

Discriminatory Data Mapping by Matrix-Based Supervised Learning Metrics

M. Strickert^{1,*}, P. Schneider², J. Keilwagen¹,
T. Villmann³, M. Biehl², and B. Hammer⁴

¹ Leibniz Institute of Plant Genetics and Crop Plant Research Gatersleben

² Institute for Mathematics and Computing Science, University of Groningen

³ Research group Computational Intelligence, University of Leipzig

⁴ Institute of Computer Science, Technical University of Clausthal

stricker@ipk-gatersleben.de

Abstract. Supervised attribute relevance detection using cross-comparisons (SARDUX), a recently proposed method for data-driven metric learning, is extended from dimension-weighted Minkowski distances to metrics induced by a data transformation matrix Ω for modeling mutual attribute dependence. Given class labels, parameters of Ω are adapted in such a manner that the inter-class distances are maximized, while the intra-class distances get minimized. This results in an approach similar to Fisher’s linear discriminant analysis (LDA), however, the involved distance matrix gets optimized, and it can be finally utilized for generating discriminatory data mappings that outperform projection pursuit methods with LDA index. The power of matrix-based metric optimization is demonstrated for spectrum data and for cancer gene expression data.

Keywords: Supervised feature characterization, adaptive matrix metrics, attribute dependence modeling, projection pursuit, LDA.

1 Introduction

Learning metrics constitute one of the most exciting topics in machine learning research [11,17,18]. The potential of metric adaptation needs exploration for facing challenges connected to the curse of dimensionality in high-throughput biomedical data sets. Mass spectra, gene expression arrays, or 2D electrophoretic gels, given as vectors of real-value measurements, are often characterized by only a low number of available experiments as compared to their huge dimensionality. Data-driven adaptation of a data metric can be used in many helpful ways. Applications of metric optimization range from attribute weighting via dimension reduction to data transformations into task-specific spaces.

The adaptive Euclidean distance $d_{\lambda}(\mathbf{x}, \mathbf{y}) = (\sum_{i=1}^q \lambda_i (x_i - y_i)^2)^{1/2}$, for example, relates attribute characterization to the choice of attribute scaling factors λ_i beneficial for the separation of labeled and unlabeled data in supervised and

* Corresponding author.

unsupervised manners. This aim is shared with projection pursuit methods for which matrix parameters of a linear projection mapping are optimized with respect to criteria of data spreading and clusterability of the low-dimensional projections [4]. In this work, evaluation takes place in a space of original dimensionality where only the comparison criterion, the metric, is changed. If desired, attribute-related parameters with low impact, for example expressed as low scaling factors, can be pruned for dimension reduction, after adaptation. For the parametric Euclidean distance, small attribute scaling factors λ_i would indicate negligible attributes. Scaling factors can be also used for transforming the data to the non-adapted Euclidean space for further utilization of standard Euclidean methods. This kind of attribute characterization is different from many other methods for feature extraction [6], such as the recently suggested Iterative Relief algorithm [16] for which the attribute weights do not coincide with a canonic rescaling of the data space.

Matrix-based metrics help to extend the view of individual attribute processing to a model of dependence between pairs of attributes. Generally, matrix methods can be used for optimizing linear data transformations aiming at criteria related to the data spreading. In the unsupervised case, interesting transforms include sphering of the data covariance matrix to the unity matrix, or the projection of data to directions of maximum variance (PCA) or to directions along maximum non-Gaussianity (ICA) [10]. The projection pursuit method [4] is a very flexible approach to extract projections of interest by optimizing a target function, called the index of the projection. Such indices exist for unsupervised cases aiming at mappings to continuous or sharply clustered views. In addition, there are supervised indices like projection entropy and class separability according to linear discriminant analysis (LDA) criteria. A good environment for the study of projection pursuit and other matrix methods is, for example, provided by the free R statistical language with rGGobi and classPP packages [2,12], an application of classPP for the visualization of gene expression data, is provided in [3].

An alternative view on seeking optimum data transformations is data-driven adaptation of the data metric or, more generally, of a data similarity measure. Learning vector quantization (LVQ), for example, can implement metric adaptation for better data classification by boosting class-separating attributes between data prototype vectors. The generalized relevance LVQ method (GRLVQ) realizes such metric adaptation by using a misclassification cost function – minimized by gradient descent – making use of data labels for attribute rescaling [8]. For Euclidean distances, large-margin optimization is realized, but also non-Euclidean similarity-measures profit from parameter adaptation [7,14]. Recently, matrix learning has been integrated into the GRLVQ framework for modeling attribute-attribute dependencies by generalized Mahalanobis distance [13]. This allows to express scalings of the data space along arbitrary directions, and very good classification accuracies are obtained on difficult classification problems ranging from spectrum classification to image segmentation.

The success of matrix metric adaptation in GRLVQ classifiers initiated the present work. Here, no classifier will be build though, but the data space will

be transformed: directions in the data space relevant to data label separation will be emphasized, while within-class variations will be damped. After all, relevant combinations of attributes, trained and expressed in form of a matrix, are identified for the discrimination task. Since only few data samples per class can be expected in costly and time-consuming biomedical studies, prototype-based data abstraction, like provided by GRLVQ, is avoided in order to keep maximum information. For the analyzed data sets it turned out that the transform matrix could be effectively compressed to only a few prominent eigenvectors, possibly only one, without significant loss of metric structure. After all, we are able to compute relatively compact discriminatory data models that allow hypotheses generation for supporting biomedical experts.

2 Method

Data. The q -dimensional row input vector $\mathbf{x} \in \mathbb{R}^{1 \times q}$ is taken from a data set containing n data vectors $\{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n\}$. The proposed metric adaptation requires that each vector \mathbf{x}^k is labeled with one class-specific index $c(k)$, assuming at least two unique classes in the whole data set.

Metric. Most essential is the definition of the matrix-based metric $d_{\Omega}^{ij} \in [0; \infty)$ between data vectors \mathbf{x}^i and \mathbf{x}^j :

$$d_{\Omega}^{ij} = d_{\Omega}(\mathbf{x}^i, \mathbf{x}^j) = (\mathbf{x}^i - \mathbf{x}^j) \cdot \mathbf{A} \cdot (\mathbf{x}^i - \mathbf{x}^j)^{\top}, \quad (\mathbf{A} = \mathbf{\Omega} \cdot \mathbf{\Omega}^{\top}) \in \mathbb{R}^{q \times q}. \quad (1)$$

Choosing the identity matrix $\mathbf{A} = \mathbf{\Omega} = \mathbf{I}$ induces the special case of the squared Euclidean distance; other diagonal matrices yield weighted squared Euclidean distances as discussed in [15]. Generally, metrics are obtained for arbitrary positive-definite matrices \mathbf{A} . Then the value expressed by $\mathbf{\Delta} \cdot \mathbf{A} \cdot \mathbf{\Delta}^{\top} \geq 0$, getting zero only for trivial difference vectors $\mathbf{\Delta} = \mathbf{0}$, is a metric. It is known that in context of metrics non-symmetric positive-definite matrices can be replaced by equivalent symmetric positive-definite matrices. Since any symmetric positive-definite matrix \mathbf{A} can be decomposed by Cholesky decomposition into a product of a lower triangular matrix and its transposed, it is in principle sufficient to learn a lower triangular matrix $\mathbf{\Omega}$ for expressing \mathbf{A} . Alternatively, symmetric positive-definite \mathbf{A} can be represented by the self-product $\mathbf{A} = \mathbf{\Omega} \cdot \mathbf{\Omega}$ of a symmetric $\mathbf{\Omega}$ [13]. Here, we consider products $\mathbf{A} = \mathbf{\Omega} \cdot \mathbf{\Omega}^{\top}$ with arbitrary $\mathbf{\Omega} \in \mathbb{R}^{q \times q}$. These full matrices $\mathbf{\Omega}$, possess more adaptive matrix elements than degrees of freedom needed for expressing the product solution space of \mathbf{A} . For the data sets discussed, the interaction of matrix element pairs Ω_{ij} and Ω_{ji} leads to a faster convergence of \mathbf{A} during optimization, compared to the convergence properties obtained for symmetric or triangular matrices $\mathbf{\Omega}$.

Note that for some the data \mathbf{A} might become positive-semidefinite during optimization, i.e. $\mathbf{\Delta} \cdot \mathbf{A} \cdot \mathbf{\Delta}^{\top} = 0$ with difference vectors $\mathbf{\Delta} \neq \mathbf{0}$. Then, the metric property gets relaxed to a mathematical distance with vanishing self-scalar

product $(\Delta \cdot \Omega) \cdot (\Delta \cdot \Omega)^\top = \langle \Delta \cdot \Omega, \Delta \cdot \Omega \rangle = 0$ becoming zero for certain configurations of Ω with $\Delta \cdot \Omega = \mathbf{0}$, else positive.

Adaptation. Driven by the goal to minimize within-class differences while maximizing between class differences, the following cost function is minimized over pairs of all n data items:

$$s(\Omega) := \frac{\sum_{i=1}^n \sum_{j=1}^n d_{\Omega}(\mathbf{x}^i, \mathbf{x}^j) \cdot \delta_{ij}}{\sum_{i=1}^n \sum_{j=1}^n d_{\Omega}(\mathbf{x}^i, \mathbf{x}^j) \cdot (1 - \delta_{ij})} = \frac{d_C}{d_D}, \quad \delta_{ij} = \begin{cases} 0 : c(i) \neq c(j) \\ 1 : c(i) = c(j) \end{cases} \quad (2)$$

Distances d_{Ω}^{ij} between data vectors \mathbf{x}^i and \mathbf{x}^j depend on the adaptive matrix parameters $\Omega = (\Omega_{kl})_{\substack{k=1 \dots q \\ l=1 \dots m}}$ of interest. The numerator represents within-class data scatter, which should be small; the denominator is related to inter-class distances, which should be large. Thus, optimization of $s(\Omega)$ handles both parts of the fraction simultaneously. Compromise solutions must be found in cases when within-class variation, potentially caused by outliers, needs compression, while inter-class separability would require inflation.

Using the chain rule, the cost function $s(\Omega)$ is iteratively optimized by gradient descent $\Omega \leftarrow \Omega - \gamma \cdot \frac{\partial s(\Omega)}{\partial \Omega}$, which requires adaptation of the matrix Ω in small steps γ into the direction of steepest gradient

$$\frac{\partial s(\Omega)}{\partial \Omega} = \sum_{i=1}^n \sum_{j=1}^n \frac{\partial s(\Omega)}{\partial d_{\Omega}^{ij}} \cdot \frac{\partial d_{\Omega}^{ij}}{\partial \Omega}. \quad (3)$$

The quotient rule applied to the fraction $s(\Omega) = d_C/d_D$ in Eqn. 2 yields

$$\frac{\partial s(\Omega)}{\partial d_{\Omega}^{ij}} = \frac{\delta_{ij} \cdot d_D}{d_D^2} + \frac{(\delta_{ij} - 1) \cdot d_C}{d_D^2} = \begin{cases} 1/d_D : c(i) = c(j) \\ -d_C/d_D^2 : c(i) \neq c(j) \end{cases}. \quad (4)$$

The right factor in Eqn. 3 is obtained by matrix derivative of Eqn. 1:

$$\frac{\partial d_{\Omega}^{ij}}{\partial \Omega} = 2 \cdot (\mathbf{x}^i - \mathbf{x}^j)^\top \cdot (\mathbf{x}^i - \mathbf{x}^j) \cdot \Omega. \quad (5)$$

If desired, adaptation can be restricted to certain structures of Ω , such as to the lower triangular elements. In that case, undesired elements must be initially masked out by zeros in Ω . Additionally, the same zero masking pattern must be applied to the matrix resulting from Eqn. 5, because the equation calculates $\partial d_{\Omega}^{ij}/\partial \Omega$ correctly only for full adaptive matrices Ω . By consistent masking operations, though, the matrix of derivatives is mathematically correct. In practice, the gradient from Eqn. 3 is computed and reused as long the cost function decreases. Potential increase of $s(\Omega)$ triggers a recomputation of the gradient. The step size γ is dynamically determined as the initial size γ_0 , being exponentially cooled down by rate η , divided by the maximum absolute element in the matrix $\partial s(\Omega)/\partial \Omega$.

Initialization. Empirically, the initial step size γ_0 can be chosen from the interval $[0.05; 1)$, such as 0.75 in the conducted experiments. The number of iterations

should be set to a value between 50 and 1000, depending on the saturation characteristics of the cost function. The exponential cooling rate should diminish the original step size by some orders of magnitude during training, for example, set to $\eta = 0.995$ for 1000 iterations.

The initialization of matrix Ω is of particular interest. If chosen as identity matrix $\Omega = \mathbf{I}$, the algorithm starts from the usual squared Euclidean distance. For data sets with strong mutual attribute dependencies, i.e. prominent non-diagonal elements, the uniform structure of the identity matrix might lead to unnecessary iterations required for the symmetry breaking, as often encountered in neural network adaptations. Therefore, the alternatively proposed method is random matrix element sampling from uniform noise in the interval $[-0.5; 0.5]$. This noise matrix $\mathbf{A} \in \mathbb{R}^{q \times q}$ is broken by QR-decomposition into $\mathbf{A} = \mathbf{Q} \cdot \mathbf{R}$, of which the \mathbf{Q} -part is known to form an orthonormal basis with $\mathbf{Q} \cdot \mathbf{Q}^\top = \mathbf{I}$. This makes $\Omega = \mathbf{Q}$ our preferred initial candidate.

Relation to LDA. At first glance, the proposed cost function looks quite similar to the inverse fraction of the LDA cost function for C classes that is maximized:

$$S_{LDA} = \frac{\mathbf{v} \cdot \left[\sum_{i=1}^C n_i \cdot (\boldsymbol{\mu}_i - \boldsymbol{\mu})^\top \cdot (\boldsymbol{\mu}_i - \boldsymbol{\mu}) \right] \cdot \mathbf{v}^\top}{\mathbf{v} \cdot \left[\sum_{i=1}^C \boldsymbol{\Sigma}_i \right] \cdot \mathbf{v}^\top}, \quad n_i = |\{\mathbf{x}^j : c(j) = i\}|. \quad (6)$$

The numerator contains the between-class variation as the squared difference between class centers $\boldsymbol{\mu}_i$ of all vectors \mathbf{x}^j belonging to class i and the overall center $\boldsymbol{\mu} = 1/n \cdot \sum_{k=1}^n \mathbf{x}^k$. The denominator describes the within-class variation over all classes i expressed by the sum of squared differences from class centers $\boldsymbol{\mu}_i$ contained in the covariance matrices $\boldsymbol{\Sigma}_i = \sum_{j:c(j)=i} (\mathbf{x}^j - \boldsymbol{\mu}_i)^\top \cdot (\mathbf{x}^j - \boldsymbol{\mu}_i)$.

LDA seeks an optimum direction vector \mathbf{v} representing a good compromise of being collinear along the class centers (numerator, separating) and orthogonal to maximum within-class variation (denominator, compressing).

If multiple directions $\mathbf{V} = (\mathbf{v}_k)^\top$ are computed simultaneously, the products in the numerator and denominator of Eqn. 6, involving the matrices in square brackets, become matrices as well. In order to circumvent the problem of valid ratio calculation with matrices, determinants of the obtained matrices can be taken, as discussed in the LDA-based projection pursuit approach [12]. As a result, the LDA ratio optimizes low-dimensional projections onto discriminatory directions.

Our approach is structurally different, because the (inverse) LDA ratio in Eqn. 2 operates in the original data space, subject to the dynamically optimized metric. This explains the higher computational demands compared to LDA for which covariance matrices and class centers can be initially computed and then reused. As a benefit of the new approach, numerator and denominator of the new ratio in Eqn. 2 naturally contain sums of real-valued distances, which avoids problems of handling singular determinants in low-rank matrices.

3 Experiments

3.1 Tecator Spectral Data Set

The benchmark spectral data set, taken from the UCI repository of machine learning [1], contains 215 samples of 100-dimensional infrared absorbance spectra recorded on a Tecator Infratec Food and Feed Analyzer working in the wavelength range 850–1050nm by the Near Infrared Transmission (NIT) principle. The original regression problem accompanying the data set is reformulated as attribute identification task for explaining the separation of 183 samples with low fat content and 77 high fat meat probes.

View 1. An exploratory data view is obtained from the left panel of Fig. 1 and from the PCA projection shown in the left scatter plot of Fig. 2. As expected, the strong spectrum overlap cannot be resolved by PCA projection. After application of the matrix learning all spectra were transformed according to $\mathbf{z} = \mathbf{x} \cdot \mathbf{\Omega}$, which realizes the left transformation part of the metric given in Eqn. 1; the right part is just \mathbf{z}^T . The result of this data transformation leads to a good separation with

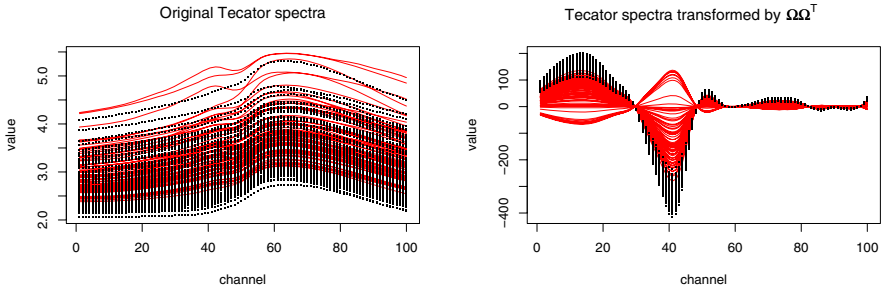


Fig. 1. Tecator spectra, raw (left) and transformed (right). Low fat content is reflected by dashed lines, high fat content by solid lines.

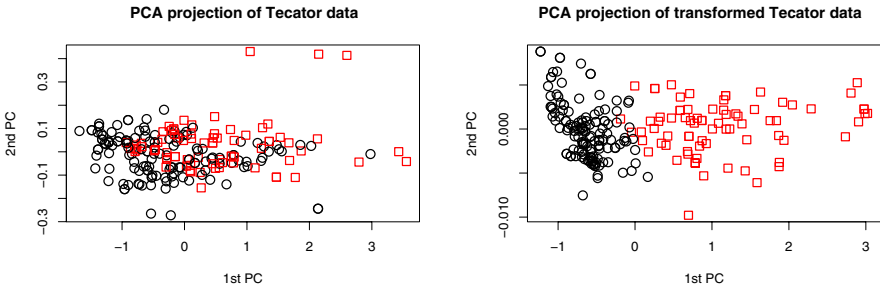


Fig. 2. Scatter plots of Tecator data. Bullets (\circ) denote low-fat samples, squares (\square) high fat content. Left: PCA projection of original data. Right: PCA projection of data transformed by $\mathbf{\Omega}$.

almost no overlap in the PCA projection. This is shown in the right panel of Fig. 2.

View 2. By reformulating the metric definition in Eqn. 1 according to

$$d_{\Omega}^{ij} = (\mathbf{x}^i - \mathbf{x}^j) \cdot \Omega \cdot \Omega^T \cdot (\mathbf{x}^i - \mathbf{x}^j)^T = \langle \mathbf{x}^i \cdot \Omega \cdot \Omega^T - \mathbf{x}^j \cdot \Omega \cdot \Omega^T, \mathbf{x}^i - \mathbf{x}^j \rangle. \quad (7)$$

another interesting perspective on the data is obtained. This is a formal metric decomposition into a static part of difference vectors of the original data (right part of the scalar product) and a dynamically adapted transformation space of the data (left argument of the scalar product). A look into this space is obtained by the transformation to $\mathbf{x}^* = \mathbf{x} \cdot \Omega \cdot \Omega^T$. The resulting transformed spectra with their amazingly separated attributes are shown in the right panel of Fig. 1.

The learned metric can be nicely presented by the matrices Ω and Λ shown in the left and right panel of Fig. 3, respectively. As displayed for Λ , attribute dependence is most prominent in the channel range 35–45. Strong emphasis of these channels around the diagonal is accompanied by simultaneous repression of the off-diagonal channels 5–30.

Matrix reduction. Since full matrices are quite big models, the study of their compressibility is important. Eigen decomposition of $\Lambda = \mathbf{S} \cdot \mathbf{W} \cdot \mathbf{W}^{-1}$ into the diagonal eigenvalue matrix \mathbf{S} and the eigenvectors matrix \mathbf{W} helps to reach substantial compressions. In the current case, the highest eigenvalue contributes an amount of 95.3%, thus most variation in the learned matrix Ω can be explained by the corresponding eigenvector \mathbf{w} , a column vector. Therefore, up to a scaling factor, a very good reconstruction of Λ by $\mathbf{w} \cdot \mathbf{w}^T$ is obtained, as confirmed in the left matrix plot of Fig. 4. If the spectra are projected onto \mathbf{w} , still a very good class separation is obtained, as demonstrated by the corresponding class-specific box plot in the right panel of Fig. 4.

The computational demands are quite high, though, requiring roughly one hour for 1000 updates of the matrix gradient. In contrast to that, the classPP [12] package is much faster, if only a class-separating projection is desired. In principle, classPP takes only several seconds or minutes, depending on the choice of

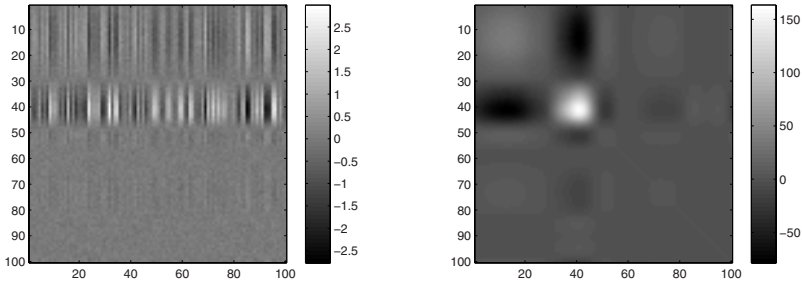


Fig. 3. Matrix representation of optimum metric for the 100-D Tecator data set. The learned matrix Ω is shown on the left, its squared counterpart $\Lambda = \Omega \cdot \Omega^T$ on the right. Interesting dependencies are found around channel index 40.

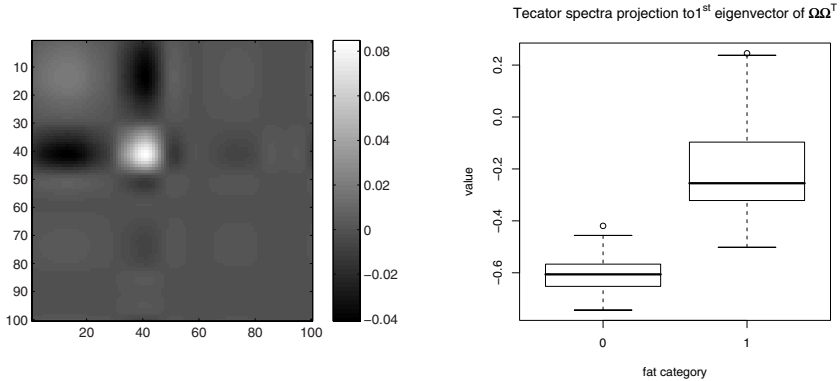


Fig. 4. Representation of \mathbf{A} by its first eigenvector. Left: plot of reconstructed matrix. Right projection of Tecator spectra to the eigenvector.

annealing parameters. However, no stable solution could be obtained, because of the degeneration of the projection vectors, probably caused by near-singular matrix determinants during the computation. Training with our proposed method showed very stable results, converging to the presented solution for different random initializations of $\mathbf{\Omega}$. Projection displays are just a by-product of our method. It is important to remember that the original dimensionality of the data space is preserved by the transformation, enabling further utilization with any classification or projection method.

3.2 Gene Expression Analysis of AML/ALL Cancer

Many well-documented and deeply investigated data sets are freely available in cancer research. Since its publication in 1999 the leukemia gene expression dataset [5] used here has become a quasi-benchmark for testing feature selection methods. The original research aimed at the identification of the most informative genes for modeling and classification of two cancer types, acute lymphoblastic leukemia (ALL) and acute myeloid leukemia (AML). The training data covers 7129 genes by 27 cases of ALL and 11 cases of AML. The test data set contains 20 cases of ALL and 14 cases of AML.

The projection of the complete set of training and test data to the first two principal components yields the scatter plot shown in the left panel of Fig. 5. It is worth noticing that a systematic difference between AML training set and test set is indicated by the unbalanced distribution of closed and open bullets. Thus, a training set specific bias is induced during training. Matrix learning is computationally very expensive, because $\mathbf{\Omega}$ is a 7129x7129 matrix. Thus, it takes roughly 40 hours on a 2.4 GHz system in order to achieve 500 gradient changes expressed by Eqn. 5. Yet, gradients are reused until first cost function degradation, creating several thousand updates, after all. Since experiment preparation in lab requires much more time, a two day calculation period is no principal problem.

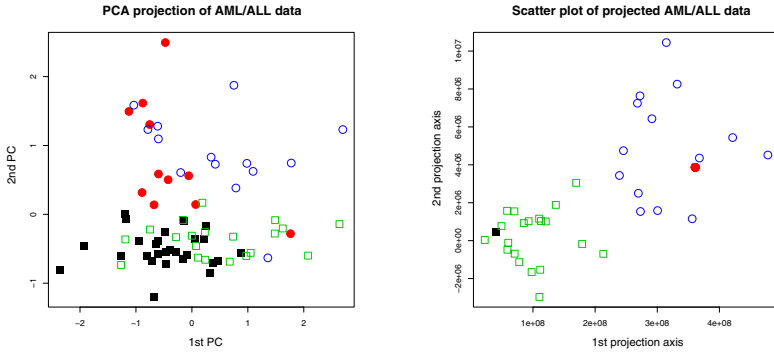


Fig. 5. Scatter plots of AML/ALL gene expression data. Bullets (\circ) denote expression samples of AML, squares (\square) are related to ALL cancer. Closed symbols indicate training data, open symbols test data. Left panel: principal component projection without distinction between training and test data. Right panel: data projected to the first two eigenvectors of the trained interaction matrix $\Omega \cdot \Omega^T$. The training data is perfectly arranged, being very distinct and almost contracted to points.

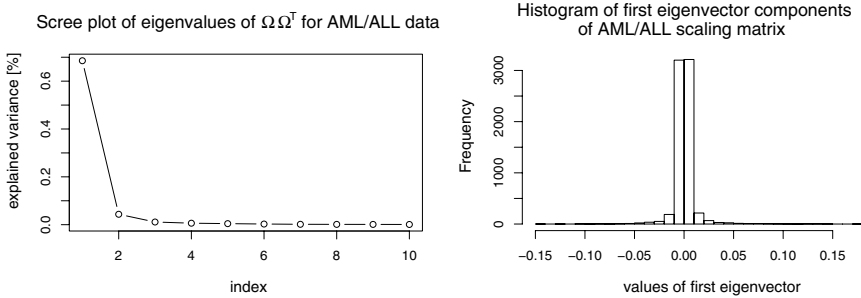


Fig. 6. Eigenvalues of the interaction matrix $\Omega \cdot \Omega^T$ (left) for the AML/ALL data, and the distribution of values in the most prominent eigenvector (right)

Several interesting results are obtained after training. The top ten eigenvalues of $\Omega \cdot \Omega^T$, displayed in the left panel of Fig. 6, point out a very strong explicative power of roughly 70% explained variance of the first eigenvector, dominating the all other eigenvectors. The right panel in Fig. 6 further indicates that only a small fraction of extreme values in that first eigenvector really contains interesting attribute magnifications. In case of data projection, many other attributes are transformed to near zero values, i.e. only few differentiating genes become emphasized.

The projection to the first two eigenvectors of the interaction matrix already yields a perfect separation of the training data into AML and ALL, as shown in the right panel of Fig. 5. Also a very strong compression of the within-class scatter almost to points is obtained. The test data, projected the same way and added to the plot, is still well-separated, but shows much larger variability. This is a clear indication of over-fitting. Sure, a huge 7129x7129 model has been

trained; yet, the displayed projection only contains several hundred effective parameters, because the first two eigenvectors are much dominated by the few most prominent entries of the first eigenvector. A simple center-of-gravity model of the data, for example, would already be much larger.

For comparison, LDA indices of the projections shown in the right panel of Fig. 5 were calculated according to Eqn. 6 using the classPP package. An almost perfect near-one value of $S_{LDA} = 1 - 7.894 \cdot 10^{-12}$ was achieved for the projected training data and good value of $S_{LDA} = 0.8298$ for the test data. Using the built-in simulated annealing strategy, the classPP package itself reported a best seen index value of $S_{LDA} = 1 - 2.825 \cdot 10^{-6}$ for the optimized projections of the training data. Since the corresponding projection matrix is not returned correctly from the classPP package, an application to the test set was not possible. The authors of classPP have identified internal rounding errors in their package.

Individual gene variances correspond to gene-specific scaling factors, compensated by the cost function by adaptation of the related matrix entries. As a consequence, components in the most prominent eigenvector show low correlation with the variance in the original data set and, accordingly, systematically separating low-variance genes were able to gain high rankings in the eigenvector.

The real benefit of the proposed method is the possibility to infer putative gene-gene interactions responsible for cancer type separation. For that purpose the indices i, j corresponding to the most extreme (high and low) values in the matrix $\Omega \cdot \Omega^T$ are extracted and associated with the genes i and j . Because of symmetry, only the lower triangular matrix, including diagonal, is considered. The top 100 pairs extracted this way are compiled in Tab. 1. After all, 14 prominent self-dependent genes are detected as individual factors on the diagonal, three of them are coinciding with the list of Golub et al. of 50 genes. Three more genes of that study are detected on non-diagonal elements as dependent.

Table 1. Table of genes specific to separation of AML/ALL cancer in alphabetic reading order. The listed genes correspond to the 100 most extreme entries in the lower triangular part of the obtained symmetric matrix $\Omega \cdot \Omega^T$. As single genes participate multiple times in combination with others, only 30 different out of 200 possible genes appear in the table. Numbers indicate the frequencies of occurrence. Underlined genes appear also on the diagonal, stressing their individual importance. For illustration purposes, bold face genes are those acting in combination with M19507 which is the overall top-ranked gene. Asterisks mark genes coinciding with top-rated genes from the study of Golub et al.

<u>D49824</u> 6	<u>HG3576-HT3779</u> 13	<u>L06797</u> 3	L20688 1	L20941 1	M11147 1
M14328 1	M17733 1	<u>M19507</u> 26	M24485 1	M27891* 3	M28130_rna1* 1
M33600 4	M69043* 1	M77232_rna1 1	<u>M91036_rna1</u> 13	M91438 1	<u>M96326_rna1*</u> 14
S73591 1	<u>U01317_cds4</u> 13	U14968 1	V00594 1	X14046 1	<u>X17042*</u> 12
<u>X78992</u> 13	<u>Y00433</u> 14	<u>Y00787*</u> 14	<u>Z19554</u> 14	<u>Z48501</u> 11	<u>Z70759</u> 13

The most prominent gene found by the new method is M19507, which is not mentioned in the Golub study. Yet, the gene is confirmed as relevant in more recent publications, such as [9]. As this gene is connected to more than 20 other top-rated genes, its central role in the discriminatory transcriptome is clearly pointed out. Yet, the whole potential of the analysis, including proper interpretation of the findings, must be thoroughly worked out together with biological experts.

4 Conclusions and Outlook

A data-driven metric in flavor of a generalized Mahalanobis distance has been proposed that makes use of label information for emphasizing or repressing class-specific attribute combinations. Similar to LDA, metric optimization of $\mathbf{A} = \mathbf{\Omega} \cdot \mathbf{\Omega}^T$ seeks improved inter-class separation with simultaneous minimization of within-class variation. In contrast to LDA, it is not the low-dimensional projection to be optimized, but a transformation in the data space. The new method is not primarily designed for visual projection or classification, but it is a first step towards, because the resulting transformed data can be used as a preprocessing step for subsequent standard methods. As illustrated, visual data exploration is easily possible by projecting the data to the most prominent eigenvectors of \mathbf{A} . No sophisticated optimization method is required, simple gradient descent works very reliably on the inverse LDA-like cost function. Both investigated data sets led to convergence to very useful label-specific metrics for different initializations of $\mathbf{\Omega}$. The main drawback of the new method is its long runtime for handling the potentially large matrices. Yet, as the discussed cases showed a strong dominance of only the first principal direction, future work will focus on the development of a sparse learning scheme for computing only the k most prominent eigenvectors instead of the whole matrix. This will help to reduce the model size and to speed up the optimization procedure. Finally, a better control of intra- and inter-class contributions to the cost function will be investigated.

Acknowledgment

The work is supported by grant XP3624HP/0606T, Ministry of Culture Saxony-Anhalt.

References

1. Blake, C., Merz, C.: UCI repository of machine learning databases (1998)
2. Cook, D., Swayne, D.: Interactive and Dynamic Graphics for Data Analysis with R and GGobi. Springer, Heidelberg (2007)
3. Faith, J., Mintram, R., Angelova, M.: Targeted projection pursuit for visualizing gene expression data classifications. *Bioinformatics* 22(21), 2667–2673 (2006)

4. Friedman, J.: Exploratory projection pursuit. *Journal of the American Statistical Association* 82, 249–266 (1987)
5. Golub, T., Slonim, D., Tamayo, P., Huard, C., Gaasenbeek, M., Mesirov, J., Coller, H., Loh, M., Downing, J., Caligiuri, M., Bloomfield, C., Lander, E.: Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *Science* 286(5439), 531–537 (1999)
6. Guyon, I., Gunn, S., Nikravesh, M., Zadeh, L.: *Feature Extraction: Foundations and Applications*. Springer, Berlin (2006)
7. Hammer, B., Strickert, M., Villmann, T.: Supervised neural gas with general similarity measure. *Neural Processing Letters* 21(1), 21–44 (2005)
8. Hammer, B., Villmann, T.: Generalized relevance learning vector quantization. *Neural Networks* 15, 1059–1068 (2002)
9. Hu, S., Rao, J.: Statistical redundancy testing for improved gene selection in cancer classification using microarray data. *Cancer Informatics* 2, 29–41 (2007)
10. Hyvärinen, A., Oja, E.: Independent component analysis: Algorithms and applications. *Neural Networks* 13(4–5), 411–430 (2000)
11. Kaski, S.: From learning metrics towards dependency exploration. In: Cottrell, M. (ed.) *Proceedings of the 5th International Workshop on Self-Organizing Maps (WSOM)*, pp. 307–314 (2005)
12. Lee, E., Cook, D., Klinke, S., Lumley, T.: Projection pursuit for exploratory supervised classification. *Journal of Computational and Graphical Statistics* 14(4), 831–846 (2005)
13. Schneider, P., Biehl, M., Hammer, B.: Adaptive relevance matrices in learning vector quantization (Submitted to *Machine Learning*) (2008)
14. Strickert, M., Seiffert, U., Sreenivasulu, N., Weschke, W., Villmann, T., Hammer, B.: Generalized relevance LVQ (GRLVQ) with correlation measures for gene expression data. *Neurocomputing* 69, 651–659 (2006)
15. Strickert, M., Witzel, K., Mock, H.-P., Schleif, F.-M., Villmann, T.: Supervised attribute relevance determination for protein identification in stress experiments. In: *Proc. of Machine Learning in Systems Biology (MLSB)*, pp. 81–86 (2007)
16. Sun, Y.: Iterative relief for feature weighting: Algorithms, theories, and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 29(6), 1035–1051 (2007)
17. Weinberger, K., Blitzer, J., Saul, L.: Distance metric learning for large margin nearest neighbor classification. In: Weiss, Y., Schölkopf, B., Platt, J. (eds.) *Advances in Neural Information Processing Systems* 18, pp. 1473–1480. MIT Press, Cambridge (2006)
18. Xing, E., Ng, A., Jordan, M., Russell, S.: Distance metric learning, with application to clustering with side-information. In: Becker, S., Thrun, S., Obermayer, K. (eds.) *Advances in Neural Information Processing Systems* 15 (NIPS), pp. 505–512. MIT Press, Cambridge (2003)