

Grzegorz Zadora,¹ Ph.D.

Classification of Glass Fragments Based on Elemental Composition and Refractive Index*

ABSTRACT: The aim of this study was to assess the efficiency of likelihood ratio (LR)-based measures when they are applied to solving various classification problems for glass objects which are described by elemental composition, and refractive index (RI) values, and compare LR-based methods to other classification methods such as support vector machines (SVM) and naïve Bayes classifiers (NBC). One hundred and fifty-three glass objects (23 building windows, 25 bulbs, 32 car windows, 57 containers, and 16 headlamps) were analyzed by scanning electron microscopy coupled with an energy dispersive X-ray spectrometer. Refractive indices for building and car windows were measured before (RI_b), and after (RI_a) an annealing process. The proposed scheme for glass fragment(s) classification demonstrates some efficiency, although the classification of car windows (*c*) and building windows (*w*) must be treated carefully. This is because of their very similar elemental content. However, a combination of elemental content and information on the change in RI during annealing ($\Delta RI = RI_a - RI_b$) gave very promising results. A LR model for the classification of glass fragments into use-type categories for forensic purposes gives slightly higher misclassification rates than SVM and NBC. However, the observed differences between results obtained by all three approaches were very similar, especially when applied to the car window and building window classification problem. Therefore, the LR model can be recommended because of the ease of interpretation of LR-based measures of certainty.

KEYWORDS: forensic sciences, glass, Naïve Bayes Classifiers, support vector machines, likelihood ratio, SEM-EDX, GRIM

Glass is a material frequently present at the scene of such events like car accidents, burglaries, assaults and other accidents, and criminal offenses. Very small glass fragments (linear dimension generally lower than ca. 0.5 mm) that arise during these events can be trapped by, and persist upon the clothes, shoes, and hair, of participants in these activities (1,2).

One of the problems of forensic glass analysis is the comparison problem. This aims to solve the question as to whether two specimens of glass, for example a glass fragment recovered from the clothes of a suspect, and a glass fragment collected from the scene of a crime, might have originated from the same glass object. The comparison problem has been dealt within the forensic literature (3–6). The other task related to glass analysis, and frequently addressed, is a determination of use-type category of a glass fragment, for example, was the fragment from a window, or from some glass container such as a bottle? This process is also called a classification problem, and becomes especially important when there may be no control sample found at the scene of crime. Some knowledge of the type of glass could help investigators focus their search for appropriate control materials.

Most glass fragments analyzed by forensic experts are smaller than 0.5 mm, so the analysis of morphological features, such as thickness and color, are of no value for the assignment of use category. Thus, it is necessary to employ the glass fragments physico-chemical features such as refractive index (RI), and elemental composition, to base judgments of use category.

When choosing an analytical method for glass analysis for forensic purposes, one should take into account the fact that the amount of material is usually very small. So, the method chosen should be non-destructive leaving the material available for re-use. The Glass

Refractive Index Measurement (GRIM) method is often used (7–9), and scanning electron microscopy coupled with an energy dispersive X-ray spectrometer (SEM-EDX), are routinely used in many forensic institutes for the investigation of glass, and other forensic problems (5,6,10–12). Other methods of elemental analysis of glass fragments are μ -X-ray fluorescence (13) and laser ablation inductively coupled plasma–mass spectrometry (LA-ICP-MS) (14). However, these methods require relatively large fragments of glass, for example: LA-ICP-MS gives good results with pieces of glass greater than 0.5 mm. SEM-EDX has the drawback that it can only provide information about major and minor elements, such as O, Na, Al, Mg, Si, K, Ca, Fe, from any glass fragment. Trace elements exist in concentrations below the detection limits of this method. It is commonly believed that trace element concentrations are essential to enable the glass investigator to effectively compare and individualize glass evidence. However, it has been shown that some headway can be made on the basis of the major and minor element concentrations (e.g., 4–6, 10, 11). Given this, is it possible to also solve the problem of classification on the basis of the result of analyses performed by the SEM-EDX method, and in the absence of information for trace elements?

Some classification work has already been conducted using physico-chemical data obtained by the analytical methods usually used in the field of glass analysis (10,11,13,15–18). For example, a non-statistical approach attempted for glass objects used SEM-EDX data, and permitted the correct classification of glass fragments which came from categories having special physico-chemical features, and which therefore also possessed special and distinct elemental compositions such as those of lighting bulbs (10). Glass objects from light bulbs were found to contain more potassium and/or barium than glass objects from other categories. Oxides of these elements improve the optical properties of these glass objects. Attempts to apply cluster analysis for solving the classification problems of glass fragments have been described (10,15,17,18). The efficiency of the proposed approaches was generally

¹Institute of Forensic Research, Westerplatte 9, 31-033 Krakow, Poland.

* The research was partially supported by a grant of the State Committee for Scientific Research, Poland (Project 0T00C 013 26 and 0T00C 028 29).

Received 21 May 2007; and in revised form 1 Dec. 2007; accepted 30 May 2008.

satisfactory. However, the number of incorrectly classified objects was occasionally high.

The determination of use-type category for glass fragments based upon the concentrations of the major and minor elements might be approached by the application of different types of classification methods (11) such as support vector machines (SVM) and naïve Bayes classifiers (NBC). Some form of decision theory is necessary to place any specific glass fragment into classes H_1 or H_2 from the posterior probabilities, $p(H_1|E)$ and $p(H_2|E)$. These probabilities might be calculated from the statistical properties of the physico-chemical observations, which can be termed E , or evidence. In the forensic sphere, posterior probabilities are the province of the court and justice system. The realm of the forensic expert is to evaluate evidence (E). This can be in the context of two propositions, which here are notated H_1 and H_2 . So, it could be suggested that the role of the forensic scientist is the estimation of the conditional probabilities $p(E|H_1)$ and $p(E|H_2)$. The likelihood ratio (LR):

$$LR = \frac{p(E|H_1)}{p(E|H_2)} \quad (1)$$

is a well-documented measure of evidence value in the forensic field (3–6), and is increasingly employed as a measure of evidential strength in the forensic comparison problem. The LR compares the probability of the measurements (E – physico-chemical features determined for compared glass samples) in the light of the assumption of a common source for these samples, termed the prosecution proposition (H_1), with the probability of the measurements (E) assuming different sources for these samples, termed the defense proposition (H_2).

In this paper, an algorithm of classification of glass fragments has been applied to a set of 153 samples, each originating from one of five known glass categories. The categories were: car windows (c), building windows (w), containers (p), car headlamps (h), and light bulbs (b). The first problem attempted was a classification of each of the glass objects into a cw category (of car and building windows) and a bhp category (bulbs, headlamps, and containers). Glass objects from c and w categories have very similar elemental compositions because they are manufactured in a very similar way (a float glass manufacturing method). Glass objects from categories b , h , and p represent other types of glass commonly met in forensic practice. Then, a refinement to the original cw/bhp classification was attempted where if a glass sample was classified to the bhp category an attempt was made to further classify it into bh or p category. It could be expected (10) that bulbs and car headlamps (bh) could have systematically different elemental composition to the containers glass (p —jars and bottles) because of their optical properties. Again, were a glass specimen to be classified as belonging to the bh category, then this could be further classified to b or h category.

Finally, an attempt to further classify the cw class into either car glass, or building glass was made. It has been found (11) that it could be done sometimes on the basis of major and minor elemental composition; however, classification is poor. So, a new feature has been examined to aid in this classification task, i.e., RI analyzed before annealing process (RI_b) or ΔRI (a difference between RI measured after annealing process [RI_a] and RI_b). Samples of glass were annealed at about 600°C (8,19,20), which removes tensions in glass. This process has already been applied to car window glass. It should be mentioned that simultaneous analysis of a glass fragment by SEM-EDX and GRIM techniques is only possible if a glass fragment is relatively large (much bigger than

0.5 mm). Then, it is possible to divide it into smaller parts which could be used separately for SEM-EDX and GRIM analysis (with or without annealing process).

Materials and Methods

SEM-EDX Analysis of Glass Fragments

One large piece of glass from each of 153 glass objects (23 building windows, 25 bulbs, 32 car windows, 57 containers, and 16 headlamps) was selected. Each of these pieces was wrapped in a sheet of gray paper and further fragmented. Four glass fragments, of linear dimension less than 0.5 mm with surfaces as smooth and flat as possible, were placed on self-adhesive carbon tabs on an aluminum stub and then carbon coated using an SCD sputter (Bal-Tech, Balzers, Liechtenstein). Analysis of the elemental content of each glass fragment was carried out using a scanning electron microscope (JSM-5800; Jeol, Tokyo, Japan), with an energy dispersive X-ray spectrometer (Link ISIS 300; Oxford Instruments Ltd., Witney, Oxfordshire, U.K.). Three replicate measurements were taken from different areas on each of the four fragments. The measurement conditions were: accelerating voltages 20 kV, life time 50 sec, magnification 1000–2000 \times . The calibration element was cobalt. The SEMQUANT option (part of the software LINK ISIS; Oxford Instruments Ltd.) was used in the process of determining the weight percentage of particular elements in a fragment. The option applied a ZAF correction procedure, which takes into account corrections for the effects of difference in the atomic number (Z), absorption (A), and X-ray fluorescence (F). The selected analytical conditions allowed the determination of all elements except lithium (Li) and boron (B). However, as only the concentrations of the main elements could be determined by SEM-EDX, only oxygen (O), sodium (Na), magnesium (Mg), aluminum (Al), silicon (Si), potassium (K), calcium (Ca), and iron (Fe) are further considered in this paper.

GRIM Analysis of Glass Fragments

The RI was determined by the thermo-immersion method, using the GRIM 2 system made by Foster & Freeman (Evesham, Worcestershire, U.K.), at 589 nm. Glass fragments from each of the 55 objects from the car and building window category were mounted onto separate clean microscope slides. Each glass fragment was covered with silicone oil (silicone oil B; Locke Scientific, Southwick, Fareham, U.K.), covered with a cover slip. The match temperature (MT), that is, the temperature when the RI of the immersion oil and the RI of the glass fragment are the equal, was determined, and the value of the RI was determined automatically from the calibration model. The calibration model ($RI = -3.74 \times 10^{-4} MT + 1.54491$) was calculated earlier from measurements of glass standards (Locke Scientific). The RI was measured five times only from the bulk, or the non-float surfaces, because the RI of the float surface differs from the other surfaces.

Glass fragments of length ca. 0.5 mm, were put into a metal sample holder, which ensured a uniform and reproducible thermal environment, and annealed in a muffle furnace (Nabertherm L3; Nabertherm GmbH, Lilienthal, Germany). The fragments were then heated to a temperature of 550°C, then cooled to 480°C at a rate of 15°C/h. The glass fragments were then cooled freely to room temperature. The relatively short schedule of annealing was long enough to reveal a comparatively large difference in RI of glass fragments.

SVM

SVM (21–30) have been considered by many researchers to be the best performing discriminative classifiers. To visualize how SVM work, we shall here consider only two variables; however, all the concepts and derivations can be extended to higher dimensional problems. The simplest situation is where two classes are linearly separable. This means that we can find a line (a plane in three dimensions, a hyper-plane in higher dimensions) that perfectly separates the two classes.

The example of Fig. 1a shows that theoretically many different solutions (separating lines) could be chosen, but an SVM selects an optimal boundary line that maximizes the distance between the classes and the objects at the border (Fig. 1b). This distance is known as the margin, while the objects at the border are called support vectors (SVs) because the optimization task involving the maximization of the margin depends only on their displacement in the dataspace.

When classifying a new object y , the decision function will simply be the equation of a separating hyper-plane explicitly expressed in terms of the SVs:

$$z = \text{sgn} \left(\sum_{i=1} \alpha_i z_i s_i^T y + b \right) \quad (2)$$

where z (equal to ± 1) is the label assigned to the object y that we wish to classify, s_i is a generic SV, z_i its label, the coefficients α_i , and b is the offset parameter found during a learning process. The sign of the function determines the class membership of y , while its absolute value represents a confidence level, because the higher this is, the stronger the SVM believes one object to belong to one class.

In the case of the non-linearly separable cases (Fig. 2a), the optimal separating hyper-plane is searched in a higher dimensional space (called a feature space) than a dimension of input space (Fig. 2). The objects are projected by means of the feature function. The optimal separating hyper-plane found in this way will have a form of a curved boundary of certain complexity when this will be projected from the so-called feature space to the input space (Fig. 2c). In practice, more than one dimension is added to the input space (Fig. 2a) when the feature space is created because there is a simple rule—the more dimensions added, the easier to find a separating hyper-plane. The kernel functions are the most often used in SVM like the feature functions, e.g., the radial basis function (RBF) (see Eq. [10]).

NBC

NBC (31–34) have been employed to address various problems in the forensic sciences (3). NBC are generative, meaning that they

represent a joint probability distribution over the data and the labels. NBC are called naïve because they make the assumption that variables are generated independently of the others. This assumption is almost always false; however, NBC are robust and generally perform well (11). The probability model for classifiers is a conditional model $p(H|X_1, \dots, X_n)$ and the application of the Bayes theory allows us to write:

$$p(H|X_1, \dots, X_n) = \frac{p(H)p(X_1, \dots, X_n|H)}{p(X_1, \dots, X_n)} \quad (3)$$

where H is the considered category, X_1, \dots, X_n — n variables which describe i th object presented in a tuning set.

The following assumption is used for the joint probability model presented in the denominator—each feature X_i is conditionally independent of every other feature X_j for $j \neq i$, that is,

$$p(H|X_1, \dots, X_n) = \frac{1}{Z} p(H) \prod_{i=1}^n p(X_i|H) \quad (4)$$

where Z is the scaling factor dependent only on X_1, \dots, X_n . This is constant if the values of the variables are known.

The learning task requires only estimation of $p(X_i|H)$ and $p(H)$ on the basis of the tuning data and it could be carried out independently because of the assumption that variables are independent. This assumption reduces a n -dimensional task to n one-dimensional tasks. Kernel density estimation is used the most often in the aim of estimation of probability density functions.

The classification of the new object y characterized by determined values of particular variables (y_1, \dots, y_n) is based on the following decision rule—the new object belongs to the class for which the hypothesis is most probably what is known as the maximum a posteriori decision rule:

$$\text{classify}(y_1, \dots, y_n) = \arg \max_c p(H) \prod_{i=1}^n p(y_i|H) \quad (5)$$

LR Approach

In this paper, a multivariate model based on the ideas which created the model for comparison problem published in (4,6) was used. The model for comparison problem proposed by authors of these papers considered two sources of variability (within glass object and between glass objects). SVM and NBC models only allow to consider a between-object variability. Therefore, the LR model which considers only this source of variability is proposed for LR calculations in the aim of determination of use-type categories of a glass fragment described by a vector \bar{y} of n variables:

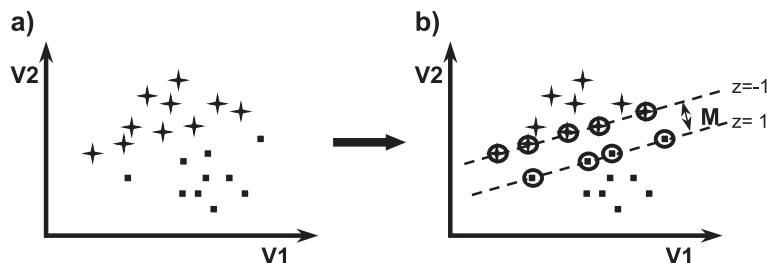


FIG. 1—A graphical illustration of support vector machines—a linearly separable case. M —the margin between two different categories. Objects (squares or stars) in circles are support vectors.

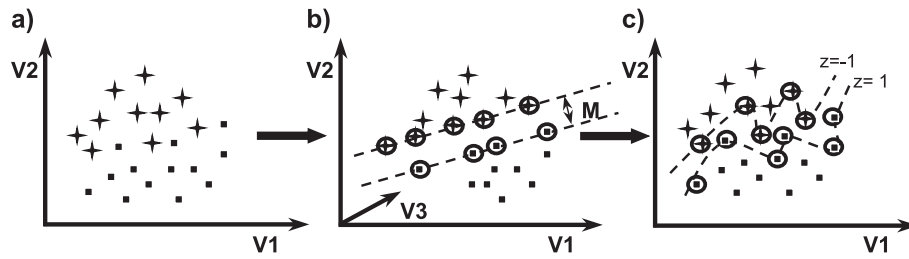


FIG. 2—A graphical illustration of support vector machines—a non-linearly separable case; (a, c) an input space, (b) a feature space. M —the margin between two different categories. Objects (squares or stars) in circles are support vectors.

$$LR = \frac{p(E|H_1)}{p(E|H_2)} = \frac{\int f(\bar{y}|\mu_1, C_1)d\mu_1}{\int f(\bar{y}|\mu_2, C_2)d\mu_2} \quad (6)$$

where H_1 —the object comes from class 1, and H_2 —the object comes from class 2, \bar{y} —vector of means of considered n variables calculated for a glass object on the base of all performed measurements, that is, 12 made by SEM-EDX and five made by GRIM technique, μ_1, μ_2 —vectors of means of considered n variables within calculated for populations considered as category H_1 and H_2 , C_1, C_2 —variance-covariance matrices (between-objects variability) calculated for considered variables which describe glass objects presented in data sets within class 1 and class 2.

The between-object distribution was modeled by a multivariate kernel density. It could be shown (Eq. [6]) that:

$$\int f(\bar{y}|\mu_1, C_1)d\mu_1 = (2\pi)^{-k/2} |h^2 C_1|^{-1/2} \frac{1}{m_1} \sum_{i=1}^{m_1} \exp\left\{-\frac{1}{2}(\bar{y}_1 - \bar{x}_{1i})^T [(h^2 C_1)]^{-1} (\bar{y}_1 - \bar{x}_{1i})\right\} \quad (7)$$

and similar for $\int f(\bar{y}|\mu_2, C_2)d\mu_2$ which gives:

$$LR = \frac{|h^2 C_1|^{-1/2} \frac{1}{m_1} \sum_{i=1}^{m_1} \exp\left\{-\frac{1}{2}(\bar{y} - \bar{x}_{1i})^T [(h^2 C_1)]^{-1} (\bar{y} - \bar{x}_{1i})\right\}}{|h^2 C_2|^{-1/2} \frac{1}{m_2} \sum_{i=1}^{m_2} \exp\left\{-\frac{1}{2}(\bar{y} - \bar{x}_{2i})^T [(h^2 C_2)]^{-1} (\bar{y} - \bar{x}_{2i})\right\}} \quad (8)$$

where $\bar{x}_{1i}, \bar{x}_{2i}$ —vectors of means of n variables calculated on the basis of all 12 or five performed measurements for an i th glass object from a population with categories H_1 and H_2 . h —a smoothing parameter $h = h_{opt} = \left(\frac{4}{2n+1}\right)^{\frac{1}{m+4}} \frac{1}{m+4}$, m_1, m_2 —numbers of samples in sets of background data used for estimation of probability density functions in nominator and denominator.

The decision rules are that values of LR above 1 support H_1 , and values of LR below 1 support H_2 . A value of LR close to 1 provides little support for either proposition. Also the larger, or lower, the value of the LR, the stronger the support of E for H_1 or H_2 .

Calculation of a full model, which takes into account all variables, requires the estimation of the probability density function which takes into account all variables under each of two propositions, H_1 and H_2 . As here glass fragments were described by at least seven variables, then it is necessary to reliably estimate seven means, seven variances, and 21 covariance, which is difficult from a sample of only 153. An approach based on graph theory has been used to factorize the joint density function into

the product of several density functions on lower dimensions which allow the parameters to be estimated with the appropriate levels of reliability (8). It can be shown that the elements of the scaled inverse correlation matrix are the negative partial correlation coefficients, and that values of partial correlation can be used to construct a decomposable graphical model of the full density into cliques representing product of several density functions in lower dimensions. The relationships between the elements in glass are not causal, and so the graphs used are undirected. The graphical model was selected by the sequential addition of edges decided by inspection of the partial correlation matrix. First, the largest magnitude partial correlation was selected, and an edge was added between the two nodes connected by this partial correlation. This process was repeated until all nodes are part of the graphical model. The factorization of the full model is given by:

$$f(C_i|S_i) = \frac{f(C_i)}{f(S_i)} \quad (9)$$

where C_i is the i th clique in the model, S_i —is the set of all separators for the i th clique calculated from a set chain of the cliques for the model.

A subset of variables in which all the nodes are connected to each other is known as a complete subgraph, and the corresponding subset of variables is known as a clique. To find a set chain, which is a particular ordering of the cliques in the model, the following algorithm was applied to the collection of cliques: select a node arbitrarily from the model graph and denote this as the lowest numbered node, number each remaining node in turn ordered by the number of edges linking it to any other already numbered node; break ties arbitrarily, assign a rank to each clique based upon the highest numbered node in the clique, if two cliques share a highest numbered node then rank arbitrary between the two nodes. Given the cliques for the model (C_i), and a suitable set chain, the sets of separators (S_i) for each clique is found. The first clique in the set chain is always a complete subgraph, and there are no separator sets. After that, the next clique presented in the set chain is added to the model. The intersection of elements between these two cliques becomes the first separator set. The process is continued until all cliques are joined to the model. A practical example is presented in the Results section.

Software

Functions written in R (35) were applied. The package *e1071* was used for data analysis by SVM and a package *klaR* was used for data analysis by NBC. Routines for LR calculations were written by the author based on functions developed previously (4,6).

Results

Descriptive Statistics

Reports from SEMQUANT software contain only data about selected elements. These were O, Na, Mg, Al, Si, K, Ca, and Fe; it was assumed that a sum of concentration of these elements was equal to 100 wt %. From the eight elements measured, seven independent variables were derived by taking the log₁₀ of oxygen concentration to concentration of each element. The following abbreviations of such ratio were used in this paper: Na', Mg', Al', Si', K', Ca', and Fe'. This normalization effectively removes stochastic fluctuations in instrumental measurement. In the case, when an element was not present, or its concentration was below a SEM-EDX detection limit (0.1 wt % for most of considered elements), then zero was substituted by a very small value (0.0001 wt %).

The mean of each of seven variables within each glass object was calculated from 12 replicate measurements (three measurements were made on each of four glass fragments collected from each glass object). This resulted in a 153 × 7 data matrix. For glass objects from *c* and *w* categories additional features were analyzed. These were RI_b and ΔRI. Each glass object was measured five times in order to determine the RI before (RI_b), and after the annealing process (RI_a). Mean values for each glass item for RI_b and RI_a were calculated, and from that ΔRI was calculated as ΔRI = RI_a - RI_b. It is reasonable to expect that variables used in a classification task should have values of similar order of magnitude, to prevent a variable having undue influence on the results of the classification problem than another variable by having systematically much larger values. The ΔRI values were three orders of magnitude lower than other considered variables. Therefore, a log transform of ΔRI was used during this analysis (|log₁₀(ΔRI)|). Descriptive statistics of RI_b and ΔRI could be found in Table 1.

Tuning and Test Set Creation

Tuning data sets which were necessary at learning process when NBC and SVM classifiers were being used, contained 75% of the

glass objects from the database (Table 2). The remaining glass objects formed a test set. Glass objects to form the tuning and test sets were chosen randomly and uniformly with the assumption that the proportion of glass objects from each glass category within the tuning and testing sets were the same as in the original database (Table 2), that is,

- (i) *cw* versus *bhp*—36% glass objects from *cw* category and 64% from *bhp* category;
- (ii) *bh* versus *p*—42% glass objects from *bh* category and 58% from *p* category;
- (iii) *b* versus *h*—61% glass objects from *b* category and 39% from *h* category;
- (iv) *c* versus *w* (all three combinations of variables)—58% glass objects from *c* category and 42% from *w* category.

Ten different tuning and test sets were created and employed. Tuning sets were also used for determination of parameters of LR models (i.e., variance-covariance matrices and graphical model).

Classification of Glass Objects by SVM

A C-classification approach with the most common used kernel function, the RBF kernel, was used. This is defined by the equation:

$$K(x_i, x_j) = \exp\left(-\sigma \|x_i - x_j\|^2\right) \tag{10}$$

for $\sigma > 0$, where x_i, x_j are two generic training objects in the data set. This function requires only the tuning of a single parameter, the radial width σ . It was also necessary to select the most suitable value of penalty error *C*. A tuning procedure was conducted separately on each of the 10 tuning sets, for each of the classification problems, to try to estimate a suitable range for \hat{C} and $\hat{\sigma}$ values. The ranges of \hat{C} and $\hat{\sigma}$ were selected by observation of suitable tuning plots such as those given as Fig. 3. It was found that the best combination of \hat{C} and $\hat{\sigma}$ were:

- (i) $\hat{\sigma} \in \langle 0.01; 0.30 \rangle$;

TABLE 1—Descriptive statistics of each variable in each considered glass category.

| Category | Parameter | Variable* | | | | | | | | |
|----------|-----------|-----------|-------|-------|-------|-------|-------|-------|-------|-------|
| | | Na' | Mg' | Al' | Si' | K' | Ca' | Fe' | RI | dRI |
| <i>b</i> | \bar{x} | 0.767 | 2.762 | 1.624 | 0.137 | 1.558 | 2.678 | 5.392 | n.a. | n.a. |
| | SD | 0.142 | 1.766 | 0.131 | 0.018 | 0.821 | 2.047 | 0.227 | n.a. | n.a. |
| | Min | 0.603 | 1.326 | 1.370 | 0.105 | 0.721 | 1.020 | 4.849 | n.a. | n.a. |
| | Max | 0.944 | 5.509 | 1.781 | 0.177 | 3.621 | 5.509 | 5.575 | n.a. | n.a. |
| <i>c</i> | \bar{x} | 0.694 | 1.364 | 2.400 | 0.152 | 4.164 | 0.911 | 4.569 | 1.519 | 2.726 |
| | SD | 0.015 | 0.151 | 0.920 | 0.032 | 1.638 | 0.064 | 1.671 | 0.002 | 0.083 |
| | Min | 0.661 | 1.181 | 1.359 | 0.082 | 2.066 | 0.789 | 1.882 | 1.515 | 2.638 |
| | Max | 0.722 | 2.119 | 5.690 | 0.209 | 5.698 | 1.019 | 5.705 | 1.524 | 2.959 |
| <i>h</i> | \bar{x} | 0.700 | 3.863 | 1.890 | 0.148 | 2.022 | 0.970 | 5.574 | n.a. | n.a. |
| | SD | 0.036 | 1.834 | 1.001 | 0.055 | 1.392 | 0.280 | 0.025 | n.a. | n.a. |
| | Min | 0.618 | 1.335 | 1.444 | 0.089 | 1.301 | 0.776 | 5.544 | n.a. | n.a. |
| | Max | 0.732 | 5.579 | 5.596 | 0.285 | 5.575 | 1.955 | 5.632 | n.a. | n.a. |
| <i>p</i> | \bar{x} | 0.711 | 2.023 | 1.827 | 0.183 | 3.098 | 0.908 | 5.101 | n.a. | n.a. |
| | SD | 0.031 | 0.928 | 0.102 | 0.033 | 1.520 | 0.065 | 1.215 | n.a. | n.a. |
| | Min | 0.625 | 1.310 | 1.698 | 0.088 | 1.721 | 0.785 | 2.282 | n.a. | n.a. |
| | Max | 0.769 | 5.705 | 2.328 | 0.227 | 5.712 | 1.111 | 5.718 | n.a. | n.a. |
| <i>w</i> | \bar{x} | 0.693 | 1.363 | 3.257 | 0.139 | 4.770 | 0.905 | 3.481 | 1.520 | 3.031 |
| | SD | 0.018 | 0.168 | 1.585 | 0.033 | 1.322 | 0.077 | 1.596 | 0.002 | 0.218 |
| | Min | 0.658 | 1.180 | 1.846 | 0.080 | 2.144 | 0.817 | 1.896 | 1.515 | 2.301 |
| | Max | 0.722 | 1.983 | 5.682 | 0.235 | 5.719 | 1.126 | 5.719 | 1.525 | 3.523 |

b, bulbs; *c*, car windows; *h*, headlamps; *p*, containers; *w*, building windows; n.a., not applicable.

*Na' is an abbreviation of log₁₀(O/Na) and so on; dRI is an abbreviation of |log₁₀(ΔRI)|.

TABLE 2—Number of samples from each of the categories considered in a particular classification problem.

| Category | No. of Samples | | |
|---|----------------|------------|-------------|
| | Primary Set | Tuning Set | Testing Set |
| Classification into car and building windows or bulbs, headlamps, and containers— <i>cw</i> versus <i>bhp</i> | | | |
| Car and building windows (<i>cw</i>) | 55 | 41 | 14 |
| Bulbs, headlamps, and containers (<i>bhp</i>) | 98 | 74 | 24 |
| Classification into bulbs and headlamps or containers— <i>bh</i> versus <i>p</i> | | | |
| Bulbs and headlamps (<i>bh</i>) | 41 | 31 | 10 |
| Containers (<i>p</i>) | 57 | 43 | 14 |
| Classification into bulbs or headlamps— <i>b</i> versus <i>h</i> | | | |
| Bulbs (<i>b</i>) | 25 | 19 | 6 |
| Headlamps (<i>h</i>) | 16 | 11 | 5 |
| Classification into car or building windows— <i>c</i> versus <i>w</i> | | | |
| Car windows (<i>c</i>) | 32 | 24 | 8 |
| Building windows (<i>w</i>) | 23 | 17 | 6 |

- (ii) $\hat{C} \in \langle 1; 10 \rangle$ —*c* versus *w* classification problem (seven variables, i.e., elemental content only) and *b* versus *h* problem;
 (iii) $\hat{C} \in \langle 10; 100 \rangle$ —other considered classification problems.

Subsequently 10 different values of \hat{C} and $\hat{\sigma}$ were taken for analysis, e.g., when $\hat{C} \in \langle 1; 10 \rangle$ then $\hat{C} = 1, 2, \dots, 10$ and 30 different values when $\hat{\sigma} \in \langle 0.01; 0.30 \rangle$ then $\hat{\sigma} = 0.01, 0.02, \dots, 0.30$. These created 300 different pairs of \hat{C} and $\hat{\sigma}$ within each of the considered classification problems. Next, each combination of \hat{C} and $\hat{\sigma}$ was applied in the classification process of samples within the tuning sets. The accuracy of classification, i.e., the percent of the total number of predictions which were corrected, as well as the number of necessary SVs were taken into account. It was assumed that the best pair of \hat{C}

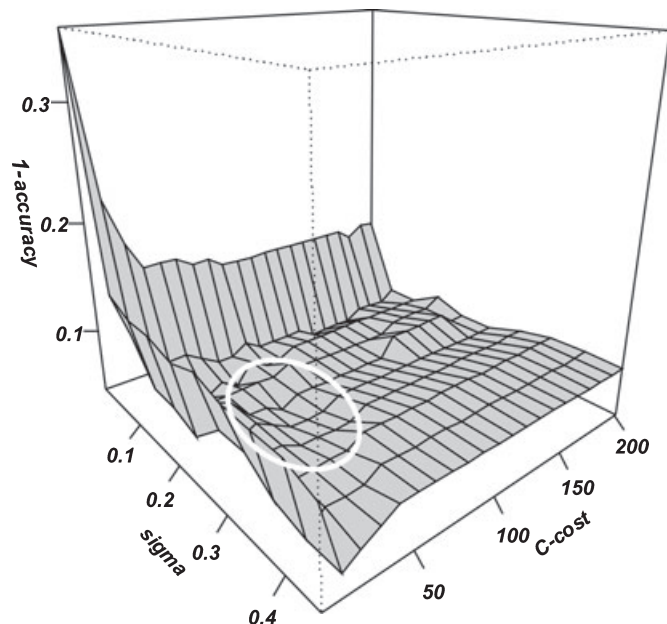


FIG. 3—An example tuning plot obtained from the first tuning set when the *cw* versus *bhp* classification problem was considered. The ellipse indicates the region of best combinations of \hat{C} and $\hat{\sigma}$ which give the best accuracy, i.e., the ratio of the total number of predictions which were correct.

and $\hat{\sigma}$ was the one which gave relatively high accuracy in tuning sets and a relatively small number of SVs necessary to construct a proper hyper-plane. The optimal sets for \hat{C} and $\hat{\sigma}$ were selected on the basis of the analysis of Table 3. These sets were:

- (i) $\hat{C} = 80$ and $\hat{\sigma} = 0.10$ when the problem was a classification into car and building windows or bulbs, headlamps, and containers (*cw* vs. *bhp*);
 (ii) $\hat{C} = 100$ and $\hat{\sigma} = 0.05$ when the problem was a classification into bulbs and headlamps or containers (*bh* vs. *p*);
 (iii) $\hat{C} = 10$ and $\hat{\sigma} = 0.05$ when the problem was a classification into bulbs or headlamps (*b* vs. *h*) and when the problem was a classification into car windows or building windows (*c* vs. *w*) on the basis of results of elemental content;
 (iv) $\hat{C} = 40$ and $\hat{\sigma} = 0.05$ when the problem was a classification into car windows or building windows (*c* vs. *w*) on the basis of results of elemental content and RI_b values or ΔRI values.

Values of accuracy, that is, the proportion of the total number of predictions which were correct, within each of classification problems and each of the test sets obtained by application of SVM as well as NBC and LR approach, are given in Table 4.

Classification of Glass Objects by NBC

The estimation of $p(X_i|H)$ was made using R. A problem with the NBC approach is that knowledge of *a priori* probabilities of objects within each category is required. The estimation of these probabilities for the wider population can be difficult. Here the sample prior probabilities $p(H)$ were used, that is: $p(H = cw) = 0.36$ and $p(H = bhp) = 0.64$ for the classification of fragments into car and building windows categories, or bulbs, headlamps, and containers categories. The analysis of the influence of various values of *a priori* probability $p(H)$ on the level of false positives and negatives, showed (11) that a number of incorrect classifications depend on the assumed value of $p(H)$ only when these values were near to 0 or 1. Therefore, it was considered sufficient to estimate $p(H)$ as a proportion of an object of particular category within test sets, and use those as prior probabilities. Results of analysis performed by NBC model are given in Table 4.

Classification of Glass Objects by LR Approach

The rescaled inverse of the variance-covariance matrices for the considered seven or eight variables in the particular classification problem are given in Table 5. Graphical models were constructed as described earlier. They were obtained on the base of data of all glass objects considered in particular classification problem. The model presented in Fig. 4a represents the final graphs obtained for the *cw* versus *bhp* problem. The cliques found and numbered in the graph are presented in Fig. 4b. The graph was factorized in a way described by (9). The clique (Mg' , Ca') has the highest numbered node Mg' , so the clique is given the same number as Mg' , which is 7. A clique (Ca' , Si') has the next highest numbered node Si' , and so has the number 6. Putting these into numerical order, the set given in Table 6, column 1, is obtained.

Given the cliques for the model, and a suitable set chain, the sets of separators (S_i in Eq. [9]) for each clique was found. The first clique (Table 6) in the set chain is (Na' , Ca'). This is a complete subgraph, and at the moment there are no other cliques added to the graph, so there can be no separator sets. The next clique in the set chain is (Na' , K'), and it is added to the model. The intersection of elements between these two cliques is (Na'), and so this

TABLE 3—Selection of the best combinations of penalty error \hat{C} and radial width $\hat{\sigma}$.

| Set | Classification Problem | | | | | | | | |
|-----|------------------------------------|------------------------|------------------|--|----------------|-----|--|----------------|-----|
| | \hat{C}^* | $\hat{\sigma}^\dagger$ | SVs [‡] | \hat{C} | $\hat{\sigma}$ | SVs | \hat{C} | $\hat{\sigma}$ | SVs |
| | <i>cw</i> versus <i>bhp</i> | | | <i>bh</i> versus <i>p</i> | | | <i>b</i> versus <i>h</i> | | |
| 1 | 80–100 | 0.18–0.22 | 28 | 90–100 | 0.08–0.10 | 11 | 9–10 | 0.09–0.10 | 8 |
| 2 | 80–100 | 0.16–0.22 | 30 | 90–100 | 0.08–0.10 | 15 | 9–10 | 0.05 | 11 |
| 3 | 70–100 | 0.08–0.10 | 29 | 90–100 | 0.10–0.12 | 12 | 8–10 | 0.05–0.08 | 12 |
| 4 | 80–100 | 0.08–0.10 | 25 | 100 | 0.05–0.07 | 13 | 8–10 | 0.05–0.08 | 10 |
| 5 | 80–100 | 0.06–0.16 | 26 | 60–100 | 0.05 | 14 | 10 | 0.07–0.08 | 9 |
| 6 | 70–100 | 0.10 | 28 | 80–100 | 0.03 | 14 | 8–10 | 0.05 | 8 |
| 7 | 80–100 | 0.08–0.16 | 28 | 100 | 0.05–0.10 | 16 | 10 | 0.05 | 8 |
| 8 | 80–100 | 0.08–0.10 | 28 | 100 | 0.01–0.08 | 12 | 7–10 | 0.06–0.07 | 8 |
| 9 | 100 | 0.10–0.22 | 28 | 70–100 | 0.03–0.08 | 19 | 9–10 | 0.05–0.06 | 12 |
| 10 | 80–100 | 0.10–0.12 | 27 | 70–100 | 0.03 | 14 | 7–10 | 0.04–0.05 | 12 |
| | <i>c</i> versus <i>w</i> —elements | | | <i>c</i> versus <i>w</i> —elements + RI _b | | | <i>c</i> versus <i>w</i> —elements + dRI | | |
| 1 | 7–8 | 0.05–0.10 | 23 | 40–50 | 0.05–0.10 | 20 | 50 | 0.05 | 21 |
| 2 | 7–10 | 0.05–0.10 | 20 | 25–50 | 0.05 | 22 | 25–50 | 0.05 | 18 |
| 3 | 8–10 | 0.10 | 23 | 20 | 0.05 | 24 | 50 | 0.05 | 20 |
| 4 | 9–10 | 0.05 | 24 | 40–50 | 0.10 | 26 | 35–50 | 0.05 | 17 |
| 5 | 10 | 0.05 | 25 | 50 | 0.10 | 24 | 40–50 | 0.05 | 19 |
| 6 | 6–10 | 0.05 | 21 | 40 | 0.05 | 20 | 40–50 | 0.05–0.10 | 21 |
| 7 | 8–10 | 0.05–0.10 | 26 | 35–50 | 0.05 | 24 | 40–50 | 0.05 | 19 |
| 8 | 10 | 0.05 | 22 | 40 | 0.05–0.10 | 26 | 30–50 | 0.05 | 18 |
| 9 | 10 | 0.10 | 25 | 40–50 | 0.05 | 23 | 40–50 | 0.05 | 21 |
| 10 | 9–10 | 0.05–0.10 | 24 | 40–50 | 0.05–0.10 | 25 | 20–50 | 0.05 | 17 |

* \hat{C} — estimated range of penalty error C.
[†] $\hat{\sigma}$ — estimated radial width (see Eq. [10]).
[‡]SV—a number of support vectors necessary to construct a proper hyper-plane.

TABLE 4—Values of accuracy—a comparison of results obtained by application of SVM, NBC, and LR models within each of the 10 test sets within each classification problem.

| Set | Classification Problem | | | | | | | | |
|-----|---|-----|----|---|-----|----|---|-----|-----|
| | SVM | NBC | LR | SVM | NBC | LR | SVM | NBC | LR |
| | <i>cw</i> versus <i>bhp</i> | | | <i>bh</i> versus <i>p</i> | | | <i>b</i> versus <i>h</i> | | |
| 1 | 90* | 90 | 79 | 100 | 92 | 67 | 73 | 73 | 73 |
| 2 | 95 | 95 | 84 | 96 | 100 | 88 | 91 | 100 | 91 |
| 3 | 92 | 92 | 79 | 88 | 92 | 63 | 100 | 100 | 82 |
| 4 | 95 | 90 | 79 | 88 | 88 | 63 | 91 | 91 | 82 |
| 5 | 87 | 90 | 79 | 96 | 100 | 71 | 100 | 100 | 82 |
| 6 | 87 | 95 | 82 | 92 | 92 | 71 | 82 | 91 | 82 |
| 7 | 95 | 92 | 71 | 92 | 92 | 71 | 100 | 73 | 100 |
| 8 | 95 | 82 | 87 | 96 | 92 | 92 | 100 | 82 | 91 |
| 9 | 90 | 97 | 84 | 100 | 96 | 83 | 100 | 100 | 82 |
| 10 | 97 | 92 | 71 | 99 | 96 | 67 | 100 | 100 | 82 |
| | <i>c</i> versus <i>w</i> —elemental content | | | <i>c</i> versus <i>w</i> —elemental content + RI _b | | | <i>c</i> versus <i>w</i> —elemental content + dRI | | |
| 1 | 64 | 64 | 64 | 64 | 57 | 64 | 93 | 79 | 93 |
| 2 | 64 | 57 | 43 | 57 | 57 | 64 | 93 | 71 | 71 |
| 3 | 71 | 79 | 43 | 64 | 57 | 50 | 93 | 86 | 86 |
| 4 | 64 | 86 | 43 | 64 | 79 | 71 | 93 | 86 | 86 |
| 5 | 71 | 79 | 57 | 64 | 71 | 57 | 100 | 93 | 93 |
| 6 | 57 | 71 | 57 | 57 | 64 | 50 | 86 | 93 | 71 |
| 7 | 71 | 64 | 64 | 64 | 71 | 71 | 100 | 86 | 100 |
| 8 | 71 | 71 | 50 | 57 | 64 | 64 | 100 | 86 | 100 |
| 9 | 79 | 71 | 86 | 79 | 71 | 57 | 93 | 93 | 93 |
| 10 | 64 | 86 | 57 | 71 | 71 | 43 | 86 | 79 | 64 |

SVM, support vector machine; NBC, naïve Bayes classifiers; LR, likelihood ratio.
 *Accuracy—i.e., the percent of the total number of predictions which were corrected.

TABLE 5—The rescaled inverse of the variance-covariance matrices for the seven or eight variables from the glass samples in each of six considered problems (only the upper right triangle of the matrix is given, the lower left triangle is given by symmetry).

| | Na' | Mg' | Al' | Si' | K' | Ca' | Fe' | RI _b /dRI |
|--|-------|--------|--------|--------|--------|--------|--------|----------------------|
| (a) <i>cw</i> versus <i>bph</i> | | | | | | | | |
| Na' | 1.000 | -0.038 | -0.108 | -0.182 | 0.261 | -0.610 | 0.015 | |
| Mg' | | 1.000 | 0.124 | 0.024 | 0.053 | -0.355 | -0.114 | |
| Al' | | | 1.000 | 0.208 | -0.406 | 0.003 | 0.345 | |
| Si' | | | | 1.000 | -0.037 | 0.212 | 0.057 | |
| K' | | | | | 1.000 | 0.046 | -0.005 | |
| Ca' | | | | | | 1.000 | -0.013 | |
| Fe' | | | | | | | 1.000 | |
| (b) <i>bh</i> versus <i>p</i> | | | | | | | | |
| Na' | 1.000 | -0.041 | 0.045 | -0.215 | 0.324 | -0.595 | -0.101 | |
| Mg' | | 1.000 | 0.121 | 0.095 | -0.056 | -0.347 | -0.067 | |
| Al' | | | 1.000 | -0.087 | -0.073 | -0.046 | -0.077 | |
| Si' | | | | 1.000 | -0.131 | 0.227 | 0.226 | |
| K' | | | | | 1.000 | 0.071 | -0.126 | |
| Ca' | | | | | | 1.000 | 0.035 | |
| Fe' | | | | | | | 1.000 | |
| (c) <i>b</i> versus <i>p</i> | | | | | | | | |
| Na' | 1.000 | -0.048 | 0.090 | 0.116 | 0.560 | -0.581 | -0.383 | |
| Mg' | | 1.000 | 0.187 | 0.063 | -0.100 | -0.352 | -0.345 | |
| Al' | | | 1.000 | -0.115 | 0.049 | -0.088 | -0.211 | |
| Si' | | | | 1.000 | 0.096 | -0.027 | -0.097 | |
| K' | | | | | 1.000 | -0.017 | -0.272 | |
| Ca' | | | | | | 1.000 | 0.273 | |
| Fe' | | | | | | | 1.000 | |
| (d) <i>c</i> versus <i>w</i> (seven variables; elemental content) | | | | | | | | |
| Na' | 1.000 | 0.032 | -0.127 | -0.414 | 0.142 | 0.231 | 0.273 | |
| Mg' | | 1.000 | 0.077 | -0.385 | -0.034 | 0.285 | -0.114 | |
| Al' | | | 1.000 | 0.132 | -0.466 | 0.025 | 0.314 | |
| Si' | | | | 1.000 | -0.167 | -0.830 | 0.099 | |
| K' | | | | | 1.000 | 0.105 | -0.088 | |
| Ca' | | | | | | 1.000 | -0.152 | |
| Fe' | | | | | | | 1.000 | |
| (e) <i>c</i> versus <i>w</i> (eight variables; elemental content and RI _b) | | | | | | | | |
| Na' | 1.000 | 0.034 | -0.127 | -0.394 | 0.142 | 0.214 | 0.262 | 0.009 |
| Mg' | | 1.000 | 0.079 | -0.424 | -0.039 | 0.340 | -0.038 | 0.208 |
| Al' | | | 1.000 | 0.119 | -0.466 | 0.030 | 0.303 | 0.018 |
| Si' | | | | 1.000 | -0.149 | -0.849 | -0.017 | -0.325 |
| K' | | | | | 1.000 | 0.085 | -0.092 | -0.027 |
| Ca' | | | | | | 1.000 | 0.003 | 0.411 |
| Fe' | | | | | | | 1.000 | 0.326 |
| Rib | | | | | | | | 1.000 |
| (f) <i>c</i> versus <i>w</i> (eight variables; elemental content and dRI) | | | | | | | | |
| Na' | 1.000 | 0.030 | -0.126 | -0.406 | 0.142 | 0.221 | 0.271 | 0.009 |
| Mg' | | 1.000 | 0.117 | -0.402 | -0.032 | 0.315 | -0.138 | -0.171 |
| Al' | | | 1.000 | 0.084 | -0.449 | 0.087 | 0.261 | -0.251 |
| Si' | | | | 1.000 | -0.166 | -0.834 | 0.123 | 0.166 |
| K' | | | | | 1.000 | 0.104 | -0.088 | -0.009 |
| Ca' | | | | | | 1.000 | -0.185 | -0.253 |
| Fe' | | | | | | | 1.000 | 0.159 |
| dRI | | | | | | | | 1.000 |

becomes the first separator set. The running union of the first two sets is now (Na', Ca', K'). Working through the entire set chain one arrives at the following factorization:

$$f(C_i|S_i) = \frac{f(\text{Na}', \text{Ca}')f(\text{Al}', \text{K}')f(\text{Al}', \text{K}')f(\text{Fe}', \text{K}')f(\text{Ca}', \text{Si}')f(\text{Mg}', \text{Ca}')}{f(\text{Ca}')^2f(\text{Al}')f(\text{K}')f(\text{Na}')} \quad (11)$$

Probability density functions presented in Eq. (11) could be expressed by (Eq. [8]) and then (Eq. [11]) become a form of LR (Eq. [12]):

$$\text{LR} = \frac{\text{LR}(\text{Na}', \text{Ca}')\text{LR}(\text{Al}', \text{K}')\text{LR}(\text{Al}', \text{K}')\text{LR}(\text{Fe}', \text{K}')\text{LR}(\text{Ca}', \text{Si}')\text{LR}(\text{Mg}', \text{Ca}')}{\text{LR}(\text{Ca}')^2\text{LR}(\text{Al}')\text{LR}(\text{K}')\text{LR}(\text{Na}')} \quad (12)$$

The model represents a minimal model that decomposes the seven variables of the full data set into six sets of two variables (see column C_i in Table 6). Parameters (means, variances, and covariance) of the two-variate distributions could also be satisfactorily estimated on the basis of information represented by a relatively small database, that is, 153 glass samples.

Graphical models for other problems considered in this paper were analyzed in the same way and the following LR formulas were obtained:

(i) *bh* versus *p*:

$$\text{LR} = \frac{\text{LR}(\text{Na}', \text{Ca}')\text{LR}(\text{K}', \text{Na}')\text{LR}(\text{Mg}', \text{Ca}')\text{LR}(\text{Si}', \text{Ca}')\text{LR}(\text{Fe}')\text{LR}(\text{Al}')}{\text{LR}(\text{Ca}')^2\text{LR}(\text{Na}')}$$

(ii) *b* versus *h*:

$$\text{LR} = \frac{\text{LR}(\text{Na}', \text{Ca}')\text{LR}(\text{K}', \text{Na}')\text{LR}(\text{Mg}', \text{Ca}')\text{LR}(\text{Na}', \text{Fe}')\text{LR}(\text{Si}')\text{LR}(\text{Al}')}{\text{LR}(\text{Ca}')\text{LR}(\text{Na}')^2}$$

(iii) *c* versus *w* (seven variables; elemental content):

$$\text{LR} = \frac{\text{LR}(\text{Mg}', \text{Si}')\text{LR}(\text{Na}', \text{Si}')\text{LR}(\text{Ca}', \text{Si}')\text{LR}(\text{K}', \text{Al}')\text{LR}(\text{Fe}')}{\text{LR}(\text{Si}')^2}$$

(iv) *c* versus *w* (eight variables; elemental content and RI_b):

$$\text{LR} = \frac{\text{LR}(\text{Ca}', \text{Si}')\text{LR}(\text{RI}_b, \text{Ca}')\text{LR}(\text{Na}', \text{Si}')\text{LR}(\text{Mg}', \text{Si}')\text{LR}(\text{K}', \text{Al}')\text{LR}(\text{Fe}')}{\text{LR}(\text{Si}')^2\text{LR}(\text{Ca}')}$$

(v) *c* versus *w* (eight variables; elemental content and ΔRI):

$$\text{LR} = \frac{\text{LR}(\text{Mg}', \text{Si}')\text{LR}(\text{Na}', \text{Si}')\text{LR}(\text{Ca}', \text{Si}')\text{LR}(\text{K}', \text{Al}')\text{LR}(\text{Fe}')\text{LR}(\Delta\text{RI})}{\text{LR}(\text{Si}')^2}$$

Discussion

Inspection of Table 4 shows that all models performed relatively well. The SVM model gave slightly better results for the problem of classifying fragments into *cw* versus *bhp* categories, *bh* versus *p* categories, as well as *b* versus *h* categories. In most of cases, the LR model gave slightly higher misclassification rates for most of the experiments conducted here.

The misclassification rate for the categorization of glass objects into *cw* or *bhp* categories is relatively low. The highest misclassification rate (71% correct; 11 incorrectly classified specimens from 38 glass objects) was observed for sets no. 7 and 10 when the LR model was used, and the lowest misclassification rate (97% correct; 1 incorrectly classified specimen) was obtained by application of SVM on set no. 10 and by the application of the NBC approach to set no. 9.

Perfect classification was obtained for test sets no. 1 and 9 when the SVM model was used, and at test sets no. 2 and 5 when the NBC model was used to classify into *bh* and *p* classes. The highest misclassification rate (63% correct, nine incorrectly classified specimens from 24) was within set no. 9 when the LR model was used.

An analysis of the misclassification rates obtained for the bulb/headlamp classification problem conducted on the test sets,

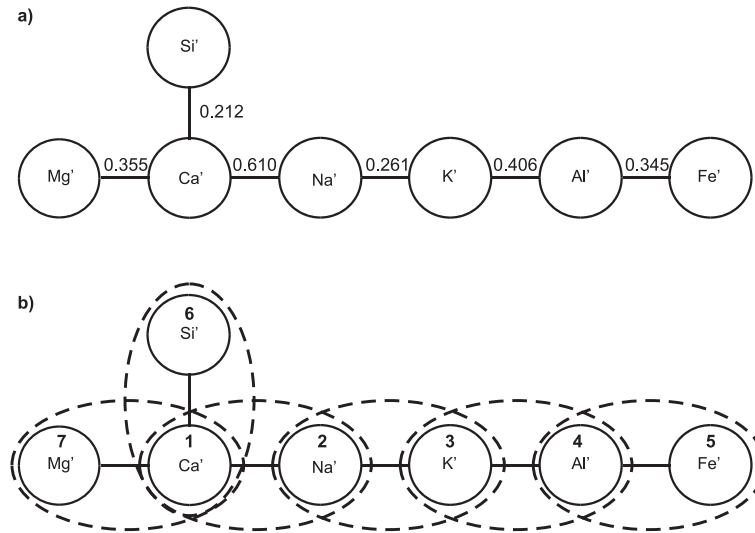


FIG. 4—Results of construction of the undirected graphical model: (a) a model calculated on the basis of the rescaled inverse of the variance-covariance matrix (Table 5a) for the seven variables from the 153 glass objects, (b) model (a) with marked cliques and numbers designated to each node. Na', Mg', Al', Si', K', Ca' and Fe'—the abbreviation of \log_{10} of oxygen concentration to concentration of each element.

TABLE 6—Results of factorization of the undirected model presented in Fig. 4.

| i | Clique (C_i) | Running Union (R_i) | Separator Set (S_i) |
|-----|------------------|------------------------------------|-------------------------|
| 1 | (Na', Ca') | (Na', Ca') | ϕ |
| 2 | (Na', K') | (Na', Ca', K') | (Na') |
| 3 | (Al', K') | (Na', Ca', K', Al') | (K') |
| 4 | (Fe', Al') | (Na', Ca', K', Al', Fe') | (Al') |
| 5 | (Ca', Si') | (Na', Ca', K', Al', Fe', Si') | (Ca') |
| 6 | (Mg', Ca') | (Na', Ca', K', Al', Fe', Si', Mg') | (Ca') |

*The abbreviation of \log_{10} of oxygen concentration to concentration of each element.

gave the highest misclassification (73% correct; three incorrectly classified specimens from 11 glass objects) from test set no. 1, when the LR model was applied. An accuracy of 100% classification was obtained 12 times, but mostly through the application of the SVM and NBC models.

When the classification problem c versus w was conducted using only elemental concentrations, the lowest misclassification was 86% correct (two incorrectly classified specimens among 14 objects) obtained when the NBC model was applied to test set no. 10. The highest misclassification rate (43% correct; sets no. 2, 3, and 4) was obtained by application of the LR model. However, the LR model gave better results than the SVM and NBC models for set no. 9. These observations might be expected because samples from car and building windows have quite similar elemental compositions. When RI_b was included as an additional feature then the misclassification rates rose, the highest misclassification rate (43% correct; eight incorrectly classified samples among 14 ones) was obtained by application of the LR model (set no. 10) and the lowest misclassification rate (79% correct; three incorrectly classified samples) was obtained for SVM (test set no. 9), and by NBC (test set no. 4) and by the LR model (set no. 9).

The application of ΔRI in form of $|\log_{10}(\Delta RI)|$ as an additional feature for the c versus w classification problem gave generally much better results. For example, the highest misclassification rate was 64% correct (five incorrectly classified specimens from 14 objects), and was obtained one time (set no. 10) when the LR model was applied. All samples were correctly classified for test

sets no. 5, 7, and 8 when the SVM model was applied, and sets no. 7 and 8 when the LR model was applied.

On no occasion did the LR model outperform the other models for the first three classification problems. It also gave the highest misclassification rates when various combinations of variables in the c versus w problem were used, however, the observed differences between the classification models were not great. These results suggest that, from a practical point of view, SEM-EDX analysis of major and minor elements contains sufficient information to reliably classify glass objects into the p , b , h , or cw categories when the following scheme is used:

- (i) Classification of the unknown specimen to cw or bhp category. When the sample is classified to cw category then halt the process because it can be concluded that the unknown specimen is most likely to belong to the cw category than to the bhp category. If a specimen is assigned to the bhp category then the classification has to be extended to assign between bh and p categories.
- (ii) If, after assignment of the glass to the bh or p categories, the specimen can be assigned to the p category, then halt the classification process because the most likely category is p . If it is found that a specimen most likely belongs to bh category then further sub-classification to either the b or h categories is possible.
- (iii) If sample most likely belongs to the cw category then either the c or the w categories can be made upon the basis of ΔRI information.

The effectiveness of the proposed algorithm for glass classification on the basis of SEM-EDX analysis has been illustrated by case-work. The clothes of a person in charge of a motor vehicle which was suspected to have been involved in a hit-and-run incident with a pedestrian were sent for analysis. The primary aim of the analysis was to recover all possible glass fragments which could help the investigators to identify the person, or persons, present at the scene of the incident. Key to this investigation was whether there were glass fragments from c (car) category on his clothes. Debris were collected by brushing, and analyzed by optical microscope. So many glass fragments of linear dimension 0.1–0.2 mm were found that it was not possible to analyze them all. Therefore, 10 fragments were taken for elemental analysis. They were only analyzed by the

TABLE 7—An application of the proposed classification algorithm—results of elemental analysis and classification process by SVM, NBC, and LR models.

| Trace | Elemental Content (% wt), Standard Deviation | | | | | | | | Results of Classification* | | |
|-------|--|-------|------|------|-------|------|------|------|----------------------------|-----------|-------------------------|
| | O | Na | Mg | Al | Si | K | Ca | Fe | SVM | NBC | LR |
| 1 | 49.56 | 9.24 | 2.15 | 0.26 | 33.07 | 0.00 | 5.68 | 0.10 | <i>cw</i> | <i>cw</i> | <i>cw</i> |
| | 0.11 | 0.09 | 0.04 | 0.08 | 0.26 | 0.00 | 0.07 | 0.08 | (0.995) | (0.899) | (184) |
| 2 | 50.07 | 9.62 | 2.17 | 0.61 | 32.67 | 0.00 | 4.80 | 0.00 | <i>cw</i> | <i>cw</i> | <i>cw</i> |
| | 1.56 | 0.10 | 0.10 | 0.05 | 1.24 | 0.00 | 0.37 | 0.00 | (0.643) | (0.943) | (106) |
| 3 | 48.50 | 9.10 | 2.15 | 0.22 | 33.79 | 0.13 | 6.01 | 0.11 | <i>cw</i> | <i>cw</i> | <i>cw</i> |
| | 0.51 | 0.07 | 0.04 | 0.04 | 0.32 | 0.04 | 0.08 | 0.09 | (0.985) | (0.825) | (1 × 10 ²⁰) |
| 4 | 45.56 | 9.00 | 2.07 | 0.40 | 35.52 | 0.20 | 7.23 | 0.00 | <i>cw</i> | <i>cw</i> | <i>cw</i> |
| | 1.04 | 0.49 | 0.09 | 0.08 | 0.76 | 0.03 | 0.74 | 0.00 | (0.911) | (0.743) | (64373) |
| 5 | 50.41 | 10.21 | 2.28 | 0.41 | 31.25 | 0.16 | 5.19 | 0.00 | <i>cw</i> | <i>cw</i> | <i>cw</i> |
| | 1.72 | 0.46 | 0.03 | 0.13 | 1.40 | 0.06 | 0.55 | 0.00 | (0.889) | (0.995) | (45999) |
| 6 | 50.39 | 9.69 | 2.23 | 0.21 | 32.06 | 0.09 | 5.27 | 0.00 | <i>cw</i> | <i>cw</i> | <i>cw</i> |
| | 0.31 | 0.08 | 0.04 | 0.03 | 0.21 | 0.01 | 0.17 | 0.00 | (0.982) | (0.683) | (621) |
| 7 | 47.67 | 9.28 | 2.19 | 0.21 | 34.34 | 0.14 | 6.07 | 0.00 | <i>cw</i> | <i>cw</i> | <i>cw</i> |
| | 0.64 | 0.12 | 0.05 | 0.04 | 0.59 | 0.02 | 0.12 | 0.00 | (0.994) | (0.899) | (91132) |
| 8 | 50.21 | 9.41 | 2.23 | 0.21 | 32.28 | 0.00 | 5.53 | 0.00 | <i>cw</i> | <i>cw</i> | <i>cw</i> |
| | 0.39 | 0.09 | 0.02 | 0.03 | 0.21 | 0.00 | 0.15 | 0.00 | (0.981) | (0.899) | (164) |
| 9 | 47.60 | 9.08 | 2.15 | 0.49 | 33.76 | 0.43 | 6.41 | 0.00 | <i>cw</i> | <i>cw</i> | <i>cw</i> |
| | 0.05 | 0.05 | 0.03 | 0.02 | 0.06 | 0.00 | 0.05 | 0.00 | (0.706) | (0.456) | (7 × 10 ⁶) |
| 10 | 51.88 | 9.83 | 2.30 | 0.36 | 30.57 | 0.11 | 4.86 | 0.00 | <i>cw</i> | <i>cw</i> | <i>cw</i> |
| | 0.13 | 0.04 | 0.03 | 0.03 | 0.10 | 0.01 | 0.09 | 0.00 | (0.714) | (0.869) | (4 × 10 ²⁰) |

SVM, support vector machine; NBC, naïve Bayes classifiers; LR, likelihood ratio.

**cw*—glass fragments originated from car and building windows. The number in brackets is a probability obtained on each of the classification problems that an object belongs to the *cw* category, i.e., $p(cw|x_1, \dots, x_n)$ when SVM and NBC were applied or it is an LR value.

SEM-EDX method because of their relatively small size. The results of SEM-EDX analysis and the classification process are presented at Table 7. Other traces which could be thought to be characteristic of car accidents (e.g., paint smears) were not found.

It was concluded that the elemental composition of all 10 glass fragments most likely belonged to the *cw* category. Nevertheless, it was thought that in this case a further classification into *c* or *w* category would not be reliable. However, the relatively large number of recovered glass fragments classified to the joint *cw* category suggested to the public prosecutor that the suspect had played a part in this hit-and-run incident.

Conclusions

The proposed scheme for glass fragment(s) classification works with relative efficiency, except that assignment to car windows (*c*) and building windows (*w*) needs to be treated with care because of the very similar elemental content of these two categories. Research on classification of glass objects to *c* and *w* categories should focus on the combination of elemental content and information on ΔRI which seems to be giving very promising results. The application of SVM and NBC gave slightly better results than application of the LR model. However, the observed differences in misclassification rates were not great (especially in *c* vs. *w* problem), and no single classification method was obviously more effective than any other. Despite the slight underperformance compared with the SVM and NBC methods, the LR model can be recommended as this framework has the advantage that an LR is easily interpretable, does not act as a “black-box” unlike the SVM method, and does not require the investigator to make some assumption about prior belief, unlike the NBC method (1,3–6).

Moreover, the LR model might be easily adapted to other forensic classification problems where the observations may be multivariate, for instance, where questions arise surrounding the differentiation between kerosene and diesel fuel traces which might be recovered from fire debris. These two organic compounds have

very similar chemical compositions, and it is difficult to distinguish between them in any simple fashion on the basis of GC/MS measurements.

Acknowledgments

The author wishes to thank Prof. Janina Zieba-Palus, Institute of Forensic Research, Krakow, Poland for delivery of refractive index data and helpful discussion and Dr. Tereza Neocleous, School of Mathematics, University of Edinburgh, Edinburgh, U.K. and Dr. David Lucy, Department of Statistics, Lancaster University, U.K. for helpful comments and language support.

References

- Curran JM, Hicks TN, Buckleton JS. Forensic interpretation of glass evidence. Boca Raton, FL: CRC Press LLC, 2000.
- Caddy B. Forensic examination of glass and paint. Boca Raton, FL: CRC Press LLC, 2001.
- Aitken CGG, Taroni F. Statistics and the evaluation of evidence for forensic scientists. Chichester, U.K.: John Wiley & Sons, 2004.
- Aitken CGG, Lucy D. Evaluation of trace evidence in the form of multivariate data. Appl Stat 2004;53:109–22.
- Aitken CGG, Lucy D, Zadora G, Curran JM. Evaluation of transfer evidence for three level multi-variate data with the use of Graphical Models. Comput Stat Data Anal 2006;50:2571–88.
- Aitken CGG, Zadora G, Lucy D. A Two level model for evidence evaluation. J Forensic Sci 2007;52(2):412–9.
- Pawluk-Kolc M, Zięba-Palus J, Parczewski A. Application of false discovery rate procedure to pairwise comparisons of refractive index of glass fragments. For Sci Int 2006;160:53–8.
- Pawluk-Kolc M, Zięba-Palus J, Parczewski A. The effect of re-annealing on the distribution of refractive index in a windscreen and a windowpane. For Sci Int 2008;174:222–8.
- Zadora G. Examination of the refractive index of selected samples of glass for forensic purposes. Probl Forensic Sci 2001;XLV:36–45.
- Zadora G. The role of statistical methods in assessing the evidential value of physico-chemical data. Probl Forensic Sci 2006;LXV:91–103.
- Zadora G. Glass analysis for forensic purposes—a comparison of classification methods. J Chemometrics 2007;21:174–86.

12. Zadora G, Brozek-Mucha Z. SEM-EDX—a useful tool for forensic examinations. *Material Chem Phys* 2003;81:345–8.
13. Hicks T, Monard Sermier F, Goldmann T, Brunelle A, Champod C, Margot P. The classification and discrimination of glass fragments using non destructive energy dispersive X-ray fluorescence. *Forensic Sci Int* 2003;137:107–18.
14. Trejos T, Almirall JR. Sampling strategies for the analysis of glass fragments LA-ICP-MS Part 1: micro-homogeneity study of glass and its application to the interpretation of forensic evidence. *Talanta* 2005;67:388–95.
15. Hickam DA. A classification scheme for glass. *Forensic Sci Int* 1981;17:265–81.
16. Hickman DA. Glass types identified by chemical analysis. *Forensic Sci Int* 1987;33:23–46.
17. Koons RD, Fiedler C, Rawalt R. Classification and discrimination of sheet and container glasses by inductively coupled plasma—atomic emission spectrometry and pattern recognition. *J Forensic Sci* 1988;33:49–67.
18. Zadora G, Piekoszewski W, Parczewski A. An application of chosen similarity measurements of objects for forensic purposes. *Bulletin of the International Statistical Institute 54th Session, LX, Book 2*. Berlin, Germany: International Statistical Institute, 2003; 364–7.
19. Locke J, Rockett L. The application of annealing to improve the discrimination between glasses. *Forensic Sci Int* 1985;29:237–45.
20. Winstanley R, Rydeard C. Concepts of annealing applied to small glass fragments. *Forensic Sci Int* 1985;29:1–10.
21. Evgeniou T, Pontil M, Poggio T. Regularization networks and support vector machines. *Adv Comput Math* 2000;13:1–50.
22. Girosi F, Jones M, Poggio T. Regularization theory and neural networks architectures. *Neural Comput* 1995;7:219–69.
23. Girosi F. An equivalence between sparse approximation and support vector machines. *Neural Comput* 1998;10:1455–80.
24. Smola A, Scholkopf B, Muller K-R. The connection between regularization operators and support vector kernels. *Neural Netw* 1998;11:637–49.
25. Vapnik VN. *The nature of statistical learning theory*. Berlin: Springer-Verlag, 1995.
26. Cortes C, Vapnik VN. Support-vector network. *Mach Learn* 1995;20:273–97.
27. Cristianini N, Shawe-Taylor J. *An introduction to support vector machines and other kernel-based learning methods*. Cambridge, U.K.: Cambridge University Press, 2000.
28. Burges CJC. A tutorial on support vector machines for pattern recognition. *Data Mining Knowl Discov* 1998;2:121–67.
29. Smola A, Scholkopf B. *A tutorial on support vector regression*. Neuro-COLT Technical Report NC-RR-98-030. U.K.: Royal Holloway College, University of London, 1998.
30. Bishop CM. *Pattern recognition and machine learning*. U.K.: Springer, 2006.
31. Domingos P, Pazzani M. Beyond independence: conditions for the optimality of the simple Bayesian classifiers. *Mach Learn* 1997;29:103–30.
32. Duda R, Hart P. *Pattern classification and scene analysis*. New York, U.S.A.: John Wiley & Sons, 1973.
33. Friedman N, Geoger D, Goldszmidt M. Bayesian network classifiers. *Mach Learn* 1997;29:131–63.
34. Hand DJ, Yu K. Idiot's Bayes—not so stupid after all? *Int Stat Rev* 2001;3:385–98.
35. The R Foundation for Statistical Computing Version 2.0.1. Available at <http://www.r-project.org>.

Additional information and reprint requests:

Grzegorz Zadora, Ph.D.
 Institute of Forensic Research
 Westerplatte 9
 31-033 Krakow
 Poland
 E-mail: gzadora@ies.krakow.pl