

Receiver operating characteristics of perceptrons: Influence of sample size and prevalence

Ansgar Freking,¹ Michael Biehl,¹ Christian Braun,² Wolfgang Kinzel,¹ and Malte Meesmann²

¹*Institut für Theoretische Physik, Universität Würzburg, Würzburg, Germany*

²*Medizinische Universitätsklinik Würzburg, Würzburg, Germany*

(Received 21 May 1999)

In many practical classification problems it is important to distinguish false positive from false negative results when evaluating the performance of the classifier. This is of particular importance for medical diagnostic tests. In this context, receiver operating characteristic (ROC) curves have become a standard tool. Here we apply this concept to characterize the performance of a simple neural network. Investigating the binary classification of a perceptron we calculate analytically the shape of the corresponding ROC curves. The influence of the size of the training set and the prevalence of the quality considered are studied by means of a statistical-mechanics analysis. [S1063-651X(99)06911-1]

PACS number(s): 87.10.+e, 07.05.-t, 05.90.+m

I. INTRODUCTION

Classification problems in general and medical diagnostic tests in particular are often well suited for the application of neural networks [1–3]. The rule how to classify an item is generally not available, but can be derived from examples, e.g., patients with a known clinical status. The task for the network is to learn from these examples, i.e., to extract the implicit information in order to classify other items. We focus on medical diagnostic tests, where the aim is to discriminate between absence and presence of a certain disease or risk, yet our analysis is not restricted to this domain.

To gain further insight into the applicability of the neural network approach in a clinical setting we calculated the required size of the training set in relation to the prevalence of the disease considered. For this purpose we provide analytical expressions for the influence of prevalence on the shape of receiver operating characteristic curves.

II. DEFINITION OF DESCRIPTIVE MEASURES

For the evaluation of a diagnostic test, it is often important to distinguish false positive from false negative results. In this case the so-called generalization error as defined in the context of learning theory is not sufficient to assess the validity of the test.

A common way to summarize the results of a medical test is to list the frequencies of positive and negative test results in a 2×2 cross table where the columns correspond to the clinical status of the patients (see Table I). In the following we use the convention that all frequencies are counted relatively to the total number of patients, i.e., $a + b + c + d = 1$.

From the cross table one defines the ratios

$$\begin{aligned} u_+ &= \frac{a}{a+c}, & v_+ &= \frac{a}{a+b}, \\ u_- &= \frac{d}{b+d}, & v_- &= \frac{d}{c+d}, \end{aligned} \tag{1}$$

which are standard measures to describe the performance of a medical diagnostic test. The *sensitivity* u_+ gives the per-

centage of correctly classified diseased persons. Thus, a sensitivity of 100% means that any occurrence of the disease is detected by the test. The *specificity* u_- gives the analog ratio of the correctly classified persons without the disease. Whereas sensitivity and specificity have a more global meaning, physicians might be more interested in the ratios taken with respect to the rows. If the test is positive, the *positive predictive value* v_+ gives the probability to have the disease. The *negative predictive value* v_- tells the reliability of a negative test result. The fraction of diseased persons in the sample is called *prevalence*,

$$\lambda = a + c. \tag{2}$$

Usually, diagnostic tests yield a continuous valued quantity which is compared to a threshold value for binary classification. By varying this threshold, the test can be made more or less stringent to meet given requirements with respect to sensitivity or specificity.

A more stringent test gives high specificity, i.e., the fraction of correctly classified healthy persons is high. On the other hand, an increase in specificity usually results in a loss of sensitivity, i.e., more persons with the disease are missed.

This trade-off between specificity and sensitivity is described by the *receiver operating characteristic* (ROC) curve. An ROC curve is the plot of a test's true positive rate, i.e., the fraction $y_{\text{ROC}} = a/(a+c)$ of correctly classified persons with the disease, as a function of its false positive rate,

TABLE I. Cross table containing the relative frequencies of positive and negative test results subdivided with respect to the clinical status (disease vs no disease or risk vs no risk). a depicts correctly classified diseased persons, d represents correctly classified normals. $b+c$ gives the fraction of false classifications, which is often referred to as *generalization error*.

		Clinical status	
		With disease	Without disease
Test result	positive	a	b
	negative	c	d

the fraction $x_{\text{ROC}}=b/(b+d)$ of misclassified healthy persons. Using the definitions (1), it is obvious that $y_{\text{ROC}}=u_+$ and $x_{\text{ROC}}=(1-u_-)$.

Each point on the ROC curve corresponds to a certain stringency level of the test. If a diagnostic test makes exclusively positive classifications, the true positive rate is 1, but the same applies to the false positive rate. The corresponding point on the ROC curve is the upper right corner, $y_{\text{ROC}}=x_{\text{ROC}}=1$. In the opposite extreme, there are no positive classifications and therefore one has $y_{\text{ROC}}=x_{\text{ROC}}=0$.

For a test without predictive power, the ratio between positive classified persons with and without disease would be the same as the corresponding ratio of the whole test population; in this case one has

$$\frac{a}{b} = \frac{a+c}{b+d}$$

and the ROC curve follows the diagonal $y_{\text{ROC}}=x_{\text{ROC}}$. Any meaningful test should yield an ROC curve above this diagonal. By changing the threshold value of the test, one obtains a curve extending between the two corners $y_{\text{ROC}}=x_{\text{ROC}}=0$ and $y_{\text{ROC}}=x_{\text{ROC}}=1$.

The given conditions in medical practice often require a certain level of specificity for a test to be useful. If one thinks of a risky or cost expensive treatment, it is clear that the fraction of misclassified healthy persons should not exceed a certain limit. On the other hand, a diagnostic test should achieve some minimal value of sensitivity in order to be effective. Since ROC curves display precisely this interplay, they can be used to determine the threshold value for a certain diagnostic test. Additionally, ROC curves allow for comparison of several tests at the same level of specificity. For a detailed discussion of the use of ROC curves see, for instance, [4] and references therein.

III. DATA MODEL

In our analysis, each patient is represented by an N -dimensional feature vector $\mathbf{S} \in \mathbf{R}^N$. Each component of this vector can be thought of as a clinically relevant measure which characterizes the patient (e.g., blood pressure, heart rate). The patient's clinical status is coded by the binary quantity S_0 , where $S_0=+1$ means *with* and $S_0=-1$ means *without* disease.

Let us now assume, that diseased constellations in feature space can be separated from the healthy ones by a $(N-1)$ -dimensional hyperplane. This means, that the actual status S_0 of a patient with feature vector \mathbf{S} is given by

$$S_0 = \text{sgn}(\mathbf{B} \cdot \mathbf{S} - \theta), \quad (3)$$

where $\mathbf{B} \in \mathbf{R}^N$, $\mathbf{B}^2=1$, is a unit-vector perpendicular to the hyperplane; for obvious reasons \mathbf{B} is called the *rule* vector. θ is related to the prevalence, i.e., the fraction of actually diseased individuals. The components of \mathbf{S} are taken to be Gaussian distributed with zero mean and unit variance, therefore the overlap $y=\mathbf{B} \cdot \mathbf{S}$ of the feature vectors with the rule is Gaussian distributed as well,

$$p(y) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}y^2\right). \quad (4)$$

Note that in case of large N , Eq. (4) holds true under more general conditions according to the central limit theorem. Given the distribution (4) of y , the threshold θ is related to the prevalence through

$$\lambda = \Phi(\theta), \quad \text{with} \quad \Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x dt \exp\left(-\frac{1}{2}t^2\right). \quad (5)$$

IV. ROC CURVES OF A PERCEPTRON

In this section we consider a perceptron with fixed weight vector $\mathbf{J} \in \mathbf{R}^N$, $\mathbf{J}^2=1$. As does the rule vector, \mathbf{J} defines a hyperplane in the feature space, and

$$\sigma = \text{sgn}(\mathbf{J} \cdot \mathbf{S} - \gamma) \quad (6)$$

is the test result for the patient \mathbf{S} . We use the overlap $R = \mathbf{J} \cdot \mathbf{B}$ as the usual measure for the perceptron's knowledge about the rule. For a given R , the projection $x = \mathbf{J} \cdot \mathbf{S}$ of a feature vector on the perceptron vector \mathbf{J} and the projection $y = \mathbf{B} \cdot \mathbf{S}$ on the rule vector are jointly Gaussian distributed according to

$$p(x,y) = \frac{1}{2\pi\sqrt{1-R^2}} \exp\left(-\frac{1}{2} \frac{x^2 - 2Rxy + y^2}{1-R^2}\right). \quad (7)$$

In order to plot the ROC curves for the perceptron classification, we need to know the entries of the respective cross table. If we want to calculate, say b , which is the relative frequency of $[(S_0=-1) \wedge (\sigma=+1)]$, we have to perform the average of $\Theta(-S_0)\Theta(+\sigma)$ over the distribution of feature vectors \mathbf{S} . Since \mathbf{S} enters only through the scalar products x and y , we arrive in an average over Eq. (7),

$$\begin{aligned} b &= \int_{\gamma}^{\infty} dx \int_{-\infty}^{\theta} dy p(x,y) \\ &= \frac{1}{\sqrt{2\pi}} \int_{\gamma}^{\infty} dx \exp\left(-\frac{1}{2}x^2\right) \Phi\left(\frac{\theta - Rx}{\sqrt{1-R^2}}\right) \\ &=: \Psi(\gamma, \theta, R). \end{aligned} \quad (8)$$

All other entries of the cross table can be calculated in the same way and hence can be expressed through the function Ψ defined in Eq. (8),

$$a = \Psi(\gamma, -\theta, -R), \quad c = \Psi(-\gamma, -\theta, R), \quad (9)$$

$$\text{and} \quad d = \Psi(-\gamma, \theta, -R).$$

This allows us to express the quantities in Eq. (1) in terms of the perceptron threshold γ , the bias θ which is related to the prevalence (5), and the overlap R ,

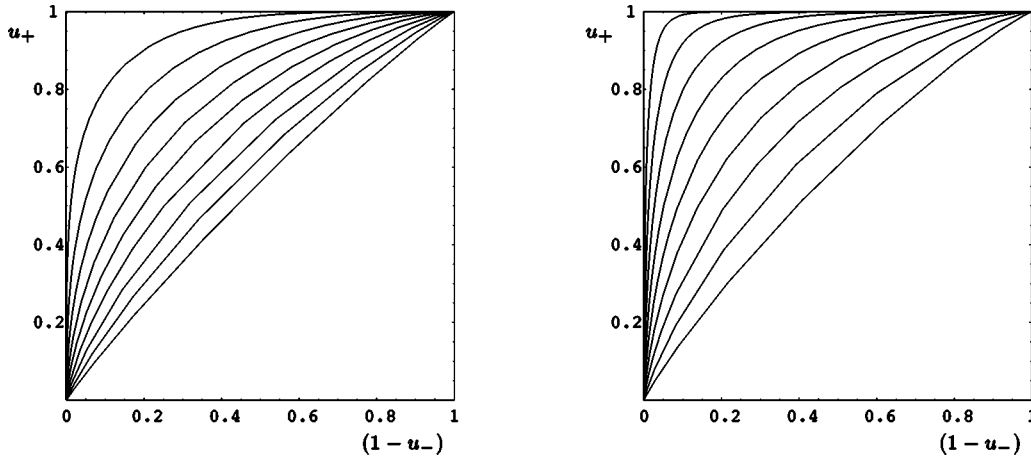


FIG. 1. ROC curves of perceptrons with overlaps $R=0.1, 0.2, \dots, 0.9$ (both panels); higher values of R give larger areas under the curves. For the left-hand panel the prevalence was set to $\lambda=50\%$, in the right-hand panel to $\lambda=1\%$. For a prevalence of 50%, the curves are symmetric with respect to the line $(0,1)-(1,0)$, which reflects the symmetry between *diseased* and *healthy* in this case.

$$\begin{aligned}
 u_+ &= \frac{\Psi(\gamma, -\theta, -R)}{\Phi(\theta)}, & v_+ &= \frac{\Psi(\gamma, -\theta, -R)}{\Phi(\gamma)}, \\
 u_- &= \frac{\Psi(-\gamma, \theta, -R)}{\Phi(-\theta)}, & v_- &= \frac{\Psi(-\gamma, \theta, -R)}{\Phi(-\gamma)}.
 \end{aligned}
 \tag{10}$$

For given external conditions, i.e., θ and R fixed, the performance measures in Eqs. (10) can only be tuned relative to each other by the choice of the perceptron threshold γ . A higher value of γ gives less positive classifications [cf. Eq. (6)] and therefore an increase of u_- with a concomitant decrease of u_+ .

The plot of u_+ , i.e., the fraction of positive classifications among the diseased persons, versus $(1-u_-)$, the fraction of positive classifications in the healthy subgroup, gives the ROC curve.

As $\gamma \rightarrow +\infty$ there are no positive test results and the respective point on the ROC curve is the lower left corner, $\gamma \rightarrow -\infty$ equivalently corresponds to the upper right corner. For finite γ the value of R determines the path in between the corners. The higher the value of R , the higher the sensitivity at a certain specificity and the larger the area under the curve.

Figure 1 presents ROC curves for two different prevalences (1% and 50%) and several values of R . In comparing the two plots, we find that a perceptron with a certain overlap R reaches higher sensitivities for all specificities when there are less patients with disease than without disease. The ROC curves for the case $\lambda=1\%$ are steeper than the ones for $\lambda=50\%$. On the other hand, it should be harder for the perceptron to achieve a certain knowledge about the rule, if there are less examples of one group at a fixed total amount of training examples.

V. LEARNING FROM A TRAINING SET

Up to now, we have considered a perceptron with fixed weights \mathbf{J} and a certain overlap R with the rule. The aim of this section is to include a learning process into the analysis. The perceptron gains knowledge about the rule by learning from examples provided in the training set. Consequently, we shall analyze the influence of the size of the training set

and the prevalence of the quality considered to the performance of the perceptron. This means we have to calculate the quantity R as a function of the total number P of training examples and the relative frequency of the label $S_0 = +1$. Finally, this gives us access to the desired accuracy measures defined in the previous sections.

We apply the standard statistical mechanics approach and consider the components of the perceptron weight vector \mathbf{J} as the N degrees of freedom of a physical system with energy

$$H(\mathbf{J}) = \sum_{\mu=1}^P \Theta(-\sigma^\mu S_0^\mu) = \sum_{\mu=1}^P \Theta[-(\mathbf{J} \cdot \mathbf{S}^\mu - \gamma)(\mathbf{B} \cdot \mathbf{S}^\mu - \theta)],
 \tag{11}$$

where the sum extends over all feature vectors \mathbf{S}^μ in the training set. The threshold γ is considered to be fixed during the training process.

In an ensemble of perceptrons at formal temperature $1/\beta$, a configuration \mathbf{J} occurs with the corresponding Gibbs-Boltzmann density, hence the term *Gibbs learning* has been coined for this scenario [5]. The quenched average over the randomness contained in the training data is performed using the replica method assuming replica symmetry. A sketch of the calculation is given in the appendix.

The limit $\beta \rightarrow \infty$ forces the system into its ground state. For $\gamma = \theta$ the energy of the ground state is $H=0$, independent of $\alpha = P/N$. This corresponds to an error-free classification of a training set of any size. In the limit $N \rightarrow \infty$ with finite normalized sample size $\alpha = P/N$, θ can be determined exactly from the prevalence and $\gamma = \theta$ is a valid choice for the perceptron threshold during training. Since we expect that this choice already gives the largest achievable overlap within the framework of Gibbs learning, we proceed without an explicit optimization of the learning strategy with respect to γ for given θ , and use $\gamma = \theta$ instead. Note that this particular choice of γ is only used to describe a specific training process and does not affect the role of the threshold as described in the preceding sections.

The overlap R defined in Sec. IV now plays the role of an order parameter. Together with the quantity q , which represents the typical overlap of two error free perceptron vectors

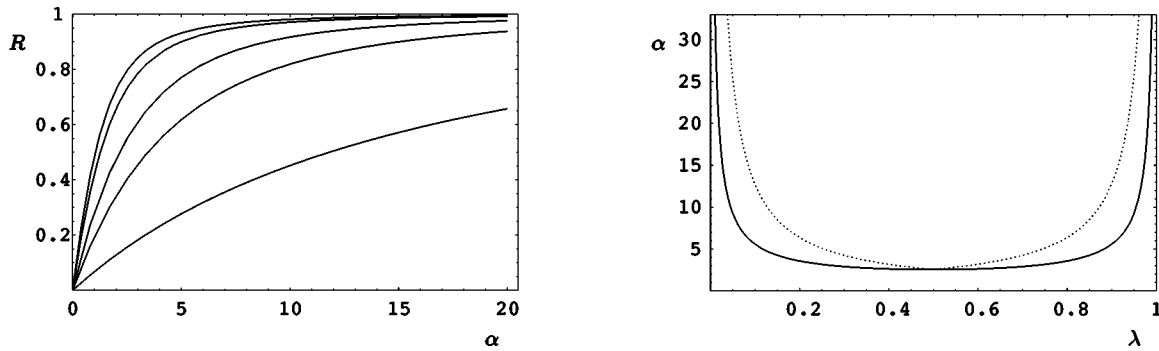


FIG. 2. On the left-hand panel, the overlap R of the perceptron vector \mathbf{J} with the rule vector \mathbf{B} is plotted as a function of the ratio $\alpha = P/N$. The individual curves correspond to the prevalences $\lambda = 0.01, 0.05, 0.10, 0.25,$ and 0.50 from bottom to top. Before having seen any example, i.e., at $\alpha=0$, \mathbf{J} is perpendicular to \mathbf{B} and the curves start at $R=0$ in any case. With increasing α , R increases and approaches asymptotically the value 1 as $\alpha \rightarrow \infty$. For a fixed α the largest value of R is achieved when both classes have equal weights, i.e., $\lambda = 0.50$. Nevertheless, only very asymmetric class weights cause pronounced reductions. This can also be seen from the right-hand part of the figure. Here the value of α required for a fixed overlap ($R=0.8$) is plotted against the prevalence (solid line). Whereas the curve is rather flat in the center region, it grows dramatically as $\lambda \rightarrow 0$ or $\lambda \rightarrow 1$. For comparison, the dotted line shows the value of α required to obtain the same number of examples in the smaller of the subgroups as there are for $\lambda=0.50$.

(A2), it is sufficient to describe the properties of the system in the thermodynamic limit $N \rightarrow \infty$.

In analogy to the case $\gamma = \theta = 0$, which was studied in [6], one can argue that the unknown vector \mathbf{B} coincides with equal probability with any error free \mathbf{J} . As a consequence the relation $q=R$ holds true, which would be violated in more general settings with $\gamma \neq \theta$.

The remaining saddle point condition $d\Gamma/dR=0$ [cf. Eq. (A4)] yields the so-called learning curve, i.e., R in dependence of α . Figure 2 displays such learning curves for different values of θ . As intuitively clear, the more examples are provided, the better the rule is captured by the perceptron. This works best if there are equal numbers of examples of both classification types. In case of very few examples for one of the classes, i.e., for prevalences far away from the balance $\lambda = 50\%$, it takes a huge total amount to gain knowledge about the rule. Nevertheless, for ratios from about $\lambda = 10\%$ up to $\lambda = 90\%$ there is a remarkably weak dependence of α on the prevalence.

A rough but widely used approximation for the dependence of the required sample size on the prevalence is given by the statement, that the smaller subgroup determines the actual sample size. This implies the dependence $\tilde{\alpha} \propto 1/\lambda$ for $\lambda < 0.50$, since the normalized number of examples of the smaller group is given by $(\alpha\lambda)$, in this case. Equivalently the dependence should be $\tilde{\alpha} \propto 1/(1-\lambda)$ in the case $\lambda > 0.50$. The right-hand graph of Fig. 2 shows the value of α necessary to obtain a certain overlap ($R=0.8$) as a function of λ ; for comparison, $\tilde{\alpha}$ is plotted as well. As expected the approximation $\tilde{\alpha}$ overestimates the actually needed sample size for any value of λ .

VI. HOW MUCH TRAINING DATA IS NECESSARY?

The relation between R , θ , and α as derived in the preceding section allows for a quantitative description of the influence of sample size and disease prevalence on the performance of a perceptron. As discussed in Sec. II, a quantity of interest might be the sensitivity of a test at a given speci-

ficity, i.e., the ordinate of the ROC curve at fixed abscissa $(1-u_-)$.

Proceeding on a given prevalence the overlap R can be calculated for any value of α (cf. Sec. V). From this, the perceptron threshold γ , which yields the desired specificity can be determined by using the ROC equations (10). Finally, by using Eq. (10) again, one obtains the sensitivity as a function of the normalized size $\alpha = P/N$ of the training set. This dependence is shown in Fig. 3. The corresponding positive predictive values, which are also presented in Fig. 3, can be obtained the same way.

From graphs like the ones in Fig. 3 one can easily read the required size of the training set for a desired test performance. Additionally, such graphs may indicate regions, where further collection of training examples gives only neglectable improvements to the test performance. The latter applies especially with respect to the positive predictive values.

VII. DISCUSSION

We applied the concept of ROC curves to describe the performance of a perceptron that realizes a threshold classification. For this, we revisited the thoroughly studied scenario where the perceptron learns from examples classified by a rule, which is of the same architecture as the perceptron itself (see, for instance, [1,5–8] and references therein).

Investigating the shape of ROC curves of a perceptron with fixed overlap R , we observed a pronounced dependence of sensitivity and specificity on the prevalence of the quality considered (cf. Fig. 1). This seems to contradict the frequently encountered statement that these quantities should be prevalence independent measures of validity, which is based on their definition by means of the cross table. It is important to realize that this predication refers to ROC curves obtained from differently composed validation sets for exactly the same classification problem. But even in this sense, the statement does not necessarily apply to realistic situations as discussed in [9].

By introducing the bias θ in the rule, we extended the

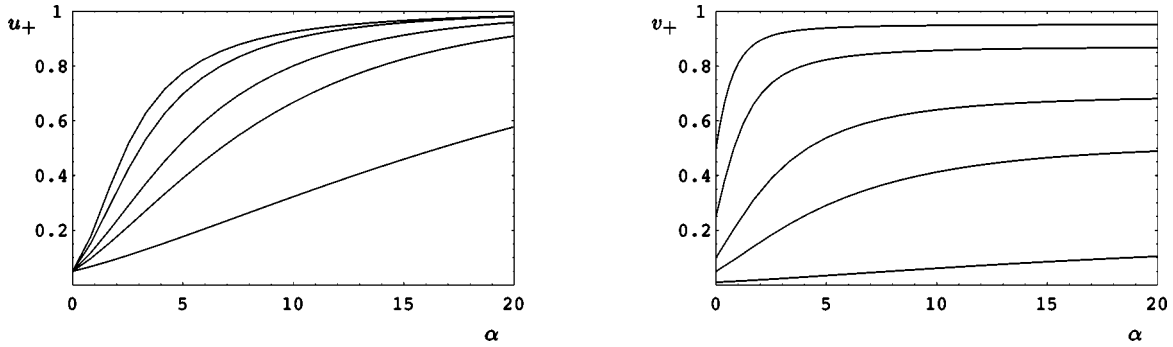


FIG. 3. Variation of the perceptron performance with increasing size of the training set. As in the left-hand panel of Fig. 2 the individual curves of both plots correspond to $\lambda = 0.01, 0.05, 0.10, 0.25,$ and 0.50 from bottom to top. The graph on the left shows the sensitivity at fixed specificity ($u_- = 0.95$). The sensitivity at $\alpha = 0$ is common to all curves. This is due to the fact that the ROC curve for a test without predictive power is given by the line of identity; the value 5% just reflects the considered specificity level of 95%. The α dependence of the sensitivity looks very similar to the shape of the learning curves shown in Fig. 2. The corresponding positive predictive values are displayed in the right-hand panel; the values at $\alpha = 0$ coincide with the respective values of the prevalence. The increase in v_+ becomes astonishingly slow, yet all curves tend to 1 as $\alpha \rightarrow \infty$.

analysis in [6] to situations where items of the two subgroups occur with different frequencies. Combining the results of this statistical-mechanics approach with standard measures of validity as defined in the context of biomedical statistics, we described analytically the influence of sample size and disease prevalence on the performance of a diagnostic test. It is hoped that this theoretical knowledge allows for more effective planning of clinical studies.

Our work is based on simple, yet rather general assumptions on the underlying data distribution and the classification scheme. Naturally, real world problems are more complicated in many respects. They are usually not completely learnable since the architecture of the rule is not known. In addition, classification tasks in practice will be affected by noise, i.e., the feature vectors can contain inaccurately measured values or the example classifications might be wrong themselves. The considered ideal training situation provides first insights and we expect the results to hold qualitatively in a wider range of settings. In particular, earlier studies of the perceptron have shown that the presence of noise does not prohibit the success ($R \rightarrow 1$) of appropriate training schemes, in principle [6,10,11]. Nevertheless, further research should incorporate such more realistic situations, including nonisotropic data distributions and more complicated classification schemes which require, for instance, the use of multilayer networks.

ACKNOWLEDGMENTS

We would like to thank E. Domany, I. Kanter, M. Opper, and K. Wegscheider for useful and stimulating discussions. This work was supported by the Deutsche Forschungsgemeinschaft (Me 799/3).

APPENDIX

The partition function corresponding to Eq. (11) reads

$$Z = \left(\prod_{j=1}^N \int dJ_j \right) \delta \left(\sum_{j=1}^N J_j^2 - 1 \right) \exp \left[-\beta \sum_{\mu} \Theta(-E^{\mu}) \right], \quad (\text{A1})$$

where

$$E^{\mu} = (\mathbf{J} \cdot \mathbf{S}^{\mu} - \gamma) \text{sgn } F^{\mu}, \quad \text{with } F^{\mu} = (\mathbf{B} \cdot \mathbf{S}^{\mu} - \theta).$$

The quenched free energy $-1/\beta \langle \ln Z \rangle$ can be performed using the identity

$$\langle \ln Z \rangle = \lim_{n \rightarrow 0} \frac{\partial}{\partial n} \langle Z^n \rangle.$$

For integer n , Z^n is the partition function of an n -fold replicated system. Its average can be calculated by means of a saddle point integration and involves the order parameters

$$R^a = \mathbf{J}^a \cdot \mathbf{B} \quad \text{and} \quad q_{ab} = \mathbf{J}^a \cdot \mathbf{J}^b \quad \text{for } a \neq b. \quad (\text{A2})$$

Here $a, b = 1, \dots, n$ denote the replica indices.

In analogy to [6] we assume replica symmetry, i.e., $R^a = R$ for all a and $q_{ab} = q$ for all $a \neq b$. Within this simplifying scheme it is straightforward to identify the solution in the limit $n \rightarrow 0$. We obtain, for $\beta \rightarrow \infty$, the quenched free energy

$$\begin{aligned} \Gamma = \frac{\partial}{\partial n} \Big|_{n=0} \frac{\langle Z^n \rangle}{N} &= \frac{q - R^2}{2(1 - q)} + \frac{1}{2} \ln(1 - q) \\ &+ \alpha \int_{-\infty}^{\infty} \frac{dx}{\sqrt{2\pi}} \exp \left(-\frac{x^2}{2} \right) \int_{-\infty}^{\infty} \frac{dF}{\sqrt{2\pi}} \\ &\times \exp \left[-\frac{(F + \theta)^2}{2} \right] \\ &\times \ln \Phi \left[\frac{x\sqrt{q - R^2} + |F|R - (\gamma - \theta R) \text{sgn } F}{\sqrt{1 - q}} \right]. \quad (\text{A3}) \end{aligned}$$

Further, we restrict the analysis to the case $\gamma = \theta$. As discussed in the text, the relation $q = R$ is satisfied in this case and one obtains the simplified free energy

$$\Gamma = \frac{1}{2} [R + \ln(1-R)] + \alpha \sqrt{\frac{1-R}{R}} \int_{-\infty}^{\infty} \frac{dt}{\sqrt{2\pi}} \left[\exp\left(t - \frac{\Theta}{\sqrt{R}}\right) + \exp\left(t + \frac{\Theta}{\sqrt{R}}\right) \right] \Phi\left(t \sqrt{\frac{R}{1-R}}\right) \ln \Phi\left(t \sqrt{\frac{R}{1-R}}\right). \quad (\text{A4})$$

-
- [1] J. A. Hertz, A. Krogh, and R. G. Palmer, *Introduction to the Theory of Neural Computation* (Addison-Wesley, Redwood City, CA, 1991).
- [2] C. M. Bishop, *Neural Networks for Pattern Recognition* (Oxford University Press, Oxford, 1995).
- [3] M. Akay, *Biol. Cybern.* **67**, 361 (1992).
- [4] L. E. Moses, D. Shapiro, and B. Littenberg, *Stat. Med.* **12**, 1293 (1993).
- [5] T. L. H. Watkin, A. Rau, and M. Biehl, *Rev. Mod. Phys.* **65**, 499 (1993).
- [6] G. Györgyi and N. Tishby, *Neural Networks and Spin Glasses* (World Scientific, Singapore, 1990), p. 3.
- [7] M. Opper and W. Kinzel, *Physics of Neural Networks* (Springer, Berlin, 1991), p. 149.
- [8] H. S. Seung, H. Sompolinsky, and N. Tishby, *Phys. Rev. A* **45**, 6056 (1992).
- [9] H. Brenner and O. Gefeller, *Stat. Med.* **16**, 981 (1997).
- [10] M. Biehl, P. Riegler, and M. Stechert, *Phys. Rev. E* **52**, R4624 (1995).
- [11] M. Opper and D. Haussler, *Proceedings of the 4th Annual Workshop on Computational Learning Theory* (Morgan Kaufman, Santa Cruz, 1991), p. 75.