

# Functional Relevance Learning in Generalized Learning Vector Quantization

M. Kästner<sup>1</sup>, B. Hammer<sup>2</sup>, M. Biehl<sup>3</sup>, and T. Villmann<sup>\*</sup>

<sup>1</sup> University of Applied Sciences Mittweida, Computational Intelligence Group, Technikumplatz 17, 09648 Mittweida, Germany

<sup>2</sup> University Bielefeld, Center of Excellence - Cognitive Interaction Technology CITEC, Universitätsstrasse 21-23, 33615 Bielefeld, Germany

<sup>3</sup> University of Groningen, Johann Bernoulli Institute for Mathematics and Computer Science, P.O. Box 407, 9700 AK Groningen, The Netherlands

## Abstract

Relevance learning in learning vector quantization is a central paradigm for classification task depending feature weighting and selection. We propose a functional approach to relevance learning for high-dimensional functional data. For this purpose we compose the relevance profile by a superposition of only a few parametrized basis functions taking into account the functional character of the data. The number of these parameters is usually significantly smaller than the number of relevance weights in standard relevance learning, which is the number of data dimensions. Thus, instabilities in learning are avoided and an inherent regularization takes place. In addition, we discuss strategies to obtain sparse relevance models for further model optimization.

**keywords:** functional vector quantization, relevance learning, feature weighting and selection, sparse models

---

\*corresponding author, email: villmann@hs-mittweida.de

# 1 Introduction

Functional data frequently occur in many fields of data analysis and processing. These data usually are high-dimensional vectors representing functions with up to hundreds or thousands dimensions [1, 2]. Depending on data, the dimensions are also called (spectral) bands, time or frequency, etc. Examples for such data are time series, spectra, or density distributions to name just a few. Frequently, smoothness of the underlying function is assumed in functional data analysis [3]. The main property distinguishing functional data vectors from other high-dimensional vectors is that the sequence of the vector dimensions carries information and cannot be changed without information loss. Hence, adequate data processing is required paying attention to this special feature.

One successful method for classification of vectorial data is the robust Learning Vector Quantization (LVQ) approach introduced by KOHONEN [4]. The idea behind it is to represent the data classes by typical prototype vectors. This aim is in contrast to the widely applied Support Vector Machines (SVM) [5, 6], which cover the class borders by so-called support vectors while maximizing the separation margin between the classes. Yet, a generalization of the standard LVQ (Generalized LVQ, GLVQ) suggested by SATO&YAMADA [7] has been shown to be a (hypothesis) margin optimizer [8]. Moreover, LVQ and GLVQ show robust behavior also for very high-dimensional data and are therefore suitable for functional data analysis [9].

Frequently, the different data dimensions do not equally contribute to class discrimination. Peaks or valleys in certain ranges of functional vectors may be characterizing class features and, therefore, important for classification. An automatic detection of class distinguishing vector dimensions can be incorporated into GLVQ [10]. This strategy, called *relevance learning*, can be seen as a kind of task specific metric adaptation weighting each data dimension according to its influence for class separation. The vector of all weighting parameters forms the so-called relevance profile. For the weighted Euclidean metric, the resulting Generalized Relevance LVQ (GRLVQ), is still a margin optimizer [11, 12]. In context of high-dimensional data this relevance learning approach leads to a large number of weighting coefficients to be adapted. However, in GRLVQ each data dimension is treated independently such that for functional data the functional information would be ignored.

The idea proposed in this paper is to explicitly take into account the functional property for relevance learning in GRLVQ. In particular, the relevance profile vector is suggested to be a superposition of only a few number of basis functions. This strategy leads to a drastically reduced number of parameters to be adapted. We denote the resulting algorithm as Generalized Functional Relevance LVQ (GFRLVQ).

The outline of the paper is as follows: First, we give a brief review of GRLVQ to clarify notations. Thereafter, we introduce the new GFRLVQ. Further, we investigate sparseness in the relevance profile of the GFRLVQ model to obtain smart models. Thereby, sparseness is considered in two different kinds: structural and feature sparseness emphasizing different aspects of relevance model

reduction. The theoretical part is followed by an experimental section demonstrating the abilities and properties of the new model compared to standard GRLVQ for different data examples.

## 2 Relevance Learning in GLVQ – GRLVQ

To clarify notation, we start with a brief repetition of the GRLVQ stated as the standard LVQ algorithm incorporating relevance learning. It is based on GLVQ which provides a cost function for LVQ.

Generally, given a set  $V \subseteq \mathbb{R}^D$  of data vectors  $\mathbf{v}$  with class labels  $c_{\mathbf{v}} \in \Upsilon = \{1, 2, \dots, C\}$ , the prototypes  $\mathbf{w} \in W \subset \mathbb{R}^D$  with class labels  $y_j$  ( $j = 1, \dots, N$ ) should be distributed in such a way that they represent the data classes as accurately as possible. The following cost function, approximating the classification error, is minimized by GLVQ:

$$E(W) = \frac{1}{2} \sum_{\mathbf{v} \in V} f(\mu(\mathbf{v})) \quad (1)$$

where the function

$$\mu(\mathbf{v}) = \frac{d^+(\mathbf{v}) - d^-(\mathbf{v})}{d^+(\mathbf{v}) + d^-(\mathbf{v})} \quad (2)$$

is the classifier function with  $d^+(\mathbf{v}) = d(\mathbf{v}, \mathbf{w}^+)$  denotes the dissimilarity measure between the data vector  $\mathbf{v}$  and the closest prototype  $\mathbf{w}^+$  with the same class label  $y_{\mathbf{w}^+} = c_{\mathbf{v}}$ , and  $d^-(\mathbf{v}) = d(\mathbf{v}, \mathbf{w}^-)$  is the dissimilarity degree for the best matching prototype  $\mathbf{w}^-$  with a class label  $y_{\mathbf{w}^-}$  different from  $c_{\mathbf{v}}$ . The transformation function  $f$  is a monotonically non-decreasing function usually chosen as sigmoidal or the identity function. A typical sigmoidal choice is the Fermi function

$$f(x) = \frac{1}{1 + a \cdot \exp\left(-\frac{(x-x_0)^2}{2\zeta^2}\right)} \quad (3)$$

with  $x_0 = 0$  and  $a = 1$  as standard parameter values. The  $\zeta$ -parameter allows a control of the sensitivity of the classifier at class borders but is generally fixed with  $\zeta = 1$ .

The dissimilarity measure  $d(\mathbf{v}, \mathbf{w})$  is supposed to be differentiable with respect to the second argument but not necessarily to be a mathematical distance. Usually, the (squared) Euclidean distance is applied. More general dissimilarity measures could also be considered. An important example is the weighted counterpart of the standard Euclidean distance

$$d_{\lambda}^E(\mathbf{v}, \mathbf{w}) = \sum_{i=1}^D \lambda_i (v_i - w_i)^2 \quad (4)$$

with relevance weights  $\lambda_i \geq 0$  and

$$\sum_i \lambda_i = 1 \quad (5)$$

as normalization constraint. The vector  $\lambda = (\lambda_1, \dots, \lambda_D)^T$  is called relevance profile.

Learning of  $\mathbf{w}^+$  and  $\mathbf{w}^-$  in GLVQ is done using the stochastic gradient

$$\frac{\partial_S E(W)}{\partial \mathbf{w}} = \frac{\partial E(W)}{\partial \mathbf{w}} \Bigg|_{\mathbf{v}|P(\mathbf{v})}$$

with respect to the cost function  $E(W)$  for a given data vector  $\mathbf{v}$  according to

$$\frac{\partial_S E(W)}{\partial \mathbf{w}^+} = \xi^+ \cdot \frac{\partial d^+(\mathbf{v})}{\partial \mathbf{w}^+} \quad (6)$$

and

$$\frac{\partial_S E(W)}{\partial \mathbf{w}^-} = \xi^- \cdot \frac{\partial d^-(\mathbf{v})}{\partial \mathbf{w}^-} \quad (7)$$

with

$$\xi^+ = f' \cdot \frac{2 \cdot d^-(\mathbf{v})}{(d^+(\mathbf{v}) + d^-(\mathbf{v}))^2} \quad (8)$$

and

$$\xi^- = -f' \cdot \frac{2 \cdot d^+(\mathbf{v})}{(d^+(\mathbf{v}) + d^-(\mathbf{v}))^2}. \quad (9)$$

In general, a parametrized dissimilarity measure  $d_\lambda(\mathbf{v}, \mathbf{w})$  can be automatically adjusted by adaptation of the parameters, again by gradient descent:

$$\frac{\partial_S E(W)}{\partial \lambda_j} = \xi^+ \cdot \frac{\partial d_\lambda^+(\mathbf{v})}{\partial \lambda_j} + \xi^- \cdot \frac{\partial d_\lambda^-(\mathbf{v})}{\partial \lambda_j}. \quad (10)$$

The respective algorithm is named Generalized Relevance LVQ – GRLVQ [10]. It should be emphasized at this point that in the GRLVQ model the relevance weights as well as the vector components are treated independently as it seems natural in the Euclidean distance or its weighted variant.

### 3 Functional Relevance for GLVQ

As we have seen, the data dimensions are handled in GRLVQ independently according to their given (natural) ordering and, hence, an information loss occurs as a consequence in case of functional data. Further, GRLVQ requires a large number of relevance weights to be adjusted, if the data vectors are really high-dimensional as it is the case in many applications of functional data analysis. For example, processing of hyperspectral data frequently requires the consideration of hundreds or thousands of spectral bands; time series may consist of a huge number of time steps. This huge dimensionality may lead to unstable behavior of relevance learning in GRLVQ. One approach to remedy this instability is proposed by Mendenhall and Merényi suggesting a modified update strategy of GRLVQ [13], without further exploiting the functional aspect.

Therefore, relevance learning should use the functional property of data, if available, to reduce the number of parameters for relevance profile adjustment. For this purpose, we assume in the following that data vectors  $\mathbf{v} = (v_1, \dots, v_D)^T$  are representations of functions  $v(t)$  with given values  $v_i = v(t_i)$ . Analogously we interpret the prototype vectors  $\mathbf{w} = (w_1, \dots, w_D)^T$  functionally as  $w_i = w(t_i)$ . Now, in *functional relevance learning* the relevance profile vector  $\lambda$  is interpreted as a function  $\lambda(t)$  with  $\lambda_j = \lambda(t_j)$ . More precisely, the relevance function  $\lambda(t)$  is supposed to be a superposition

$$\lambda(t) = \sum_{l=1}^K \beta_l \mathcal{K}_l(\omega_l, t) \quad (11)$$

of  $K$  simple basis functions,  $\mathcal{K}_l$  depending on only a few parameters  $\omega_l = (\omega_{l,1}, \dots, \omega_{l,p})^T$  and with the common restriction  $\sum_{l=1}^K \beta_l = 1$ . Doing so, each parameter  $\beta_l$  describes the influence of a certain basis function. The normalization constraint (5) reads now as

$$\int \lambda(t) dt = 1. \quad (12)$$

In consequence, the scaled Euclidean distance (4) now results in

$$d_{\lambda(t)}^E(\mathbf{v}, \mathbf{w}) = \int \lambda(t) (v(t) - w(t))^2 dt \quad (13)$$

which can be also written as

$$d_{\lambda(t)}^E(\mathbf{v}, \mathbf{w}) = \sum_{l=1}^K \beta_l \int \mathcal{K}_l(\omega_l, t) (v(t) - w(t))^2 dt \quad (14)$$

Obviously, if the basis functions  $\mathcal{K}_l(\omega_l, t)$  form an orthogonal basis system in the Hilbert space  $\mathcal{L}_2$  of quadratic integrable functions, an arbitrary approximation precision can be achieved for sufficiently large  $K$ -values [14, 15]. To ensure the normalization constraint (12) the orthogonal functions have to be normalized. However, an arbitrary mixture of any kind of basis functions is possible with, maybe, reduced quality and additional requirements concerning the normalization.

Famous examples of normalized basis functions are standard Gaussians or Lorentzians:

$$\mathcal{K}_l^G(\omega_l, t) = \frac{1}{\sigma_l \sqrt{2\pi}} \exp\left(-\frac{(t - \Theta_l)^2}{2\sigma_l^2}\right) \quad (15)$$

with  $\omega_l = (\sigma_l, \Theta_l)$  and

$$\mathcal{K}_l^L(\omega_l, t) = \frac{1}{\eta_l \pi} \frac{\eta_l^2}{\eta_l^2 + (t - \Theta_l)^2}, \quad (16)$$

with  $\omega_l = (\eta_l, \Theta_l)$ , respectively.

The adaptation of the relevance profile  $\lambda(t)$  is now achieved by a gradient descent of the cost function with respect to the basis functions' parameters  $\omega_l$  as well as the weighting coefficients  $\beta_l$ . In particular, we get for the weighting coefficients of the GFRLVQ cost function

$$E_{GFRLVQ} = \frac{1}{2} \sum_{\mathbf{v} \in V} f(\mu_{\lambda(t)}(\mathbf{v})) \quad (17)$$

with classifier function

$$\mu_{\lambda(t)}(\mathbf{v}) = \frac{d_{\lambda(t)}^+(\mathbf{v}) - d_{\lambda(t)}^-(\mathbf{v})}{d_{\lambda(t)}^+(\mathbf{v}) + d_{\lambda(t)}^-(\mathbf{v})} \quad (18)$$

the stochastic gradient as

$$\frac{\partial_S E_{GFRLVQ}}{\partial \beta_l} = \xi^+ \cdot \frac{\partial d_{\lambda(t)}^+(\mathbf{v})}{\partial \beta_l} + \xi^- \cdot \frac{\partial d_{\lambda(t)}^-(\mathbf{v})}{\partial \beta_l} \quad (19)$$

with  $\xi^+$ ,  $\xi^-$  as in (8) and (9), respectively. In case of the scaled Euclidean dissimilarity (14) this yields

$$\frac{\partial d_{\lambda(t)}^E}{\partial \beta_l} = \int \mathcal{K}_l(\omega_l, t) (v(t) - w(t))^2 dt. \quad (20)$$

For the parameters  $\omega_{l,i}$  the derivatives  $\frac{\partial_S E}{\partial \omega_{l,i}}$  are calculated analogously to (19) with

$$\frac{\partial d_{\lambda(t)}^E}{\partial \omega_{l,i}} = \beta_l \int \frac{\partial \mathcal{K}_l(\omega_l, t)}{\partial \omega_{l,i}} (v(t) - w(t))^2 dt \quad (21)$$

for the scaled Euclidean (14), where we assumed that integration and differentiation can be interchanged. This operation is allowed if the partial derivative  $|\frac{\partial \mathcal{K}_l(\omega_l, t)}{\partial \omega_{l,i}}|$  can be majorized by an integrable function  $\phi$  in the sense of Lebesgue and  $\mathcal{K}_l(\omega_l, t)$  is itself Lebesgue-integrable for each  $\omega_l$  [16]. For the introduced Gaussians  $\mathcal{K}_l^G(\omega_l, t)$  we simply obtain

$$\frac{\partial \mathcal{K}_l^G(\omega_l, t)}{\partial \sigma_l} = \frac{1}{\sigma_l} \left( \frac{(t - \Theta_l)^2}{\sigma_l^2} - 1 \right) \cdot \mathcal{K}_l^G(\omega_l, t) \quad (22)$$

and

$$\frac{\partial \mathcal{K}_l^G(\omega_l, t)}{\partial \Theta_l} = \frac{1}{\sigma_l^2} (t - \Theta_l) \cdot \mathcal{K}_l^G(\omega_l, t) \quad (23)$$

whereas for the Lorentzians  $\mathcal{K}_l^L(\omega_l, t)$  we have

$$\frac{\partial \mathcal{K}_l^L(\omega_l, t)}{\partial \eta_l} = \frac{1}{\eta_l} \frac{(t - \Theta_l)^2 - \eta_l^2}{\eta_l^2 + (t - \Theta_l)^2} \cdot \mathcal{K}_l^L(\omega_l, t) \quad (24)$$

and

$$\frac{\partial \mathcal{K}_l^L(\omega_l, t)}{\partial \Theta_l} = \frac{1}{\eta_l} \frac{2(t - \Theta_l)}{\eta_l^2 + (t - \Theta_l)^2} \cdot \mathcal{K}_l^L(\omega_l, t). \quad (25)$$

The choice of basis function influences the achievable relevance profile. If smooth basis functions are applied, the resulting profile is smooth, too. Sharply peaked or non-differentiable basis functions would lead to roughly structured profiles. However, a mixture of different types is also possible to reflect different aspects. Yet, this may lead to a reduced capability in terms of minimum approximation error compared to orthogonal basis systems for fixed maximum  $K$ -value.

It should be mentioned here that we discussed so far a global metric valid for all prototypes. Obviously, this is not a real restriction: each prototype in GFRLVQ may be equipped with its own local metric, as it is also known from standard GRLVQ.

The learning of the prototypes takes place as in standard GLVQ/GRLVQ according to the formulae (6)-(9) applying the relevance profile  $\lambda(t)$ , i.e. replacing there the distance measure  $d$  by  $d_{\lambda(t)}^E$ . Thus, the learning is structurally kept as vector shifts.

## 4 Sparsity in GFRLVQ

One important question is the optimum number  $K_{opt}$  of basis functions needed for maximum classification performance. According to the paradigm of Occam's razor this can be related to a requirement of sparsity in the relevance model of GFRLVQ. We have to distinguish at least two different kinds of sparsity. The first one is *structural sparsity* emphasizing the sparsity of the generative model of the relevance profile with respect to the selection of number of basis functions. The second one we call *feature sparsity* reflecting the sparsity in terms of data dimensions, which are taken into account for classification, i.e. the number of  $t$ -values for which  $\lambda(t) \approx 0$ . Of course, feature sparsity is also influenced by structural sparsity. However, in feature sparsity the sparseness in data dimensions used for classification is explicitly demanded.

Generally, the sparsity requirement leads to a regularization effect but needs additional terms for learning the parameters of the relevance profile  $\lambda(t)$  as it is explained for both kinds of sparsity in the following subsections.

### 4.1 Structural Sparsity

In the GFRLVQ model the number  $K$  of basis functions to be used can be chosen arbitrarily so far. Obviously, if  $K$  is too small, an appropriate relevance weighting is impossible. Otherwise, a  $K$ -value too large complicates the problem more than necessary. Hence, a good compromise is desirable. This problem can be seen as a structural sparseness requirement in functional relevance learning model. Hence, structural sparsity controls the complexity of the generating model.

A suitable methodology to judge sparsity is based on information theory. In particular, the Shannon entropy

$$H_S = - \sum_{i=1}^K \beta_i \log(\beta_i) \quad (26)$$

of the weighting coefficients  $\beta = (\beta_1, \dots, \beta_K)$  can be applied to quantify structural sparsity. Maximum sparseness, i.e., minimum entropy is obtained, iff  $\beta_l = 1$  for exactly one  $l$  and all the other  $\beta_m$  are equal to zero. However, maximum sparseness may be accompanied by a decrease of accuracy in classification and/or increased cost function value  $E_{GFRLVQ}$ .

To achieve an optimal balance, we propose the following strategy: The cost function  $E_{GFRLVQ}$  from (17) is extended to

$$E_{GFRLVQ}^S = E_{GFRLVQ} + \gamma_S(\tau) \cdot H_S(\beta) \quad (27)$$

with  $\tau$  counting the adaptation steps and  $\gamma_S(\tau) \geq 0$  is the weighting factor which controls the sparsity. Let  $\tau_0$  be the final time step of the GFRLVQ-learning. In the GFRLVQ with sparsity this can be interpreted such that  $\gamma_S(\tau) = 0$  for  $\tau < \tau_0$  holds. Thereafter,  $\gamma_S(\tau)$  is slowly increased in an adiabatic manner [17], such that all parameters can persistently follow the drift of the system. An additional term for  $\beta_l$ -adaptation occurs for non-vanishing  $\gamma_S(\tau)$ -values according to this new cost function (27):

$$\frac{\partial E_{GFRLVQ}^S}{\partial \beta_l} = \frac{\partial E_{GFRLVQ}}{\partial \beta_l} + \gamma_S(\tau) \frac{\partial H_S}{\partial \beta_l} \quad (28)$$

with

$$\frac{\partial H_S}{\partial \beta_l} = -(\log(\beta_l) + 1). \quad (29)$$

This last term causes the vector  $\beta$  to become sparse. The adaptation process is stopped if the  $E_{GFRLVQ}$ -value or the classification error shows a significant increase compared to the time  $\tau_0$ .

## 4.2 Feature Sparsity

A different sparsity requirement concerns the data dimensions  $t$  contributing to the classification decision. In GFRLVQ this feature selection can be controlled by an additional additive entropy term

$$H_F(\lambda) = - \int \lambda(t) \ln(\lambda(t)) dt \quad (30)$$

enforcing the sparsity in the relevance profile  $\lambda(t)$ . This explicite sparsity control of dimensions needed for classification is substantially different from the structural sparsity. Here, the complexity of the model is not optimized but sparseness is judged in terms of non-vanishing parts of the relevance profile  $\lambda(t)$ .

Maximum feature sparseness would lead to so-called line spectra as known from MALDI-TOF spectra [18], for example, with spikes at the centers of the  $K$  basis functions in case of Gaussians or Lorentzians. The resulting cost function is

$$E_{GFRLVQ}^F = E_{GFRLVQ} + \gamma_F(\tau) \cdot H_F(\lambda(t)) \quad (31)$$

and the parameter  $\gamma_F(\tau)$  plays the same role as  $\gamma_S(\tau)$  in the case of structural sparsity.

According to the functional profile model (11) using the basis functions  $\mathcal{K}_l(t)$  we obtain for the derivatives

$$\frac{\partial H_F(\lambda)}{\partial \beta_j} = - \int \mathcal{K}_j(\omega_j, t) [1 + \ln(\lambda(t))] dt \quad (32)$$

and

$$\frac{\partial H_F(\lambda)}{\partial \omega_{j,i}} = - \int \beta_j \frac{\partial \mathcal{K}_j(\omega_j, t)}{\partial \omega_{j,i}} [1 + \ln(\lambda(t))] dt, \quad (33)$$

where we have made the same assumptions about the reversing order of integration and differentiation as before Sec.3.

Obviously, the two sparsity models can be combined. However, in that case, a careful balancing between the two sparsity constraints is important to avoid instabilities and unexpected behavior.

## 5 Experiments and Results

We tested the GFRLVQ for the classification of two well known real world spectral data sets obtained from StatLib and UCI: the *Tecator data set* [19], and the *Wine data set* [20].

The *Tecator data set* consists of 215 spectra measured for several meat probes, see Fig.1. The spectral range is 850nm – 1050nm with 100 spectral bands. The data are labeled according to their fat levels into two classes (low/high). Further, the data are split randomly into 144 training and 71 test spectra.

The *Wine data set* contains 121 absorbing infrared spectra of wine between wavenumbers  $4000\text{cm}^{-1} - 400\text{cm}^{-1}$  (256 bands) split into 91 training and 30 test data, see Fig.2. These data are classified according to their two alcohol levels (low/high) as given in [21].

According to the general shape of the data, we applied GFRLVQ using Gaussian basis functions for the Tecator and Lorentzians for the Wine, because the latter one is more sharply peaked. In the first standard experiments we investigated the performance of the GFRLVQ for different numbers of prototypes and number of basis functions to show the general capability of the algorithm and its basic behavior. In the second step we run simulation to underlay the sparsity methodology.

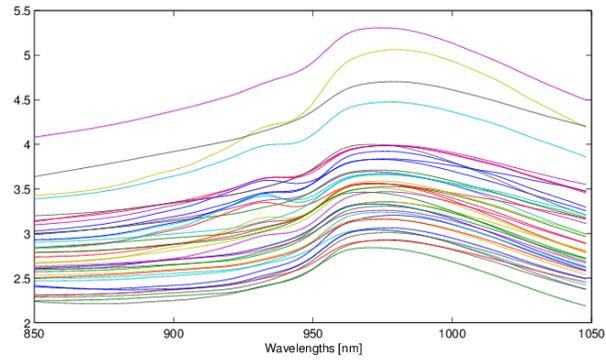


Figure 1: The *Tecator data set*

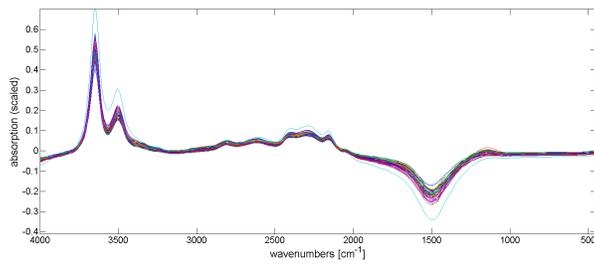


Figure 2: The *Wine data set*

$K \setminus  W $	Tecator			Wine set		
	10	20	40	4	8	12
1	70.8	84.1	90.0	90.1	91.2	93.4
	70.5	70.5	85.3	73.3	80.0	83.3
5	71.1	81.7	90.0	91.2	89.0	93.4
	71.6	75.8	83.2	76.7	73.3	83.3
10	75.0	83.0	90.8	89.0	90.1	93.4
	76.8	80.0	84.2	80.0	80.0	86.7
GRLVQ	71.7	87.5	94.2	93.4	91.2	93.4
	70.5	77.9	83.2	83.3	86.7	80.0

Table 1: Correct classification rate (in %) of the training (1st value) and test (2nd value) sets with different numbers of basis functions  $K$  and prototypes  $|W|$

## 5.1 Standard experiments

In these experiments we varied the number of prototypes as well as the number of basis functions. In particular we used  $K \in \{1, 5, 10\}$  basis functions for the experiments for each data set. Hence, the parameters for relevance learning are drastically reduced to 3, 15, and 30 from 100 and 256 for Tecator and Wine, respectively. The obtained classification accuracies are depicted in the Tab.1. As we can see, these results are comparable to those of standard GRLVQ (see Tab.1) or other approaches, see [21]. However, the results are achieved using a considerably lower number of parameters to be adapted for relevance learning compared to GRLVQ

The resulting relevance profiles for the *Tecator data set* using 20 prototypes are depicted in Fig.3. For the case of  $K = 10$  Gaussian basis functions the respective decomposition is displayed in Fig.4

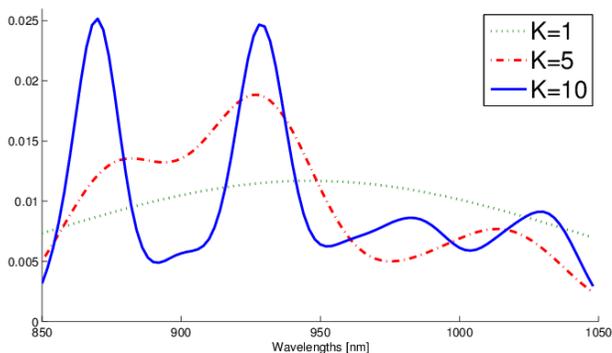


Figure 3: Global relevance profiles obtained using different number  $K$  of Gaussian basis functions for functional relevance learning for the *Tecator data set* using 20 prototypes for learning. The richness in shape increases with  $K$  while keeping the relative smoothness.

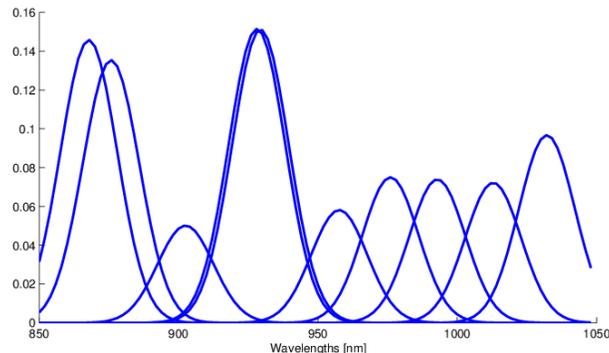


Figure 4: Distribution of the  $K = 10$  adapted Gaussians for the functional relevance profile for the *Tecator data set* using 20 prototypes for learning with global metric adaptation. The resulting relevance profile is shown in Fig.3.

For the *Wine data set* the results are depicted in Fig.5 and Fig.6, accordingly, using 8 prototypes.

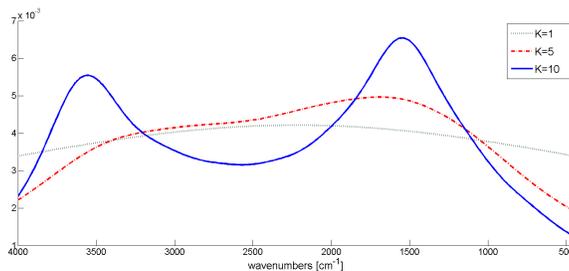


Figure 5: Global relevance profiles obtained using different number  $K$  of Lorentzian basis functions for functional relevance learning for the *Wine data set* using 12 prototypes for learning. As for the *Tecator data set*, the richness in shape increases with  $K$  while keeping the relative smoothness.

As one can expect, the shape of the relevance profile becomes richer with increasing  $K$ -values, while keeping the smooth character according to the smooth basis functions applied. Further, one can observe that heights, widths as well as the centers of the basis functions were properly adapted during learning for both experiments.

## 5.2 Sparsity experiments

According to the introduced kinds of sparsity we performed two experiments for each data set: In the beginning, models with 20 and 12 prototypes were trained using global relevance learning and  $K = 10$  basis functions for each

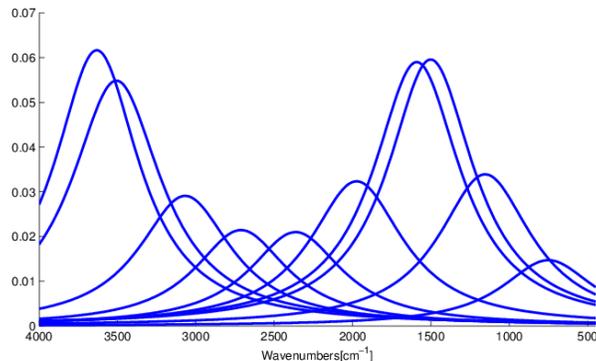


Figure 6: Distribution of the  $K = 10$  adapted Lorentzians for the functional relevance profile for the *Wine data set* using 8 prototypes for learning with global metric adaptation.

data set, respectively. There, Gaussians and Lorentzians were applied as before, accordingly. This training was based on the standard GFRLVQ cost function (17). After this basis training the influence of the sparsity constraint was slowly increased by linear growth of the weighting factor  $\gamma_S(\tau)$  and  $\gamma_F(\tau)$  for structural and feature sparsity, respectively.

In case of structural sparsity judged by the entropy  $H_S$  of the weighting coefficients  $\beta_i$ , this leads to a subsequent fading out of the several basis functions, see Fig.7 for the *Tecator data set*. In the beginning the accuracy level is kept. Above a critical sparsity, a drastic decrease of accuracy can be observed indicating the transition from sparseness optimum to a very low level, see Fig.8. The related relevance profiles at the beginning of the structural sparsity optimization, just before the critical transition and in the final phase are depicted in Fig.9. As one can see, the loss of the relevance peak around wavelength  $870nm$  leads to the breakdown.

For the *Wine data set* the results for structural sparsity optimization are depicted in Figures 10-12, respectively, showing similar behavior as discussed for the *Tecator data set*. However, here the information loss immediately begins with model reduction. This is accompanied by the disappearance of the small broad plateau/peak around wavenumber  $3400cm^{-1}$ .

In the second experiment the feature sparsity was investigated. The experimental setting was as before for structural sparseness optimization. For the *Tecator data set* the results are depicted in Fig.13-15. Again, we observe a drastic loss of accuracy after approximately 800 time cycles corresponding to vanishing  $\sigma_i$ -values for 5 basis functions. The remaining model is too poor to keep the information.

For the *Wine data set* again 10 Lorentzians were taken. As one can observe in Fig.16, the width parameters  $\eta_i$  immediately decreases. After vanishing of 4 basis functions after 4000 learning cycles a destabilization takes places but

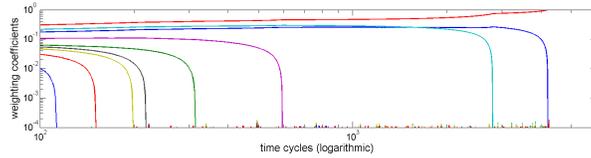


Figure 7: Development in time of weighting coefficients  $\beta_l$  for structural sparsity for the *Tecator data set*. The influence of the sparsity constraint  $H_S$  was linearly increased. Global relevance learning for 20 prototypes was applied with  $K = 10$  Gaussians in the beginning. The weighting coefficients  $\beta_l$  vanish with growing sparsity pressure except only a single remaining for maximum sparseness.

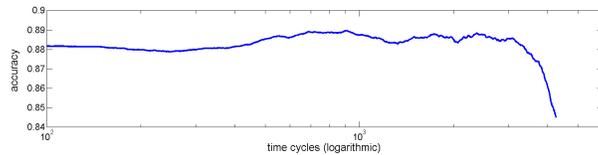


Figure 8: Development in time of the accuracy for structural sparseness optimization for the *Tecator data set*. The accuracy is still high if at least 3 basis functions are weighted to be active. After one more function becomes inactive, the accuracy decreases significantly indicating an substantial information loss.

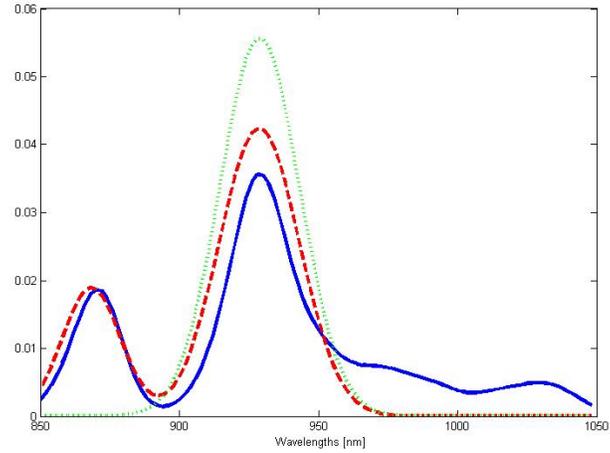


Figure 9: Relevance profiles during structural sparseness optimization process for the *Tecator data set*: at the beginning of the structural sparsity optimization (blue -solid), just before the critical transition (red - dashed) and in the final phase (green - dotted).

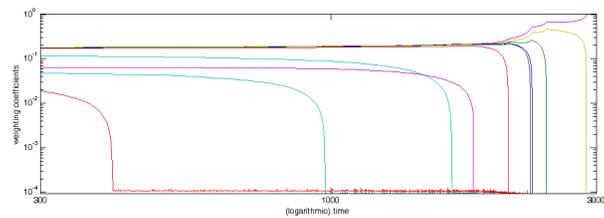


Figure 10: Development in time of weighting coefficients  $\beta_i$  for structural sparsity for the *Wine data set*. The influence of the sparsity constraint  $H_S$  was linearly increased. Global relevance learning for 12 prototypes was applied with  $K = 10$  Lorentzians in the beginning. The weighting coefficients  $\beta_i$  vanish with growing sparsity pressure except only a single remaining for maximum sparseness.

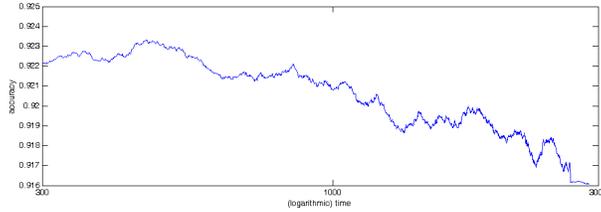


Figure 11: Development in time of the accuracy for structural sparseness optimization for the *Wine data set*. The accuracy slowly decreases after vanishing the first basis function weight. Hence, the  $K = 10$  basis functions are structurally optimal.

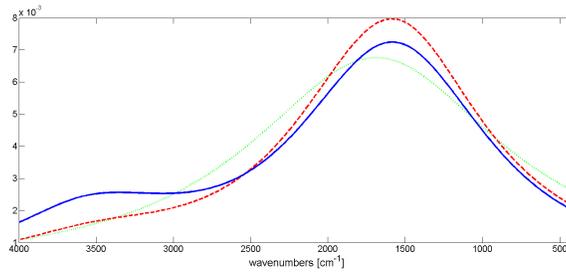


Figure 12: Relevance profiles during structural sparseness optimization process for the *Wine data set*: at the beginning of the structural sparsity optimization (blue -solid), in the middle phase (red - dashed) and in the final phase (green - dotted).

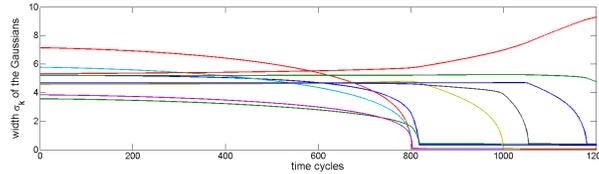


Figure 13: Development in time of the width  $\sigma_l$  of the Gaussians for feature sparsity for the *Tecator data set*. The influence of the sparsity constraint  $H_S$  was linearly increased. Global relevance learning for 20 prototypes was applied with  $K = 10$  basis functions in the beginning. The coefficients vanish with growing sparsity pressure except only a single remaining for maximum sparseness.

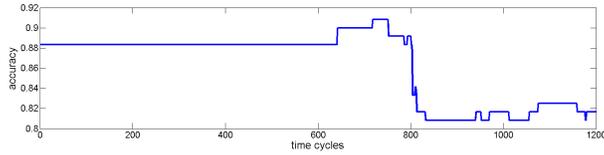


Figure 14: Development in time of the accuracy for feature sparseness optimization for the *Tecator data set*. The accuracy is still high for at least 800 time cycles. After these, the accuracy decreases significantly indicating an substantial information loss while 5 basis functions vanish at this time.

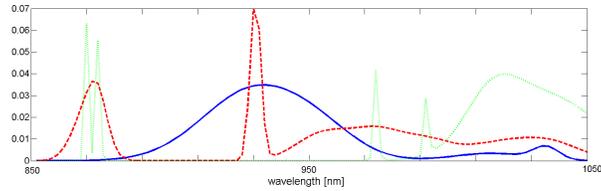


Figure 15: Relevance profiles during the feature sparseness optimization process for the *Tecator data set*: at the beginning of the feature sparsity optimization (blue -solid), just before the critical phase transition (red - dashed) and in the final phase (green - dotted).

still with keeping the high accuracy after a short stabilization phase, see Fig.17. This process is accompanied by substantial thinning out of the relevance profile, see Fig.18. If the pressure of the entropic term is further increased the accuracy significantly drops down and stabilizes at a lower degree leading to a very sparse relevance profile. In consequence, only a few spectral bands are sufficient to distinguish the wine alcoholic classes. This result justifies earlier findings: The influence of the data bands 130 – 230, corresponding to wavenumbers between  $2000\text{cm}^{-1}$  and  $1000\text{cm}^{-1}$ , seems to be class differentiating [21, 22]. However, our simulations show that more or less the observation of single bands are sufficient and, therefore, model complexity could be drastically reduced.

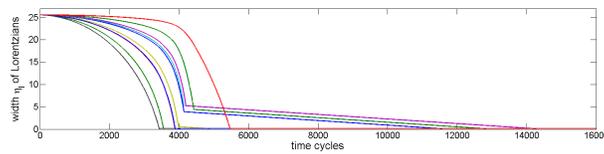


Figure 16: Development in time of the width  $\eta_l$  of the Lorentzians for feature sparsity for the *Wine data set*. The influence of the sparsity constraint  $H_S$  was linearly increased. Global relevance learning for 12 prototypes was applied with  $K = 10$  basis functions in the beginning. The coefficients vanish with growing sparsity pressure leading to a very sparse model.

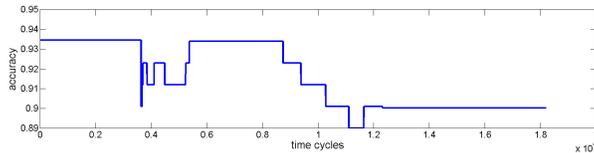


Figure 17: Development in time of the accuracy for feature sparseness optimization for the *Wine data set*. After destabilization the accuracy significantly decreases with increasing sparseness. However, the accuracy remains still high.

## 6 Conclusion

In this paper we propose a *functional* relevance learning for generalized learning vector quantization. Functional learning supposes that the data vectors are

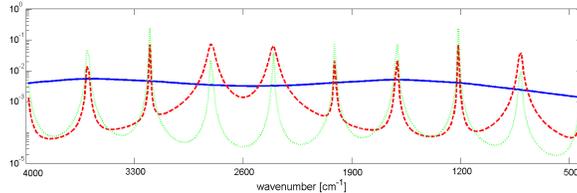


Figure 18: Relevance profiles during the feature sparseness optimization process for the *Wine data set*: at the beginning of the feature sparsity optimization (blue -solid), just before the drop down of hte accuracy after 9000 time cycles (red -dashed) and in the final phase (green - dotted).

representations of functions. In consequence, the relevance profiles are supposed to be functions, too, such that the profiles can be generated by a superposition of basis functions of simple shape depending on only a few parameters. Hence, the number of parameters to be adapted during functional relevance learning is drastically decreased compared to the huge number of relevance weights to be adjusted in standard relevance learning. To obtain an optimal number of basis functions for the superposition sparsity constraints are suggested dealing with different kinds of sparsity - structural and feature sparsity. The sparsity is judged in terms of the entropy of the respective sparsity model: structural sparsity prunes the superposition of the weights of the basis functions used in the model, whereas feature sparsity leads to a reduced number of input dimensions based on width adaptation of the basis functions.

We demonstrated the capabilities of the functional relevance learning algorithm for different data sets and parameter settings achieving good results compared to other models. Further, using sparsity constraints more compact models are obtained. The GFRLVQ approach is here exemplified in the experiments for the weighted Euclidean distance based, for simplicity. Obviously, the Euclidean distance is not based on a functional norm [2, 3, 23]. Yet, the transfer to real functional norms and distances like Sobolev norms [24, 25], the Lee-norm [23, 1], kernel based LVQ-approaches [26] or divergence based similarity measures [27],[28], which carry the functional aspect inherently, is straightforward and topic of future investigations.

Apparently, this functional relevance learning approach can easily extended to matrix learning and limited rank matrix learning, which are proposed in [29, 30]. This would offer a greater flexibility. First ideas are reported in [31] but are still under consideration and, therefore, subject of a forthcoming article.

## References

- [1] J. Lee, M. Verleysen, *Nonlinear Dimensionality Reduction*, Information Sciences and Statistics, Springer Science+Business Media, New York, 2007.
- [2] F. Rossi, N. Delannay, B. Conan-Gueza, M. Verleysen, Representation of functional data in neural networks, *Neurocomputing* 64 (2005) 183–210.
- [3] J. Ramsay, B. Silverman, *Functional Data Analysis*, 2nd Edition, Springer Science+Media, New York, 2006.
- [4] T. Kohonen, *Self-Organizing Maps*, Vol. 30 of Springer Series in Information Sciences, Springer, Berlin, Heidelberg, 1995, (Second Extended Edition 1997).
- [5] B. Schölkopf, A. Smola, *Learning with Kernels*, MIT Press, 2002.
- [6] J. Shawe-Taylor, N. Cristianini, *Kernel Methods for Pattern Analysis and Discovery*, Cambridge University Press, 2004.
- [7] A. Sato, K. Yamada, Generalized learning vector quantization, in: D. S. Touretzky, M. C. Mozer, M. E. Hasselmo (Eds.), *Advances in Neural Information Processing Systems 8. Proceedings of the 1995 Conference*, MIT Press, Cambridge, MA, USA, 1996, pp. 423–9.
- [8] K. Crammer, R. Gilad-Bachrach, A. Navot, A. Tishby, Margin analysis of the LVQ algorithm, in: S. Becker, S. Thrun, K. Obermayer (Eds.), *Advances in Neural Information Processing (Proc. NIPS 2002)*, Vol. 15, MIT Press, Cambridge, MA, 2003, pp. 462–469.
- [9] A. Witoelar, M. Biehl, B. Hammer, Equilibrium properties of offline LVQ, in: M. Verleysen (Ed.), *Proc. Of European Symposium on Artificial Neural Networks (ESANN'2009)*, d-side publications, Evere, Belgium, 2009, pp. 535–540.
- [10] B. Hammer, T. Villmann, Generalized relevance learning vector quantization, *Neural Networks* 15 (8-9) (2002) 1059–1068.
- [11] B. Hammer, M. Strickert, T. Villmann, On the generalization ability of GRLVQ networks, *Neural Processing Letters* 21 (2) (2005) 109–120.
- [12] B. Hammer, M. Strickert, T. Villmann, Relevance LVQ versus SVM, in: L. Rutkowski, J. Siekmann, R. Tadeusiewicz, L. Zadeh (Eds.), *Artificial Intelligence and Soft Computing (ICAISC 2004)*, Lecture Notes in Artificial Intelligence 3070, Springer Verlag, Berlin-Heidelberg, 2004, pp. 592–597.
- [13] M. Mendenhall, E. Merényi, Relevance-based feature extraction for hyperspectral images, *IEEE Transactions on Neural Networks* 19 (4) (2008) 658–672.

- [14] A. Kolmogorov, S. Fomin, Reelle Funktionen und Funktionalanalysis, VEB Deutscher Verlag der Wissenschaften, Berlin, 1975.
- [15] E. Pekalska, R. Duin, The Dissimilarity Representation for Pattern Recognition: Foundations and Applications, World Scientific, 2006.
- [16] I. Kantorowitsch, G. Akilow, Funktionalanalysis in normierten Räumen, 2nd Edition, Akademie-Verlag, Berlin, 1978.
- [17] T. Kato, On the adiabatic theorem of quantum mechanics, Journal of the Physical Society of Japan 5 (6) (1950) 435–439.
- [18] T. Villmann, F.-M. Schleif, M. Kostrzewa, A. Walch, B. Hammer, Classification of mass-spectrometric data in clinical proteomics using learning vector quantization methods, Briefings in Bioinformatics 9 (2) (2008) 129–143.
- [19] TECATOR: Submitted by Hans Henrik Thodberg, Tecator meat sample dataset, available on: <http://lib.stat.cmu.edu/datasets/tecator>.
- [20] Wine data set provided by Prof. Marc Meurens, Available on <http://www.ucl.ac.be/mlg/index.php?page=databases>, [meurens@bnut.ucl.ac.be](mailto:meurens@bnut.ucl.ac.be).
- [21] C. Krier, M. Verleysen, Michel, Rossi, Fabrice, D. François, Supervised variable clustering for classification of nir spectra, in: Proceedings of XVIIth European Symposium on Artificial Neural Networks (ESANN 2009), Bruges, Belgique, 2009, pp. 263–268.
- [22] C. Krier, F. Rossi, D. François, M. Verleysen, A data-driven functional projection approach for the selection of feature ranges in spectra with ICA or cluster analysis, Chemometrics and Intelligent Laboratory Systems 91 (2008) 43–53.
- [23] J. Lee, M. Verleysen, Generalization of the  $l_p$  norm for time series and its application to self-organizing maps, in: M. Cottrell (Ed.), Proc. of Workshop on Self-Organizing Maps (WSOM) 2005, Paris, Sorbonne, 2005, pp. 733–740.
- [24] B. Silverman, Smoothed functional principal components analysis by the choice of norm, The Annals of Statistics 24 (1) (1996) 1–24.
- [25] T. Villmann, F.-M. Schleif, Functional vector quantization by neural maps, in: J. Chanussot (Ed.), Proceedings of First Workshop on Hyperspectral Image and Signal Processing: Evolution in Remote Sensing (WHISPERS 2009), IEEE Press, 2009, pp. 1–4, ISBN 978-1-4244-4948-4.
- [26] T. Villmann, B. Hammer, Theoretical aspects of kernel GLVQ with differentiable kernel, IfI Technical Report Series (IfI-09-12) (2009) 133–141.

- [27] T. Villmann, S. Haase, Divergence based vector quantization, *Neural Computation* 23 (5) (2011) 1343–1392.
- [28] E. Mwebaze, P. Schneider, F.-M. Schleif, J. Aduwo, J. Quinn, S. Haase, T. Villmann, M. Biehl, Divergence based classification in learning vector quantization, *Neurocomputing* 74 (9) (2011) 1429–1435.
- [29] P. Schneider, B. Hammer, M. Biehl, Adaptive relevance matrices in learning vector quantization, *Neural Computation* 21 (2009) 3532–3561.
- [30] K. Bunte, B. Hammer, A. Wismüller, M. Biehl, Adaptive local dissimilarity measures for discriminative dimension reduction of labeled data, *Neurocomputing* 73 (2010) 1074–1092.
- [31] M. Kästner, T. Villmann, M. Biehl, About sparsity in functional relevance learning in generalized learning vector quantization, *Machine Learning Reports* 5 (MLR-03-2011) (2011) 1–12, ISSN:1865-3960, [http://www.techfak.uni-bielefeld.de/~fschleif/mlr/mlr\\_03\\_2011.pdf](http://www.techfak.uni-bielefeld.de/~fschleif/mlr/mlr_03_2011.pdf).