

Bayes-optimal performance in a discrete space

This article has been downloaded from IOPscience. Please scroll down to see the full text article.

1999 J. Phys. A: Math. Gen. 32 L555

(<http://iopscience.iop.org/0305-4470/32/50/104>)

View [the table of contents for this issue](#), or go to the [journal homepage](#) for more

Download details:

IP Address: 171.66.16.118

The article was downloaded on 02/06/2010 at 07:57

Please note that [terms and conditions apply](#).

LETTER TO THE EDITOR

Bayes-optimal performance in a discrete spaceM Copelli[†], C Van den Broeck[†] and M Opper[‡]

E-mail: mauro.copelli@luc.ac.be, christian.vandenbroeck@luc.ac.be and opperm@aston.ac.uk

[†] Limburgs Universitair Centrum, B-3590 Diepenbeek, Belgium[‡] Neural Computing Research Group, Aston University, Birmingham B4 7ET, UK

Received 1 October 1999

Abstract. We study a simple model of unsupervised learning where the single symmetry breaking vector has binary components ± 1 . We calculate exactly the Bayes-optimal performance of an estimator which is required to lie in the same discrete space. We also show that, except for very special cases, such an estimator cannot be obtained by minimization of a class of variationally optimal potentials.

Statistical mechanics techniques have been used with success to study and understand key properties of inferential learning [1,2]. This approach provides explicit and detailed results that are in many ways complementary to the more general results obtained by statistics. The case of non-smooth problems, in which the parameters that have to be estimated take discrete values, is of particular interest. On the one hand, many of the results from statistics can no longer be applied, while on the other, the estimation of these parameters is often a computationally hard problem. In this paper, we present a detailed analysis of a simple model of unsupervised learning [3–8], involving a single symmetry breaking vector with binary components ± 1 and highlight the differences with the case of smooth components. In particular, we compare the results from Gibbs learning and Bayes learning with the ones for the best binary vector and a vector which minimizes a variationally optimal potential.

The problem is as follows: p N -dimensional real patterns $\{\xi^\mu, \mu = 1, \dots, p\}$ are sampled independently from a distribution $P(\xi^\mu | \mathbf{B}) \sim \delta(\xi^\mu \cdot \mathbf{B} - N) \exp[-U(\mathbf{B} \cdot \xi^\mu / \sqrt{N})]$ with a single symmetry breaking direction \mathbf{B} . The function U modulates the distribution of the patterns along \mathbf{B} . We will focus on the properties in the thermodynamic limit $N \rightarrow \infty$, $p \rightarrow \infty$ with $\alpha = p/N$ finite. One then finds that the normalized projection $t \equiv \mathbf{B} \cdot \xi / \sqrt{N}$ is distributed according to (\mathcal{N} being a normalization constant)

$$P^*(t) = \frac{\mathcal{N}}{\sqrt{2\pi}} \exp\left\{-\frac{t^2}{2} - U(t)\right\} \quad (1)$$

while projections on any direction orthogonal to \mathbf{B} are normal. The case of a so-called spherical prior, in which \mathbf{B} is chosen at random on the sphere with radius \sqrt{N} , was discussed in [6–8]. As announced earlier, we focus here on the more complicated situation in which the components of \mathbf{B} take binary values ± 1 . The *prior distribution* is now given by

$$P(\mathbf{B}) \equiv P_b(\mathbf{B}) = \prod_{j=1}^N \left[\frac{1}{2} \delta(B_j - 1) + \frac{1}{2} \delta(B_j + 1) \right]. \quad (2)$$

The goal of unsupervised learning is to give an estimate \mathbf{J} of \mathbf{B} . One way to do so is to sample \mathbf{J} from a Boltzmann distribution with Hamiltonian $\mathcal{H}(\mathbf{J}) = \sum_{\mu=1}^p V(\lambda^\mu)$, with $\lambda^\mu \equiv \mathbf{J} \cdot \boldsymbol{\xi}^\mu / \sqrt{N}$, at temperature $T = \beta^{-1}$, for an appropriate choice of the *ad hoc* potential V [9]. The properties of such a \mathbf{J} -vector can be extracted from the partition function:

$$Z = \int d\mathbf{J} P_b(\mathbf{J}) e^{-\beta \mathcal{H}(\mathbf{J})}. \quad (3)$$

The latter is a fluctuating quantity due to the random choice of the patterns, but the free energy per component $f = -(N\beta)^{-1} \ln Z$ is expected to be self-averaging in the thermodynamic limit and can therefore be calculated by averaging over the pattern distribution with the aid of the replica trick [10]. Assuming replica symmetry (RS), one finds

$$f = \frac{1}{\beta} \text{Extr}_{R, q, \hat{R}, \hat{q}} \left\{ \frac{1}{2} (1-q) \hat{q} + \hat{R} R - \int \mathcal{D}z \ln \cosh(z\sqrt{\hat{q}} + \hat{R}) - \alpha \int \mathcal{D}^*t \int \mathcal{D}t' \right. \\ \left. \times \ln \int \frac{d\lambda}{\sqrt{2\pi(1-q)}} \exp \left(-\beta V(\lambda) - \frac{(\lambda - t'\sqrt{q - R^2} - tR)^2}{2(1-q)} \right) \right\} \quad (4)$$

where $\mathcal{D}^*t = dt P^*(t)$ and $\mathcal{D}t' = dt' (2\pi)^{-1/2} \exp(-t'^2/2)$. The extremum operator gives saddle point equations which determine the self-averaging value of the order parameters. As usual q can be interpreted as the typical mutual overlap between two samples \mathbf{J} and \mathbf{J}' , $q = \mathbf{J} \cdot \mathbf{J}'/N$, while the performance R measures the proximity between the estimate \mathbf{J} and the 'true' direction \mathbf{B} , $R = \mathbf{J} \cdot \mathbf{B}/N$. For even functions U , there is no distinction between \mathbf{B} and $-\mathbf{B}$, and a symmetry $R \rightarrow -R$ arises. In the following, only $R \geq 0$ will be considered.

As a first application of equation (4), we turn to Gibbs learning [5, 11, 12]. It corresponds to sampling from the posterior distribution and is realized by taking $\beta = 1$ and $V = U$ in equation (4) (for more details, see [6]; for the estimation of U , see [13]). In agreement with the fact that one cannot make a statistical distinction between \mathbf{B} and its Gibbsian estimate \mathbf{J} , one finds that the order parameters satisfy $q_G = R_G$ and $\hat{q}_G = \hat{R}_G$, where the subscript G refers to Gibbs learning. This observation allows us to simplify the saddle point equations further, and the Gibbs overlap is found to obey the following equation:

$$R_G = F_B^2 \left(\mathcal{F} \left(\sqrt{R_G} \right) \right) \quad (5)$$

with

$$F_B(x) = \sqrt{\int \mathcal{D}z \tanh(zx + x^2)} \quad \text{and} \quad \mathcal{F}(R) = \sqrt{\alpha \int \mathcal{D}t \frac{Y^2(t; R)}{X(t; R)}} \quad (6)$$

and

$$X(t; R) = \int \mathcal{D}t' \mathcal{N} e^{-U(Rt + \sqrt{1-R^2}t')} \quad Y(t; R) = \frac{1}{R} \frac{\partial}{\partial t} X(t; R). \quad (7)$$

Note that F_B comes from the entropic term of the free energy and does not depend on U , as opposed to \mathcal{F} , X and Y .

For R_G small, one obtains from equations (5)–(7), upon assuming a smooth behaviour as a function of α , that $(\int \mathcal{D}^*t f(t) = \langle f(t) \rangle_*)$:

$$\langle t \rangle_* \neq 0 \Rightarrow R_G \simeq \alpha \langle t \rangle_*^2 \quad (8)$$

$$\langle t \rangle_* = 0 \Rightarrow R_G \begin{cases} = 0 & \alpha \leq \alpha_G \\ \simeq C(\alpha - \alpha_G) & \alpha \geq \alpha_G \end{cases} \quad (9)$$

with critical load $\alpha_G = (1 - \langle t^2 \rangle_*)^{-2}$. These results are identical to those for a spherical prior [8]. In particular, one observes the appearance of *retarded learning* when the distribution

has a zero mean along the symmetry breaking axis. In the regime $R_G \rightarrow 1$, on the other hand, one finds an exponential approach:

$$1 - R_G(\alpha) \stackrel{\alpha \rightarrow \infty}{\simeq} \sqrt{\frac{\pi}{2\alpha \langle (U')^2 \rangle_*}} \exp\left(\frac{-\alpha \langle (U')^2 \rangle_*}{2}\right) \quad (10)$$

where $U' \equiv dU(t)/dt$. This is now different from the case of a spherical prior, where the approach is following an inverse power law $1 - R_G \sim \alpha^{-1}$ [8]. The difference becomes even more pronounced when U has singular derivatives, as is typically the case when a supervised problem is mapped onto an unsupervised version [6]. Then one finds that $R_G = 1$ is attained at a finite value of α while $1 - R_G \sim \alpha^{-2}$ for a spherical prior, see [6, 11] for an explicit example.

Apart from its intrinsic interest, Gibbs learning is also directly related to the Bayes optimal overlap by $R_B = \sqrt{R_G}$, see [5, 12, 14]. This overlap is realized by the centre of mass \mathbf{J}_B of the Gibbs ensemble. A simple reasoning [5, 12] shows that \mathbf{J}_B maximizes the overlap R averaged over the posterior distribution of \mathbf{B} . In order to exclude the case $\mathbf{J}_B = 0$ (which would follow in the presence of the symmetry $\mathbf{B} \rightarrow -\mathbf{B}$), we will implicitly assume an infinitesimally small symmetry breaking field in the Gibbs distribution.

Using the self-averaging of the mutual overlap, with $q_G = R_G$, the explicit form of \mathbf{J}_B is found to be $\mathbf{J}_B = R_G^{-1/2} Z^{-1} \int d\mathbf{J} P_b(\mathbf{J}) \mathbf{J} \exp\{-\sum_\mu U(\lambda^\mu)\}$. In general, the components of this centre of mass are continuous, while our prime interest here is in the optimal performance attainable by a binary vector. The latter vector, which we will denote by \mathbf{J}_{bb} (for best binary), can fortunately be easily obtained [1]: it is the clipped version of the centre of mass \mathbf{J}_B , with components $(\mathbf{J}_{bb})_j = \text{sign}((\mathbf{J}_B)_j)$.

To evaluate the overlap between \mathbf{J}_{bb} and \mathbf{B} , we recall the following general result for the overlap $\tilde{R} = \tilde{\mathbf{J}} \cdot \mathbf{B} / N$ of a vector $\tilde{\mathbf{J}}$ with transformed components $\tilde{J}_i = \sqrt{N} g(J_i) / \sqrt{\sum_i g^2(J_i)}$ (with g odd and \mathbf{B} binary) as a function of the overlap R of \mathbf{J} with \mathbf{B} (see [9] for details):

$$\tilde{R} = \frac{\int P(x) g(x) dx}{[\int P(x) g^2(x) dx]^{1/2}} \quad (11)$$

where $P(x)$ is the probability density for $x \equiv J_1 B_1$, which for the prior distribution equation (2) is independent of the index due to the permutation symmetry among the axes. If \mathbf{J} is sampled from a spherical distribution (with \mathbf{B} binary), then $P(x)$ is found to be a Gaussian [9] with mean R and variance $1 - R^2$.

In order to obtain $P(x)$ corresponding to the centre of mass \mathbf{J}_B , we evaluate the quenched moments of $y = x \sqrt{R_G}$:

$$\langle y^m \rangle = \left\langle \left(Z^{-1} \int d\mathbf{J} P_b(\mathbf{J}) e^{-\sum_\mu U(\lambda^\mu)} J_1 B_1 \right)^m \right\rangle. \quad (12)$$

The average $\langle \dots \rangle$ over the quenched pattern set can be performed by the replica trick with the following replica symmetric result:

$$\langle y^m \rangle = \int \mathcal{D}z \left[\tanh \left(z \sqrt{\hat{R}_G} + \hat{R}_G \right) \right]^m \quad (13)$$

where \hat{R}_G , which is determined by the saddle point equations of Gibbs learning, cf equation (5), is found to be $\hat{R}_G = \mathcal{F}^2(\sqrt{R_G})$. Recognizing equation (13) as a transformation of variables $y = \tanh(z \sqrt{\hat{R}_G} + \hat{R}_G)$, with z normally distributed, one concludes[†] that

$$P(x) = \frac{\sqrt{R_G}}{\sqrt{2\pi \hat{R}_G (1 - R_G x^2)}} \exp \left\{ \frac{-1}{2\hat{R}_G} \left[\frac{1}{2} \ln \left(\frac{1 + \sqrt{R_G} x}{1 - \sqrt{R_G} x} \right) - \hat{R}_G \right]^2 \right\}. \quad (14)$$

[†] Equation (14) is consistent with the fact that the overlap for the centre of mass cannot be improved by a transformation of its components. Indeed, the transformed overlap \tilde{R} in equation (11) is maximized [9] by setting $g_{\text{opt}}(x) = (P(x) - P(-x)) / (P(x) + P(-x))$, reducing for (14) to $g_{\text{opt}}(x) \sim x$.

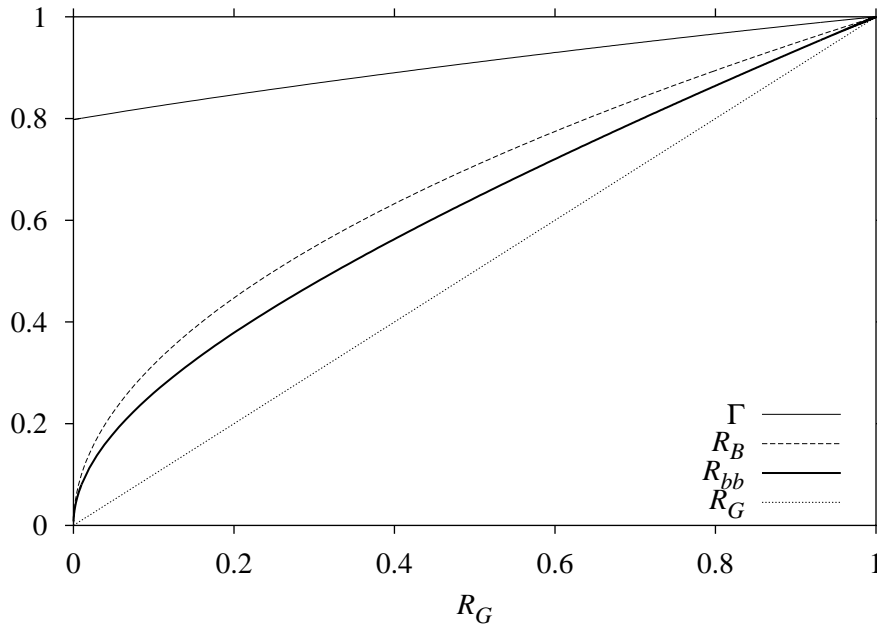


Figure 1. Γ , R_B and R_{bb} parametrized by R_G , according to equations (5) and (15).

By applying equation (11), for $g(x) = \text{sign}(x)$, with $P(x)$ given by equation (14), one finally obtains the following overlap $R_{bb} \equiv \mathbf{J}_{bb} \cdot \mathbf{B}/N$ of the best binary vector:

$$R_{bb} = 1 - 2H(F_B^{-1}(R_B)) = 1 - 2H(\mathcal{F}(R_B)) \tag{15}$$

where $H(x) \equiv \int_x^\infty Dt$. Equation (15) is a central result of this paper, providing an upper bound for the performance of any binary vector. The asymptotics of R_{bb} can be obtained from those of $R_G = R_B^2$, yielding

$$R_{bb} \stackrel{R_G \rightarrow 0}{\simeq} \sqrt{\frac{2R_G}{\pi}} \tag{16}$$

in the poor performance regime, and an exponential behaviour in the limit of $R_G \rightarrow 1$:

$$1 - R_{bb} \simeq \frac{2}{\pi}(1 - R_G). \tag{17}$$

We note that another quantity of interest, the mutual overlap $\Gamma \equiv \mathbf{J}_B \cdot \mathbf{J}_{bb}/N$ between the centre of mass and best binary, can also be evaluated quite easily, leading to the simple result $\Gamma = R_{bb}/R_B$. In the limit $R_G \rightarrow 0$ one recovers $\Gamma \rightarrow \sqrt{2/\pi}$, which is the result for the overlap between a vector sampled at random from the N -sphere and its clipped counterpart. Γ , R_B and R_{bb} are plotted as functions of R_G in figure 1.

We finally turn to the problem of a variationally optimized potential. In the case of a spherical prior, it was shown that the Bayes-optimal performance can indeed be attained by a vector that minimizes this potential [8, 15–17]. We now address the question of whether the same procedure is successful in discrete space, a problem which has been also studied in [18] for the supervised scenario. Since \mathbf{J}_{bb} is a unique optimal binary vector, one would like the desired potential to satisfy both $R = R_{bb}$ and $q = 1$. Proceeding again from the free

energy equation (4) for a general potential V , taking the limits $q \rightarrow 1, \beta \rightarrow \infty$ with finite $c \equiv \beta(1 - q)$, and rescaling the conjugate parameters $\hat{c} \equiv \hat{q}/\beta^2, \hat{y} \equiv \hat{R}/\beta$, one obtains the following saddle point equations:

$$R = 1 - 2H\left(\frac{\hat{y}}{\sqrt{\hat{c}}}\right) \quad c = \sqrt{\frac{2}{\pi\hat{c}}} \exp\left(-\frac{\hat{y}^2}{2\hat{c}}\right) \tag{18}$$

$$\hat{c} = \frac{\alpha}{c^2} \int \mathcal{D}t X(t; R)[\lambda_0(t, c) - t]^2 \quad \hat{y} = \frac{\alpha}{c} \int \mathcal{D}t Y(t; R)[\lambda_0(t, c) - t]$$

where $\lambda_0(t, c) \equiv \text{Argmin}_\lambda [V(\lambda) + (\lambda - t)^2/2c]$. The variational optimization of R with respect to the choice of V can now be performed as in [8, 15–17] invoking the Schwarz inequality. We only quote the final result for the resulting overlap R_{opt} at the minimum of this optimal potential:

$$R_{\text{opt}} = 1 - 2H(\mathcal{F}(R_{\text{opt}})). \tag{19}$$

The important issue to be examined is whether or not $R_{\text{opt}}(\alpha)$ saturates the bound given by the best binary. By comparison of equation (19) with (15), one immediately concludes that this is *not possible*, as long as \mathcal{F} is not a constant nor singular, since $R_{\text{opt}} = R_{bb}$ would imply that $\mathcal{F}(R_{bb}) = \mathcal{F}(R_B)$, and $R_{bb} = R_B$ is excluded by the first equality in equation (15). In general one thus has that $R_{\text{opt}} \leq R_{bb}$, since $\partial\mathcal{F}/\partial R \geq 0$. The equality is reached in asymptotic limits and for a special case (see below). For $R_{\text{opt}} \sim 0$ one has

$$\langle t \rangle_* \neq 0 \Rightarrow R_{\text{opt}} \simeq |\langle t \rangle_*| \sqrt{\frac{2\alpha}{\pi}} \tag{20}$$

$$\langle t \rangle_* = 0 \Rightarrow R \begin{cases} = 0 & \alpha \leq \alpha_c \\ \simeq \sqrt{C'(\alpha - \alpha_c)} & \alpha \geq \alpha_c \end{cases} \tag{21}$$

where the critical value now is $\alpha_c \equiv \pi\alpha_G/2$. Furthermore, the approach $R_{\text{opt}} \rightarrow 1$ is identical to that of R_{bb} , $1 - R_{\text{opt}} \simeq 1 - R_{bb}$. Therefore V_{opt} is successful only in the asymptotic limits $\alpha \rightarrow 0$ and $\alpha \rightarrow \infty$. Note that the second-order phase transition in equation (21) occurs at a larger value of α than for Gibbs learning.

The case $\mathcal{F}(R)$ independent of R , implying $R_{\text{opt}} = R_{bb}, \forall\alpha$, arises in a simple Gaussian scenario with a linear function U [19]. In this case, the best binary corresponds to clipped Hebbian learning. This seems to be the only case in which minimization of an optimal potential reproduces the best binary vector. We conclude that an optimal potential saturating the R_{bb} bound with $q \rightarrow 1$ cannot be constructed, in general. It motivates the search for alternative methods in discrete optimization. The main issue is to find new ways to incorporate information about the binary nature of the symmetry breaking vector, other than simply imposing the same binary constraint in the solution space. An interesting approach would be to try to construct a suitable potential for the continuous centre of mass J_B from which the best binary could be obtained by clipping. Whether such an approach is possible will be answered in future work.

The authors would like to thank Nestor Caticha for useful discussions and the organizers of the International Seminar on ‘Statistical Physics of Neural Networks’ (1999), held at the Max-Planck Institut für Physik komplexer Systeme (Dresden), during which part of this work was accomplished. We also acknowledge support from the FWO Vlaanderen and the Belgian IUAP programme (Prime Minister’s Office).

References

[1] Watkin T L H, Rau A and Biehl M 1993 *Rev. Mod. Phys.* **65**

- [2] Opper M and Kinzel W 1996 *Models of Neural Networks* vol III, ed E Domany *et al* (Berlin: Springer)
- [3] Biehl M and Mietzner A 1993 *Europhys. Lett.* **24** 421–6
- [4] Biehl M and Mietzner A 1994 *J. Phys. A: Math. Gen.* **27** 1885
- [5] Watkin T L H and Nadal J-P 1994 *J. Phys. A: Math. Gen.* **27** 1899–915
- [6] Reimann P and Van den Broeck C 1996 *Phys. Rev. E* **53** 3989
- [7] Reimann P, Van den Broeck C and Bex G J 1996 *J. Phys. A: Math. Gen.* **29** 3521
- [8] Van den Broeck C and Reimann P 1996 *Phys. Rev. Lett.* **76** 2188
- [9] Schietse J, Bouten M and Van den Broeck C 1995 *Europhys. Lett.* **32** 279–84
- [10] Mézard M, Parisi G and Virasoro M A 1987 *Spin Glass Theory and Beyond* (Singapore: World Scientific)
- [11] Györgyi G 1990 *Phys. Rev. A* **41** 7097
- [12] Watkin T L H 1993 *Europhys. Lett.* **21** 871
- [13] Reimann P 1997 *Europhys. Lett.* **40** 251–6
- [14] Opper M and Haussler D 1991 *Phys. Rev. Lett.* **66** 2677–80
- [15] Kinouchi O and Caticha N 1996 *Phys. Rev. E* **54** R54
- [16] Buhot A, Torres Moreno J-M and Gordon M B 1997 *Phys. Rev. E* **55** 7434–40
- [17] Buhot A and Gordon M B 1998 *Phys. Rev. E* **57** 3326–33
- [18] de Mattos C R 1997 *Aplicações de Mecânica Estatística ao Perceptron Binário e ao Processamento de Imagens*
PhD Thesis Universidade de São Paulo (in Portuguese)
Caticha N 1999 Private communication
- [19] Copelli M and Van den Broeck C 1999 *Preprint* available at <http://xxx.lanl.gov/abs/cond-mat/9910365>