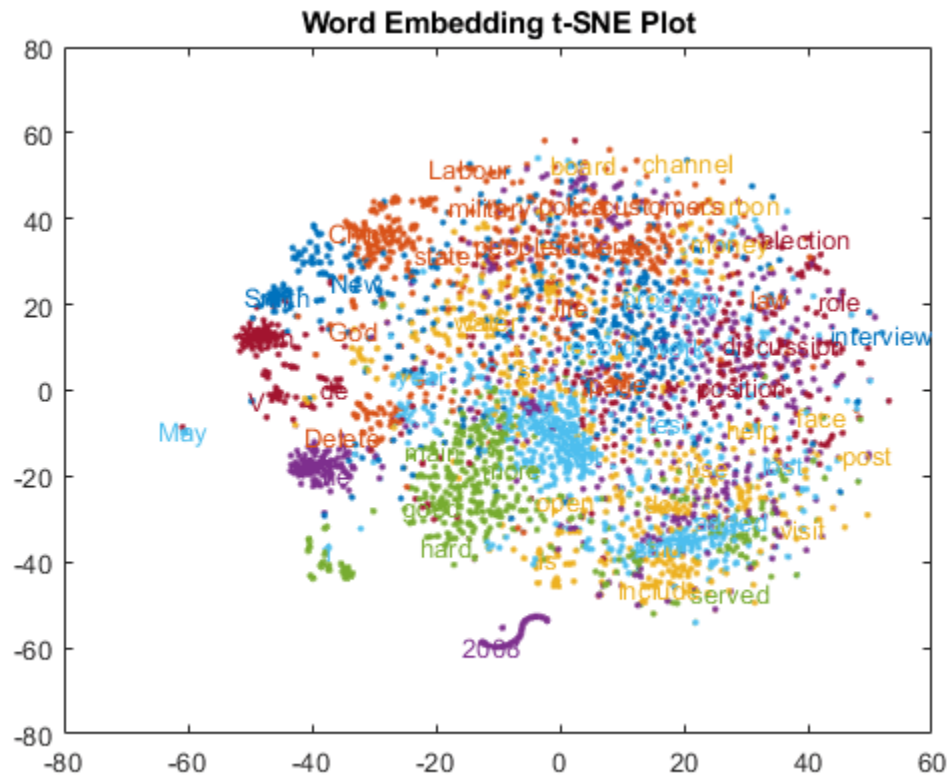
```
5000    300
```

Discover 25 clusters using `kmeans`.

```
cidx = kmeans(V,25,'dist','sqeuclidean');
```

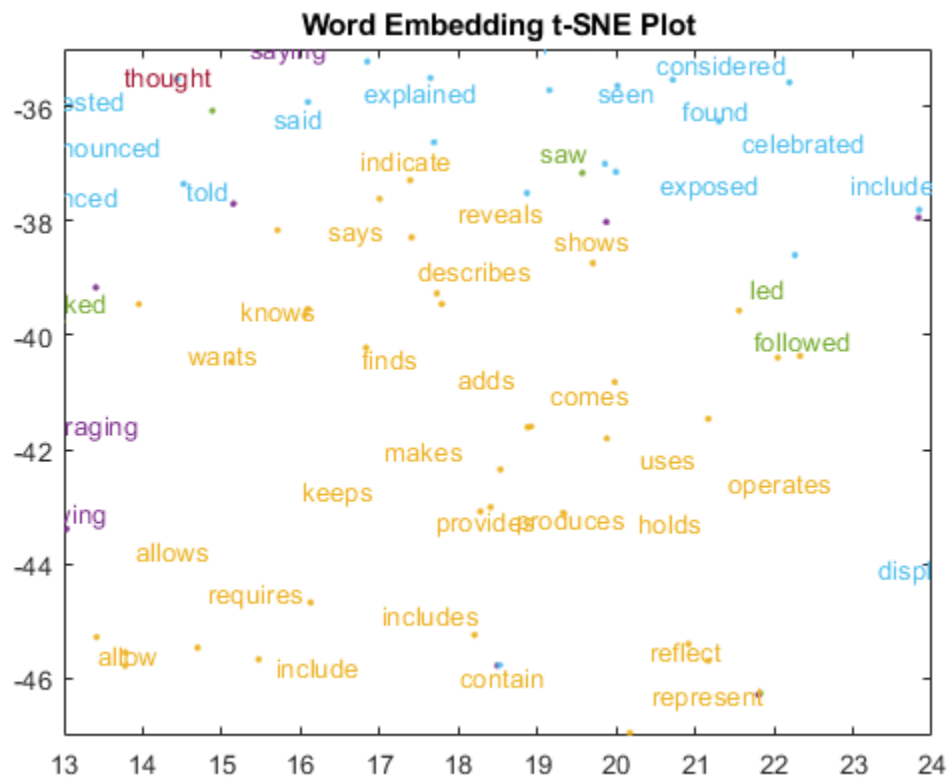
Visualize the clusters in a text scatter plot using the 2-D t-SNE data coordinates calculated earlier.

```
figure
textscatter(XY,words,'ColorData',categorical(cidx));
title("Word Embedding t-SNE Plot")
```



Zoom in on a section of the plot.

```
xlim([13 24])  
ylim([-47 -35])
```



### See Also

`readWordEmbedding` | `textscatter` | `textscatter3` | `word2vec` | `vec2word` | `wordEmbedding` | `tokenizedDocument`

### Related Examples

- “Extract Text Data from Files” on page 1-2
- “Prepare Text Data for Analysis” on page 1-11
- “Visualize Text Data Using Word Clouds” on page 3-2
- “Classify Text Data Using Deep Learning” on page 2-85

# Language Support

---

- “Language Considerations” on page 4-2
- “Japanese Language Support” on page 4-6
- “Analyze Japanese Text Data” on page 4-11
- “German Language Support” on page 4-21
- “Analyze German Text Data” on page 4-26
- “Korean Language Support” on page 4-37
- “Language-Independent Features” on page 4-39

## Language Considerations

Text Analytics Toolbox supports the languages English, Japanese, German, and Korean. Most Text Analytics Toolbox functions also work with text in other languages. This table summarizes how to use Text Analytics Toolbox features for other languages.

Feature	Language Consideration	Workaround
Tokenization	The <code>tokenizedDocument</code> function has built-in rules for English, Japanese, German, and Korean only. For English and German text, the 'unicode' tokenization method of <code>tokenizedDocument</code> detects tokens using rules based on Unicode® Standard Annex #29 [1] and the ICU tokenizer [2], modified to better detect complex tokens such as hashtags and URLs. For Japanese and Korean text, the 'mecab' tokenization method detects tokens using rules based on the MeCab tokenizer [3].	For other languages, you can still try using <code>tokenizedDocument</code> . If <code>tokenizedDocument</code> does not produce useful results, then try tokenizing the text manually. To create a <code>tokenizedDocument</code> array from manually tokenized text, set the 'TokenizeMethod' option to 'none'.  For more information, see <code>tokenizedDocument</code> .
Stop word removal	The <code>stopWords</code> and <code>removeStopWords</code> functions support English, Japanese, German, and Korean stop words only.	To remove stop words from other languages, use <code>removeWords</code> and specify your own stop words to remove.
Sentence detection	The <code>addSentenceDetails</code> function detects sentence boundaries based on punctuation characters and line number information. For English and German text, the function also uses a list of abbreviations passed to the function.	For other languages, you might need to specify your own list of abbreviations for sentence detection. To do this, use the 'Abbreviations' option of <code>addSentenceDetails</code> .  For more information, see <code>addSentenceDetails</code> .

Feature	Language Consideration	Workaround
Word clouds	For string input, the <code>wordcloud</code> and <code>wordCloudCounts</code> functions use English, Japanese, German, and Korean tokenization, stop word removal, and word normalization.	<p>For other languages, you might need to manually preprocess your text data and specify unique words and corresponding sizes in <code>wordcloud</code>.</p> <p>To specify word sizes in <code>wordcloud</code>, input your data as a table or arrays containing the unique words and corresponding sizes.</p> <p>For more information, see <code>wordcloud</code>.</p>
Word embeddings	File input to the <code>trainWordEmbedding</code> function requires words separated by whitespace.	<p>For files containing non-English text, you might need to input a <code>tokenizedDocument</code> array to <code>trainWordEmbedding</code>.</p> <p>To create a <code>tokenizedDocument</code> array from pretokenized text, use the <code>tokenizedDocument</code> function and set the <code>'TokenizeMethod'</code> option to <code>'none'</code>.</p> <p>For more information, see <code>trainWordEmbedding</code>.</p>
Keyword extraction	The <code>rakeKeywords</code> function supports English, Japanese, German, and Korean text only.	<p>The <code>rakeKeywords</code> function extracts keywords using a delimiter-based approach to identify candidate keywords. The function, by default, uses punctuation characters and the stop words given by the <code>stopWords</code> with language given by the language details of the input documents as delimiters.</p> <p>For other languages, specify an appropriate set of delimiters using the <code>'Delimiters'</code> and <code>'MergingDelimiters'</code> options.</p> <p>For more information, see <code>rakeKeywords</code>.</p>

Feature	Language Consideration	Workaround
	<p>The <code>textRankKeywords</code> function supports English, Japanese, German, and Korean text only.</p>	<p>The <code>textRankKeywords</code> function extracts keywords by identifying candidate keywords based on their part-of-speech tag. The function uses part-of-speech tags given by the <code>addPartOfSpeechDetails</code> function which supports English, Japanese, German, and Korean text only.</p> <p>For other languages, try using the <code>rakeKeywords</code> instead and specify an appropriate set of delimiters using the <code>'Delimiters'</code> and <code>'MergingDelimiters'</code> options.</p> <p>For more information, see <code>textRankKeywords</code>.</p>

## Language-Independent Features

### Word and N-Gram Counting

The `bagOfWords` and `bagOfNgrams` functions support `tokenizedDocument` input regardless of language. If you have a `tokenizedDocument` array containing your data, then you can use these functions.

### Modeling and Prediction

The `fitlda` and `fitlsa` functions support `bagOfWords` and `bagOfNgrams` input regardless of language. If you have a `bagOfWords` or `bagOfNgrams` object containing your data, then you can use these functions.

The `trainWordEmbedding` function supports `tokenizedDocument` or file input regardless of language. If you have a `tokenizedDocument` array or a file containing your data in the correct format, then you can use this function.

## References

- [1] *Unicode Text Segmentation*. <https://www.unicode.org/reports/tr29/>
- [2] *Boundary Analysis*. <https://unicode-org.github.io/icu/userguide/boundaryanalysis/>
- [3] *MeCab: Yet Another Part-of-Speech and Morphological Analyzer*. <https://taku910.github.io/mecab/>



**See Also**

`stopWords` | `removeWords` | `normalizeWords` | `bagOfWords` | `bagOfNgrams` |  
`tokenizedDocument` | `fitlda` | `fitlsa` | `wordcloud` | `addSentenceDetails` |  
`addLanguageDetails`

**More About**

- “Text Data Preparation”
- “Modeling and Prediction”
- “Display and Presentation”
- “Japanese Language Support” on page 4-6
- “Analyze Japanese Text Data” on page 4-11
- “German Language Support” on page 4-21
- “Analyze German Text Data” on page 4-26

## Japanese Language Support

This topic summarizes the Text Analytics Toolbox features that support Japanese text. For an example showing how to analyze Japanese text data, see “Analyze Japanese Text Data” on page 4-11.

### Tokenization

The `tokenizedDocument` function automatically detects Japanese input. Alternatively, set the 'Language' option in `tokenizedDocument` to 'ja'. This option specifies the language details of the tokens. To view the language details of the tokens, use `tokenDetails`. These language details determine the behavior of the `removeStopWords`, `addPartOfSpeechDetails`, `normalizeWords`, `addSentenceDetails`, and `addEntityDetails` functions on the tokens.

To specify additional MeCab options for tokenization, create a `mecabOptions` object. To tokenize using the specified MeCab tokenization options, use the 'TokenizeMethod' option of `tokenizedDocument`.

### Tokenize Japanese Text

Tokenize Japanese text using `tokenizedDocument`. The function automatically detects Japanese text.

```
str = [
    "恋に悩み、苦しむ。"
    "恋の悩みで苦しむ。"
    "空に星が輝き、瞬いている。"
    "空の星が輝きを増している。"];
documents = tokenizedDocument(str)

documents =
    4x1 tokenizedDocument:

    6 tokens: 恋 に 悩み 、 苦しむ 。
    6 tokens: 恋 の 悩み で 苦しむ 。
    10 tokens: 空 に 星 が 輝き 、 瞬い て い る 。
    10 tokens: 空 の 星 が 輝き を 増し て い る 。
```

### Part of Speech Details

The `tokenDetails` function, by default, includes part of speech details with the token details.

### Get Part of Speech Details of Japanese Text

Tokenize Japanese text using `tokenizedDocument`.

```
str = [
    "恋に悩み、苦しむ。"
    "恋の悩みで 苦しむ。"
    "空に星が輝き、瞬いている。"
    "空の星が輝きを増している。"
    "駅までは遠くて、歩けない。"
    "遠くの駅まで歩けない。"]
```

```
"すもももももものうち。"];
documents = tokenizedDocument(str);
```

For Japanese text, you can get the part-of-speech details using `tokenDetails`. For English text, you must first use `addPartOfSpeechDetails`.

```
tdetails = tokenDetails(documents);
head(tdetails)
```

ans=8×8 table

Token	DocumentNumber	LineNumber	Type	Language	PartOfSpeech	Lemma
"恋"	1	1	letters	ja	noun	"恋"
"に"	1	1	letters	ja	adposition	"に"
"悩み"	1	1	letters	ja	verb	"悩む"
"、"	1	1	punctuation	ja	punctuation	"、"
"苦しむ"	1	1	letters	ja	verb	"苦しむ"
"。"	1	1	punctuation	ja	punctuation	"。"
"恋"	2	1	letters	ja	noun	"恋"
"の"	2	1	letters	ja	adposition	"の"

## Named Entity Recognition

The `tokenDetails` function, by default, includes entity details with the token details.

### Add Named Entity Tags to Japanese Text

Tokenize Japanese text using `tokenizedDocument`.

```
str = [
  "マリーさんはボストンからニューヨークに引っ越しました。"
  "駅で鈴木さんに迎えに行きます。"
  "東京は大阪より大きいですか？"
  "東京に行った時、新宿や渋谷などいろいろな所を訪れました。"];
documents = tokenizedDocument(str);
```

For Japanese text, the software automatically adds named entity tags, so you do not need to use the `addEntityDetails` function. This software detects person names, locations, organizations, and other named entities. To view the entity details, use the `tokenDetails` function.

```
tdetails = tokenDetails(documents);
head(tdetails)
```

ans=8×8 table

Token	DocumentNumber	LineNumber	Type	Language	PartOfSpeech	Lemma
"マリー"	1	1	letters	ja	proper-noun	"マリー"
"さん"	1	1	letters	ja	noun	"さん"
"は"	1	1	letters	ja	adposition	"は"
"ボストン"	1	1	letters	ja	proper-noun	"ボストン"
"から"	1	1	letters	ja	adposition	"から"
"ニューヨーク"	1	1	letters	ja	proper-noun	"ニューヨーク"
"に"	1	1	letters	ja	adposition	"に"

```
"引っ越し"          1          1          letters          ja          verb          "引っ越
```

View the words tagged with entity "person", "location", "organization", or "other". These words are the words not tagged "non-entity".

```
idx = tdetails.Entity ~= "non-entity";
tdetails(idx,:).Token
```

```
ans = 11x1 string
    "マリー"
    "さん"
    "ボストン"
    "ニューヨーク"
    "鈴木"
    "さん"
    "東京"
    "大阪"
    "東京"
    "新宿"
    "渋谷"
```

## Stop Words

To remove stop words from documents according to the token language details, use `removeStopWords`. For a list of Japanese stop words set the 'Language' option in `stopWords` to 'ja'.

### Remove Japanese Stop Words

Tokenize Japanese text using `tokenizedDocument`. The function automatically detects Japanese text.

```
str = [
    "ここは静かなので、とても穏やかです"
    "企業内の顧客データを利用し、今年の売り上げを調べることが出来た。"
    "私は先生です。私は英語を教えています。"];
documents = tokenizedDocument(str);
```

Remove stop words using `removeStopWords`. The function uses the language details from `documents` to determine which language stop words to remove.

```
documents = removeStopWords(documents)
```

```
documents =
    3x1 tokenizedDocument:
```

```
    4 tokens: 静か 、 とても 穏やか
   10 tokens: 企業 顧客 データ 利用 、 今年 売り上げ 調べる 出来 。
    5 tokens: 先生 。 英語 教え 。
```

## Lemmatization

To lemmatize tokens according to the token language details, use `normalizeWords` and set the 'Style' option to 'lemma'.

### Lemmatize Japanese Text

Tokenize Japanese text using the `tokenizedDocument` function. The function automatically detects Japanese text.

```
str = [
  "空に星が輝き、瞬いている。"
  "空の星が輝きを増している。"
  "駅までは遠くて、歩けない。"
  "遠くの駅まで歩けない。"];
documents = tokenizedDocument(str);
```

Lemmatize the tokens using `normalizeWords`.

```
documents = normalizeWords(documents)

documents =
  4x1 tokenizedDocument:

  10 tokens: 空 に 星 が 輝 く 、 瞬 く て い る 。
  10 tokens: 空 の 星 が 輝 き を 増 す て い る 。
  9 tokens: 駅 ま で は 遠 い て 、 歩 け る な い 。
  7 tokens: 遠 く の 駅 ま で 歩 け る な い 。
```

## Language-Independent Features

### Word and N-Gram Counting

The `bagOfWords` and `bagOfNgrams` functions support `tokenizedDocument` input regardless of language. If you have a `tokenizedDocument` array containing your data, then you can use these functions.

### Modeling and Prediction

The `fitlda` and `fitlsa` functions support `bagOfWords` and `bagOfNgrams` input regardless of language. If you have a `bagOfWords` or `bagOfNgrams` object containing your data, then you can use these functions.

The `trainWordEmbedding` function supports `tokenizedDocument` or file input regardless of language. If you have a `tokenizedDocument` array or a file containing your data in the correct format, then you can use this function.

### See Also

`tokenizedDocument` | `removeStopWords` | `stopWords` | `addPartOfSpeechDetails` | `tokenDetails` | `normalizeWords` | `addLanguageDetails` | `addEntityDetails`

### **More About**

- “Language Considerations” on page 4-2
- “Analyze Japanese Text Data” on page 4-11



- Fit a topic model and visualize the results.

### Import Data

Read the data from "吾輩は猫である" by 夏目漱石 from [https://www.aozora.gr.jp/cards/000148/files/789\\_14547.html](https://www.aozora.gr.jp/cards/000148/files/789_14547.html) using the `webread` function.

Specify the character encoding of the text using the `weboptions` function. To find the correct character encoding for an HTML, look in the header of the HTML code. For this file, specify the character encoding to be "Shift\_JIS".

```
url = "https://www.aozora.gr.jp/cards/000148/files/789_14547.html";
options = weboptions('CharacterEncoding', 'Shift_JIS');
code = webread(url, options);
```

View the first few lines of the HTML code.

```
extractBefore(code, "<script")

ans =
'<?xml version="1.0" encoding="Shift_JIS"?>
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.1//EN"
"http://www.w3.org/TR/xhtml11/DTD/xhtml11.dtd">
<html xml:lang="ja" >
<head>
  <meta http-equiv="Content-Type" content="text/html; charset=Shift_JIS" />
  <meta http-equiv="content-style-type" content="text/css" />
  <link rel="stylesheet" type="text/css" href="../../aozora.css" />
  <title>夏目漱石 吾輩は猫である</title>
```

Extract the text data from the HTML using `extractHTMLText`. Split the text by newline characters.

```
textData = extractHTMLText(code);
textData = string(split(textData, newline));
textData(1:10)
```

```
ans = 10×1 string array
    "吾輩は猫である"
    ""
    "夏目漱石"
    ""
    ""
    ""
    ""
    ""
    " 吾輩は猫である。名前はまだ無い。"
    " どこで生れたかとうと見当がつかぬ。何でも薄暗いじめじめした所でニャーニャー泣いていた事だけは記憶している。吾輩
```

Remove the empty lines of text.

```
idx = textData == "";
textData(idx) = [];
textData(1:10)

ans = 10×1 string array
    "吾輩は猫である"
```





```

0 tokens:
1 tokens: 一
11 tokens: 吾輩は猫である。名前はまだ無い。
264 tokens: どこで生れたか とんと見当がつかぬ。何でも薄暗いじめじめした所でニャーニャー
100 tokens: この書生の掌の裏でしばらくはよい心持に坐っておったが、しばらくすると非常
92 tokens: ふと気が付いて見ると書生はいない。たくさんおった兄弟が一疋も見えぬ。肝
693 tokens: ようやくの思いで笹原を這い出すと向うに大きな池がある。吾輩は池の前に坐っ
276 tokens: 吾輩の主人は滅多に吾輩と顔を合せる事がない。職業は教師だそうだ。学校か

```

## Get Part-of-Speech Tags

Get the token details and then view the details of the first few tokens.

```

tdetails = tokenDetails(documents);
head(tdetails)

```

ans=8x8 table

Token	DocumentNumber	LineNumber	Type	Language	PartOfSpeech	Lemma
"吾輩"	1	1	letters	ja	pronoun	"吾輩"
"は"	1	1	letters	ja	adposition	"は"
"猫"	1	1	letters	ja	noun	"猫"
"で"	1	1	letters	ja	auxiliary-verb	"だ"
"ある"	1	1	letters	ja	auxiliary-verb	"ある"
"夏目"	2	1	letters	ja	proper-noun	"夏目"
"漱石"	2	1	letters	ja	proper-noun	"漱石"
"一"	4	1	letters	ja	numeral	"一"

The PartOfSpeech variable in the table contains the part-of-speech tags of the tokens. Create word clouds of all the nouns and adjectives, respectively.

```

figure
idx = tdetails.PartOfSpeech == "noun";
tokens = tdetails.Token(idx);
subplot(1,2,1)
wordcloud(tokens);
title("Nouns")

idx = tdetails.PartOfSpeech == "adjective";
tokens = tdetails.Token(idx);
subplot(1,2,2)
wordcloud(tokens);
title("Adjectives")

```



```

2 tokens: 吾輩 猫
2 tokens: 夏目 漱石
0 tokens:
0 tokens:
4 tokens: 吾輩 猫 まだ 無い
102 tokens: 生れ とんと 見当 つか ぬ 薄暗い じめじめ ニャーニャー 泣い いた事 記憶 吾輩 始め 人間 という
36 tokens: 書生 掌 裏 しばらく よい 心持 坐っ おっ しばらく 非常 速力 運転 始め 書生 動く 動く 分ら 無
38 tokens: ふと 付い 見る 書生 おっ 兄弟 一疋 見え ぬ 肝心 母親 姿 隠し しまっ 上今 違っ 無 暗に 明る
274 tokens: ようやく 思い 笹原 這い出す 向う 大きな 池 吾輩 池 坐っ たら よかる 考え 別に という 分別 出
101 tokens: 吾輩 主人 滅多 吾輩 顔 合せる 職業 教師 学校 帰る 終日 書齋 這入っ ぎりほとんど 出 来る 大変

```

Lemmatize the text using `normalizeWords`.

```
documents = normalizeWords(documents);
documents(1:10)
```

```
ans =
10×1 tokenizedDocument:

2 tokens: 吾輩 猫
2 tokens: 夏目 漱石
0 tokens:
0 tokens:
4 tokens: 吾輩 猫 まだ 無い
102 tokens: 生れる とんと 見当 つく ぬ 薄暗い じめじめ ニャーニャー 泣く いた事 記憶 吾輩 始める 人間 とい
36 tokens: 書生 掌 裏 しばらく よい 心持 坐る おる しばらく 非常 速力 運転 始める 書生 動く 動く 分る 無
38 tokens: ふと 付く 見る 書生 おる 兄弟 一疋 見える ぬ 肝心 母親 姿 隠す しまう 上今 違う 無 暗に 明る
274 tokens: ようやく 思い 笹原 這い出す 向う 大きな 池 吾輩 池 坐る た よい 考える 別に という 分別 出る
101 tokens: 吾輩 主人 滅多 吾輩 顔 合せる 職業 教師 学校 帰る 終日 書齋 這入る ぎりほとんど 出る 来る 大変

```

Some preprocessing steps, such as removing stop words and erasing punctuation, return empty documents. Remove the empty documents using the `removeEmptyDocuments` function.

```
documents = removeEmptyDocuments(documents);
```

### Create Preprocessing Function

Creating a function that performs preprocessing can be useful to prepare different collections of text data in the same way. For example, you can use a function to preprocess new data using the same steps as the training data.

Create a function which tokenizes and preprocesses the text data to use for analysis. The function `preprocessJapaneseText`, performs these steps:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Erase punctuation using `erasePunctuation`.
- 3 Remove a list of stop words (such as "あそこ", "あたり", and "あちら") using `removeStopWords`.
- 4 Lemmatize the words using `normalizeWords`.

Remove the empty documents after preprocessing using the `removeEmptyDocuments` function. Removing documents after using a preprocessing function makes it easier to remove corresponding data such as labels from other sources.

In this example, use the preprocessing function `preprocessJapaneseText`, listed at the end of the example, to prepare the text data.

```
documents = preprocessJapaneseText(textData);
documents(1:5)
```

```
ans =
  5×1 tokenizedDocument:

    2 tokens: 吾輩 猫
    2 tokens: 夏目 漱石
    0 tokens:
    0 tokens:
    4 tokens: 吾輩 猫 まだ 無い
```

Remove the empty documents.

```
documents = removeEmptyDocuments(documents);
```

### Fit Topic Model

Fit a latent Dirichlet allocation (LDA) topic model to the data. An LDA model discovers underlying topics in a collection of documents and infers word probabilities in topics.

To fit an LDA model to the data, you first must create a bag-of-words model. A bag-of-words model (also known as a term-frequency counter) records the number of times that words appear in each document of a collection. Create a bag-of-words model using `bagOfWords`.

```
bag = bagOfWords(documents);
```

Remove the empty documents from the bag-of-words model.

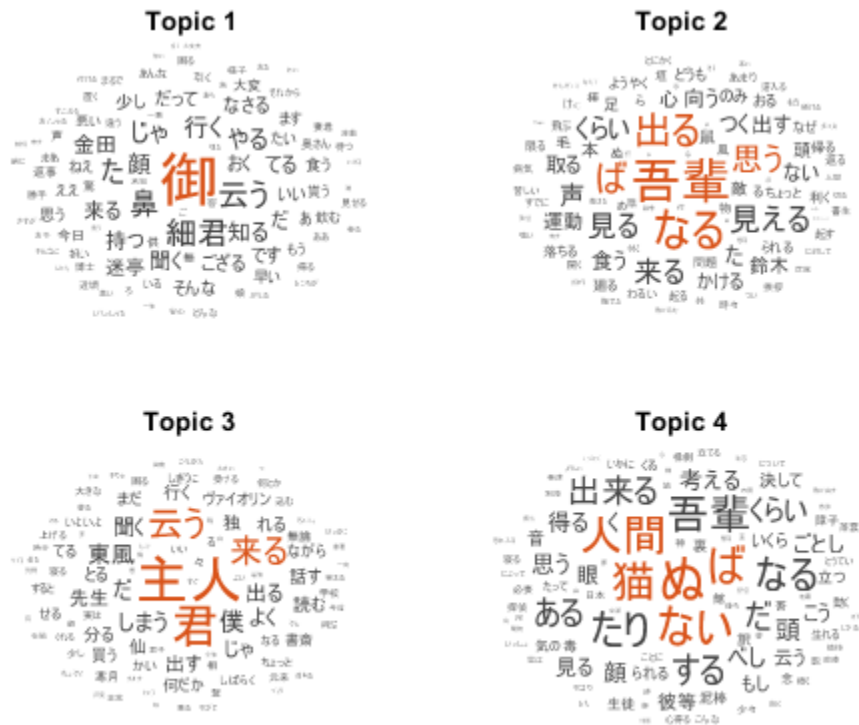
```
bag = removeEmptyDocuments(bag);
```

Fit an LDA model with seven topics using `fitlda`. To suppress the verbose output, set `'Verbose'` to `0`.

```
numTopics = 7;
mdl = fitlda(bag,numTopics,'Verbose',0);
```

Visualize the first four topics using word clouds.

```
figure
for i = 1:4
    subplot(2,2,i)
    wordcloud(mdl,i);
    title("Topic " + i)
end
```



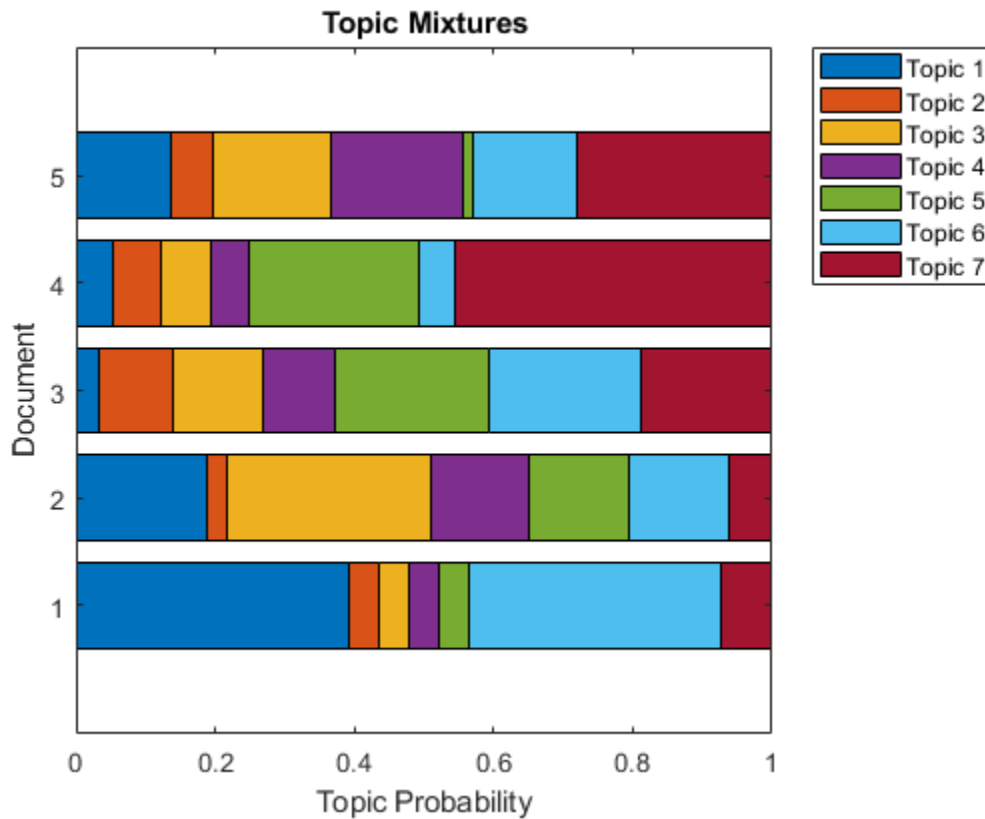
Visualize multiple topic mixtures using stacked bar charts. View five input documents at random and visualize the corresponding topic mixtures.

```
numDocuments = numel(documents);
idx = randperm(numDocuments,5);
documents(idx)
```

```
ans =
    5×1 tokenizedDocument:
```

```
    4 tokens: 細君 細君 なかなか さばける
    7 tokens: 進行 せる 山々 どうしても 暮れる くれる 困る
   13 tokens: 来る そんな 仙骨 相手 少々 骨 折れる 過ぎる 宛然 たり 仙 伝 人物
    3 tokens: 先生 譜 下さる
   23 tokens: 立つ 月給 上がる いくら 勉強 褒める くれる 郎 君 独 寂寞 中学 時代 覚える 詩 句 細君 朗吟 細君
```

```
topicMixtures = transform mdl,documents(idx));
figure
barh(topicMixtures(1:5,:), 'stacked')
xlim([0 1])
title("Topic Mixtures")
xlabel("Topic Probability")
ylabel("Document")
legend("Topic " + string(1:numTopics), 'Location', 'northeastoutside')
```



### Example Preprocessing Function

The function `preprocessJapaneseText`, performs these steps:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Erase punctuation using `erasePunctuation`.
- 3 Remove a list of stop words (such as "あそこ", "あたり", and "あちら") using `removeStopWords`.
- 4 Lemmatize the words using `normalizeWords`.

```
function documents = preprocessJapaneseText(textData)
```

```
% Tokenize the text.
```

```
documents = tokenizedDocument(textData);
```

```
% Erase the punctuation.
```

```
documents = erasePunctuation(documents);
```

```
% Remove a list of stop words.
```

```
documents = removeStopWords(documents);
```

```
% Lemmatize the words.
```

```
documents = normalizeWords(documents, 'Style', 'lemma');  
end
```

### See Also

[tokenizedDocument](#) | [removeStopWords](#) | [stopWords](#) | [addPartOfSpeechDetails](#) | [tokenDetails](#) | [normalizeWords](#)

### More About

- “Language Considerations” on page 4-2
- “Create Simple Text Model for Classification” on page 2-2
- “Analyze Text Data Using Topic Models” on page 2-13
- “Analyze Text Data Using Multiword Phrases” on page 2-7
- “Analyze Text Data Containing Emojis” on page 2-52
- “Train a Sentiment Classifier” on page 2-71
- “Classify Text Data Using Deep Learning” on page 2-85
- “Generate Text Using Deep Learning” (Deep Learning Toolbox)

### See Also



## German Language Support

This topic summarizes the Text Analytics Toolbox features that support German text. For an example showing how to analyze German text data, see “Analyze German Text Data” on page 4-26.

### Tokenization

The `tokenizedDocument` function automatically detects German input. Alternatively, set the 'Language' option in `tokenizedDocument` to 'de'. This option specifies the language details of the tokens. To view the language details of the tokens, use `tokenDetails`. These language details determine the behavior of the `removeStopWords`, `addPartOfSpeechDetails`, `normalizeWords`, `addSentenceDetails`, and `addEntityDetails` functions on the tokens.

#### Tokenize German Text

Tokenize German text using `tokenizedDocument`. The function automatically detects German text.

```
str = [
    "Guten Morgen. Wie geht es dir?"
    "Heute wird ein guter Tag."];
documents = tokenizedDocument(str)

documents =
    2x1 tokenizedDocument:

    8 tokens: Guten Morgen . Wie geht es dir ?
    6 tokens: Heute wird ein guter Tag .
```

### Sentence Detection

To detect sentence structure in documents, use the `addSentenceDetails`. You can use the `abbreviations` function to help create custom lists of abbreviations to detect.

#### Add Sentence Details to German Documents

Tokenize German text using `tokenizedDocument`.

```
str = [
    "Guten Morgen, Dr. Schmidt. Geht es Ihnen wieder besser?"
    "Heute wird ein guter Tag."];
documents = tokenizedDocument(str);
```

Add sentence details to the documents using `addSentenceDetails`. This function adds the sentence numbers to the table returned by `tokenDetails`. View the updated token details of the first few tokens.

```
documents = addSentenceDetails(documents);
tdetails = tokenDetails(documents);
head(tdetails,10)
```

ans=10x6 table

Token	DocumentNumber	SentenceNumber	LineNumber	Type	Language
-------	----------------	----------------	------------	------	----------

"Guten"	1	1	1	letters	de
"Morgen"	1	1	1	letters	de
","	1	1	1	punctuation	de
"Dr"	1	1	1	letters	de
". "	1	1	1	punctuation	de
"Schmidt"	1	1	1	letters	de
". "	1	1	1	punctuation	de
"Geht"	1	2	1	letters	de
"es"	1	2	1	letters	de
"Ihnen"	1	2	1	letters	de

### Table of German Abbreviations

View a table of German abbreviations. Use this table to help create custom tables of abbreviations for sentence detection when using `addSentenceDetails`.

```
tbl = abbreviations('Language', 'de');
head(tbl)
```

ans=8x2 table

Abbreviation	Usage
"A.T"	regular
"ABl"	regular
"Abb"	regular
"Abdr"	regular
"Abf"	regular
"Abfl"	regular
"Abh"	regular
"Abk"	regular

### Part of Speech Details

To add German part of speech details to documents, use the `addPartOfSpeechDetails` function.

#### Get Part of Speech Details of German Text

Tokenize German text using `tokenizedDocument`.

```
str = [
    "Guten Morgen. Wie geht es dir?"
    "Heute wird ein guter Tag. "];
documents = tokenizedDocument(str)

documents =
    2x1 tokenizedDocument:

    8 tokens: Guten Morgen . Wie geht es dir ?
    6 tokens: Heute wird ein guter Tag .
```

To get the part of speech details for German text, first use `addPartOfSpeechDetails`.

```
documents = addPartOfSpeechDetails(documents);
```

To view the part of speech details, use the `tokenDetails` function.

```
tdetails = tokenDetails(documents);
head(tdetails)
```

ans=8x7 table

Token	DocumentNumber	SentenceNumber	LineNumber	Type	Language	Part
"Guten"	1	1	1	letters	de	adj
"Morgen"	1	1	1	letters	de	noun
". "	1	1	1	punctuation	de	punc
"Wie"	1	2	1	letters	de	adv
"geht"	1	2	1	letters	de	verb
"es"	1	2	1	letters	de	pron
"dir"	1	2	1	letters	de	pron
"?"	1	2	1	punctuation	de	punc

## Named Entity Recognition

To add entity tags to documents, use the `addEntityDetails` function.

### Add Named Entity Tags to German Text

Tokenize German text using `tokenizedDocument`.

```
str = [
    "Ernst zog von Frankfurt nach Berlin."
    "Besuchen Sie Volkswagen in Wolfsburg."];
documents = tokenizedDocument(str);
```

To add entity tags to German text, use the `addEntityDetails` function. This function detects person names, locations, organizations, and other named entities.

```
documents = addEntityDetails(documents);
```

To view the entity details, use the `tokenDetails` function.

```
tdetails = tokenDetails(documents);
head(tdetails)
```

ans=8x8 table

Token	DocumentNumber	SentenceNumber	LineNumber	Type	Language	Part
"Ernst"	1	1	1	letters	de	per
"zog"	1	1	1	letters	de	v
"von"	1	1	1	letters	de	a
"Frankfurt"	1	1	1	letters	de	pl
"nach"	1	1	1	letters	de	a
"Berlin"	1	1	1	letters	de	pl
". "	1	1	1	punctuation	de	p
"Besuchen"	2	1	1	letters	de	v

View the words tagged with entity "person", "location", "organization", or "other". These words are the words not tagged with "non-entity".

```
idx = tdetails.Entity ~= "non-entity";
tdetails(idx,:)
```

ans=5x8 table

Token	DocumentNumber	SentenceNumber	LineNumber	Type	Language	Part
"Ernst"	1	1	1	letters	de	prop
"Frankfurt"	1	1	1	letters	de	prop
"Berlin"	1	1	1	letters	de	prop
"Volkswagen"	2	1	1	letters	de	noun
"Wolfsburg"	2	1	1	letters	de	prop

## Stop Words

To remove stop words from documents according to the token language details, use `removeStopWords`. For a list of German stop words set the 'Language' option in `stopWords` to 'de'.

### Remove German Stop Words from Documents

Tokenize German text using `tokenizedDocument`. The function automatically detects German text.

```
str = [
    "Guten Morgen. Wie geht es dir?"
    "Heute wird ein guter Tag."];
documents = tokenizedDocument(str)

documents =
    2x1 tokenizedDocument:

    8 tokens: Guten Morgen . Wie geht es dir ?
    6 tokens: Heute wird ein guter Tag .
```

Remove stop words using the `removeStopWords` function. The function uses the language details from documents to determine which language stop words to remove.

```
documents = removeStopWords(documents)

documents =
    2x1 tokenizedDocument:

    5 tokens: Guten Morgen . geht ?
    5 tokens: Heute wird guter Tag .
```

## Stemming

To stem tokens according to the token language details, use `normalizeWords`.

## Stem German Text

Tokenize German text using the `tokenizedDocument` function. The function automatically detects German text.

```
str = [
    "Guten Morgen. Wie geht es dir?"
    "Heute wird ein guter Tag."];
documents = tokenizedDocument(str);
```

Stem the tokens using `normalizeWords`.

```
documents = normalizeWords(documents)
```

```
documents =
  2x1 tokenizedDocument:

    8 tokens: gut morg . wie geht es dir ?
    6 tokens: heut wird ein gut tag .
```

## Language-Independent Features

### Word and N-Gram Counting

The `bagOfWords` and `bagOfNgrams` functions support `tokenizedDocument` input regardless of language. If you have a `tokenizedDocument` array containing your data, then you can use these functions.

### Modeling and Prediction

The `fitlda` and `fitlsa` functions support `bagOfWords` and `bagOfNgrams` input regardless of language. If you have a `bagOfWords` or `bagOfNgrams` object containing your data, then you can use these functions.

The `trainWordEmbedding` function supports `tokenizedDocument` or file input regardless of language. If you have a `tokenizedDocument` array or a file containing your data in the correct format, then you can use this function.

## See Also

[tokenizedDocument](#) | [removeStopWords](#) | [stopWords](#) | [addPartOfSpeechDetails](#) | [tokenDetails](#) | [normalizeWords](#) | [addLanguageDetails](#)

## More About

- “Language Considerations” on page 4-2
- “Analyze German Text Data” on page 4-26

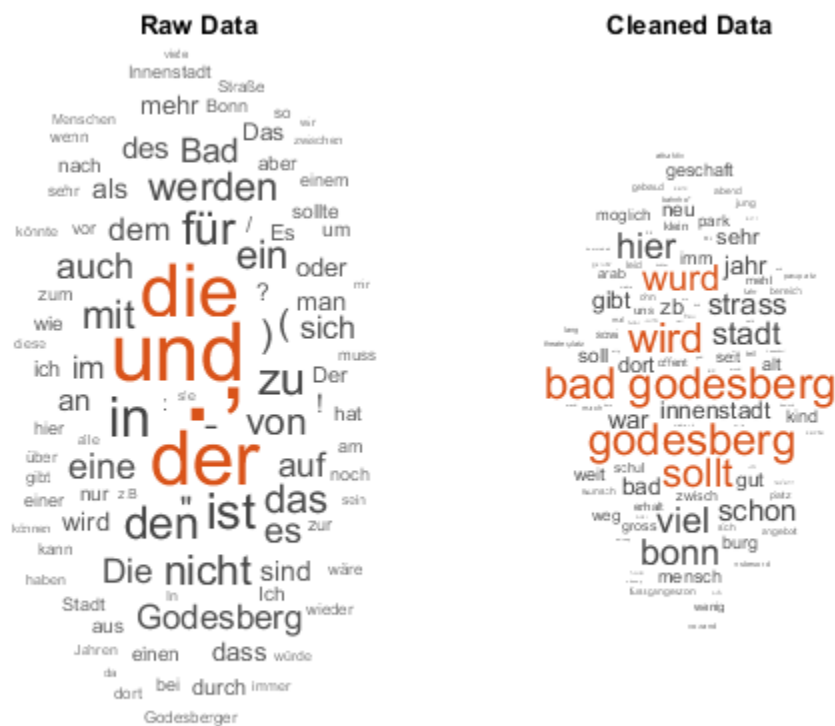
## Analyze German Text Data

This example shows how to import, prepare, and analyze German text data using a topic model.

German text data can be large and can contain lots of noise that negatively affects statistical analysis. For example, the text data can contain the following:

- Variations in word forms. For example, „rot“, „rote“, and „roten“.
- Words that add noise. For example, stop words such as „der“, „die“, and „das“.
- Punctuation and special characters.

These word clouds illustrate word frequency analysis applied to some raw text data and a preprocessed version of the same text data.



This example first shows how to import and prepare German text data, and then it shows how to analyze the text data using a Latent Dirichlet Allocation (LDA) model. An LDA model is a topic model that discovers underlying topics in a collection of documents and infers the word probabilities in topics. Use these steps in preparing the text data and fitting the model:

- Import the text data from a CSV file and extract the relevant data.
- Prepare the text data for analysis using standard preprocessing techniques.
- Fit a topic model and visualize the results.

## Import Data

Download the data `vorhaben.csv` from <https://opendata.bonn.de/dataset/vorhabenliste-b%C3%BCrgerbeteiligungen-planungen-und-projekte>. This file can change over time, so the results in the example can vary.

Use `detectImportOptions` to determine the format of the CSV file and set the text type to string. Set the 'Encoding' option to 'ISO-8859-15'. Read the data using the `readtable` function and view the first few rows.

```
filename = "vorhaben.csv";
options = detectImportOptions(filename, 'TextType', 'string', 'Encoding', 'ISO-8859-15');
data = readtable(filename, options);
head(data)
```

`ans=8x19 table`

Titel

---

```
"Bauleitplanverfahren zur Aufstellung des vorhabenbezogenen-Bebauungsplans Nr. 6620-1 ?Bundes
"Bauleitplanverfahren zur Aufstellung des vorhabenbezogenen-Bebauungsplans Nr. 6522-1 "Didin
"Bauleitplanverfahren zur Aufstellung des Bebauungsplans-Nr. 7621-56 ?Sebastianstraße?"
"EPICURO - European Partnership for Innovative Cities within and Urban Resilience Outlook"
"Bauleitplanverfahren zur Aufstellung des Bebauungsplanes Nr. 6719-3 "Schwimmbad Wasserland"
"Bürgerbeteiligung an der Konzepterstellung für den Neubau eines Schwimmbades in Bonn-Dotten
"Integriertes Handlungskonzept Grüne Infrastruktur (InHK GI) zur-zukünftigen Freiraumsicheru
"Verlängerung des Teufelsbachweges bis zur L 83n"
```

Extract the text data from the variable `InhaltlicheBeschreibungUndZielsetzung` (the description of the content and the goal).

```
textData = data.InhaltlicheBeschreibungUndZielsetzung;
```

Visualize the text data in a word cloud.

```
figure
wordcloud(textData);
```





```
tdetails = tokenDetails(documents);
head(tdetails)
```

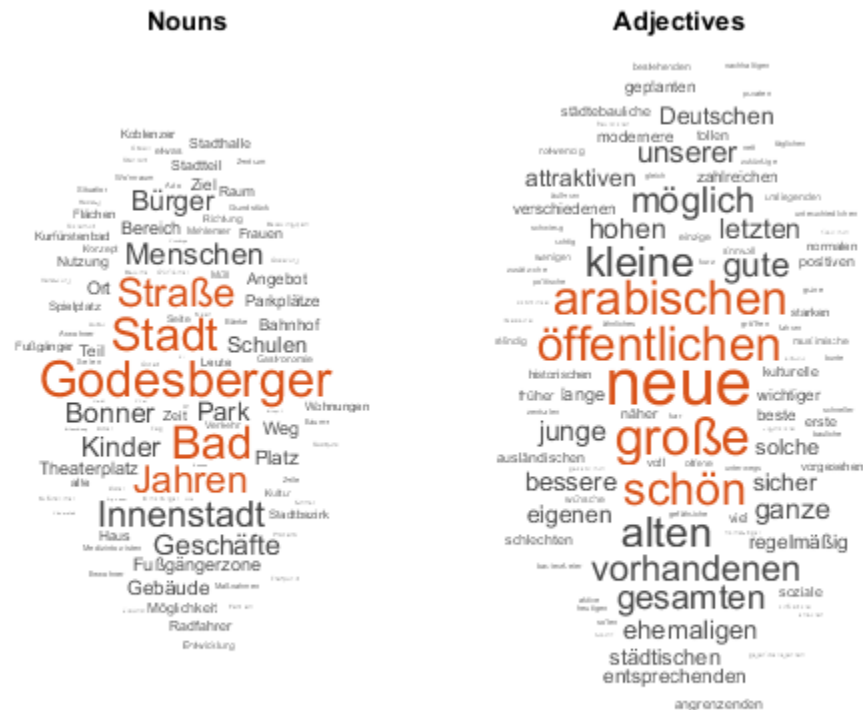
```
ans=8x7 table
```

Token	DocumentNumber	SentenceNumber	LineNumber	Type	Language
"Für"	1	1	1	letters	de
"das"	1	1	1	letters	de
"Gebiet"	1	1	1	letters	de
"zwischen"	1	1	1	letters	de
"Reuterstraße"	1	1	1	letters	de
", "	1	1	1	punctuation	de
"Bundeskanzlerplatz"	1	1	1	letters	de
", "	1	1	1	punctuation	de

The `PartOfSpeech` variable in the table contains the part-of-speech tags of the tokens. Create word clouds of all the nouns and adjectives, respectively.

```
figure
idx = tdetails.PartOfSpeech == "noun";
tokens = tdetails.Token(idx);
subplot(1,2,1)
wordcloud(tokens);
title("Nouns")

idx = tdetails.PartOfSpeech == "adjective";
tokens = tdetails.Token(idx);
subplot(1,2,2)
wordcloud(tokens);
title("Adjectives")
```



### Prepare Text Data for Analysis

Tokenize the text using `tokenizedDocument` and view the first few documents.

```
documentsRaw = tokenizedDocument(textData);
documents = documentsRaw;
documents(1:10)
```

```
ans =
    10×1 tokenizedDocument:
```

```
50 tokens: Für das Gebiet zwischen Reuterstraße , Bundeskanzlerplatz , Willy-Brandt-Allee ,
46 tokens: Für den vorhabenbezogenen Bebauungsplan Nr . 6522-1 ? Didinkirica ? der Bundessta
41 tokens: Für das Gebiet zwischen Alfred-Bucherer-Straße , Sebastianstraße und dem Fußweg
134 tokens: In den vergangenen Jahren führte der Klimawandel zu einer Vielzahl von Folgen für
24 tokens: Schaffung von Planungsrecht für den Bau eines neuen Familien - , Schul - und Spor
80 tokens: Für die begleitende Bürgerbeteiligung bei der Konzepterstellung , für das neue Sch
60 tokens: In der Gebietskulisse des Grünen C sollen die Freiräume auch zukünftig im Sinne v
51 tokens: Zur Entlastung von Pützchen / Bechlinghoven vor Durchgangsverkehr und in Verbind
29 tokens: Für das Areal der ehemaligen Landwirtschaftskammer sowie einer angrenzenden städ
37 tokens: Für das Areal Herbert-Rabius-Straße im Stadtbezirk Beuel , Ortsteil Beuel-Mitte :
```

Replace common phrases (n-grams) with a single token and remove the stop words.

```
old = ["Bad" "Godesberg"];
new = "Bad Godesberg";
documents = replaceNgrams(documents,old,new);
```

```
documents = removeStopWords(documents);
documents(1:10)
```

```
ans =
```

```
10x1 tokenizedDocument:
```

```
35 tokens: Gebiet zwischen Reuterstraße , Bundeskanzlerplatz , Willy-Brandt-Allee , Eduard-P
33 tokens: vorhabenbezogenen Bebauungsplan Nr . 6522-1 ? Didinkirica ? Bundesstadt Bonn , Sta
27 tokens: Gebiet zwischen Alfred-Bucherer-Straße , Sebastianstraße Fußweg zwischen Röckumst
81 tokens: vergangenen Jahren führte Klimawandel Vielzahl Folgen Umwelt , Wirtschaft Menschen
15 tokens: Schaffung Planungsrecht Bau neuen Familien - , Schul - Sportschwimmbades Flächen n
57 tokens: begleitende Bürgerbeteiligung Konzepterstellung neue Schwimmbad soll folgenden The
40 tokens: Gebietskulisse Grünen C sollen Freiräume zukünftig Sinne Naherholung , Landwirtsch
32 tokens: Entlastung Pützchen / Bechlinghoven Durchgangsverkehr Verbindung geplanten Anschl
19 tokens: Areal ehemaligen Landwirtschaftskammer sowie angrenzenden städtischen Fläche Stad
25 tokens: Areal Herbert-Rabius-Straße Stadtbezirk Beuel , Ortsteil Beuel-Mitte soll vorhaben
```

Normalize the text using the `normalizeWords` function.

```
documents = normalizeWords(documents);
documents(1:10)
```

```
ans =
```

```
10x1 tokenizedDocument:
```

```
35 tokens: gebiet zwisch reuterstrass , bundeskanzlerplatz , willy-brandt-alle , eduard-pflug
33 tokens: vorhabenbezog bebauungsplan nr . 6522-1 ? didinkirica ? bundesstadt bonn , stadtbe
27 tokens: gebiet zwisch alfred-bucherer-strass , sebastianstrass fussweg zwisch rockumstrass
81 tokens: vergang jahr fuhr klimawandel vielzahl folg umwelt , wirtschaft mensch . stadt ge
15 tokens: schaffung planungsrecht bau neu famili - , schul - sportschwimmbad flach nordlich
57 tokens: begleit burgerbeteil konzepterstell neu schwimmbad soll folgend them beteil geb :
40 tokens: gebietskuliss grun c soll freiraum zukunft sinn naherhol , landwirtschaft natursch
32 tokens: entlast putzch / bechlinghov durchgangsverkehr verbind geplant anschlussstell maar
19 tokens: areal ehemal landwirtschaftskamm sowi angrenz stadtisch flach stadtbezirk beuel ,
25 tokens: areal herbert-rabius-strass stadtbezirk beuel , ortsteil beuel-mitt soll vorhaben
```

Erase the punctuation using the `erasePunctuation` function.

```
documents = erasePunctuation(documents);
documents(1:10)
```

```
ans =
```

```
10x1 tokenizedDocument:
```

```
27 tokens: gebiet zwisch reuterstrass bundeskanzlerplatz willybrandtalle eduardpflugerstrass
25 tokens: vorhabenbezog bebauungsplan nr 65221 didinkirica bundesstadt bonn stadtbezirk bonn
22 tokens: gebiet zwisch alfredbuchererstrass sebastianstrass fussweg zwisch rockumstrass en
64 tokens: vergang jahr fuhr klimawandel vielzahl folg umwelt wirtschaft mensch stadt gemein
11 tokens: schaffung planungsrecht bau neu famili schul sportschwimmbad flach nordlich heizk
41 tokens: begleit burgerbeteil konzepterstell neu schwimmbad soll folgend them beteil geb d
31 tokens: gebietskuliss grun c soll freiraum zukunft sinn naherhol landwirtschaft naturschur
27 tokens: entlast putzch bechlinghov durchgangsverkehr verbind geplant anschlussstell maars
18 tokens: areal ehemal landwirtschaftskamm sowi angrenz stadtisch flach stadtbezirk beuel o
19 tokens: areal herbertrabiusstrass stadtbezirk beuel ortsteil beuelmitt soll vorhabenbezog
```

Visualize the raw and cleaned data in word clouds.

```

figure
subplot(1,2,1)
wordcloud(documentsRaw);
title("Raw Data")

subplot(1,2,2)
wordcloud(documents);
title("Cleaned Data")

```



### Create Preprocessing Function

Creating a function that performs preprocessing can be useful to prepare different collections of text data in the same way. For example, you can use a function to preprocess new data using the same steps as the training data.

Create a function which tokenizes and preprocesses the text data to use for analysis. The function `preprocessGermanText`, listed at the end of the example, performs these steps:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Replace the multiword phrase ["Bad " "Godesberg "] with "Bad Godesberg".
- 3 Remove a list of stop words (such as „der“, „die“, and „das“) using `removeStopWords`.
- 4 Normalize the words using `normalizeWords`.
- 5 Erase punctuation using `erasePunctuation`.

Remove the empty documents after preprocessing using the `removeEmptyDocuments` function. Removing documents after using a preprocessing function makes it easier to remove corresponding data such as labels from other sources.

In this example, use the preprocessing function `preprocessGermanText`, listed at the end of the example, to prepare the text data.

```
documents = preprocessGermanText(textData);
documents(1:5)
```

```
ans =
    5×1 tokenizedDocument:
```

```
    27 tokens: gebiet zwisch reuterstrass bundeskanzlerplatz willybrandtalle eduardpflugerstrass
    25 tokens: vorhabenbezog bebauungsplan nr 65221 didinkirica bundesstadt bonn stadtbezirk bonn
    22 tokens: gebiet zwisch alfredbuchererstrass sebastianstrass fussweg zwisch rockumstrass end
    64 tokens: vergang jahr fuhr klimawandel vielzahl folg umwelt wirtschaft mensch stadt gemein
    11 tokens: schaffung planungsrecht bau neu famili schul sportschwimmbad flach nordlich heizk
```

Remove the empty documents using the `removeEmptyDocuments` function.

```
documents = removeEmptyDocuments(documents);
```

### Fit Topic Model

Fit a latent Dirichlet allocation (LDA) topic model to the data. An LDA model discovers underlying topics in a collection of documents and infers word probabilities in topics.

To fit an LDA model to the data, you first must create a bag-of-words model. A bag-of-words model (also known as a term-frequency counter) records the number of times that words appear in each document of a collection. Create a bag-of-words model using `bagOfWords`.

```
bag = bagOfWords(documents);
```

Remove the empty documents from the bag-of-words model.

```
bag = removeEmptyDocuments(bag);
```

Fit an LDA model with seven topics using `fitlda`. To suppress the verbose output, set `'Verbose'` to `0`.

```
numTopics = 7;
mdl = fitlda(bag,numTopics,'Verbose',0);
```

Visualize the first four topics using word clouds.

```
figure
for i = 1:4
    subplot(2,2,i)
    wordcloud(mdl,i);
    title("Topic " + i)
end
```



Visualize multiple topic mixtures using stacked bar charts. View five input documents at random and visualize the corresponding topic mixtures.

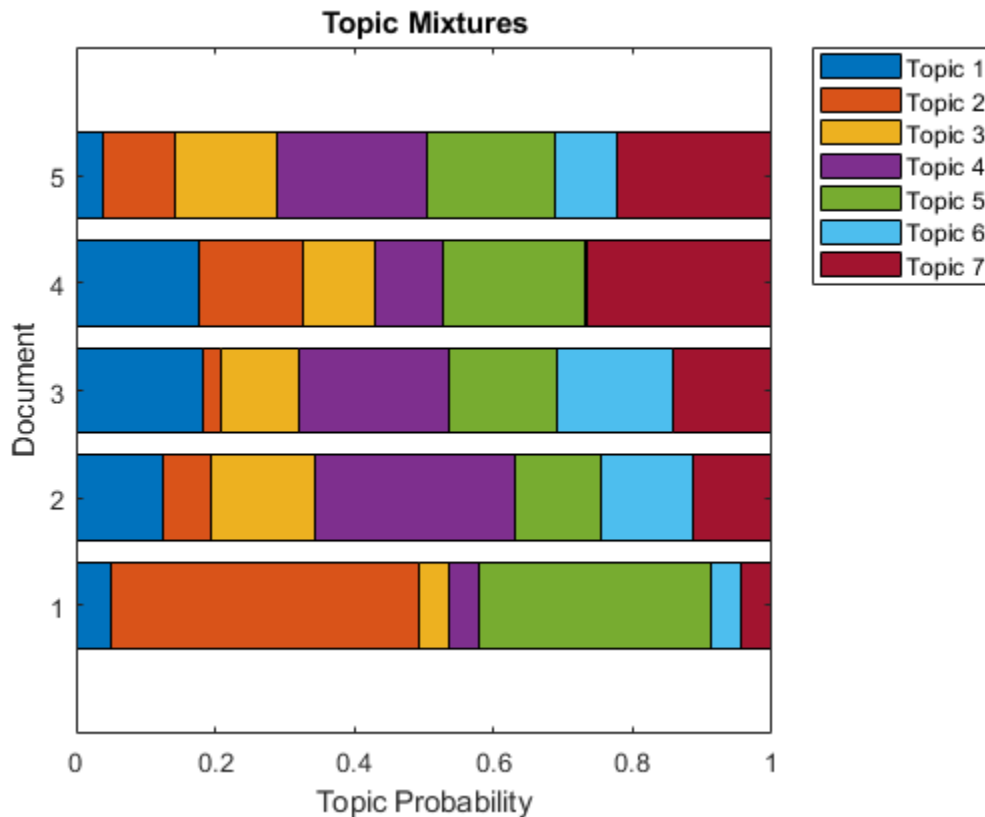
```
numDocuments = numel(documents);
idx = randperm(numDocuments,5);
documents(idx)
```

```
ans =
    5×1 tokenizedDocument:
```

```
    4 tokens: gastronom offer sollt verbessert
   82 tokens: grunflach dietrichglaunerstrass rand dorfplatz entlang fussweg mehlem bach entlar
  116 tokens: sportplatz plittersdorf kommt leid regelmass unschon vorfall einsehbar umfeld her
   64 tokens: mainz strass bereich geschaft kirch uberwieg beidseit zugeparkt unschon sond fuss
   50 tokens: "1" "bezirksverodnet" "sollt" "kulturburgermeist" "gewahlt" "hatt" "aufgab" "ver"
```

```
topicMixtures = transform mdl,documents(idx));
```

```
figure
barh(topicMixtures(1:5,:), 'stacked')
xlim([0 1])
title("Topic Mixtures")
xlabel("Topic Probability")
ylabel("Document")
legend("Topic " + string(1:numTopics), 'Location', 'northeastoutside')
```



### Example Preprocessing Function

The function `preprocessGermanText`, performs these steps:

- 1 Tokenize the text using `tokenizedDocument`.
- 2 Replace the multiword phrase ["Bad" "Godesberg"] with "Bad Godesberg".
- 3 Remove a list of stop words (such as „der“, „die“, and „das“) using `removeStopWords`.
- 4 Normalize the words using `normalizeWords`.
- 5 Erase punctuation using `erasePunctuation`.

```
function documents = preprocessGermanText(textData)
```

```
% Tokenize the text.
```

```
documents = tokenizedDocument(textData);
```

```
% Replace multiword phrases
```

```
old = ["Bad" "Godesberg"];
```

```
new = "Bad Godesberg";
```

```
documents = replaceNgrams(documents,old,new);
```

```
% Remove a list of stop words.
```

```
documents = removeStopWords(documents);
```

```
% Normalize the words.
```

```
documents = normalizeWords(documents);
```

```
% Erase the punctuation.  
documents = erasePunctuation(documents);  
  
end
```

### See Also

[tokenizedDocument](#) | [removeStopWords](#) | [stopWords](#) | [addPartOfSpeechDetails](#) | [tokenDetails](#) | [normalizeWords](#)

### More About

- “Language Considerations” on page 4-2
- “Create Simple Text Model for Classification” on page 2-2
- “Analyze Text Data Using Topic Models” on page 2-13
- “Analyze Text Data Using Multiword Phrases” on page 2-7
- “Analyze Text Data Containing Emojis” on page 2-52
- “Train a Sentiment Classifier” on page 2-71
- “Classify Text Data Using Deep Learning” on page 2-85
- “Generate Text Using Deep Learning” (Deep Learning Toolbox)



# Korean Language Support

This topic summarizes the Text Analytics Toolbox features that support Korean text.

## Tokenization

The `tokenizedDocument` function automatically detects Korean input. Alternatively, set the 'Language' option in `tokenizedDocument` to 'ko'. This option specifies the language details of the tokens. To view the language details of the tokens, use `tokenDetails`. These language details determine the behavior of the `removeStopWords`, `addPartOfSpeechDetails`, `normalizeWords`, `addSentenceDetails`, and `addEntityDetails` functions on the tokens.

To specify additional MeCab options for tokenization, create a `mecabOptions` object. To tokenize using the specified MeCab tokenization options, use the 'TokenizeMethod' option of `tokenizedDocument`.

## Part of Speech Details

The `tokenDetails` function, by default, includes part of speech details with the token details.

## Named Entity Recognition

The `tokenDetails` function, by default, includes entity details with the token details.

## Stop Words

To remove stop words from documents according to the token language details, use `removeStopWords`. For a list of Korean stop words set the 'Language' option in `stopWords` to 'ko'.

## Lemmatization

To lemmatize tokens according to the token language details, use `normalizeWords` and set the 'Style' option to 'lemma'.

## Language-Independent Features

### Word and N-Gram Counting

The `bagOfWords` and `bagOfNgrams` functions support `tokenizedDocument` input regardless of language. If you have a `tokenizedDocument` array containing your data, then you can use these functions.

### Modeling and Prediction

The `fitlda` and `fitlsa` functions support `bagOfWords` and `bagOfNgrams` input regardless of language. If you have a `bagOfWords` or `bagOfNgrams` object containing your data, then you can use these functions.

The `trainWordEmbedding` function supports `tokenizedDocument` or file input regardless of language. If you have a `tokenizedDocument` array or a file containing your data in the correct format, then you can use this function.

### See Also

`tokenizedDocument` | `removeStopWords` | `stopWords` | `addPartOfSpeechDetails` | `tokenDetails` | `normalizeWords` | `addLanguageDetails` | `addEntityDetails`

### More About

- “Language Considerations” on page 4-2

## Language-Independent Features

### Word and N-Gram Counting

The `bagOfWords` and `bagOfNgrams` functions support `tokenizedDocument` input regardless of language. If you have a `tokenizedDocument` array containing your data, then you can use these functions.

### Modeling and Prediction

The `fitlda` and `fitlsa` functions support `bagOfWords` and `bagOfNgrams` input regardless of language. If you have a `bagOfWords` or `bagOfNgrams` object containing your data, then you can use these functions.

The `trainWordEmbedding` function supports `tokenizedDocument` or file input regardless of language. If you have a `tokenizedDocument` array or a file containing your data in the correct format, then you can use this function.

### See Also

`stopWords` | `removeWords` | `normalizeWords` | `bagOfWords` | `bagOfNgrams` | `tokenizedDocument` | `fitlda` | `fitlsa` | `wordcloud` | `addSentenceDetails` | `addLanguageDetails`

### More About

- “Text Data Preparation”
- “Modeling and Prediction”
- “Display and Presentation”
- “Japanese Language Support” on page 4-6
- “Analyze Japanese Text Data” on page 4-11
- “German Language Support” on page 4-21
- “Analyze German Text Data” on page 4-26



# Glossary

---

## Text Analytics Glossary

This section provides a list of terms used in text analytics.

### Documents and Tokens

Term	Definition	More Information
Bigram	Two tokens in succession. For example, ["New" "York"].	bagOfNgrams
Complex token	A token with complex structure. For example, an email address or a hash tag.	tokenDetails
Context	Tokens or characters that surround a given token.	context
Corpus	A collection of documents.	tokenizedDocument
Document	A single observation of text data. For example, a report, a tweet, or an article.	tokenizedDocument
Grapheme	A human readable character. A grapheme can consist of multiple Unicode code points. For example, "a", "□□" or "語".	splitGraphemes
N-gram	<i>N</i> tokens in succession.	bagOfNgrams
Part of speech	Categories of words used in grammatical structure. For example, "noun", "verb", and "adjective".	addPartOfSpeechDetails
Token	A string of characters representing a unit of text data, also known as a "unigram". For example, a word, number, or email address.	tokenizedDocument
Token details	Information about the token. For example, type, language, or part-of-speech details.	tokenDetails
Token types	The category of the token. For example, "letters", "punctuation", or "email address".	tokenDetails
Tokenized document	A document split into tokens.	tokenizedDocument
Trigram	Three tokens in succession. For example, ["The" "United" "States"]	bagOfNgrams
Vocabulary	Unique words or tokens in a corpus or model.	tokenizedDocument

## Preprocessing

Term	Definition	More Information
Normalize	Reduce words to a root form. For example, reduce the word "walking" to "walk" using stemming or lemmatization.	normalizeWords
Lemmatize	Reduce words to a dictionary word (the lemma form). For example, reduce the words "running" and "ran" to "run".	normalizeWords
Stem	Reduce words by removing inflections. The reduced word is not necessarily a real word. For example, the Porter stemmer reduces the words "happy" and "happiest" to "happi".	normalizeWords
Stop words	Words commonly removed before analysis. For example "and", "of", and "the".	removeStopWords

## Modeling and Prediction

### Bag-of-Words

Term	Definition	More Information
Bag-of-n-grams model	A model that records the number of times that n-grams appear in each document of a corpus.	bagOfNgrams
Bag-of-words model	A model that records the number of times that words appear in each document of a collection.	bagOfWords
Term frequency count matrix	A matrix of the frequency counts of words occurring in a collection of documents corresponding to a given vocabulary. This matrix is the underlying data of a bag-of-words model.	bagOfWords
Term Frequency-Inverse Document Frequency (tf-idf) matrix	A statistical measure based on the word frequency counts in documents and the proportion of documents containing the words in the corpus.	tfidf

**Latent Dirichlet Allocation**

<b>Term</b>	<b>Definition</b>	<b>More Information</b>
Corpus topic probabilities	The probabilities of observing each topic in the corpus used to fit the LDA model.	ldaModel
Document topic probabilities	The probabilities of observing each topic in each document used to fit the LDA model. Equivalently, the topic mixtures of the training documents.	ldaModel
Latent Dirichlet allocation (LDA)	A generative statistical topic model that infers topic probabilities in documents and word probabilities in topics.	fitlda
Perplexity	A statistical measure of how well a model describes the given data. A lower perplexity indicates a better fit.	logp
Topic	A distribution of words, characterized by the "topic word probabilities".	ldaModel
Topic concentration	The concentration parameter of the underlying Dirichlet distribution of the corpus topics mixtures.	ldaModel
Topic mixture	The probabilities of topics in a given document.	transform
Topic word probabilities	The probabilities of words in a given topic.	ldaModel
Word concentration	The concentration parameter of the underlying Dirichlet distribution of the topics.	ldaModel

**Latent Semantic Analysis**

<b>Term</b>	<b>Definition</b>	<b>More Information</b>
Component weights	The singular values of the decomposition, squared.	lsaModel
Document scores	The score vectors in lower dimensional space of the documents used to fit the LSA model.	transform
Latent semantic analysis (LSA)	A dimension reducing technique based on principal component analysis (PCA).	fitlsa



Term	Definition	More Information
Word scores	The scores of each word in each component of the LSA model.	<code>lsaModel</code>

### Word Embeddings

Term	Definition	More Information
Word embedding	A model, popularized by the <code>word2vec</code> , <code>GloVe</code> , and <code>fastText</code> libraries, that maps words in a vocabulary to real vectors.	<code>wordEmbedding</code>
Word embedding layer	A deep learning network layer that learns a word embedding during training.	<code>wordEmbeddingLayer</code>
Word encoding	A model that maps words to numeric indices.	<code>wordEncoding</code>

### Visualization

Term	Definition	More Information
Text scatter plot	A scatter plot with words plotted at specified coordinates instead of markers.	<code>textscatter</code>
Word cloud	A chart that displays words with sizes corresponding to numeric data, usually frequency counts.	<code>wordcloud</code>

### See Also

`tokenizedDocument` | `tokenDetails` | `addPartOfSpeechDetails` | `removeStopWords` | `normalizeWords` | `bagOfWords` | `fitlda` | `wordEmbedding` | `wordEncoding` | `wordEmbeddingLayer` | `bagOfNgrams` | `wordcloud` | `textscatter`

### More About

- “Try Text Analytics in 10 Lines of Code”
- “Import Text Data into MATLAB”
- “Create Simple Preprocessing Function”
- “Get Started with Topic Modeling”
- “Visualize Text Data Using Word Clouds” on page 3-2

