

# Classification and Feature Selection on Matrix Data with Application to Gene-Expression Analysis

Sepp Hochreiter  
hochreit@cs.tu-berlin.de

Klaus Obermayer  
Fakultät für Elektrotechnik und Informatik  
Technische Universität Berlin  
Franklinstr. 28/29, 10587 Berlin, Germany  
oby@cs.tu-berlin.de

We consider the classification task for datasets which are described by matrices. Rows and columns of these matrices correspond to objects where row and column objects may be from different sets and column objects are labeled. Data matrix entries express relationships between row and column objects and are produced by an unknown kernel. These kernels represent dot products in some (unknown) feature space. In this feature space a linear column object classifier should be constructed. However the dot products between column objects are not available. Therefore standard support vector techniques cannot be utilized. We derive a new objective function for model selection in such a feature space according to the principle of structural risk minimization. The new objective allows the analysis of matrices which are not positive definite, and not even symmetric or square. Because row objects can be interpreted as features and our method assigns support vector weights to the row objects can be used for feature selection. An additional constraint, which imposes sparseness on the row objects resulting in few selected features. We analyse data obtained from DNA microarrays, where “column” objects correspond to samples, “row” objects correspond to genes and matrix elements correspond to expression levels. Benchmarks are conducted using standard one-gene classification and support vector machines and K-nearest neighbors after standard feature selection. Our new method extracts a sparse set of genes and provides superior classification results.

## Introduction

Many data in the real world are characterized by matrices, e.g. data in Microsoft Excel tables, transaction data or micro arrays. Rows and columns of these matrices indicate the relationship between objects. One typical case are so-called pairwise data, where rows as well as columns of the matrix represent the objects of the dataset (Fig. 1a) and where the entries of the matrix denote similarity values. Another typical case occurs, if objects are described by a set of features (Fig. 1b). In this case, the column objects are the objects to be characterized, the row objects correspond to their features and the matrix elements denote the strength with which a feature is expressed in a particular object.

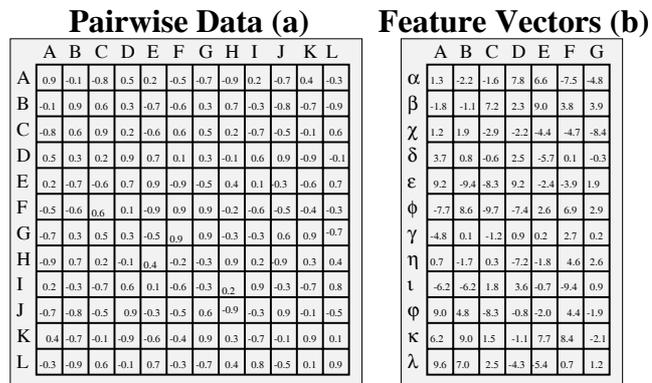


Figure 1: Two matrix data examples: pairwise data and feature vectors (see text).

In the following we consider the task of learning a classification problem on matrix data. We consider the case that class labels are assigned to the column objects of the training set. Given the matrix and the class labels we then want to construct a classifier with good generalization properties. From all the possible choices we select classifiers from the support vector

machine (SVM) family [8, 4] and we use the principle of structural risk minimization [8] for model selection - because of its recent success [4] and its theoretical properties [8].

A serious problem arises when the number of features becomes large and comparable to the number of objects: Without feature selection, SVMs are prone to overfitting, despite the complexity regularization which is implicit in the learning method [1]. Rather than being sparse in the number of support vectors, the classifier should be sparse in the number of features used for classification. This relates to the result [8] that the number of features provide an upper bound on the number of “essential” support vectors. We show how to limit the number of features and, therefore, how to increase the generalization capability, i.e. the performance, of support vector approaches.

## The New Objective And The Resulting Optimization Formulation

Given are two sets  $\mathcal{X}$  and  $\mathcal{Z}$  of objects, which are described by feature vectors  $\mathbf{x}$  and  $\mathbf{z}$ . Based on the feature vectors  $\mathbf{x}$  we construct a linear classifier defined as  $f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$ , where  $\langle \cdot, \cdot \rangle$  denotes a dot product. The zero isoline of  $f$  is a hyperplane with unit normal vector  $\hat{\mathbf{w}}$  and perpendicular distance  $b/\|\mathbf{w}\|_2$  from the origin. The hyperplane’s margin  $\gamma$  is given by  $\gamma = \min_{\mathbf{x} \in \mathcal{X}} |\langle \hat{\mathbf{w}}, \mathbf{x} \rangle + b/\|\mathbf{w}\|_2|$ . Setting  $\gamma = \|\mathbf{w}\|_2^{-1}$  allows us to treat normal vectors  $\mathbf{w}$  which are not normalized, if the margin is normalized to 1. The hyperplane with largest margin is then obtained by minimizing  $\|\mathbf{w}\|_2^2$  for a margin which equals 1. It has been shown [7, 6, 5] that the generalization error of a linear classifier can be bounded from above with probability  $1 - \delta$  by the bound  $\mathcal{B}$ ,

$$(1) \quad \mathcal{B}(L, a/\gamma, \delta) = \frac{2}{L} \left( \log_2 \left( EN \left( \frac{\gamma}{2a}, \mathcal{F}, 2L \right) \right) + \log_2 \left( \frac{4L a}{\delta \gamma} \right) \right),$$

provided that the training classification error is zero and  $f(\mathbf{x})$  is bounded by  $-a \leq f(\mathbf{x}) \leq a$  for all  $\mathbf{x}$  drawn iid.  $L$  denotes the number of training objects  $\mathbf{x}$  and  $EN(\epsilon, \mathcal{F}, L)$  the expected  $\epsilon$ -covering number of a class  $\mathcal{F}$  of functions that map data objects from  $T$  to  $[0, 1]$  (see [5]). To select a classifier with good generalization properties we suggest to minimize  $a/\gamma$  under proper constraints.  $a$  is not known in general. Therefore we approximate  $a$  by the range  $m = 0.5 (\max_i \langle \hat{\mathbf{w}}, \mathbf{x}^i \rangle - \min_i \langle \hat{\mathbf{w}}, \mathbf{x}^i \rangle)$  of values in the training set and minimize  $\mathcal{B}(L, m/\gamma, \delta)$  instead of eq. (1).

Let  $\mathbf{X} := (\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^L)$  be the matrix of feature vectors of  $L$  objects from the set  $\mathcal{X}$  and  $\mathbf{Z} := (\mathbf{z}^1, \mathbf{z}^2, \dots, \mathbf{z}^P)$  be the matrix of feature vectors of  $P$  objects from the set  $\mathcal{Z}$ . The objects of set  $\mathcal{X}$  are labeled, and we summarize all labels using a diagonal label matrix  $\mathbf{Y}$ . The feature vectors  $\mathbf{X}$  and  $\mathbf{Z}$  are unknown, but the matrix  $\mathbf{K} := \mathbf{X}^T \mathbf{Z}$  of the corresponding scalar products is given. The training set is  $\mathbf{K}, \mathbf{Y}$ . The principle of structural risk minimization is implemented by minimizing an upper bound on  $(m/\gamma)^2$  given by  $\|\mathbf{X}^T \mathbf{w}\|_2^2$ , as can be seen from  $m/\gamma \leq \|\mathbf{w}\|_2 \max_i |\langle \hat{\mathbf{w}}, \mathbf{x}^i \rangle| \leq \sqrt{\sum_i (\langle \mathbf{w}, \mathbf{x}^i \rangle)^2} = \|\mathbf{X}^T \mathbf{w}\|_2$ . In the following optimization formulation the constraints  $f(\mathbf{x}^i) = y^i$  imposed by the training set are taken into account including slack variables  $\xi_i^+, \xi_i^- \geq 0$ :

$$(2) \quad \begin{aligned} \min_{\mathbf{w}, b, \xi^+, \xi^-} \quad & \frac{1}{2} \|\mathbf{X}^T \mathbf{w}\|_2^2 + M^+ \mathbf{1}^T \xi^+ + M^- \mathbf{1}^T \xi^- \\ \text{s.t.} \quad & \mathbf{Y}^{-1} (\mathbf{X}^T \mathbf{w} + b\mathbf{1}) - \mathbf{1} + \xi^+ \geq \mathbf{0} \\ & \mathbf{Y}^{-1} (\mathbf{X}^T \mathbf{w} + b\mathbf{1}) - \mathbf{1} - \xi^- \leq \mathbf{0}, \xi^+, \xi^- \geq \mathbf{0}. \end{aligned}$$

$M^+$  penalizes wrong classification and for classification  $M^-$  may be set to zero.

Let  $\tilde{\alpha}^+, \tilde{\alpha}^-$  be the dual variables for the constraints imposed by the training set,  $\tilde{\alpha} := \tilde{\alpha}^+ - \tilde{\alpha}^-$ , and  $\alpha$  a vector with  $\tilde{\alpha} = \mathbf{Y} (\mathbf{X}^T \mathbf{Z}) \alpha$ . Two cases must be treated:  $\alpha$  is not unique

and  $\alpha$  does not exist. If  $\alpha$  is not unique we choose  $\alpha$  which is most sparse in its components (see below). If  $\alpha$  does not exist we set  $\alpha = (\mathbf{Z}^T \mathbf{X} \mathbf{Y}^{-T} \mathbf{Y}^{-1} \mathbf{X}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{X} \mathbf{Y}^{-T} \tilde{\alpha}$ , where  $\mathbf{Y}^{-T} \mathbf{Y}^{-1}$  is the identity. We derive for the optimal values:  $\mathbf{w} = \mathbf{Z} \alpha$ ,  $0 = \mathbf{1}^T (\mathbf{X}^T \mathbf{Z}) \alpha$ ,  $\tilde{\alpha}_i \leq M^+$ , and  $-\tilde{\alpha}_i \leq M^-$ . For  $M^+ = M^- = M$  and  $C := M \|\mathbf{Y} (\mathbf{X}^T \mathbf{Z})\|_{row}^{-1}$  we get  $\|\alpha\|_\infty \leq C$  and  $\|\tilde{\alpha}\|_\infty \leq M$ , where  $\|\cdot\|_{row}$  is the row-sum norm. With an additional sparseness term  $\epsilon \|\alpha\|_1$ , we obtain the dual problem:

$$(3) \quad \begin{aligned} \min_{\alpha} \quad & \frac{1}{2} (\alpha^+ - \alpha^-)^T \mathbf{K}^T \mathbf{K} (\alpha^+ - \alpha^-) - \mathbf{1}^T \mathbf{Y} \mathbf{K} (\alpha^+ - \alpha^-) + \epsilon \mathbf{1}^T (\alpha^+ + \alpha^-) \\ \text{s.t.} \quad & \mathbf{1}^T \mathbf{K} (\alpha^+ - \alpha^-) = 0, C \mathbf{1} \geq \alpha^+, \alpha^- \geq \mathbf{0} . \end{aligned}$$

If a classifier has been selected according to eq. (3), a new example  $\mathbf{u}$  is classified according to the sign of  $f(\mathbf{u}) = \langle \mathbf{w}, \mathbf{u} \rangle + b = \sum_{i=1}^P \alpha_i \langle \mathbf{z}^i, \mathbf{u} \rangle + b$ .

The optimal classifier is selected by optimizing eq. (3), and as long as  $a = m$  holds true for all possible objects  $\mathbf{x}$  (which are assumed to be drawn iid), the generalization error is bounded by eq. (1). If outliers are rejected, condition  $a = m$  can always be enforced. For large training sets the number of rejections is small: The probability  $P\{|\langle \mathbf{w}, \mathbf{x} \rangle| > m\}$  that an outlier occurs can be bounded with confidence  $1 - \delta$  using the additive Chernoff bounds (e.g. [8]):  $P\{|\langle \mathbf{w}, \mathbf{x} \rangle| > m\} \leq \sqrt{\frac{-\log \delta}{2L}}$ . But note, that not all outliers are misclassified, and the trivial bound on the generalization error is still of the order  $L^{-1}$ .

## Kernel Functions Are Dot Products

In the last section we have assumed that the matrix  $\mathbf{K}$  is derived from scalar products between the feature vectors  $\mathbf{x}$  and  $\mathbf{z}$  which describe the objects from the sets  $\mathcal{X}$  and  $\mathcal{Z}$ . For all practical purposes, however, the only information available is summarized in the matrices  $\mathbf{K}$  and  $\mathbf{Y}$ . The feature vectors are not known and it is even unclear whether they exist. In order to apply the results of previous section to practical problems the following question remains to be answered: What are the conditions under which the measurement operator  $k(\cdot, \mathbf{z})$  can indeed be interpreted as a scalar product between feature vectors and under which the matrix  $\mathbf{K}$  can be interpreted as a Gram matrix?

The following theorem states under which assumptions on  $k(\cdot, \mathbf{z})$  this kernel represents a dot product. Let  $L^2(H)$  denote the set of functions  $h$  from  $H$  with  $\int h^2(\mathbf{x}) d\mathbf{x} < \infty$ .

**Theorem 1 (Singular Value Expansion)** *Let  $H_1$  and  $H_2$  be Hilbert spaces. Let  $\alpha$  be from  $L^2(H_1)$  and let  $k$  be a kernel from  $L^2(H_2, H_1)$  which defines a Hilbert-Schmidt operator  $T_k : H_1 \rightarrow H_2$  ( $T_k \alpha$ )( $\mathbf{x}$ ) =  $f(\mathbf{x}) = \int k(\mathbf{x}, \mathbf{z}) \alpha(\mathbf{z}) d\mathbf{z}$ . Then there exists an expansion  $k(\mathbf{x}, \mathbf{z}) = \sum_n s_n e_n(\mathbf{z}) g_n(\mathbf{x})$  which converges in the  $L^2$ -sense. The  $s_n \geq 0$  are the singular values of  $T_k$ , and  $e_n \in H_1$ ,  $g_n \in H_2$  are the corresponding orthonormal functions.*

## DNA Microarray Data Analysis

We apply our new method eq. (3) to the DNA microarray data published in [3]. Column objects are samples from different brain tumors of the medullablastoma kind. The samples were obtained from 60 patients, which were treated in a similar way and the samples were labeled according to whether a patient responded well to chemo- or radiation therapy. Row objects correspond to genes. For every sample-gene pair a snapshot of the level of gene expression was measured. This gave rise to a  $60 \times 7,129$  real valued sample-gene matrix.

The task is now to construct a classifier which predicts therapy outcome on the basis of samples taken from new patients. The major problem of this classification task is the

limited number of samples - given the large number of genes. Therefore, feature selection is a prerequisite for good generalization [2]. For feature selection we apply our new method on a  $59 \times 7,129$  matrix, where one column object was withheld to avoid biased feature selection.

**Table 1: Prediction of the outcome of chemo therapy based on gene-expression. Our method (P-SVM) is compared to standard methods via the number of wrong classifications (“E”). Standard method results are taken from [3].**

Meth	F	E	Meth	C	F	E	Meth	C	F	E
TrkC	1	<b>20</b>	KNN		8	<b>13</b>	P-SVM	0.1	25/30/35	<b>8/5/8</b>
SVM		<b>15</b>	C2			<b>12</b>	P-SVM	0.15	25/30/35	<b>8/5/7</b>
C1		<b>14</b>	P-SVM	0.05	25/30/35	<b>8/5/8</b>	P-SVM	0.2	25/30/35	<b>6/4/6</b>

Table 1 shows benchmark results for DNA microarray data where the classification error given by the number of wrong classifications (“E”) for different numbers of selected features (“F”) and for different values of the parameter  $C$ . Data are provided for “TrkC”-Gene classification, standard SVMs, weighted “TrkC”/SVM (C1), K nearest neighbor (KNN), combined SVM/TrkC/KNN (C2), and our procedure (P-SVM). Standard methods feature selection was based on the correlation of features with classes using signal-to-noise-statistics and  $t$ -statistics according to [1]. The feature selection procedure (also a classifier) had its lowest misclassification rate between 20 and 40 features. Our method clearly outperforms standard methods — the number of misclassification is down by a factor of 2 (for 30 selected genes) using p-SVM.

## REFERENCES

- [1] R. Golub, D. K. Slonim, P. Tamayo, C. Huard, M. Gaasenbeek, J. P. Mesirov, H. Coller, M. Loh, J. R. Downing, M. A. Caligiuri, C. D. Bloomfield, and E. S. Lander. Molecular classification of cancer: Class discovery and class prediction by gene expression monitoring. *Science*, 286(5439):531–537, 1999.
- [2] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik. Gene selection for cancer classification using support vector machines. *Mach. Learn.*, 46:389–422, 2002.
- [3] S. L. Pomeroy, P. Tamayo, M. Gaasenbeek, L. M. Sturla, M. Angelo, M. E. McLaughlin, J. Y. H. Kim, L. C. Goumnerova, P. M. Black, C. Lau, J. C. Allen, D. Zagzag, J. M. Olson, T. Curran, C. Wetmore, J. A. Biegel, T. Poggio, S. Mukherjee, R. Rifkin, A. Califano, G. Stolovitzky, D. N. Louis, J. P. Mesirov, E. S. Lander, and T. R. Golub. Prediction of central nervous system embryonal tumour outcome based on gene expression. *Nature*, 415(6870):436–442, 2002.
- [4] B. Schölkopf and A. J. Smola. *Learning with kernels — Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, Cambridge, 2002.
- [5] J. Shawe-Taylor, P. L. Bartlett, R. C. Williamson, and M. Anhtony. A framework for structural risk minimisation. In *Comp. Learn. Th.*, pages 68–76, 1996.
- [6] J. Shawe-Taylor, P. L. Bartlett, R. C. Williamson, and M. Anhtony. Structural risk minimization over data-dependent hierarchies. *IEEE Trans. Inf. Theory*, 44:1926–1940, 1998.
- [7] J. Shawe-Taylor and N. Cristianini. On the generalisation of soft margin algorithms. Technical Report NC2-TR-2000-082, Dep. Computer Science, Royal Holloway, London, 2000.
- [8] V. Vapnik. *The nature of statistical learning theory*. Springer, NY, 1995.