

Adaptive relevance matrices in Learning Vector Quantization

Petra Schneider¹, Michael Biehl¹, Barbara Hammer²

¹Institute for Mathematics and Computing Science, University of Groningen
P.O. Box 407, 9700 AK Groningen, The Netherlands
{p.schneider,m.biehl}@rug.nl

²Institute of Computing Science, Clausthal University of Technology
Julius Albert Strasse 4, 38678 Clausthal-Zellerfeld, Germany
hammer@in.tu-clausthal.de

Abstract

We propose a new matrix learning scheme to extend relevance learning vector quantization (RLVQ), an efficient prototype-based classification algorithm, towards a general adaptive metric. By introducing a full matrix of relevance factors in the distance measure, correlations between different features and their importance for the classification scheme can be taken into account and automated, general metric adaptation takes place during training. In comparison to the weighted Euclidean metric used in RLVQ and its variations, a full matrix is more powerful to represent the internal structure of the data appropriately. Large margin generalization bounds can be transferred to this case leading to bounds which are independent of the input dimensionality. This also holds for local metrics attached to each prototype which corresponds to piecewise quadratic decision boundaries. The algorithm is tested in comparison to alternative LVQ schemes using an artificial data set, a benchmark multi-class problem from the UCI repository, and a problem from bioinformatics, the recognition of splice sites for *C. elegans*.

Keywords: learning vector quantization, generalized LVQ, metric adaptation, generalization bounds

1 Introduction

Learning vector quantization (LVQ) as introduced by Kohonen is a particularly intuitive and simple though powerful classification scheme (Kohonen, 1997) which is very appealing for several reasons: the method is easy to implement; the complexity of the resulting classifier can be controlled by the user; the classifier can naturally deal with multi-class problems; and, unlike many alternative classification schemes such as feedforward networks or the support vector machine

(SVM), the LVQ system is straightforward to interpret because of the intuitive assignment of data to the class of the closest prototype. For these reasons, LVQ has been used in a variety of academic and commercial applications such as image analysis, bioinformatics, telecommunication, robotics, etc. (*Bibliography on the Self-Organizing Map (SOM) and Learning Vector Quantization (LVQ)*, 2002).

Original LVQ, however, suffers from several drawbacks such as potentially slow convergence and instable behavior because of which various alternatives have been proposed, see for instance Kohonen (1997). Still, there are two major drawbacks of these methods, which have been tackled only recently.

On the one hand, LVQ relies on heuristics and a full mathematical investigation of the algorithm is lacking. This problem relates to unexpected behavior and instabilities of training. It has been shown that already slight variations of the basic LVQ learning scheme yield quite different results (Biehl m. fl., 2006, 2007). Variants of LVQ which can be derived from an explicit cost function are particularly interesting. Several proposals for cost functions can be found in the literature, one example being generalized LVQ (GLVQ) which forms the basis for the method we will consider in this article (Sato & Yamada, 1996). Two alternatives which implement soft relaxations of the original learning rule are presented in Seo & Obermayer (2003); Seo m. fl. (2003). These two approaches, however, have the drawback that the original crisp limit case does not exist (for Seo & Obermayer, 2003) resp. the discrete limit case shows poor results also in simple settings (see Ghosh m. fl., 2006 for Seo m. fl., 2003). The cost function as proposed in Sato & Yamada (1996) displays stable behavior and aims at a good generalization ability already during training as pointed out in Hammer m. fl. (2005a). On the other hand, LVQ and variants severely rely on the standard Euclidean metric and they are not appropriate for situations where the Euclidean metric does not represent the underlying semantic. This is the case, e.g., for high dimensional data where noise accumulates and likely disrupts the classification, for heterogeneous data where the importance and nature of the dimensions differs, and for data which involves correlations of the dimensions. In these cases, which are quite common in practice, simple LVQ may fail. Recently, a cost function based generalization of LVQ has been proposed which allows to incorporate general differentiable similarity measures (Hammer m. fl., 2005b). The specific choice of the similarity measure as a simple weighted diagonal metric with adaptive relevance terms has turned out particularly suitable in many practical applications since it can account for irrelevant or inadequately scaled dimensions. At the same time, it allows for straightforward interpretation of the result because the relevance profile can directly be interpreted as the contribution of the dimensions to the classification (Hammer & Villmann, 2002). For an adaptive diagonal metric, dimensionality independent large margin generalization bounds can be derived (Hammer m. fl., 2005a). This fact is remarkable since it accompanies the good experimental classification results for high dimensional data by a theoretical counterpart. The same bounds also hold for kernelized versions, but not for arbitrary choices of the metric. Often, different features are correlated in classification tasks. In unsupervised cluster-

ing, correlations of data are accounted for, e.g., by the classical Mahalanobis distance (Duda m. fl., 2000) or fuzzy-covariance matrices as proposed, e.g., in the approaches Gath & Geva (1989); Gustafson & Kessel (1997). It has been shown recently that general metric learning based on large margin principles can greatly improve the results obtained by distance-based schemes such as k-nearest neighbor classifier (Shalev-Schwartz m. fl., 2004; Weinberger m. fl., 2006). For supervised LVQ classification tasks, however, an explicit metric which takes into account correlations has not yet been proposed. Based on the general framework as presented in Hammer m. fl. (2005a), we develop an extension of LVQ to an adaptive matrix of relevances which parameterizes a general Euclidean metric and accounts for pairwise correlations of features. By means of an implicit scaling and rotation of the data, the algorithm yields a discriminative distance measure which is particularly suitable for the given classification task. It can be parameterized in terms of a single, global matrix or by individual matrices attached to the prototypes. Interestingly, one can derive generalization bounds which are similar to the case of a simple diagonal metric for this more complex case. Apart from this theoretical guarantee, we demonstrate the usefulness of the novel scheme in the context of several classification problems.

Note that extensions of LVQ schemes towards general kernels exist which take into account a more general metric than the Euclidean one, see e.g. (Hammer m. fl., 2005c). Learning rules can easily be derived, provided the LVQ scheme follows a cost function and the kernel is differentiable. However, although this idea equips LVQ schemes with a larger capacity, the kernel has to be fixed a priori. In this contribution, we consider a specific metric which is *adaptive* according to the given classification task and corresponding generalization bounds.

2 Review of LVQ

LVQ aims at approximating a classification scheme by prototypes. Assume training data $(\boldsymbol{\xi}_i, y_i) \in \mathbb{R}^N \times \{1, \dots, C\}$ are given, N denoting the data dimensionality and C the number of different classes. An LVQ network consists of a number of prototypes which are characterized by their location in the weight space $\boldsymbol{w}_i \in \mathbb{R}^N$ and their class label $c(\boldsymbol{w}_i) \in \{1, \dots, C\}$. Classification is implemented as a winner takes all scheme. For this purpose, a possibly parameterized similarity measure d^λ is fixed for \mathbb{R}^N , where λ specifies the metric parameters which can be adapted during training. Often, the standard Euclidean metric is chosen. A data point $\boldsymbol{\xi} \in \mathbb{R}^N$ is mapped to the class label $c(\boldsymbol{\xi}) = c(\boldsymbol{w}_i)$ of the prototype i for which $d^\lambda(\boldsymbol{w}_i, \boldsymbol{\xi}) \leq d^\lambda(\boldsymbol{w}_j, \boldsymbol{\xi})$ holds for every $j \neq i$, breaking ties arbitrarily. Hence, it is mapped to the class of the closest prototype, the so-called winner.

Learning aims at determining weight locations for the prototypes such that the given training data are mapped to their corresponding class labels. This is usually achieved by a modification of Hebbian learning, which moves prototypes closer to the data points of their respective class. A very flexible learning

approach has been introduced in Sato & Yamada (1996): training is derived as a minimization of the cost function

$$\sum_i \Phi(\mu_i) \text{ where } \mu_i = \frac{d_J^\lambda(\xi_i) - d_K^\lambda(\xi_i)}{d_J^\lambda(\xi_i) + d_K^\lambda(\xi_i)} \quad (1)$$

based on the steepest descent method. Φ is a monotonic function, e.g. the logistic function or the identity $\Phi(x) = x$ which we use throughout the following; $d_J^\lambda(\xi) = d^\lambda(\mathbf{w}_J, \xi)$ is the distance of data point ξ from the closest prototype \mathbf{w}_J with the same class label y , and $d_K^\lambda(\xi) = d^\lambda(\mathbf{w}_K, \xi)$ is the distance from the closest prototype \mathbf{w}_K with a different class label than y . Note that the numerator is smaller than 0 iff the classification of the data point is correct. The smaller the numerator, the greater the security of classification, i.e. the difference of the distance from a correct and wrong prototype. The denominator scales the argument of Φ such that it satisfies $-1 < \mu(\xi) < 1$. A further possibly nonlinear scaling by Φ might be beneficial for applications. This formulation can be seen as a kernelized version of so-called generalized LVQ introduced in Sato & Yamada (1996).

The learning rule can be derived from this cost function by taking derivatives. We assume that the similarity measure $d^\lambda(\mathbf{w}, \xi)$ is differentiable with respect to the parameters \mathbf{w} and λ . As shown in Hammer m. fl. (2005a), for a given pattern ξ the derivatives yield

$$\Delta \mathbf{w}_J = -\epsilon \cdot \Phi'(\mu(\xi)) \cdot \mu^+(\xi) \cdot \nabla_{\mathbf{w}_J} d_J^\lambda(\xi), \quad (2)$$

where $\epsilon > 0$ is the learning rate, the derivative of Φ is taken at position $\mu(\xi)$, and $\mu^+(\xi) = 2 \cdot d_K^\lambda(\xi) / (d_J^\lambda(\xi) + d_K^\lambda(\xi))^2$. Further,

$$\Delta \mathbf{w}_K = \epsilon \cdot \Phi'(\mu(\xi)) \cdot \mu^-(\xi) \cdot \nabla_{\mathbf{w}_K} d_K^\lambda(\xi), \quad (3)$$

where $\mu^-(\xi) = 2 \cdot d_J^\lambda(\xi) / (d_J^\lambda(\xi) + d_K^\lambda(\xi))^2$. The derivative with respect to the parameters λ yields the update

$$\Delta \lambda = \epsilon \cdot \Phi'(\mu(\xi)) \cdot (\mu^+(\xi) \cdot \nabla_\lambda d_J^\lambda(\xi) - \mu^-(\xi) \cdot \nabla_\lambda d_K^\lambda(\xi)) . \quad (4)$$

The adaptation of λ is often followed by normalization during training, e.g. enforcing $\sum_i \lambda_i = 1$ to prevent degeneration of the metric. It has been shown in Hammer m. fl. (2005a) that these update rules are valid whenever the metric is differentiable. Using delta-functions, the computation can also be done for distributions P with continuous support, (Hammer m. fl., 2005a), i.e. a valid cost function results for every differentiable metric and reasonable distribution P .

It has been demonstrated in Hammer m. fl. (2005a) that the squared weighted Euclidean metric $d^\lambda(\mathbf{w}, \xi) = \sum_i \lambda_i (w_i - \xi_i)^2$ with $\lambda_i \geq 0$ and $\sum_i \lambda_i = 1$ is a simple and powerful choice which allows to use prototype based learning also in the presence of high dimensional data with features of different, yet a priori unknown, relevance. This measure has the advantage that the relevance factors λ_i can be interpreted directly and provide insight into the classification

task: dimensions with large λ_i are most important for the classification while very small or zero relevances indicate that the corresponding feature could be omitted. We refer to this method as generalized relevance learning vector quantization (GRLVQ, see Hammer & Villmann (2002)). Very similar schemes have been motivated heuristically which lack an interpretation of the algorithm as stochastic gradient descent with respect to a cost function (Bojer m. fl., 2001). Alternative choices have been introduced in Hammer m. fl. (2005a), including, for example, metrics which take local windows into account, e.g., for time series data.

Note that the relevance factors, i.e. the choice of the metric need not be global, but can be attached to single prototypes, locally. In this case, individual updates take place for the relevance factors λ^j for each prototype j , and the distance of a data point $\boldsymbol{\xi}$ from prototype \boldsymbol{w}_j , $d^{\lambda^j}(\boldsymbol{w}_j, \boldsymbol{\xi})$ is computed based on λ^j . This allows a local relevance adaptation, taking into account that the relevance might change within the data space. This method has been investigated e.g. in Schleif & Villmann (2005). We refer to this version as localized GRLVQ (LGRLVQ).

3 The GMLVQ Algorithm

Here, we introduce an important extension of the above concept, which employs a full matrix of relevances in the similarity measure. We consider a generalized distance of the form

$$d^\Lambda(\boldsymbol{w}, \boldsymbol{\xi}) = (\boldsymbol{\xi} - \boldsymbol{w})^T \Lambda (\boldsymbol{\xi} - \boldsymbol{w}), \quad (5)$$

where Λ is a full $N \times N$ matrix which can account for correlations between the features. For Λ to define a valid metric, symmetry and positive definiteness has to be enforced (we will discuss in a moment how this property can be guaranteed). This way, arbitrary Euclidean metrics can be realized by an appropriate choice of the parameters. In particular, correlations of dimensions and rotation of the axes can be accounted for. Such choices have already successfully been introduced in unsupervised clustering methods such as fuzzy clustering (Gath & Geva, 1989; Gustafson & Kessel, 1997), however, at the expense of increased computational costs, since these methods require a matrix inversion at each adaptation step. For the metric as introduced above, a variant which costs $\mathcal{O}(N^2)$ can be derived.

Note that, as already stated, the above similarity measure defines a general squared Euclidean distance in an appropriately transformed space only if Λ is positive definite and symmetric. We can achieve this by substituting

$$\Lambda = \Omega^T \Omega \quad (6)$$

which yields $\boldsymbol{u}^T \Lambda \boldsymbol{u} = \boldsymbol{u}^T \Omega^T \Omega \boldsymbol{u} = (\Omega^T \boldsymbol{u})^2 \geq 0$ for all \boldsymbol{u} , where Ω is an arbitrary real $N \times N$ matrix. In addition, $\det \Lambda \neq 0$ has to be enforced to guarantee that Λ is positive definite. However, in practice, positive semi-definiteness of

the matrix is sufficient, since data often only populates a sub-manifold of the full data space and definiteness has to hold only with regard to the relevant subspace of data. Therefore, we do not enforce $\det \Lambda \neq 0$. Using the relation in Eq. (6), the squared distance reads

$$d^\Lambda(\mathbf{w}, \boldsymbol{\xi}) = \sum_{ijk} (\xi_i - w_i) \Omega_{ki} \Omega_{kj} (\xi_j - w_j). \quad (7)$$

To obtain the adaptation formulas we need to compute the derivatives with respect to \mathbf{w} and Ω . The derivative of d^Λ with respect to \mathbf{w} yields

$$\nabla_{\mathbf{w}} d^\Lambda(\mathbf{w}, \boldsymbol{\xi}) = -2\Lambda (\boldsymbol{\xi} - \mathbf{w}) = -2\Omega^T \Omega (\boldsymbol{\xi} - \mathbf{w}). \quad (8)$$

Derivatives with respect to a single element Ω_{lm} give

$$\begin{aligned} \frac{\partial d^\Lambda(\mathbf{w}, \boldsymbol{\xi})}{\partial \Omega_{lm}} &= \sum_j (\xi_m - w_m) \Omega_{lj} (\xi_j - w_j) + \sum_i (\xi_i - w_i) \Omega_{li} (\xi_m - w_m) \\ &= 2 \cdot (\xi_m - w_m) [\Omega(\boldsymbol{\xi} - \mathbf{w})]_l, \end{aligned} \quad (9)$$

where subscripts l, m specify components of vectors. Thus, we get the update equations

$$\begin{aligned} \Delta \mathbf{w}_J &= \epsilon \cdot 2 \cdot \Phi'(\mu(\boldsymbol{\xi})) \cdot \mu^+(\boldsymbol{\xi}) \cdot \Lambda \cdot (\boldsymbol{\xi} - \mathbf{w}_J), \\ \Delta \mathbf{w}_K &= -\epsilon \cdot 2 \cdot \Phi'(\mu(\boldsymbol{\xi})) \cdot \mu^-(\boldsymbol{\xi}) \cdot \Lambda \cdot (\boldsymbol{\xi} - \mathbf{w}_K). \end{aligned} \quad (10)$$

Note that these updates correspond to the standard Hebb terms of LVQ, pushing the closest correct prototype towards the considered data point and the closest wrong prototype away from the considered data point. For the update of the matrix elements Ω_{lm} we get

$$\begin{aligned} \Delta \Omega_{lm} &= -\epsilon \cdot 2 \cdot \Phi'(\mu(\boldsymbol{\xi})) \cdot \\ &\quad \left(\mu^+(\boldsymbol{\xi}) \cdot \left((\xi_m - w_{J,m}) [\Omega(\boldsymbol{\xi} - \mathbf{w}_J)]_l \right) - \right. \\ &\quad \left. \mu^-(\boldsymbol{\xi}) \cdot \left((\xi_m - w_{K,m}) [\Omega(\boldsymbol{\xi} - \mathbf{w}_K)]_l \right) \right). \end{aligned} \quad (11)$$

This update also corresponds to a Hebbian term, since the driving force consists of the derivative of the distance from the closest correct prototype (scaled with -1) and the closest incorrect prototype. Thus, the parameters of the matrix are changed in such a way that the distance from the closest correct prototype becomes smaller, whereas the distance from the closest wrong prototype is increased. Similar to the case of diagonal relevances (Bojer m. fl., 2001), matrix updates can be formulated on heuristic grounds, as well. Such variants will be discussed in forthcoming publications.

The learning rate for the metric can be chosen independently of the learning rate of the prototypes. We set it an order of magnitude smaller in order

to achieve a slower time-scale of metric learning compared to the weight updates. After each update Λ should be normalized to prevent the algorithm from degeneration. One possibility is to enforce

$$\sum_i \Lambda_{ii} = 1 \quad (12)$$

by dividing all elements of Λ by the row value of $\sum_i \Lambda_{ii}$ after each step. In this way we fix the sum of diagonal elements which coincides with the sum of eigenvalues. This generalizes the normalization of relevances $\sum_i \lambda_i = 1$ for a simple diagonal metric. One can interpret the eigendirections of Λ as the temporary coordinate system with the relevances corresponding to the eigenvalues. Since

$$\Lambda_{ii} = \sum_k \Omega_{ki} \Omega_{ki} = \sum_k (\Omega_{ki})^2, \quad (13)$$

normalization can be done by dividing all elements of Ω by $(\sum_{ki} (\Omega_{ki})^2)^{1/2} = (\sum_i [\Omega^T \Omega]_{ii})^{1/2}$ after every update step. Note that for one global matrix, normalization of the matrix does not affect the cost function such that this procedure can still be interpreted as a standard gradient technique. The explicit normalization of Λ changes the learning dynamics in the general case of individual matrices because of the explicit projection of the metric tensor to 1 after every adaptation step. In this case, for small step size, we observed a reasonable behavior since the method can be interpreted as an approximate gradient projection (albeit with non-convex constraints, see (Bertsekas, 1976)).

We term the algorithm specified by Eq.s (10) and (11) generalized matrix LVQ (GMLVQ). The complexity of one adaptation step is determined by the computation of the closest correct and incorrect prototypes ($\mathcal{O}(N^2 \cdot N_w)$, N_w being the number of prototypes), and the adaptation ($\mathcal{O}(N^2)$). Usually, this procedure is repeated a number of time steps which is linear in the number of patterns to achieve convergence. Thus, this procedure is faster than unsupervised fuzzy-clustering variants which use a similar form of the metric but which require a matrix inversion in each step. Apart from this improved efficiency, the metric is determined in a supervised way in this approach, such that the parameters are optimized with respect to the given classification task.

We can work with one full matrix which accounts for a transformation of the whole input space, or, alternatively, with local matrices attached to the individual prototypes. In the latter case, the squared distance of data point $\boldsymbol{\xi}$ from a prototype \boldsymbol{w}_j is computed as $d^{\Lambda^j}(\boldsymbol{w}_j, \boldsymbol{\xi}) = (\boldsymbol{\xi} - \boldsymbol{w}_j)^T \Lambda^j (\boldsymbol{\xi} - \boldsymbol{w}_j)$. Each matrix is adapted individually in the following way: given $\boldsymbol{\xi}$ with closest correct prototype \boldsymbol{w}_J and closest incorrect prototype \boldsymbol{w}_K , we get the update equations

$$\begin{aligned} \Delta \Omega_{lm}^J &= -\epsilon \cdot 2 \cdot \Phi'(\mu(\boldsymbol{\xi})) \cdot \\ &\quad \mu^+(\boldsymbol{\xi}) \cdot \left((\xi_m - w_{J,m}) [\Omega^J(\boldsymbol{\xi} - \boldsymbol{w}_J)]_l \right), \\ \Delta \Omega_{lm}^K &= +\epsilon \cdot 2 \cdot \Phi'(\mu(\boldsymbol{\xi})) \cdot \\ &\quad \mu^-(\boldsymbol{\xi}) \cdot \left((\xi_m - w_{K,m}) [\Omega^K(\boldsymbol{\xi} - \boldsymbol{w}_K)]_l \right). \end{aligned} \quad (14)$$

Localized matrices have the potential to take into account correlations which can vary between different classes or regions in feature space. For instance, clusters with ellipsoidal shape and different orientation could be present in the data. We refer to this general version as localized GMLVQ (LGMLVQ).

Note that LGMLVQ leads to nonlinear decision boundaries which are composed of quadratic pieces, unlike GMLVQ which is characterized by piecewise linear decision boundaries. This way, the receptive fields of the prototypes need no longer be convex or even connected for LGMLVQ, as we will see in the experiments. Depending on the data at hand, this effect can largely increase the capacity of the system.

4 Generalization ability

One of the benefits of prototype-based learning algorithms consists in the fact that they show very good generalization ability also for high dimensional data. This observation can be accompanied by theoretical guarantees. It has been proved in Crammer *et al.* (2003) that basic LVQ networks equipped with the Euclidean metric possess dimensionality independent large-margin generalization bounds, whereby the margin refers to the security of the classification, i.e. the distance of a given data point to the classification boundary. A similar result has been derived in Hammer *et al.* (2005b) for LVQ networks as considered above which possess an adaptive diagonal metric. Remarkably, the margin is directly correlated to the numerator of the cost function as introduced above, i.e. these learning algorithms inherently aim at margin optimization during training. As pointed out in Hammer *et al.* (2005a), these results transfer immediately to kernelized versions of the algorithm where the similarity measure can be interpreted as the composition of the standard scaled Euclidean metric and a fixed kernel map. In the case of an adaptive full matrix, however, these results are not applicable, because the matrix is changed during training. Hence a large number of additional free parameters is introduced since the kernel is optimized according to the given classification task.

Here, we directly derive a large margin generalization bound for LGMLVQ networks with a full adaptive matrix attached to every prototype, whereby we use the ideas of Hammer *et al.* (2005b). It turns out that the so-called hypothesis margin dominates generalization bounds, while the input dimensionality and thus, the number of free parameters of the GMLVQ network does not occur explicitly in the bounds. Thus a good generalization ability of trained networks can be expected independent of the input dimensionality provided a large margin can be achieved. Since the function class implemented by GMLVQ networks is contained in the class of LGMLVQ networks, the bounds also hold for the simpler case.

We consider a LGMLVQ network given by N_w prototypes \mathbf{w}_i . We assume that all inputs $\boldsymbol{\xi}$ fulfill the condition $|\boldsymbol{\xi}| \leq B$ for some $B > 0$, and we assume that weights are also restricted by $|\mathbf{w}_i| \leq B$. As beforehand, we assume that Λ^i is a symmetric positive semidefinite matrix such that the trace is normalized to

1. We consider the case of a binary classification, i.e. only two classes, -1 and 1 , are present. We refer to prototypes labeled with $S \in \{+, -\}$ by \mathbf{w}_i^S .

Classification takes place by a winner takes all rule, i.e.

$$\boldsymbol{\xi} \mapsto \operatorname{sgn} \left(\min_{\mathbf{w}_i^-} \{d^{\Lambda^i}(\mathbf{w}_i^-, \boldsymbol{\xi})\} - \min_{\mathbf{w}_j^+} \{d^{\Lambda^j}(\mathbf{w}_j^+, \boldsymbol{\xi})\} \right), \quad (15)$$

where $d^{\Lambda^i}(\mathbf{w}_i, \boldsymbol{\xi}) = (\boldsymbol{\xi} - \mathbf{w}_i)^T \Lambda^i (\boldsymbol{\xi} - \mathbf{w}_i)$ as beforehand and sgn selects the sign ± 1 of the real number. A trainable LGMLVQ network corresponds to a function in the class

$$\mathcal{F} := \{f : \mathbb{R}^N \rightarrow \{-1, 1\} \mid \exists |\mathbf{w}_i| \leq B, \exists \Lambda^i \text{ such that } \Lambda^i \text{ is symmetric and positive semidefinite with trace 1, such that } f \text{ is given by (15)}\}. \quad (16)$$

Assume some unknown underlying probability measure P is given on $\mathbb{R}^N \times \{-1, 1\}$ according to which training examples are drawn. The goal of learning is to find a function $f \in \mathcal{F}$ such that the generalization error

$$E_P(f) := P(y \neq f(\boldsymbol{\xi})) \quad (17)$$

is as small as possible. However, P is not known during training; instead, examples for the distribution $(\boldsymbol{\xi}_i, y_i)$, $i = 1, \dots, m$ are available, which are independent and identically distributed according to P . Training aims at minimizing the empirical error on the given training data

$$\hat{E}_m(f) := \sum_{i=1}^m |y_i - f(\boldsymbol{\xi}_i)| / m. \quad (18)$$

Thus, the learning algorithm generalizes to unseen data if $\hat{E}_m(f)$ becomes representative for $E_P(f)$ for an increasing number of examples m with high probability, i.e. if we can automatically guarantee a small error on *any* possible input to the learned function, given the *trained* inputs are correct. This bound should hold simultaneously for any function f of the class, in particular for the network trained according to the given sample set.

We will not derive bounds which are directly based on the empirical error $\hat{E}_m(f)$, rather, we incorporate the security of a classification in terms of the classification margin. For a function f as given by (15), we consider the related real-valued function

$$M_f : \boldsymbol{\xi} \mapsto \left(\min_{\mathbf{w}_i^-} \{d^{\Lambda^i}(\mathbf{w}_i^-, \boldsymbol{\xi})\} - \min_{\mathbf{w}_j^+} \{d^{\Lambda^j}(\mathbf{w}_j^+, \boldsymbol{\xi})\} \right) \quad (19)$$

which is obtained by dropping the function sgn . The sign of this real value determines the output class and the size of its absolute value indicates the security of the classification, i.e. the margin of the classifier with respect to input $\boldsymbol{\xi}$ around the decision boundary. The larger this margin, the more robust is the

classification of ξ with respect to noise in the input or function parameters. We refer to the resulting class of real-valued functions implemented by LGMLVQ networks by $M_{\mathcal{F}}$.

Assume $\rho > 0$ estimates the minimum security of the classification. For the moment, we assume that $\rho \in (0, 1)$ is a fixed number which is chosen a priori independent of the training set and the final LGMLVQ function. Following the approach (Bartlett & Mendelson, 2003), we define the loss function

$$L : \mathbb{R} \rightarrow \mathbb{R}, t \mapsto \begin{cases} 1 & \text{if } t \leq 0 \\ 1 - t/\rho & \text{if } 0 < t \leq \rho \\ 0 & \text{otherwise} \end{cases} . \quad (20)$$

The term

$$\hat{E}_m^L(f) := \sum_{i=1}^m L(y_i \cdot M_f(\xi_i))/m \quad (21)$$

is the so called empirical risk. It accumulates the number of errors for a given data set, and, in addition, also punishes all correct classifications if their margin is smaller than ρ , i.e. it measures the classification accuracy and its robustness with respect to noise. The empirical risk is small iff the number of misclassifications is small and almost all correctly classified points have margin larger than ρ .

It is possible to correlate the generalization error of LGMLVQ networks and this modified empirical error by a dimensionality independent bound: according to Bartlett & Mendelson (2003, Theorem 7) the inequality

$$E_P(f) \leq \hat{E}_m^L(f) + \frac{2}{\rho} \cdot R_m(M_{\mathcal{F}}) + \sqrt{\frac{\ln(4/\delta)}{2m}} \quad (22)$$

holds simultaneously for all functions in $M_{\mathcal{F}}$ with probability at least $1 - \delta/2$, whereby $R_m(M_{\mathcal{F}})$ is the so-called Rademacher complexity of the function class $M_{\mathcal{F}}$. The term $2/\rho \cdot R_m(M_{\mathcal{F}}) + \sqrt{\ln(4/\delta)/2m}$ specifies the so called structural risk, denoted by $\hat{E}_m(M_{\mathcal{F}})$ in the following.

The Rademacher complexity of a function class measures the complexity of the class by considering the correlation of outputs of a function of the class on a given set of points and random variables. The empirical Rademacher complexity of the function class $M_{\mathcal{F}}$ given m samples ξ_i , is defined as the expectation

$$\hat{R}_m(M_{\mathcal{F}}) = E_{\sigma_1, \dots, \sigma_m} \left(\sup_{M_f \in M_{\mathcal{F}}} \left| \frac{2}{m} \sum_{i=1}^m \sigma_i \cdot M_f(\xi_i) \right| \right), \quad (23)$$

where σ_i are independent $\{-1, 1\}$ -valued random variables with zero mean. The Rademacher complexity is defined as the expectation over the samples

$$R_m(M_{\mathcal{F}}) = E_{\xi_1, \dots, \xi_m} \hat{R}_m(M_{\mathcal{F}}), \quad (24)$$

where ξ_i are independent and identically distributed according to the marginal distribution of P on the input space.

Using techniques of Bartlett & Mendelson (2003), we will show in the appendix, that the Rademacher complexity of functions given by a LGMLVQ architecture, $R_m(M_{\mathcal{F}})$, can be limited by the term

$$\mathcal{O}\left(\frac{N_w^{3/2}B^3 + \sqrt{\ln(1/\delta)}}{\sqrt{m}}\right) \quad (25)$$

with probability at least $1 - \delta/2$ where N_w denotes the number of prototypes and B the size restriction of inputs and prototypes. Thus, the overall inequality

$$E_P(f) \leq \hat{E}_m^L(f) + \frac{1}{\sqrt{m}}\mathcal{O}\left(\frac{N_w^{3/2}B^3}{\rho} + \frac{\sqrt{\ln(1/\delta)}}{\min\{1, \rho\}}\right) \quad (26)$$

results which holds simultaneously for all $M_f \in M_{\mathcal{F}}$ and training data with probability at least $1 - \delta$. Since this is valid for all $M_f \in M_{\mathcal{F}}$, matrix parameters as well as prototypes can be adaptive. This bound allows to estimate the deviation of the generalization error of a LGMLVQ network and its result on a given training set. Obviously, the bound is small for a small number of errors and a large margin ρ for almost all correctly classified points. Note that the cost function of LGMLVQ is correlated to the classification error, but it also contains the margin of a data point as denominator of the summands. Thus, LGMLVQ aims at an optimization of the margin during training and at a corresponding simplification of the classifier, such that good generalization can be expected. Note that the bound is independent of the dimensionality of the data. Thus, excellent generalization can be expected also for high dimensional settings.

So far, we assumed that the bound ρ for the margin (points with smaller margin contribute to the error) is fixed a priori. In applications, it is reasonable to choose ρ based on the outcome of a training algorithm such that almost all training examples have a margin larger than ρ . For this case, a generalization of the argumentation is possible which only assumes some prior about a reasonable range of the margin: Assume the empirical margin can be upper bounded by $C > 0$, a naive bound being e.g. the maximum distance of data in the given training set. We define $\rho_i = C/i$ for $i \geq 1$, and we choose prior probabilities $p_i \geq 0$ with $\sum p_i = 1$ which indicate the confidence in achieving an empirical margin of size at least ρ_i for almost all training data. We define the cost function L_i as above associated to margin ρ_i and the corresponding empirical error as $\hat{E}_m^{L_i}(f)$. We are interested in the probability

$$P\left(\exists i E_P(f) \geq \hat{E}_m^{L_i}(f) + \epsilon(i)\right), \quad (27)$$

where the bound

$$\epsilon(i) = \frac{1}{\sqrt{m}}\mathcal{O}\left(\frac{N_w^{3/2}B^3}{\rho_i} + \frac{\sqrt{\ln(1/(p_i\delta))}}{\min\{1, \rho_i\}}\right) \quad (28)$$

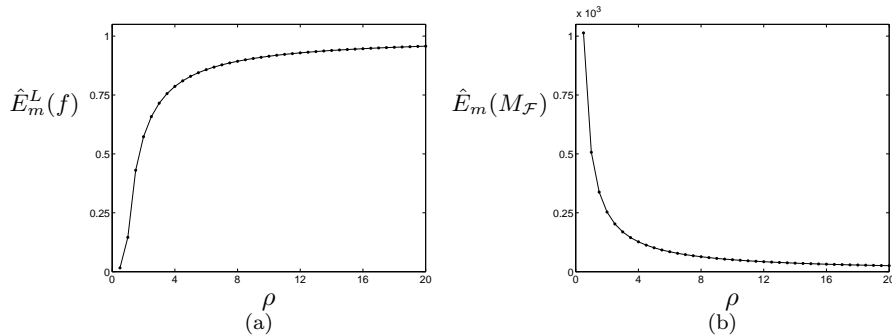


Figure 1: Empirical risk (a) and structural risk (b) vs. ρ for an LGMLVQ network as derived for the Splice Site Recognition data set in Sec. 5 with $\delta = 0.7$.

is chosen according to the inequality (26). We can argue

$$\begin{aligned}
 P\left(\exists i E_P(f) \geq \hat{E}_m^{L_i}(f) + \epsilon(i)\right) &\leq \sum_i P\left(E_P(f) \geq \hat{E}_m^{L_i}(f) + \epsilon(i)\right) \\
 &\leq \sum_i p_i \cdot \delta = \delta,
 \end{aligned}
 \tag{29}$$

because the bounds $\epsilon(i)$ are chosen according to equation (26). Thus, posterior bounds depending on the empirical margin and the prior confidence in achieving this margin can be derived.

The theoretical bound as derived in Eq. (22) is illustrated in an example: Fig. 1 shows the empirical risk and the structural risk of Eq. (22) for different values of ρ for an LGMLVQ network derived in Sec. 5. One can observe the characteristic monotonicity of these two risks with respect to ρ such that the combination of both terms leads to a minimum of the overall bound for an appropriate choice of ρ . However, the constants of the overall bound (Eq. (25)) as derived in this article are still too weak in the sense that the structural risk is dominating for relevant domains of the parameters, i.e. $E_P(f)$ is monotonically decreasing with ρ in the relevant domain. This bound shows the qualitative dependencies for large values of m rather than leading to tight bounds for realistic small-sample settings.

5 Experiments

In the following experiments, we study the performance of matrix relevance adaptation in the context of several learning problems. We compare global and local matrix schemes with the corresponding schemes for relevance vectors as used in GRLVQ, for instance. To this end, we restrict our GMLVQ algorithm to the adaptation of diagonal matrices. Note that we implement gradient steps

in Ω , while the original GRLVQ scheme corresponds to steepest descent w.r.t. diagonal elements of $\Lambda = \Omega^T \Omega$, directly. This modification was necessary in order to allow for a fair comparison of vector and matrix adaptation with the same learning rates.

Note that GMLVQ optimization involves number of hyperparameters, including the number of prototypes per class, the learning rates, and their annealing schedule. As for LVQ itself, the number of prototypes per class is a critical parameter and it depends on the number of modes of the underlying class distribution. In our experiments, one prototype per class turned out to be sufficient to obtain reasonable results; in general, more than one prototype per class can improve the classification accuracy depending on the underlying data distribution. To initialize the prototypes, we choose the mean values of random subsets of data points selected from each class. The learning rates have been optimized by means of crossvalidation within a reasonable range. Learning rate annealing is performed according to the following schedule:

$$\epsilon_{p,m}(t) = \frac{\epsilon_{p,m}(0)}{1 + \tau \cdot (t - 1)}, \quad (30)$$

where $\epsilon_p = 4 \cdot \epsilon$ in Eq. (10) and $\epsilon_m = 4 \cdot \epsilon$ in Eq.s (11) and (14), respectively. t denotes the number of sweeps through the training set. The parameter τ determines the speed of annealing and is selected for every application individually. Prior to learning, the matrix Λ is set to be diagonal with $\Lambda_{ii} = 1/N, \forall i$. The same holds for local relevance matrices, respectively.

Artificial data

In a first illustrative experiment, the algorithms are applied to two-dimensional artificial data in a binary classification problem. Each class corresponds to a cigar-shaped cluster with equal prior weights. Raw data is generated according to axis-aligned Gaussians with mean $\mu_1 = [1.5, 0.0]$ for class 1 and $\mu_2 = [-1.5, 0.0]$ for class 2 data, respectively. In both classes the standard deviations are $\sigma_{11} = 0.5$ and $\sigma_{22} = 3.0$. These clusters are rotated independently by the angles $\varphi_1 = \pi/4$ and $\varphi_2 = -\pi/6$ so that the two clusters intersect. Training and test set consist of 600 data points per class, respectively. In order to reduce the influence of *lucky set* compositions, the experiments are performed on ten statistically independent data sets. One of these data sets is visualized in Fig. 2. It will be used for demonstration purposes in the following.

For training, we use one prototype per class and the following settings: we use the squared Euclidean metric (GLVQ), an adaptive diagonal metric (GRLVQ), individual adaptive diagonal metrics for each prototype (LGRLVQ), a global adaptive matrix (GMLVQ), and individual adaptive matrices for every prototype (LGMLVQ). Training is done for 1000 epochs in total. In all experiments, the learning rates are chosen differently for prototypes and metric parameters and are annealed during training. The initial learning rate $\epsilon_p(0)$ for prototypes is chosen as 0.005, the initial learning rate for the metric parameters

$\epsilon_m(0)$ is set to 0.0005 and $\tau = 0.001$. The mean classification accuracies on the training and test sets are summarized in the left panel of Tab. 1. The position of the resulting prototypes and decision boundaries for the example data set are shown in Fig. 2 (a)-(e).

The relevance matrix

$$\Lambda \approx \begin{pmatrix} 0.7783 & 0.4154 \\ 0.4154 & 0.2217 \end{pmatrix}$$

which results from GMLVQ training on the example data set has the eigenvalues one and zero. The same eigenvalue spectrum is obtained in all runs, i.e. for all independent data sets. It implies that the algorithm determines only one new feature to discriminate the data. The respective direction in feature space is defined by the first eigenvector of Λ . The corresponding matrix Ω projects the data onto this 1-dimensional subspace as depicted in Fig. 2(g). Furthermore, this figure displays that class 2 samples spread only slightly around their prototype in the new feature space. The opposite holds for class 1 samples, implying large distances of these data points to their prototype. Accordingly, the classification performance is much better for class 2 samples. This can also be seen in the receptive fields in Fig. 2(d). On this data set, almost 98% of the training error goes back to class 1 data.

For local matrix adaptation, the algorithm also tends towards a state with eigenvalues one and zero for both matrices Λ_1 and Λ_2 . The resulting local relevance matrices on the example data set are

$$\Lambda_1 \approx \begin{pmatrix} 0.5155 & -0.4993 \\ -0.4993 & 0.4845 \end{pmatrix} \quad \Lambda_2 \approx \begin{pmatrix} 0.7498 & 0.4331 \\ 0.4331 & 0.2502 \end{pmatrix}.$$

Figures 2(h) and 2(i) denote the projections of the training set to the feature spaces which are determined for the two prototypes individually. One can clearly observe the benefit of individual matrix adaptation: it allows each prototype to shape its distance measure according to the local ellipsoidal form of the class.

Table 1: Percentage of correctly classified patterns for the artificial data and the image segmentation data using different LVQ algorithms.

Artificial data			Image data		
Algorithm	Training	Test	Algorithm	Training	Test
GLVQ	74.5	73.4	GLVQ	84.8	83.0
GRLVQ	74.9	73.7	GRLVQ	89.0	88.9
GMLVQ	79.9	79.0	GMLVQ	90.4	90.2
LGRLVQ	80.0	79.5	LGRLVQ	91.4	90.0
LGMLVQ	91.4	92.2	LGMLVQ	98.8	94.4

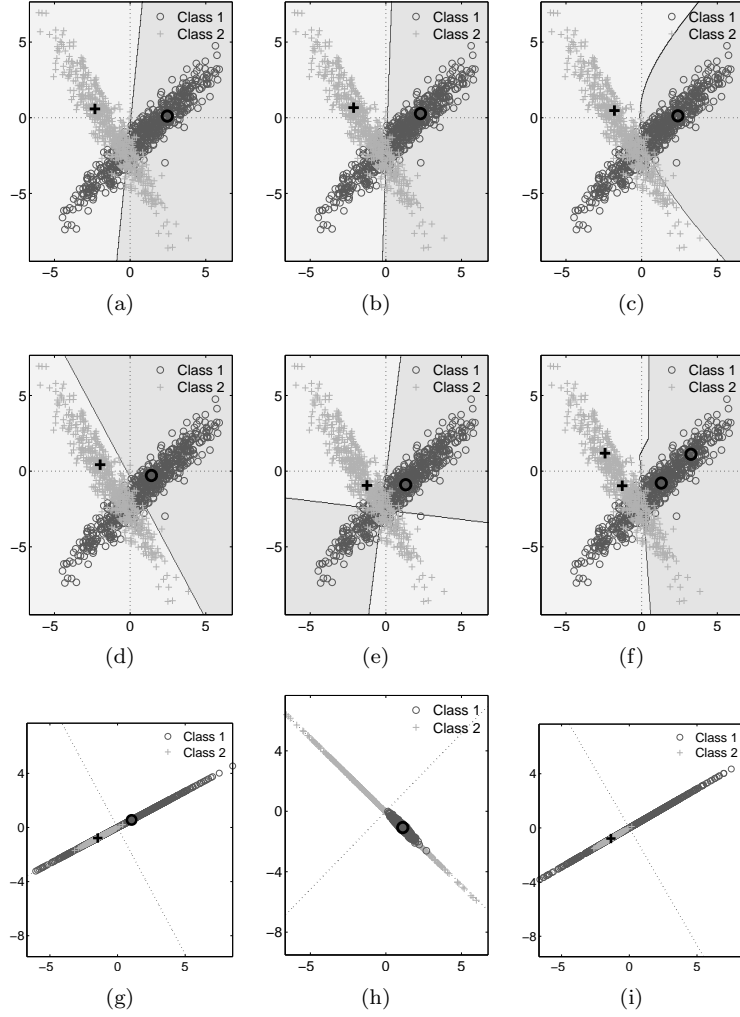


Figure 2: Artificial data set. (a)-(f) Prototypes and receptive fields, (a) GLVQ, (b) GRLVQ, (c) LGRLVQ, (d) GMLVQ, (e) LGMLVQ, (f) GLVQ using two prototypes per class, (g) Training set transformed by global matrix Ω , (h) Training set transformed by local matrix Ω_1 , (i) Training set transformed by local matrix Ω_2 . In (g), (h), (i) the dotted lines correspond to the eigendirections of Λ , Λ_1 and Λ_2 , respectively.

This way, the data points of both ellipsoidal clusters can be classified correctly except for the tiny region where the classes overlap. A partitioning of the feature space into four receptive fields could also be realized based on the Euclidean distance, if two prototypes per class are adapted to the data. However, we

observe that GLVQ additionally requires a proper initialization in order to find the respective prototype constellation. Using the same initialization strategy as in the previous experiments, GLVQ yields the prototype constellation visible in Fig. 2(f). The classification performance obtained by this set of prototypes is nearly the same as in Fig. 2(a). Hence, the crucial factor to influence the classification performance of LVQ-classifiers on this data set is not the number of prototypes, but the choice of an appropriate distance measure. Note that, for local metric parameter adaptation, the receptive fields of the prototypes are no longer separated by straight lines (Fig. 2(c)) and need no longer be convex (Fig. 2(e)).

Image Segmentation Data

In a second experiment, we apply the algorithms to the image segmentation data set provided by the UCI repository of Machine Learning (Newman m. fl., 1998). The data set contains 19-dimensional feature vectors, which encode different attributes of 3×3 pixel regions extracted from outdoor images. Each such region is assigned to one of seven classes (brickface, sky, foliage, cement, window, path, grass). The features 3-5 are (nearly) constant and are eliminated for these experiments. As a further preprocessing step, the features are normalized to zero mean and unit variance. The training set contains 210 data points (30 samples per class), the test data consists of 300 samples per class. We split the test data into a test and a validation set of equal size in order to optimize the model settings.

In our first set of experiments, each class is approximated by one prototype, respectively. We use the parameters settings $\epsilon_p(0) = 0.005$, $\epsilon_m(0) = 0.0001$ and $\tau = 0.0001$ and continue learning until the validation error remains constant or starts indicating overfitting effects. In order to reduce the influence of random fluctuations, we average our results over ten runs with varying initializations.

The mean classification accuracies are summarized in the right panel of Tab. 1. The algorithms based on adaptive distance measures show a better performance than GLVQ. Remarkably, using different metrics influences the final location of the prototypes in feature space only slightly. The prototypes saturate close to the class conditional means in all experiments. Clear differences affect only a small number of features in certain classes (see Fig. 3).

Figure 4(a) displays the averaged classification errors on training- and validation sets in the course of GMLVQ-Training. The validation error starts degrading slightly after approximately 1500 epochs. One of the matrices after 1500 epochs is visualized in Figure 5. The eigenvalue spectrum shows that the classifier uses a ten dimensional space to classify the data. The dimension weighted as most relevant in the original space is feature 16 (hue-mean, see Newman m. fl., 1998). GRLVQ training with identical initializations and learning parameters also weights the same dimension as most important. However, the relevance profile is much more pronounced ($\lambda_{16}^{\text{GRLVQ}} \approx 0.9$). The additional consideration of correlations for the computation of distance values causes a less distinct relevance profile.

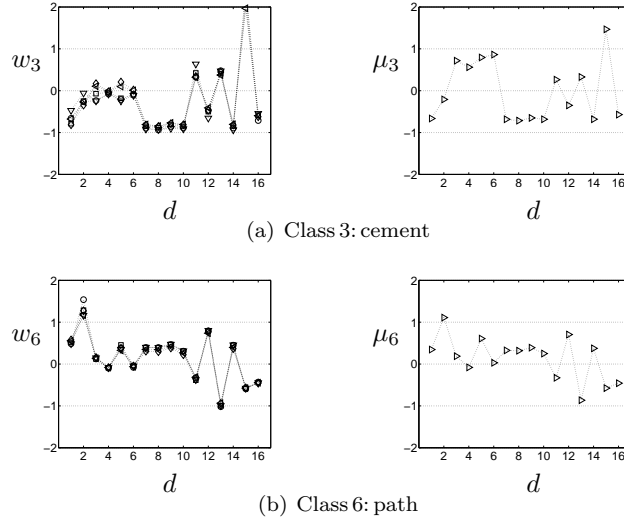


Figure 3: Class conditional means (μ) of 2 classes in comparison to the prototypes identified by the different algorithms. GLVQ: circle, GRLVQ: square, GMLVQ: diamond, LGRLVQ: triangle (down), LGMLVQ: triangle (left).

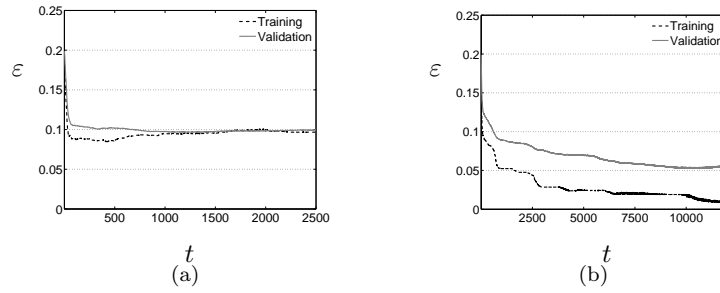


Figure 4: Evolution of the mean training and test error in the course of (a) GMLVQ-Training. (b) LGMLVQ-Training.

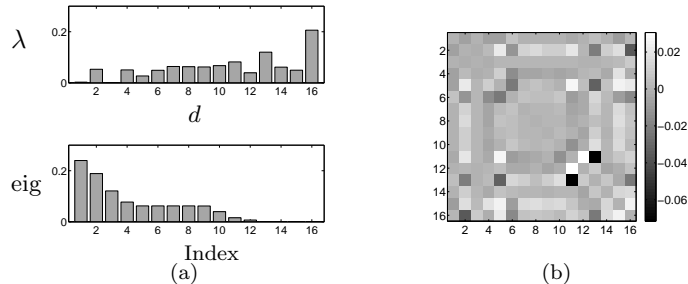


Figure 5: Visualization of the global relevance matrix Λ after 1000 epochs GMLVQ-Training (a) Diagonal elements and eigenvalues (b) Off-diagonal elements. The diagonal elements are set to zero for the plot.

In a further sequence of GMLVQ experiments, we analyse the algorithms' sensitivity with respect to the number of prototypes. We determine the class-wise rates of misclassification after each experiment and repeat the training providing one further prototype for the class contributing most to the overall rate of misclassification. A system consisting of 10 prototypes achieves 90.9% mean test accuracy, using 13 prototypes reaches 91.2% mean accuracy on the test sets. The slight improvement in classification performance is accompanied by an increasing training time until convergence.

The training error during LGMLVQ-training remains constant after approximately 12 000 sweeps through the training set (Fig. 4(b)). The validation error shows slight over-fitting effects. It reaches a minimum after approximately 10 000 epochs and increases in the further course of training. In the following we present the results obtained after 10 000 epochs. At this point, training of individual matrices per prototype achieves a test accuracy of 94.4%, an improvement of approximately 4.9% compared to GMLVQ and LGRLVQ. We are aware of only one SVM result in the literature which is applicable for comparing the performance. In Prudent & Ennaji (2005), the authors achieve 93.95% accuracy on the test set.

Figure 6 shows the diagonal elements and eigenvalue spectra of all local matrices we obtain in one run which are also representative for the other experiments. Matrices with a clear preference for certain dimensions on the diagonal also display a distinct eigenvalue profile (e.g. Λ_1, Λ_5). Similarly, matrices with almost balanced relevance values on the diagonal exhibit only a weak decay from the first to the second eigenvalue (e.g. Λ_2, Λ_7). This observation for diagonal elements and eigenvalues coincides with a similar one for the off-diagonal elements. Figure 7 visualizes the off-diagonal elements of the local matrices Λ_1, Λ_2 and Λ_5 . Corresponding to the balanced relevance- and eigenvalue profile of matrix Λ_2 , the off-diagonal elements are only slightly different from zero. This may indicate diffuse data without a pronounced, hidden structure. There are obviously no other directions in feature space which could be used to significantly minimize distances within this class. On the contrary, the matrices

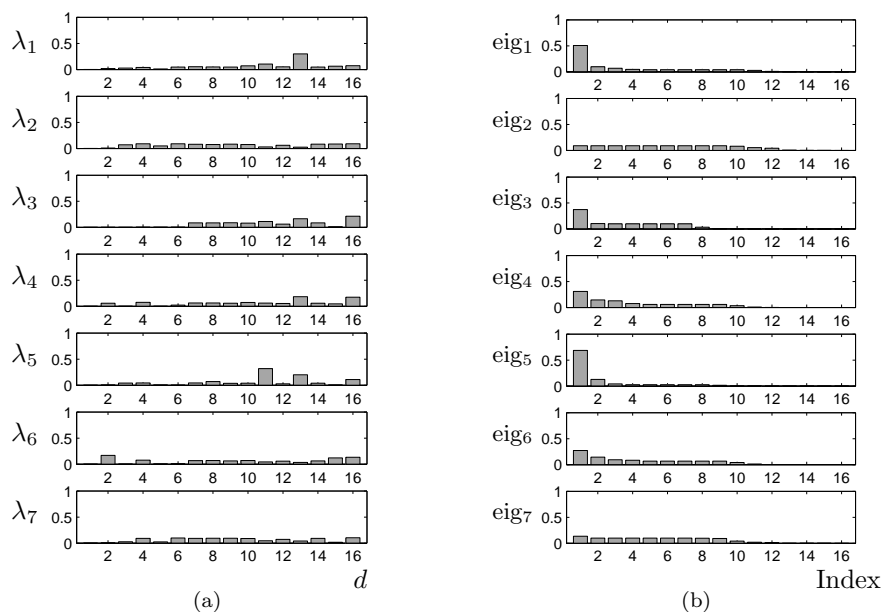


Figure 6: (a) Diagonal elements of local relevance matrices Λ_{1-7} . (b) Eigenvalue spectra of local relevance matrices Λ_{1-7} .

Λ_1 and Λ_5 show a clearer structure. The off-diagonal elements cover a much wider range and there is a clearer emphasis on particular dimensions. This implies that class-specific correlations between certain features have significant influence. The most distinct weights for correlations with other dimensions are obtained for features, which also gain high relevance values on the diagonal. It is visible that especially relations between the dimensions encoding color information are emphasized. The dimensions weighted as most important are features 11: exred-mean ($2R - (G + B)$) and 13: exgreen-mean ($2G - (R + B)$) in both classes. Furthermore, the off-diagonal elements highlight correlations with e.g. feature 8: rawred-mean (average over the regions red values), feature 9: rawblue-mean (average over the regions green values), feature 10: rawgreen-mean (average over the regions green values). For a description of the features, see Newman m. fl. (1998).

Splice Site Recognition

As a second benchmark test, we apply the algorithms to the publicly available *C. elegans* data set for the detection of splice sites. The data can be downloaded at <http://www2.fml.tuebingen.mpg.de/raetsch/projects/AnuSplice>. The feature vectors encode a sequence of 50 nucleotides with a potential splice site in the center, in between the characteristic dinucleotide AG. The classification task

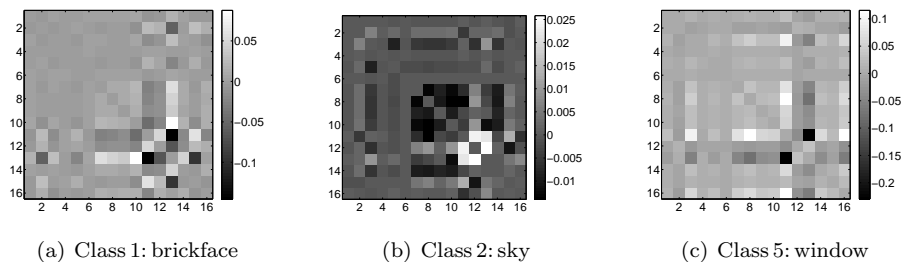


Figure 7: Visualization of the off-diagonal elements of the local matrices $\Lambda_{1,2,5}$. The diagonal elements are set to zero for the plot.

consists in separating sequences containing a splice site from sequences without a splice site, i.e. a two-class problem is defined accordingly. The 4 nucleotides are encoded as corners of a tetrahedron in a 3-dimensional vector space. This realizes equal pairwise distances between the nucleotides. The redundant dinucleotide AG in the center of all feature vectors is removed. Accordingly, the sequence information is represented by a vector in \mathbb{R}^N with $N = 144$. The data consists of 5 data sets containing 1000 data points for training and 10000 data points for testing, respectively. The sets are not balanced and their composition varies slightly. They contain approximately two times more non-splice site-samples than examples of splice sites.

We would like to stress that our main interest in this experiment is not related to the biological aspects of the classification problem. We will put emphasis on the analysis of our method and the comparison of the new algorithm to the adaptation of relevance vectors.

We choose the simplest setting and approximate each class with one prototype respectively. The initial learning parameters are chosen as $\eta_p(0) = 1 \cdot 10^{-3}$ and $\eta_m(0) = 1 \cdot 10^{-5}$. Equation (30) is used for annealing the values during training with $\tau = 1 \cdot 10^{-4}$.

In all experiments, GMLVQ turns out to be more robust than GRLVQ. The learning curves of GRLVQ show strong fluctuations until they finally saturate at a constant level. The mean test set accuracy in this limit is $86.8\% \pm 0.03\%$. In earlier states of training the system shows better classification accuracy of above 90%. But GRLVQ performs a very strong feature selection in the further course of learning and the performance degrades in response to this oversimplification. Fig. 8(a) displays the evolution of the weight values on one of the five data sets which is also representative for the other experiments. When the error finally converges, only three factors remain significantly different from zero. Typical relevance factors are

$$\lambda_{\text{GRLVQ}}^{61,64,72} \approx (0.04, 0.63, 0.33).$$

GMLVQ shows a larger stability during learning. The error curves display only small oscillations, but indicate slight over-fitting effects. In the course of training, we observe an immediate focus on a single linear combination of the

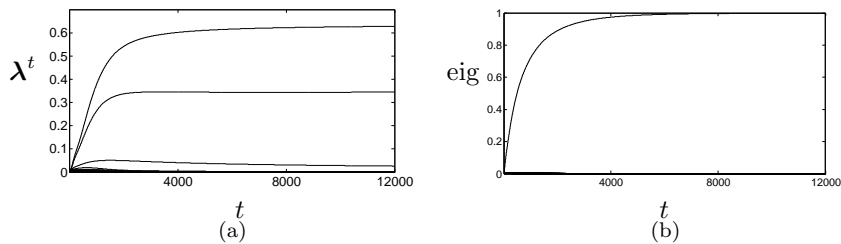
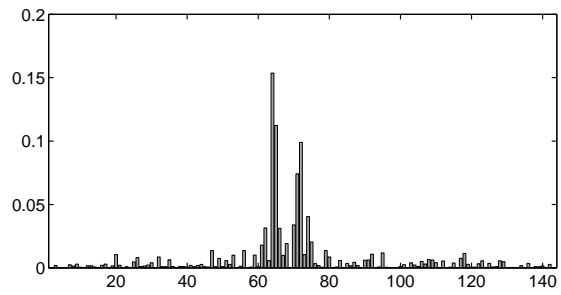


Figure 8: (a) Evolution of the elements of relevance vector λ in the course of GRLVQ-Training. (b) Evolution of the eigenvalues of relevance matrix Λ in the course of GMLVQ training.

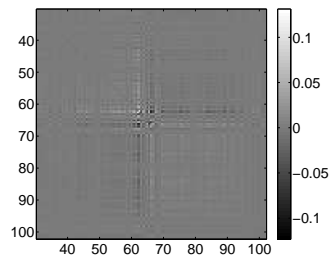
original features. Fig. 8(b) displays the eigenvalues of Λ as a function of training time. Except for the first eigenvalue, all other factors begin to decrease to zero immediately after starting metric adaptation. After approximately 10 000 epochs, the system finally reaches a state with only one eigenvalue remaining. At this point, the mean classification accuracy is $93.2\% \pm 7 \cdot 10^{-4}\%$ on the test sets. Due to these extreme configurations of the relevance matrix, the same accuracy can be achieved in the 1-dimensional subspace defined by the first eigenvector of Λ . Accordingly, our method allows to reduce the number of features dramatically, without losing classification performance.

Fig. 9 visualizes one of the final global matrices Λ . Features in the center display the highest relevances on the diagonal. This implies that the region around the potential splice site is of particular importance and mirrors biological knowledge. Additionally, the classifier considers correlations between different features to evaluate the similarity between the prototypes and new feature vectors. Similar to the previous experiment, the most significant off-diagonal Λ_{ij} relate to the features with the highest diagonal relevances. These correlations result in a cross-like structure in the visualization of the off-diagonal elements, see Fig. 9(b),(c).

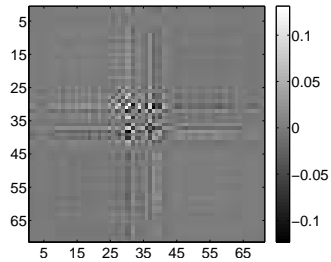
The prototypes can be interpreted as a sequence of 48 nontrivial combinations of the four bases. They converge after approximately 4000 epochs, independent of the additional adaptation of a relevance vector or a relevance matrix in the distance measure. The representatives detected by GMLVQ approximate the data more appropriately compared to the GRLVQ-prototypes. GRLVQ slightly pushes the prototypes away from the data, several components leave the boundaries of the tetrahedron. This effect is even stronger, if we train GLVQ with the fixed Euclidean metric. Fig. 10 visualizes the prototypes identified by GMLVQ on one data set. All subcomponents of the class 1 prototype are located close to the origin, the tetrahedron’s center of mass. In this position they have almost equal distance to the four vertices which represent the bases. On the contrary, the splice site prototype exhibits a more specific structure. Especially the components with high relevance values are located close to one of the corners, i.e. one of the nucleotides, and allow for a better semantic



(a) Diagonal elements



(b) Off-diagonal elements



(c) Off-diagonal elements dimensions 38-108

Figure 9: Visualization of the resulting global matrix Λ after GMLVQ-Training. The diagonal elements are set to zero in (b) and (c).

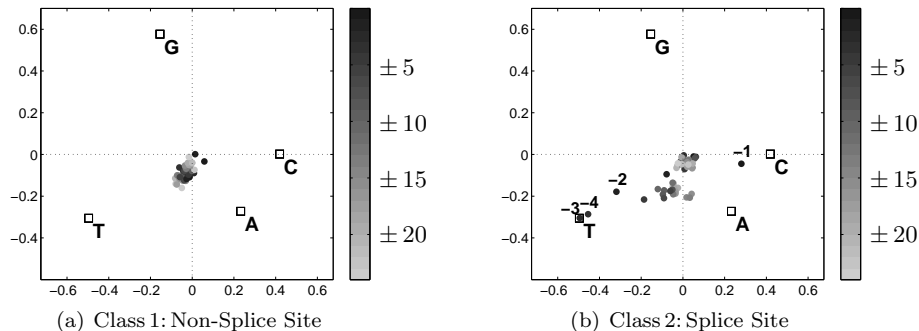


Figure 10: Visualization of the resulting prototypes after GMLVQ training. The plots show the projections of the 3-dimensional elements encoding the separate components of the sequences, onto the x-y plane. The gray values display the subcomponent’s position in the sequence, i.e. the distance to the potential splice site. In the right plot, we labeled the positions relative to the center which are lying closest to one of the bases.

interpretation.

In accordance with our findings for GMLVQ, localized matrix learning detects one dominating feature per class. As in our GMLVQ experiments, we train the LGMLVQ system for 12 000 epochs. Note, that both local matrices are updated in each learning step. The resulting matrices Λ_1 and Λ_2 do not show the extreme eigenvalue settings like the global relevance matrix Λ . But the error curves indicate over-fitting and we do not continue training. The largest eigenvalues range from 0.76 to 0.82 (Λ_1) and from 0.96 to 0.97 (Λ_2) in the different experiments. The diagonal and the off-diagonal elements of the local matrices show the same characteristic patterns as the global matrix. However, the values of matrix Λ_2 are more distinct. The mean classification accuracy is slightly better compared to GMLVQ ($93.5\% \pm 5 \cdot 10^{-4}\%$). When we perform the classification based on the two features defined by the first eigendirections of Λ_1 and Λ_2 , we loose almost no performance and still achieve $93.45\% \pm 3 \cdot 10^{-4}\%$ test accuracy. SVM results reported in the literature even lie above 96% (Hammer m. fl., 2004; Rätsch & Sonnenburg, 2004) test accuracy. Note, however, that our classifier is extremely sparse and simple and still achieves a performance which is only slightly worse.

6 Conclusion

We have proposed a new metric learning scheme for LVQ classifiers which allows to adapt a full matrix according to the given classification task. This scheme extends the successful relevance learning vector quantization algorithm such that correlations of dimensions can be accounted for during training. The

learning scheme can be derived directly as a stochastic gradient of the GLVQ cost function such that convergence and flexibility of the original GLVQ is preserved. Since the resulting classifier is represented by prototype locations and matrix parameters, the results can be interpreted by humans: prototypes show typical class representatives and matrix parameters reveal the importance of input dimensions for the diagonal elements and the importance of correlations for the off-diagonal elements. Local as well as global parameters can be used, i.e. relevance terms which contribute to a good description of single classes or the global classification, respectively, can be identified. The efficiency of the model has been demonstrated in several application scenarios, demonstrating impressively the increased capacity of local matrix adaptation schemes. Interestingly, local matrix learning obtains a classification accuracy which is similar to the performance of the SVM in several cases, while it employs less complex classification schemes and maintains intuitive interpretability of the results.

The new class of algorithms drastically increases the number of free parameters of training, since full $N \times N$ matrices are updated. For a global metric, this corresponds to an adaptive linear transformation of the space according to the given classification task. For local metrics, every prototype uses its own transformation to emphasize characteristics of the respective classes. In this case, the receptive fields are no longer separated by planes but quadratic surfaces. Furthermore, they need not be convex, such that more complex settings can easily be accounted for. Clearly, straightforward modifications can be considered which employ class-wise relevance matrices or other intermediate schemes.

Interestingly, only very mild over-fitting is observed in our experiments, and matrix adaptation leads to excellent generalization despite the increased number of free parameters. This effect can be explained by an inherent regularization which is present in GLVQ adaptation schemes: the margin of the classifier with respect to training points, i.e. the difference of their distance to the closest correct versus the closest wrong prototype is optimized. We have rigorously shown that generalization bounds which do include the margin, but which are independent of the dimensionality of the input space and the dimensionality of the adaptive matrices, can be derived. Thus, the extended classification scheme provides increased capacity without diminishing the excellent generalization capability of LVQ classifiers.

Albeit the generalization ability of GMLVQ does not decrease regarding the largely increased number of free parameters, one drawback of the method consists in computational costs which scale quadratically with the data dimensionality. Thus, quadratic instead of linear effort can be observed in every update step. Obviously, the method becomes computationally infeasible for very high dimensional data, i.e. 50 or more dimensions. This issue can be avoided using standard preprocessing such as e.g. a global PCA of the data. As an alternative, matrix adaptation can directly be adapted to this setting by reducing number of free parameters of a given matrix by enforcing e.g. a limited rank of the matrix. Formally, this can be achieved by setting $\Lambda_j = \Omega_j^T \Omega_j$ with $M \times N$ matrices Ω_j where $M \ll N$, which leads to positive-semidefinite matrices with limited rank. This way, the rank M of Λ_j constitutes another metaparameter which has to be

determined using e.g. crossvalidation. First promising results using this method can be found in (Bunte m. fl., 2008).

Appendix

Here we derive upper bounds for the Rademacher complexity of LGMLVQ networks. Assume a function class implemented by LGMLVQ networks is given as above, $M_{\mathcal{F}}$. In analogy to the Rademacher complexity, one can define the empirical Gaussian complexity, given m samples ξ_i , as the expectation

$$\hat{G}_m(M_{\mathcal{F}}) = E_{g_1, \dots, g_m} \left(\sup_{M_f \in M_{\mathcal{F}}} \left| \frac{2}{m} \sum_{i=1}^m g_i \cdot M_f(\xi_i) \right| \right) \quad (31)$$

where g_i are independent Gaussian variables with zero mean and unit variance. The Gaussian complexity is defined as the expectation over the samples

$$G_m(M_{\mathcal{F}}) = E_{\xi_1, \dots, \xi_m} \hat{G}_m(M_{\mathcal{F}}) \quad (32)$$

where ξ_i are independent and identically distributed according to the marginal of P . These quantities are closely related to the Rademacher complexity. According to Bartlett & Mendelson (2003, Lemma 4) and Ledoux & Talagrand (1991), respectively, the inequality

$$\sqrt{\pi/2} \cdot R_m(M_{\mathcal{F}}) \leq G_m(M_{\mathcal{F}}) \quad (33)$$

holds.

Our aim is to upper bound the Rademacher complexity $R_m(M_{\mathcal{F}})$ with probability at least $1 - \delta/2$ whereby M_f has the form

$$M_f(\xi) = \left(\min_{\mathbf{w}_i^-} \{d^{\Lambda^i}(\mathbf{w}_i^-, \xi)\} - \min_{\mathbf{w}_j^+} \{d^{\Lambda^j}(\mathbf{w}_j^+, \xi)\} \right) \quad (34)$$

As beforehand, we assume that $|\xi| \leq B$, $|\mathbf{w}_i| \leq B$, Λ^i is symmetric and positive semidefinite with trace 1, and we assume that N_w prototypes are present. Because of equation (33), it holds

$$G_m(M_{\mathcal{F}}) \leq C_2 \Rightarrow R_m(M_{\mathcal{F}}) \leq C_1 \text{ where } C_1 = C_2 \cdot \sqrt{2/\pi} \quad (35)$$

for a bound C_2 on the Gaussian complexity. Therefore, we will limit the Gaussian complexity. The empirical Gaussian complexity and the Gaussian complexity differ by more than ϵ with probability at most $2 \exp(-\epsilon^2 m/8)$ according to Bartlett & Mendelson (2003, Theorem 11), i.e. they differ by no more than $\sqrt{8/m \cdot \ln(4/\delta)}$ with probability at least $1 - \delta/2$. Thus, it is sufficient to upper bound the empirical Gaussian complexity of LGMLVQ networks, since

$$\hat{G}_m(M_{\mathcal{F}}) \leq C_3 \Rightarrow G_m(M_{\mathcal{F}}) \leq C_2 \text{ where } C_2 = C_3 + \sqrt{8/m \cdot \ln(4/\delta)} \quad (36)$$

for a bound C_3 on the empirical Gaussian complexity and where the right hand side holds with probability at least $1 - \delta/2$.

Note that

$$\hat{G}_m \left(\sum_i \mathcal{F}_i \right) \leq \sum_i \hat{G}_m(\mathcal{F}_i) \quad (37)$$

holds for all function classes \mathcal{F}_i due to the triangle inequality. Further, the empirical Gaussian complexity does obviously not change when multiplying a function class by -1 . The function class $M_{\mathcal{F}}$ constitutes of functions which can be expressed as the sum of two terms of identical form, one multiplied by -1 . Thus, we can upper bound $\hat{G}_m(M_{\mathcal{F}})$ by twice the complexity of a function class of functions of the form

$$\xi \rightarrow \min_{\mathbf{w}_i} \{d^{\Lambda^i}(\mathbf{w}_i, \xi)\} \quad (38)$$

where the minimum is taken over at most N_w terms, where N_w denotes the number of prototypes.

As next step, we try to obtain a bound resulting from the minimum in this term. The function which computes the minimum of N_w values,

$$(a_1, \dots, a_{N_w}) \mapsto \min\{a_1, \dots, a_{N_w}\} \quad (39)$$

is Lipschitz continuous with constant $\sqrt{8N_w}$, as can be seen as follows: We find

$$|\min\{a, 0\} - \min\{a', 0\}| = \left\{ \begin{array}{ll} |a - a'| & \text{if } a, a' \leq 0 \\ |a| = -a \leq -a + a' & \text{if } a \leq 0, a' > 0 \\ |a'| = -a' \leq -a' + a & \text{if } a' \leq 0, a > 0 \\ 0 & \text{if } a, a' > 0 \end{array} \right\} \leq |a - a'| \quad (40)$$

Further,

$$\min\{a, b\} = \left\{ \begin{array}{ll} a = a - b + b & \text{if } a \leq b \\ b = 0 + b & \text{if } a > b \end{array} \right\} = \min\{a - b, 0\} + b \quad (41)$$

Hence, by induction,

$$\begin{aligned} & |\min\{a_1, \dots, a_{N_w}\} - \min\{a'_1, \dots, a'_{N_w}\}| \\ &= |\min\{\min\{a_1, \dots, a_{N_w-1}\}, a_{N_w}\} - \min\{\min\{a'_1, \dots, a'_{N_w-1}\}, a'_{N_w}\}| \\ &= |\min\{\min\{a_1, \dots, a_{N_w-1}\} - a_{N_w}, 0\} + a_{N_w} - \min\{\min\{a'_1, \dots, a'_{N_w-1}\} - a'_{N_w}, 0\} + a'_{N_w}| \\ &\leq |a_{N_w} - a'_{N_w}| + |\min\{a_1, \dots, a_{N_w-1}\} - a_{N_w} - \min\{a'_1, \dots, a'_{N_w-1}\} + a'_{N_w}| \\ &\leq 2|a_{N_w} - a'_{N_w}| + |\min\{a_1, \dots, a_{N_w-1}\} - \min\{a'_1, \dots, a'_{N_w-1}\}| \\ &\leq \dots \leq 2|a_1 - a'_1| + \dots + 2|a_{N_w} - a'_{N_w}| \end{aligned} \quad (42)$$

Thus,

$$\begin{aligned}
|\min\{a_1, \dots, a_{N_w}\} - \min\{a'_1, \dots, a'_{N_w}\}|^2 &= 4 \sum_{ij} |a_i - a'_i| |a_j - a'_j| \\
&\leq 4 \sum_{ij} (|a_i - a'_i|^2 + |a_j - a'_j|^2) \\
&\leq 8N_w \sum_i |a_i - a'_i|^2. \tag{43}
\end{aligned}$$

Therefore, we obtain the Lipschitz constant $\sqrt{8N_w}$.

Because of Bartlett & Mendelson (2003, Theorem 14), we find

$$\hat{G}_m(\Phi \circ \mathcal{F}) \leq 2L \sum_i \hat{G}_m(\mathcal{F}_i) \tag{44}$$

for every Lipschitz continuous function Φ on a real vector space with Lipschitz constant L and a function class \mathcal{F} contained in the direct sum of the classes \mathcal{F}_i . Thus, because of the Lipschitz continuity of the function \min , it is sufficient to upper bound the empirical Gaussian complexity of function classes of the form,

$$\boldsymbol{\xi} \mapsto (\boldsymbol{\xi} - \mathbf{w})^t \Lambda (\boldsymbol{\xi} - \mathbf{w}) = \boldsymbol{\xi}^t \Lambda \boldsymbol{\xi} - 2\boldsymbol{\xi}^t \Lambda \mathbf{w} + \mathbf{w}^t \Lambda \mathbf{w}. \tag{45}$$

This decomposes into a linear function

$$\boldsymbol{\xi} \mapsto -2\boldsymbol{\xi}^t \Lambda \mathbf{w} + \mathbf{w}^t \Lambda \mathbf{w} \tag{46}$$

and a quadratic form

$$\boldsymbol{\xi} \mapsto \boldsymbol{\xi}^t \Lambda \boldsymbol{\xi} = \sum_{ij} \Lambda_{ij} \xi_i \xi_j. \tag{47}$$

According to Bartlett & Mendelson (2003, Lemma 22), the empirical Gaussian complexity of linear forms $\mathbf{x} \mapsto \mathbf{a}^t \mathbf{x}$ can be upper bounded by

$$\frac{2K_1 K_2}{\sqrt{m}} \tag{48}$$

where inputs are restricted to $|\mathbf{x}| \leq K_1$ and weights are restricted to $|\mathbf{w}| \leq K_2$. Note that inputs to the LGMLVQ network and prototypes have length at most B , further, the sum of eigenvalues of every matrix Λ^i of the LGMLVQ network is 1. Thus, functions of the form (46) correspond to linear functions with inputs restricted to $B + 1$ and weights restricted to $2B + B^2$. Functions of the form (47) can be interpreted as linear functions with enlarged inputs which size is restricted by B^2 , and weights restricted by 1, since the Frobenius-norm of the matrix is given by the sum of squared eigenvalues in this case. Therefore, we can limit the empirical Gaussian complexity of quadratic forms of this type by $2/\sqrt{m} \cdot ((B+1)(2B+B^2) + B^2)$, hence the empirical Gaussian complexity of $M_{\mathcal{F}}$ can be limited by $\hat{G}_m(M_{\mathcal{F}}) \leq 2 \cdot 2\sqrt{8N_w} N_w \cdot 2/\sqrt{m} \cdot ((B+1)(2B+B^2) + B^2)$.

Collecting all inequalities, we can finally upper bound the Rademacher complexity of LGMLVQ networks by

$$\begin{aligned} & \sqrt{2/\pi} \left(\sqrt{8/m \cdot \ln(4/\delta)} + 2 \cdot 2\sqrt{8N_w} \cdot N_w \cdot 2/\sqrt{m} \cdot ((B+1)(2B+B^2) + B^2) \right) \\ = & \frac{1}{\sqrt{m}} \cdot \mathcal{O} \left(\sqrt{\ln(1/\delta)} + N_w^{3/2} B^3 \right). \end{aligned} \quad (49)$$

Acknowledgement

The authors would like to thank the Max Planck Institute for the Physics of Complex Systems, Dresden, for hospitality during a seminar at which this work was finalized.

References

- Bartlett, P. L., & Mendelson, S. (2003). Rademacher and gaussian complexities: risk bounds and structural results. *Journal of Machine Learning Research*, 3, 463–482.
- Bertsekas, D. (1976). On the goldstein - levitin - polyak gradient projection method. *IEEE Transactions on Automatic Control*, 21(2), 174–184.
- Biehl, M., Ghosh, A., & Hammer, B. (2006). Learning vector quantization: The dynamics of winner-takes-all algorithms. *Neurocomputing*, 69(7-9), 660-670.
- Biehl, M., Ghosh, A., & Hammer, B. (2007). Dynamics and generalization ability of LVQ algorithms. *Journal of Machine Learning Research*, 8, 323-360.
- Bojer, T., Hammer, B., Schunk, D., & Toschanowitz, K. T. von. (2001). Relevance determination in learning vector quantization. I M. Verleysen (red.), *European symposium on artificial neural networks* (s. 271-276).
- Bunte, K., Schneider, P., Hammer, B., Schleif, F.-M., Villmann, T., & Biehl, M. (2008). *Discriminative visualization by limited rank matrix learning* (teknisk rapport nr. MLR-03-2008). Machine Learning Reports. (ISSN:1865-3960 http://www.uni-leipzig.de/compint/mlr/mlr_03_2008.pdf)
- Crammer, K., Gilad-Bachrach, R., Navot, A., & Tishby, A. (2003). Margin analysis of the lvq algorithm. I *Advances in neural information processing systems* (vol. 15, s. 462-469). MIT Press, Cambridge, MA.
- Duda, R. O., Hart, P. E., & Stork, D. G. (2000). *Pattern classification* (Second upplagan). Wiley-Interscience.
- Gath, I., & Geva, A. B. (1989). Unsupervised optimal fuzzy clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 11(7), 773–780.

- Ghosh, A., Biehl, M., & Hammer, B. (2006). Performance analysis of lvq algorithms: a statistical physics approach. *Neural Networks*, 19(6), 817–829.
- Gustafson, E., & Kessel, W. (1997). Fuzzy clustering with a fuzzy covariance matrix. I *Ieee conference on decision and control* (ss. 761–766).
- Hammer, B., Strickert, M., & Villmann, T. (2004). Prototype based recognition of splice sites. I *Bioinformatic using computational in intelligence paradigms* (ss. 25–56). Springer.
- Hammer, B., Strickert, M., & Villmann, T. (2005a). On the generalization ability of GRLVQ networks. *Neural Processing Letters*, 21(2), 109–120.
- Hammer, B., Strickert, M., & Villmann, T. (2005b). Supervised neural gas with general similarity measure. *Neural Processing Letters*, 21(1), 21–44.
- Hammer, B., Strickert, M., & Villmann, T. (2005c, February). Supervised neural gas with general similarity measure. *Neural Processing Letters*, 21(1), 21–44.
- Hammer, B., & Villmann, T. (2002). Generalized relevance learning vector quantization. *Neural Networks*, 15(8-9), 1059–1068.
- Kohonen, T. (1997). *Self-organizing maps* (Second upplagan). Berlin, Heidelberg: Springer.
- Ledoux, M., & Talagrand, M. (1991). *Probability in banach spaces*. Springer.
- Newman, D. J., Hettich, S., Blake, C. L., & Merz, C. J. (1998). *Uci repository of machine learning databases*. <http://archive.ics.uci.edu/ml/>. University of California, Department of Information and Computer Science.
- Prudent, Y., & Ennaji, A. (2005). A k nearest classifier design. *Electronic Letters on Computer Vision and Image Analysis*, 5(2).
- Rätsch, G., & Sonnenburg, S. (2004). Accurate splice site prediction for caenorhabditis elegans. I *Kernel methods in computational biology* (s. 277–298). MIT Press.
- Sato, A., & Yamada, K. (1996). Generalized learning vector quantization. I M. C. M. D. S. Touretzky & M. E. Hasselmo (red:er), *Advances in neural information processing systems 8. proceedings of the 1995 conference* (s. 423–9). Cambridge, MA, USA: MIT Press.
- Schleif, B. H. F.-M., & Villmann, T. (2005). *On the generalization ability of prototype-based classifiers with local relevance determination* (teknisk rapport nr. Ifi-05-14). Clausthal University of Technology.
- Schneider, P., Bunte, K., Stiekema, H., Hammer, B., Villmann, T., & Biehl, M. (2008). *Regularization in matrix relevance learning* (teknisk rapport nr. MLR-02-2008). Machine Learning Reports. (ISSN:1865-3960 http://www.uni-leipzig.de/compint/mlr/mlr_02.2008.pdf)

- Seo, S., Bode, M., & Obermayer, K. (2003). Soft nearest prototype classification. *IEEE Transactions on Neural Networks*, *14*, 390-398.
- Seo, S., & Obermayer, K. (2003). Soft learning vector quantization. *Neural Computation*, *15*(7), 1589–1604.
- Shalev-Schwartz, S., Singer, Y., & Ng, A. (2004). Online and batch learning of pseudo-metrics. I *Proceedings of the 21st international conference on machine learning* (s. 94). New York, USA: ACM.
- Bibliography on the Self-Organizing Map (SOM) and Learning Vector Quantization (LVQ)*. (2002). Neural Networks Research Centre, Helsinki University of Technology.
- Weinberger, K., Blitzer, J., & Saul, L. (2006). Distance metric learning for large margin nearest neighbor classification. I Y. Weiss, B. Schölkopf, & J. Platt (red:er), *Advances in neural information processing systems 18* (ss. 1473–1480). Cambridge, MA: MIT Press.