

# Ordinal regression based on learning vector quantization



Fengzhen Tang<sup>a,b,\*</sup>, Peter Tiño<sup>b</sup>

<sup>a</sup> Shenyang Institute of Automation, Chinese Academy of Sciences, No.114, Nanta Street, Shenyang, Liaoning Province, 110016, China

<sup>b</sup> School of Computer Science, The University of Birmingham, Edgbaston, Birmingham B15 2TT, UK

## ARTICLE INFO

### Article history:

Received 19 November 2015

Received in revised form 27 April 2017

Accepted 4 May 2017

Available online 9 May 2017

### Keywords:

Ordinal regression

Learning vector quantization

Generalized matrix learning vector quantization

## ABSTRACT

Recently, ordinal regression, which predicts categories of ordinal scale, has received considerable attention. In this paper, we propose a new approach to solve ordinal regression problems within the learning vector quantization framework. It extends the previous approach termed ordinal generalized matrix learning vector quantization with a more suitable and natural cost function, leading to more intuitive parameter update rules. Moreover, in our approach the bandwidth of the prototype weights is automatically adapted. Empirical investigation on a number of datasets reveals that overall the proposed approach tends to have superior out-of-sample performance, when compared to alternative ordinal regression methods.

© 2017 Elsevier Ltd. All rights reserved.

## 1. Introduction

Recently a new learning setting that predicts categories of ordinal scale, referred to as ordinal regression or ranking learning, has received considerable attention (Fouad & Tiño, 2012; Gutiérrez & Tiño, 2014; Lin & Li, 2012; Sánchez-Monedero, Gutiérrez, Tiño, & Hervás-Martínez, 2013; Seah, Tsang, & Ong, 2012). In this setting, the training examples are labeled by categories (ranks) exhibiting a natural order. Consequently, ordinal regression bears resemblance to both regression and classification. However, in contrast to regression, the ranks are of discrete and finite type, and the difference between the consecutive ranks is not quantified. The existence of order information in the class labels also makes ordinal regression different from classification. Ordinal regression problems arise commonly in real world applications, for example, financial movements grading (Tino, Schittenkopf, & Dorffner, 2001), credit rating (Dikkers & Rothkrantz, 2012), life quality assessments (Lall, Campbell, Walters, & Morgan, 2012), facial beauty assessment (Yan, 2014), etc.

A variety of algorithms have been developed for ordinal regression problems in the literatures. The simplest idea, presented in Kramer, Widmer, Pfahringer, and de Groeve (2001), is to map the ordinal scales into numeric values and then solve the

problem as a regular regression problem (e.g. by a regression tree). However, this may be suboptimal since the true distances between ordinal scales are unknown in most cases. Another idea is to transform the ordinal regression problem into a set of binary classification problems. For instance, Frank and Hall (2001) transforms an  $C$ -class ordinal problem into  $C - 1$  binary classification problems and then combines the  $C - 1$  model predictions to estimate the probabilities of the  $C$  original ordinal classes for test instances. This method requires no modification of the standard binary classification algorithm. However, it cannot capture the overall structure among the ordinal classes. It is also possible to formulate the ordinal regression problem as a large augmented binary classification problem by adding additional constraints. Herbrich, Graepel, and Obermayer (2000) applied the principle of structural risk minimization to ordinal regression, leading to a new distribution-independent learning algorithm based on a loss function between pairs of ranks. Crammer and Singer (2001) generalized the online perception algorithm with multiple thresholds to seek the direction and thresholds for ordinal regression. Shashua and Levin (2003) generalized the support vector machine formulation for ordinal regression by maximizing the margin of  $C - 1$  parallel hyperplanes separating adjacent classes in the feature space. The parallel separation hyperplanes share the same normal vector and are defined by a set of  $C - 1$  thresholds  $b_1, \dots, b_{C-1}$ . The margin can either be the same for all separating hyperplanes, or can be different for each adjacent class separation, in which case the sum of the margins is maximized. However, in this approach it cannot be guaranteed that the optimized thresholds will preserve the category order. To address this problem, Chu and Keerthi (2005, 2007) proposed explicit and

\* Corresponding author at: Shenyang Institute of Automation, Chinese Academy of Sciences, No.114, Nanta Street, Shenyang, Liaoning Province, 110016, China.

E-mail addresses: [tangfengzhen87@hotmail.com](mailto:tangfengzhen87@hotmail.com) (F. Tang),

[P.Tino@cs.bham.ac.uk](mailto:P.Tino@cs.bham.ac.uk) (P. Tiño).

implicit constraints enforcing the inequalities on the thresholds, i.e.  $b_1 < b_2 < \dots < b_{C-1}$ . The explicit constraints directly enforce order on the adjacent thresholds, while the implicit constraints ensure the threshold order implicitly through the data by stipulating that the  $j$ th hyperplane (corresponding to threshold  $b_j$ ) separates all points from classes  $\leq j$  from all points of classes  $> j$ . Li and Lin (2006) introduced a unified reduction framework from ordinal regression to binary classification. The explicit and implicit approaches of Chu and Keerthi (2007) can be regarded as special instances of this framework. Chu and Ghahramani (2005) proposed a simple non-parametric Bayesian approach to ordinal regression based on a generalization of the probit likelihood function for Gaussian process, achieving competitive (often better) generalization performance to the support vector approaches. Sun, Li, Wu, Zhang, and Li (2010) augmented kernel discriminate learning for ordinal regression to take class distribution into consideration.

Recently, Fouad and Tiño (2012) extended generalized matrix learning vector quantization approach (Schneider, Biehl, & Hammer, 2009) to ordinal regression, where the order information among different categories is utilized in the selection of the prototypes to be adapted, as well as in updating of the selected prototypes. Given an (*input, target class*) training example, generalized matrix learning vector quantization identifies two prototypes – the closest prototype of the target class and the closest prototype among the prototypes with a different label. The positions of this prototype pair are updated according to a well-defined cost function. In the extended generalized matrix learning vector quantization approach to ordinal regression, spatially close prototypes to the training input from neighboring classes (in terms of the class order) of the target class are selected as “correct” prototypes, while prototypes from classes far away in the class order from the target class, but lying in the neighborhood of the training input are selected as “incorrect” prototypes. Correct and incorrect prototypes are then paired into a set of (*correct prototype, incorrect prototype*) couples. Updating rules for each such prototype couple are formulated in a principled manner through a weighted cost function. The weights for correct and incorrect prototypes will decrease and increase, respectively, with growing class difference from the target class. In addition, for incorrect prototypes, the weights will diminish with increasing distance from the input data.

However, the cost function of Fouad and Tiño (2012) based on the pairing strategy does not wholistically take into account the global order on classes. While prototype pairing, given an input vector, of correct and incorrect prototypes naturally extends the generalized matrix learning vector quantization approach of Schneider et al. (2009), the price to be paid is the need to brake the global linear class order into ordered pairs that are treated independently of each other. This is unnatural and the updating rules of the prototypes derived from this cost function cannot consistently guarantee the ordering relation among the prototypes.

In this paper, we propose a new cost function that is not only more natural and intuitive, but crucially yields prototype updates that explicitly express the global prototype ordering. This is in contrast to the approach of Fouad and Tiño (2012) that employs partial pair-wise comparisons reflecting the global prototype order only implicitly. We also show that the derivatives of our cost function vanish at boundaries of the prototype receptive fields and thus gradient based optimization will constitute a valid gradient descent. Furthermore, we propose to automatically adapt parameters in the weighting functions employed in the cost functional, eliminating the need for costly grid search of Fouad and Tiño (2012).

Using both artificial and real ordinal regression datasets, we demonstrate that the generalization performance of the proposed approach is competitive and often better than the approach presented in Fouad and Tiño (2012). We will refer to the original

approach of Fouad and Tiño (2012) and the proposed approach as pair ordinal generalized matrix learning vector quantization ( $p$ -OGMLVQ) and accumulative ordinal generalized matrix learning vector quantization ( $a$ -OGMLVQ), respectively.

The rest of the paper is organized as follows. In Section 2, we briefly introduce learning vector quantization and its extension to ordinal regression. In Section 3, we present the proposed approach along with automatic hyperparameter updates in the weighting function. Experimental results are given in Section 4. Main findings and conclusions are presented in Section 5.

## 2. Background

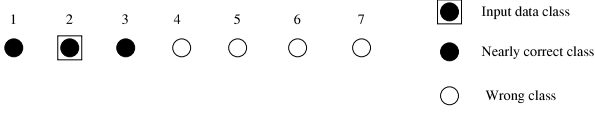
Our approach is developed within the framework of learning vector quantization (LVQ) (Kohonen, 1986, 1995). LVQ is a prototype-based supervised classification algorithm where the classifier is parametrized by a set of labeled prototypes which live in the same space as the input data. The approach has enjoyed popularity because of its simplicity, intuitive nature, and natural accommodation of multi-class classification problems. In this section, we will briefly introduce LVQ and its extension to ordinal regression.

Consider a training dataset  $(\mathbf{x}_i, y_i) \in \mathbb{R}^m \times \{1, \dots, C\}$ ,  $i = 1, \dots, n$ , where  $m$  is the dimension of the inputs,  $C$  is the number of different classes and  $n$  is the number of training examples. A typical LVQ classifier consists  $M$  ( $M \geq C$ ) prototypes  $\mathbf{w}_i \in \mathbb{R}^m$ , which are labeled by  $c(\mathbf{w}_i) \in \{1, \dots, C\}$ . The classification is implemented as a winner-takes-all scheme. For a data point  $\mathbf{x} \in \mathbb{R}^m$ , the output class is determined by the class label of its closest prototype: i.e.  $\hat{y}(\mathbf{x}) := c(\mathbf{w}_i)$  such that  $i = \arg \min_j d(\mathbf{x}, \mathbf{w}_j)$ , where  $d(\cdot, \cdot)$  is a distance measure in  $\mathbb{R}^m$ . Each labeled prototype  $\mathbf{w}_i$  with label  $c(\mathbf{w}_i)$  defines a receptive field in the input space – a set of points which pick this prototype as their winner. The goal of learning the LVQ classifier is to adapt prototypes automatically such that the class labels of data points within the receptive field coincide with the label of the respective prototype as much as possible.

A generalization of LVQ, termed generalized learning vector quantization (GLVQ), was introduced in Sato and Yamada (1996). In GLVQ the prototypes are updated based on the steepest descent method on a well defined cost function. The cost function is determined so that the obtained learning rule satisfies the convergence condition. In the training phase of GLVQ, for each labeled input  $\mathbf{x}_i$ , a pair of prototypes will be updated. The closest prototype to  $\mathbf{x}_i$  with the same label  $y_i$  (correct prototype) is rewarded by dragging it closer to  $\mathbf{x}_i$ , while the closest prototype with a different label (incorrect prototype) is penalized by pushing it away from  $\mathbf{x}_i$ .

GLVQ has been extended towards metric learning, where a general adaptive full metric tensor of the distance measure is learned. In this way important dependencies between different input features as well as their individual importance for the classification task can be revealed and utilized (Schneider et al., 2009). In this approach, termed generalized matrix learning vector quantization (GMLVQ), again at each training step a *single* pair of prototypes is updated (the closest correct and the closest incorrect prototype to the current training input). Unlike in GLVQ, the distance function is also updated so that the distance between the example  $\mathbf{x}_i$  and its closest correct prototype is shortened, while the distance between the example  $\mathbf{x}_i$  and its closest incorrect prototype is enlarged.

The GMLVQ approach has been further extended to ordinal regression (ordinal GMLVQ–OGMLVQ, which is referred to as pair-OGMLVQ ( $p$ -OGMLVQ) in this paper in order to ease the presentation) in Fouad and Tiño (2012). In the  $p$ -OGMLVQ approach, at each training step *several* pairs of “correct” and “incorrect” prototypes are updated. First, correct and incorrect



**Fig. 1.** Illustration of ordinal learning vector quantization (Fouad & Tiño, 2012), with threshold  $\tau = 1$ . The class  $y_i$  of the training input  $\mathbf{x}_i$  is 2, indicated with a square. Prototype classes indicated by black circles (i.e. 1, 2, and 3) will be regarded as correct ones, while other prototype classes denoted by white circles (i.e. 4–7) will be seen as incorrect ones.

prototype classes are identified according to the label order. If the prototype class is “close” to the class  $y_i$  of the training input  $\mathbf{x}_i$  (label distance smaller than or equal to a threshold  $\tau$ ), the prototype class will be viewed as “tolerably correct”, otherwise it will be considered as incorrect (see Fig. 1). The correct prototypes are the spatially closest prototypes to  $\mathbf{x}_i$  (one for each class) from the correct classes. The incorrect prototypes are the prototypes with incorrect class lying in the neighborhood of  $\mathbf{x}_i$ .

As in GMLVQ, the correct prototypes will be dragged towards  $\mathbf{x}_i$ , while the incorrect prototypes will be pushed away from  $\mathbf{x}_i$ . However, they will be not all dragged or pushed away to the same extent – a weighting scheme reflecting how close the class of the prototype is to the class of  $\mathbf{x}_i$  is introduced. The weights for the correct prototypes are given as follows:

$$\alpha_j^+ = \exp \left\{ -\frac{\kappa(y_i, c(\mathbf{w}_j^+))^2}{2\sigma^2} \right\}, \quad (1)$$

where  $\mathbf{w}_j^+$  denotes the  $j$ th correct prototype and  $\kappa(\cdot, \cdot)$  is the loss function on the class labels. Commonly absolute error loss is adopted, i.e.  $\kappa(y_i, c(\mathbf{w}_j)) = |y_i - c(\mathbf{w}_j)|$ .  $\sigma$  is the bandwidth parameter of the Gaussian weighting function. The weight for the  $j$ th incorrect prototype  $\mathbf{w}_j^-$  is defined as follows:

$$\alpha_j^- = \exp \left\{ -\frac{(T - \kappa(y_i, c(\mathbf{w}_j^-)))^2}{2\sigma^2} \right\} \cdot \exp \left\{ -\frac{d^A(\mathbf{x}_i, \mathbf{w}_j^-)}{2\sigma'^2} \right\}, \quad (2)$$

where  $T = \max_{\mathbf{u} \in W^-} \kappa(y_i, c(\mathbf{u}))$  and  $W^-$  is the set of incorrect prototypes of data point  $\mathbf{x}_i$ . In other words,  $T$  is the maximum class difference between the data point and the identified incorrect prototypes.  $d^A(\cdot, \cdot)$  is the squared distance between data point and prototype,<sup>1</sup> defined as  $d^A(\mathbf{x}, \mathbf{w}) = (\mathbf{x} - \mathbf{w})^T \cdot \Lambda \cdot (\mathbf{x} - \mathbf{w})$ , where  $\Lambda$  is an  $m \times m$  symmetric positive semi-definite matrix (metric tensor). The positive semi-definiteness of  $\Lambda$  can be enforced by substituting  $\Lambda = \Omega^T \Omega$ , where  $\Omega$  is an unconstrained  $m \times m$  matrix.  $\Lambda$  needs to be normalized since the (global) scale is irrelevant in GMLVQ.  $\sigma'$  is the bandwidth parameter of the Gaussian weighting function corresponding to distance between data point and prototype. These hyper-parameters can be tuned via cross-validation.

Motivated by GMLVQ, the cost function of  $p$ -OGMLVQ is defined through pairing correct and incorrect prototypes into the corresponding couples. Suppose there are  $r$  couples, for each couple  $j$ , a cost function is defined as follows:

$$f_j(\mathbf{x}_i) = \frac{\alpha_j^+ \cdot d^A(\mathbf{x}_i, \mathbf{w}_j^+) - \alpha_j^- \cdot d^A(\mathbf{x}_i, \mathbf{w}_j^-)}{\alpha_j^+ \cdot d^A(\mathbf{x}_i, \mathbf{w}_j^+) + \alpha_j^- \cdot d^A(\mathbf{x}_i, \mathbf{w}_j^-)}. \quad (3)$$

Consequently, the overall cost function is given by

$$C_{p\text{-OGMLVQ}} = \sum_{i=1}^n \sum_{j=1}^r \Phi(f_j(\mathbf{x}_i)), \quad (4)$$

<sup>1</sup> In this weight definition, the distance  $d^A(\mathbf{x}_i, \mathbf{w}_j^-)$  is involved in order to take the following aspect into consideration: if incorrect prototype is already far away from the input point, then we do not need a large weight to push it away that hard.

where  $\Phi(\cdot)$  is a monotonically increasing function. The core of the formulation in Eq. (3) has been well studied both theoretically and empirically in the past as a large margin inducing cost function for prototype based classifiers (e.g. Sato & Yamada, 1996). Moreover, dimensionality independent generalization bound can be derived by the LVQ cost function. Monotonically increasing function  $\Phi$  plays a modulating role for the margin quantification and can be simply the identity  $\Phi(x) = x$ , or the logistic function  $\Phi(x) = 1/(1 + e^{-x})$  (if the large values need to be squashed – e.g. to dampen the influence of class-conditional outliers).

The updating rules of the prototypes and the metric can be derived by minimizing the cost function:

$$\mathbf{w}_j^+ := \mathbf{w}_j^+ + \epsilon \cdot 2 \cdot \mu_j^+ \cdot \Lambda \cdot (\mathbf{x}_i - \mathbf{w}_j^+), \quad (5)$$

$$\mathbf{w}_j^- := \mathbf{w}_j^- - \epsilon \cdot 2 \cdot \mu_j^- \cdot \Lambda \cdot (\mathbf{x}_i - \mathbf{w}_j^-), \quad (6)$$

$$\Omega := \Omega - \eta \cdot \left\{ 2 \cdot \mu_j^+ \cdot \Omega \cdot (\mathbf{x}_i - \mathbf{w}_j^+) (\mathbf{x}_i - \mathbf{w}_j^+)^T - 2 \cdot \mu_j^- \cdot \Omega \cdot (\mathbf{x}_i - \mathbf{w}_j^-) (\mathbf{x}_i - \mathbf{w}_j^-)^T \right\}, \quad (7)$$

where

$$\mu_j^+ = \frac{2\alpha_j^+ \cdot \alpha_j^- \cdot d^A(\mathbf{x}_i, \mathbf{w}_j^-)}{(\alpha_j^+ \cdot d^A(\mathbf{x}_i, \mathbf{w}_j^+) + \alpha_j^- \cdot d^A(\mathbf{x}_i, \mathbf{w}_j^-))^2}, \quad (8)$$

$$\mu_j^- = \frac{2\alpha_j^+ \cdot \alpha_j^- \cdot d^A(\mathbf{x}_i, \mathbf{w}_j^+)}{(\alpha_j^+ \cdot d^A(\mathbf{x}_i, \mathbf{w}_j^+) + \alpha_j^- \cdot d^A(\mathbf{x}_i, \mathbf{w}_j^-))^2}, \quad (9)$$

$0 < \epsilon < 1$  and  $0 < \eta < 1$  are learning rates for prototypes and metric, respectively. The learning rates are monotonically decreasing with time (Fouad & Tiño, 2012; Schneider et al., 2009):

$$\epsilon(t) = \frac{\epsilon(0)}{1 + \nu(t-1)}, \quad (10)$$

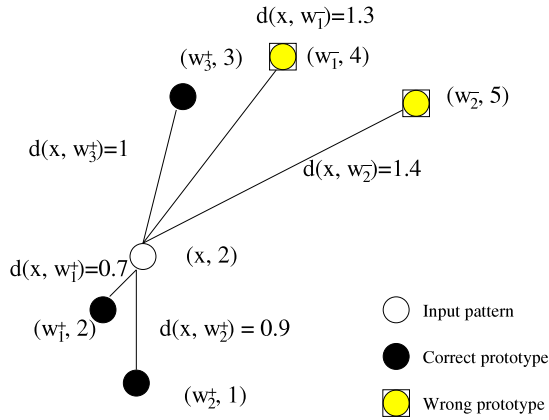
$$\eta(t) = \frac{\eta(0)}{1 + \nu(t-1)}, \quad (11)$$

where  $\nu > 0$  is a parameter determining the decay speed,  $t$  denotes the number of training epochs. Commonly the magnitude of  $\eta$  is set smaller than that of  $\epsilon$ , in order to achieve a slower time-scale of metric learning compared to the weight updates (Schneider et al., 2009).

### 3. Proposed algorithm

In the  $p$ -OGMLVQ approach,  $\mu_j^+$  and  $\mu_j^-$  control the magnitude of “forces” when updating the identified prototypes. Naturally  $\mu_j^+$  should be monotonically decreasing as a function of label distance between the prototype  $\mathbf{w}_j^+$  and training input  $\mathbf{x}_i$ . In other words, if  $\mathbf{w}_j^+$  is “less correct” than  $\mathbf{w}_k^+$ , i.e.  $\kappa(y_i, c(\mathbf{w}_j^+)) > \kappa(y_i, c(\mathbf{w}_k^+))$ , then  $\mathbf{w}_j^+$  should be dragged towards  $\mathbf{x}_i$  to a lesser extent than  $\mathbf{w}_k^+$ , i.e.  $\mu_j^+ < \mu_k^+$ . However, this property might be violated as illustrated in Fig. 2. The input pattern ( $\mathbf{x}$ , 2) represents a training input  $\mathbf{x}$  with class label 2. There are several labeled prototypes ( $\mathbf{w}_j^+, c(\mathbf{w}_j^+)$ ) in the neighborhood of  $\mathbf{x}$ . Assume  $\sigma = \sigma' = 0.5$ . Based on the prototype pairing and prototype distances  $d(\mathbf{x}, \mathbf{w}_j^+)$  from  $\mathbf{x}$  shown in the figure, we get  $\mu_1^+ = 0.0514$  and  $\mu_2^+ = 0.5381$ . Obviously,  $\mathbf{w}_2^+$  with label 1 is “less correct” than  $\mathbf{w}_1^+$  with label 2 matching the label of  $\mathbf{x}$ . However, counterintuitively, we have  $\mu_2^+ > \mu_1^+$  and the less correct prototype  $\mathbf{w}_2^+$  will be updated more than the correct one  $\mathbf{w}_1^+$ . Similar problems may occur for  $\mu^-$ .

To solve this problem mentioned above, we propose a new cost function such that the update rules derived from it will explicitly guarantee that the attractive force decreases as the class difference between the correct prototype and the input increases. It will also be the case that the repulsive force increases



**Fig. 2.** Illustration of the problem of breaking the global class linear order into a set of prototype pairs with correct and incorrect labels, as done in Fouad and Tiño (2012). The input pattern  $(\mathbf{x}, 2)$  represents a training input  $\mathbf{x}$  with class label  $y(\mathbf{x}) = 2$ . There are several labeled prototypes  $(\mathbf{w}_j^\pm, c(\mathbf{w}_j^\pm))$  in the neighborhood of  $\mathbf{x}$ . Assume  $\sigma = \sigma' = 0.5$ . Based on the prototype pairing and prototype distances  $d(\mathbf{x}, \mathbf{w}_j^\pm)$  from  $\mathbf{x}$  shown in the figure, we get  $\mu_1^+ = 0.0514$  and  $\mu_2^+ = 0.5381$ . Obviously,  $\mathbf{w}_2^+$  with label 1 is “less correct” than  $\mathbf{w}_1^+$  with label 2 matching the label of  $\mathbf{x}$ . However, counterintuitively, we have  $\mu_2^+ > \mu_1^+$  and the less correct prototype  $\mathbf{w}_2^+$  will be updated more than the correct one  $\mathbf{w}_1^+$ .

as the class difference between the incorrect prototype and the input increases. Consequently, the ordering of prototype classes is ensured. Since incorrect prototypes far from the input are of lesser importance, the repulsive force will decrease with increasing prototype distance from  $\mathbf{x}_i$ .

Equally importantly, our approach is derived in a principled manner from a single unified cost function, in contrast to somewhat ad-hoc unnatural decoupling of the  $p$ -OGMLVQ parameter updates into separate updates of binary classifiers.

As in Section 2, we assume a training dataset  $(\mathbf{x}_i, y_i) \in \mathbb{R}^m \times \{1, \dots, C\}$ ,  $i = 1, \dots, n$ , is given, where  $m$  is the dimension of the inputs,  $C$  is the number of classes and  $n$  is the number of training examples. As we are dealing with ordinal regression problems, we add an additional assumption that classes are ordered as  $1 < 2, \dots, < C$ , where  $<$  denotes the order relation on labels. As in LVQ network, the proposed classifier is defined by  $M$  labeled prototypes  $(\mathbf{w}_i, c(\mathbf{w}_i))$ , where  $\mathbf{w}_i \in \mathbb{R}^m$  and  $c(\mathbf{w}_i) \in \{1, \dots, C\}$ . The classification scheme is implemented by means of the winner-takes-all rule, where the output class for a given input  $\mathbf{x} \in \mathbb{R}^m$  is determined by its closest prototype. Formally, the output class is determined by  $\hat{y}(\mathbf{x}) := c(\mathbf{w}_j)$  such that  $j = \arg \min_i d^A(\mathbf{x}, \mathbf{w}_i)$ , where  $d^A(\mathbf{x}, \mathbf{w}_i)$  is the squared distance between input  $\mathbf{x}$  and the prototype  $\mathbf{w}_i$ , as defined in Section 2.

As in Fouad and Tiño (2012), our algorithm aims to minimize the average absolute error of class mislabeling. This is reflected by the utilization of absolute label distance  $\kappa(y_i, c(\mathbf{w})) = |y_i - c(\mathbf{w})|$ . Following Fouad and Tiño (2012), the ordering information in class label is utilized to identify the sets of correct and incorrect prototypes. The classes of prototypes within a neighborhood (in the label space) of the given input class will be regarded as “correct”. Formally, the set of correct classes is defined as follows:

$$C^+ = \{c(\mathbf{w}_j) \mid \kappa(y_i, c(\mathbf{w}_j)) \leq \tau\}, \quad (12)$$

where  $\tau > 0$  is a rank loss threshold. The classes of prototypes which are outside the neighborhood of the given input class will be treated as “incorrect” ones. The set of incorrect classes is given by:

$$C^- = \{c(\mathbf{w}_j) \mid \kappa(y_i, c(\mathbf{w}_j)) > \tau\}. \quad (13)$$

The spatially closest prototype in each correct class to the input  $\mathbf{x}_i$  will be treated as a correct prototype. The set of correct prototypes

is then

$$W^+ = \{\mathbf{w}_j \mid j = \arg \min_{c(\mathbf{w}_i)=k} d^A(\mathbf{x}_i, \mathbf{w}_i), k \in C^+\}. \quad (14)$$

Prototypes from the identified incorrect classes which lie in the neighborhood of  $\mathbf{x}_i$  will be treated as incorrect ones that need to be pushed away from  $\mathbf{x}_i$ . The set of incorrect prototypes is given as follows:

$$W^- = \{\mathbf{w}_j \mid d^A(\mathbf{x}_i, \mathbf{w}_j) < D, c(\mathbf{w}_j) \in C^-\}, \quad (15)$$

where  $D > 0$ . We chose  $D$  to be the median distance of the example  $\mathbf{x}$  to all prototypes from the incorrect classes.<sup>2</sup>

We rearrange the set  $W^+$  and  $W^-$  in increasing order in terms of distance to the input  $\mathbf{x}_i$ . Let  $r = \min(|W^+|, |W^-|)$ , where  $|W|$  is the cardinality of the set  $W$ . The first  $r$  prototypes of each prototype set will be selected for updating. However, instead of updating the prototypes in pairs as in Fouad and Tiño (2012), we propose to update all the selected prototypes collectively with a weighting scheme defined by the ordering information of the labels.

Weights for correct prototypes with respect to  $c(\mathbf{x}_i)$  are given by a Gaussian kernel:

$$\alpha_j^+ = \exp \left\{ -\frac{\kappa(y_i, c(\mathbf{w}_j^+))^2}{2\sigma_1^2} \right\}. \quad (16)$$

Label of a correct prototype  $\mathbf{w}_j^+$  further away from the label of the input pattern will result in a smaller weight  $\alpha_j^+$ . Weights for incorrect prototypes with respect to  $y_i$  are determined as:

$$\alpha_j^- = \exp \left\{ -\frac{T - \kappa(y_i, c(\mathbf{w}_j^-))^2}{2\sigma_2^2} \right\} \cdot \exp \left\{ -\frac{(d^A(\mathbf{x}_i, \mathbf{w}_j^-))^2}{2\sigma_3^2} \right\}, \quad (17)$$

where  $T$  is the maximum label difference between input pattern and wrong prototypes, i.e.  $T = \max_{\mathbf{w}_j \in W^-} \kappa(y_i, c(\mathbf{w}_j))$ . Wrong prototypes spatially closer to the input pattern and with labels further away from that of the input pattern will have larger weight  $\alpha_j^-$ . We share the formulation of the weight functions with Fouad and Tiño (2012). However, in contrast to Fouad and Tiño (2012), the scale parameters  $\sigma_1, \sigma_2$  and  $\sigma_3$  are learned as inherent part of the parameter updates, as explained in Section 3.2.<sup>3</sup> Moreover, instead of unnatural pair-wise cost functions of  $p$ -OGMLVQ, we formulate a single unified cost function:

$$C_{a\text{-OGMLVQ}} = \sum_{i=1}^n \Phi(f(\mathbf{x}_i)), \quad (18)$$

where

$$f(\mathbf{x}_i) = \frac{\sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) - \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-)}{\sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-)}. \quad (19)$$

Note that instead of having  $f_j$  for each  $j$ th pair of correct and incorrect prototypes, we only have one  $f$  for all the identified prototypes. The  $p$ -OGMLVQ cost function can be seen as sum of  $r$  weighted versions of the GMLVQ cost function (Schneider et al., 2009). The correct/incorrect prototypes are not treated equally in the  $p$ -OGMLVQ cost function. In contrast, our cost function treats the correct and incorrect prototypes equally by accumulating costs of all correct and all incorrect prototypes in a single cost  $f(\mathbf{x}_i)$ . Our approach does not take advantages of binary classifier as Lin

<sup>2</sup>  $D$  can be set as mean distance as well.

<sup>3</sup> In  $p$ -OGMLVQ,  $\sigma_1 = \sigma_2$ .

and Li (2012) does, but directly incorporates the class ordering information in a natural multi-class classification algorithm – GMLVQ.

Following Fouad and Tiño (2012), we choose  $\Phi$  as identity function. The learning rule can be derived from the cost function by differentiating with respect to prototypes:

$$\frac{\partial f(\mathbf{x}_i)}{\partial \mathbf{w}_j^+} = \frac{\partial f(\mathbf{x}_i)}{\partial d^A(\mathbf{x}_i, \mathbf{w}_j^+)} \cdot \frac{\partial d^A(\mathbf{x}_i, \mathbf{w}_j^+)}{\partial \mathbf{w}_j^+}, \quad (20)$$

where

$$\begin{aligned} \frac{\partial f(\mathbf{x}_i)}{\partial d^A(\mathbf{x}_i, \mathbf{w}_j^+)} &= \frac{\alpha_j^+ \left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)}{\left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)^2} \\ &\quad - \frac{\alpha_j^+ \left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) - \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)}{\left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)^2} \\ &= \frac{2\alpha_j^+ \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-)}{\left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)^2}, \quad (21) \end{aligned}$$

and

$$\frac{\partial d^A(\mathbf{x}_i, \mathbf{w}_j^+)}{\partial \mathbf{w}_j^+} = -2\Omega^T \Omega (\mathbf{x}_i - \mathbf{w}_j^+) = -2\Lambda (\mathbf{x}_i - \mathbf{w}_j^+). \quad (22)$$

Partial derivative of the cost function with respect to  $\mathbf{w}_j^-$  can be written as:

$$\frac{\partial f(\mathbf{x}_i)}{\partial \mathbf{w}_j^-} = \frac{\partial f(\mathbf{x}_i)}{\partial d^A(\mathbf{x}_i, \mathbf{w}_j^-)} \cdot \frac{\partial d^A(\mathbf{x}_i, \mathbf{w}_j^-)}{\partial \mathbf{w}_j^-}, \quad (23)$$

where

$$\begin{aligned} \frac{\partial f(\mathbf{x}_i)}{\partial d^A(\mathbf{x}_i, \mathbf{w}_j^-)} &= - \frac{\frac{\partial \alpha_j^-}{\partial d^A(\mathbf{x}_i, \mathbf{w}_j^-)} d^A(\mathbf{x}_i, \mathbf{w}_j^+) + \alpha_j^-}{\left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)^2} \\ &\quad \cdot \left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right) \\ &\quad - \frac{\frac{\partial \alpha_j^-}{\partial d^A(\mathbf{x}_i, \mathbf{w}_j^-)} d^A(\mathbf{x}_i, \mathbf{w}_j^+) + \alpha_j^-}{\left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)^2} \\ &\quad \cdot \left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) - \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right) \\ &= - \frac{2(1 - d^A(\mathbf{x}_i, \mathbf{w}_j^-)/(2\sigma_3^2)) \alpha_j^- \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+)}{\left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)^2}, \quad (24) \end{aligned}$$

and

$$\frac{\partial d^A(\mathbf{x}_i, \mathbf{w}_j^-)}{\partial \mathbf{w}_j^-} = -2\Omega^T \Omega (\mathbf{x}_i - \mathbf{w}_j^-) = -2\Lambda (\mathbf{x}_i - \mathbf{w}_j^-). \quad (25)$$

Denoting  $\mu_j^+ = \frac{\partial f(\mathbf{x}_i)}{\partial d^A(\mathbf{x}_i, \mathbf{w}_j^+)}$  and  $\mu_j^- = -\frac{\partial f(\mathbf{x}_i)}{\partial d^A(\mathbf{x}_i, \mathbf{w}_j^-)}$ , derivative of the cost function with respect to  $\Omega$  can be computed as follows:

$$\begin{aligned} \frac{\partial f(\mathbf{x}_i)}{\partial \Omega} &= \sum_{k=1}^r \frac{\partial f(\mathbf{x}_i)}{\partial d^A(\mathbf{x}_i, \mathbf{w}_k^+)} \cdot \frac{\partial d^A(\mathbf{x}_i, \mathbf{w}_k^+)}{\partial \Omega} \\ &\quad + \sum_{k=1}^r \frac{\partial f(\mathbf{x}_i)}{\partial d^A(\mathbf{x}_i, \mathbf{w}_k^-)} \cdot \frac{\partial d^A(\mathbf{x}_i, \mathbf{w}_k^-)}{\partial \Omega} \\ &= \sum_{k=1}^r \mu_k^+ \cdot \frac{\partial d^A(\mathbf{x}_i, \mathbf{w}_k^+)}{\partial \Omega} - \sum_{k=1}^r \mu_k^- \cdot \frac{\partial d^A(\mathbf{x}_i, \mathbf{w}_k^-)}{\partial \Omega}, \quad (26) \end{aligned}$$

where

$$\frac{\partial d^A(\mathbf{x}_i, \mathbf{w}_k^+)}{\partial \Omega} = 2\Omega (\mathbf{x}_i - \mathbf{w}_k^+) (\mathbf{x}_i - \mathbf{w}_k^+)^T, \quad (27)$$

and

$$\frac{\partial d^A(\mathbf{x}_i, \mathbf{w}_k^-)}{\partial \Omega} = 2\Omega (\mathbf{x}_i - \mathbf{w}_k^-) (\mathbf{x}_i - \mathbf{w}_k^-)^T. \quad (28)$$

Consequently, updating rules can be expressed as:

$$\mathbf{w}_j^+ := \mathbf{w}_j^+ + \epsilon \cdot 2 \cdot \mu_j^+ \cdot \Lambda \cdot (\mathbf{x}_i - \mathbf{w}_j^+), j = 1, \dots, r \quad (29)$$

$$\mathbf{w}_j^- := \mathbf{w}_j^- - \epsilon \cdot 2 \cdot \mu_j^- \cdot \Lambda \cdot (\mathbf{x}_i - \mathbf{w}_j^-), j = 1, \dots, r \quad (30)$$

$$\begin{aligned} \Omega &:= \Omega - \eta \cdot \left( 2 \cdot \sum_{k=1}^r \mu_k^+ \cdot \Omega \cdot (\mathbf{x}_i - \mathbf{w}_k^+) (\mathbf{x}_i - \mathbf{w}_k^+)^T \right. \\ &\quad \left. - 2 \cdot \sum_{k=1}^r \mu_k^- \cdot \Omega \cdot (\mathbf{x}_i - \mathbf{w}_k^-) (\mathbf{x}_i - \mathbf{w}_k^-)^T \right). \quad (31) \end{aligned}$$

Comparing the prototype updating rules in the original  $p$ -OGMLVQ approach ((5)–(6)), with those of our approach ((29)–(30)), we observe that the prototype updates differ mainly in the definition of  $\mu_j^+$  and  $\mu_j^-$ . In  $p$ -OGMLVQ, the updating coefficient  $\mu_j^+$ ,  $j = 1, \dots, r$ , corresponding to different correct prototypes are not directly comparable (the same holds for  $\mu_j^-$ ). However, in our approach  $\mu_j^+$  can be rewritten as

$$\mu_j^+ = \alpha_j^+ \cdot \gamma^+, \quad (32)$$

where

$$\gamma^+ = \frac{2 \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-)}{\left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)^2}. \quad (33)$$

The differences among  $\mu_j^+$ ,  $j = 1, \dots, r$  are solely due to the weights  $\alpha_j^+$ . Recall that  $\alpha_j^+$  will be greater if the class of the corresponding prototype is closer to the class of the input data  $\mathbf{x}_i$  in the label space. Consequently, greater  $\mu_j^+$  will be assigned to the prototype whose class is closer to the input data.

Similarly, the weight  $\mu_j^-$  can be rewritten as:

$$\mu_j^- = (1 - d^A(\mathbf{x}_i, \mathbf{w}_j^-)/(2\sigma_3^2)) \alpha_j^- \cdot \gamma^-, \quad (34)$$

where

$$\gamma^- = \frac{2 \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+)}{\left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)^2}.$$

We can see that  $\mu_j^-$  increases with  $\alpha_j^-$  but decreases with the distance  $d^A(\mathbf{x}_i, \mathbf{w}_j^-)$ , which means that if the incorrect prototype is closer to the input, with large class difference, the updating push-force  $\mu_j^-$  will be greater.

### 3.1. Gradient structure

Following Villmann (2002), we will now show that the derivatives of our cost function vanish at boundaries of the prototype receptive fields and hence gradient based optimization will constitute a valid gradient descent.

For convenience, we denote the squared distance between input  $\mathbf{x}$  and prototype  $\mathbf{w}_i$ ,  $(\mathbf{x} - \mathbf{w}_i)^T \cdot \Lambda \cdot (\mathbf{x} - \mathbf{w}_i)$ , by  $d_i(\mathbf{x})$ . Let  $J(\mathbf{x}) = \{J_1, J_2, \dots, J_{r(\mathbf{x})}\}$  and  $K(\mathbf{x}) = \{K_1, K_2, \dots, K_{r(\mathbf{x})}\}$  be the sets of  $r(\mathbf{x})$  indices of correct and incorrect prototypes, respectively, that are selected for updating, given the input  $\mathbf{x}$ . Note that  $r(\mathbf{x})$  can vary, depending on the input location  $\mathbf{x}$ . Our scaled general cost function (18)–(19) has the following form (taking  $\Phi$  as identity function):

$$\hat{S} = \frac{1}{n} \sum_{i=1}^n \frac{\sum_{k=1}^{r(\mathbf{x}_i)} \alpha_{J_k}^+(\mathbf{x}_i) \cdot d_{J_k}(\mathbf{x}_i) - \sum_{k=1}^{r(\mathbf{x}_i)} \alpha_{K_k}^-(\mathbf{x}_i) \cdot d_{K_k}(\mathbf{x}_i)}{\sum_{k=1}^{r(\mathbf{x}_i)} \alpha_{J_k}^+(\mathbf{x}_i) \cdot d_{J_k}(\mathbf{x}_i) + \sum_{k=1}^{r(\mathbf{x}_i)} \alpha_{K_k}^-(\mathbf{x}_i) \cdot d_{K_k}(\mathbf{x}_i)}. \quad (35)$$

As in Villmann (2002), suppose the input data comes from a distribution  $p$  on the input space  $\mathbb{R}^m$  and there is a labeling function  $y: \mathbb{R}^m \rightarrow \{1, \dots, C\}$  assigning each input  $\mathbf{x}$  its class label  $y(\mathbf{x})$ . The continuous version of the error function  $\hat{S}$  reads

$$S = \int_{\mathbb{R}^m} \frac{\sum_{k=1}^{r(\mathbf{x})} \alpha_{J_k}^+(\mathbf{x}) \cdot d_{J_k}(\mathbf{x}) - \sum_{k=1}^{r(\mathbf{x})} \alpha_{K_k}^-(\mathbf{x}) \cdot d_{K_k}(\mathbf{x})}{\sum_{k=1}^{r(\mathbf{x})} \alpha_{J_k}^+(\mathbf{x}) \cdot d_{J_k}(\mathbf{x}) + \sum_{k=1}^{r(\mathbf{x})} \alpha_{K_k}^-(\mathbf{x}) \cdot d_{K_k}(\mathbf{x})} dp(\mathbf{x}). \quad (36)$$

To simplify the presentation, in what follows we do not explicitly denote dependence on  $\mathbf{x}$ , with the understanding that quantities such as  $r$ ,  $\alpha_i^\pm$ ,  $d_i$  etc. are all functions of  $\mathbf{x}$ . We partition the input space and the set of prototypes into class sets  $X_\ell = \{\mathbf{x} | y(\mathbf{x}) = \ell\}$  and  $W_\ell = \{\mathbf{w}_j | c(\mathbf{w}_j) = \ell\}$ , respectively,  $\ell = 1, 2, \dots, C$ . Assuming  $X_\ell$  are measurable with respect to  $p$ , we decompose  $S$  as follows:

$$\sum_{\ell=1}^C \int_{X_\ell} \frac{\sum_{||-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_l} \alpha_j^+ d_j h(j, l, \ell) - \sum_{||-\ell| > \tau} \sum_{\mathbf{w}_k \in W_l} \alpha_k^- d_k \tilde{h}(k, \ell)}{\sum_{||-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_l} \alpha_j^+ d_j h(j, l, \ell) + \sum_{||-\ell| > \tau} \sum_{\mathbf{w}_k \in W_l} \alpha_k^- d_k \tilde{h}(k, \ell)} dp(\mathbf{x}) \quad (37)$$

where

- $h(j, l, \ell)$  is an indicator function attaining value 1 if and only if the prototype  $\mathbf{w}_j \in W_l$  is the closest to  $\mathbf{x}$  among the prototypes in  $W_l$  and the class  $l$  is among the  $r(\mathbf{x})$  correct classes w.r.t. class  $\ell$  that are selected for prototype update, given input  $\mathbf{x}$ ; otherwise  $h(j, l, \ell)$  is zero.
- $\tilde{h}(k, \ell)$  is an indicator function indicating that  $\mathbf{w}_k$  is an incorrect prototype w.r.t. class  $\ell$  and is among the  $r(\mathbf{x})$  incorrect prototypes closest to  $\mathbf{x}$ .

To find a closed-form expressions for the indicator functions  $h(j, l, \ell)$  and  $\tilde{h}(k, \ell)$  we first introduce two further indicator functions:

- $\mathcal{I}(j, l)$  indicates that the prototype  $\mathbf{w}_j$  is the closest prototype to  $\mathbf{x}$  among those labeled with class  $l$ ,
- $\tilde{\mathcal{I}}(k, D)$  indicates that the prototype  $\mathbf{w}_k$  is spatially within the  $D$ -neighborhood of  $\mathbf{x}$ . Recall that  $D$  is the median distance of the example  $\mathbf{x}$  of class  $\ell = y(\mathbf{x})$  to the incorrect prototypes, i.e. prototypes labeled with  $l$ , such that  $||-\ell| > \tau$ .

Denoting by  $H$  the Heaviside function, we have

$$\mathcal{I}(j, l) = H \left( \sum_{\mathbf{w}_q \in W_l} H(d_q - d_j) - |W_l| \right), \quad (38)$$

where  $|W_l|$  is the size of  $W_l$  (and hence the number of prototypes labeled with class  $l$ ) and

$$\tilde{\mathcal{I}}(k, D) = H(D - d_k). \quad (39)$$

Denote by  $N(\ell)$  and  $\tilde{N}(\ell)$  the number of correct and incorrect prototypes, respectively, given the input  $\mathbf{x}$  and its class  $\ell = y(\mathbf{x})$ . We have,

$$N(\ell) = \sum_{||-\ell| \leq \tau} 1, \quad \tilde{N}(\ell) = \sum_{||-\ell| > \tau} \sum_{\mathbf{w}_k \in W_l} \tilde{\mathcal{I}}(k, D).$$

Hence,

$$\begin{aligned} r(\mathbf{x}) &= \min \{N(\ell), \tilde{N}(\ell)\} \\ &= N(\ell) - (N(\ell) - \tilde{N}(\ell)) \cdot H(N(\ell) - \tilde{N}(\ell)) \\ &= \tilde{N}(\ell) - (\tilde{N}(\ell) - N(\ell)) \cdot H(\tilde{N}(\ell) - N(\ell)). \end{aligned}$$

We are now ready to formally express  $h(j, l, \ell)$  and  $\tilde{h}(k, \ell)$ :

$$h(j, l, \ell) = \mathcal{I}(j, l) \cdot H(g(j, \ell) + f(\ell)), \quad (40)$$

$$\tilde{h}(k, \ell) = \tilde{\mathcal{I}}(k, D) \cdot H(\tilde{g}(j, \ell) + \tilde{f}(\ell)), \quad (41)$$

where

$$g(j, \ell) = \sum_{||-\ell| \leq \tau} \sum_{\mathbf{w}_q \in W_l} \mathcal{I}(q, l) \cdot H(d_q - d_j), \quad (42)$$

$$\begin{aligned} f(\ell) &= -(N(\ell) - r(\mathbf{x}) + 1) \\ &= -(N(\ell) - \tilde{N}(\ell)) \cdot H(N(\ell) - \tilde{N}(\ell)) - 1, \end{aligned} \quad (43)$$

$$\tilde{g}(k, \ell) = \sum_{||-\ell| > \tau} \sum_{\mathbf{w}_q \in W_l} \tilde{\mathcal{I}}(q, D) \cdot H(d_q - d_k), \quad (44)$$

$$\begin{aligned} \tilde{f}(\ell) &= -(\tilde{N}(\ell) - r(\mathbf{x}) + 1) \\ &= -(\tilde{N}(\ell) - N(\ell)) \cdot H(\tilde{N}(\ell) - N(\ell)) - 1. \end{aligned} \quad (45)$$

The first factor of  $h(j, l, \ell)$  verifies whether the prototype  $\mathbf{w}_j$  is one of the ‘‘correct’’ prototypes, while the second factor further checks whether this prototype is one of the first  $r$  closest ones among all the ‘‘correct’’ prototypes. This is realized by ordering the ‘‘correct’’ prototypes through  $g(j, \ell)$  and selecting the  $r$  first closest ones by  $f(\ell)$ . The same trick is used in  $\tilde{h}(k, \ell)$  to pick the first  $r$  closest prototypes among those ‘‘incorrect’’ ones.

Let  $c_i = c(\mathbf{w}_i)$  be the label of prototype  $\mathbf{w}_i$ . Denoting

$$\begin{aligned} D_i &= \Lambda \cdot (\mathbf{x} - \mathbf{w}_i), \\ E(\ell) &= \sum_{||-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_l} \alpha_j^+ d_j h(j, l, \ell), \end{aligned}$$

and

$$\tilde{E}(\ell) = \sum_{||-\ell| > \tau} \sum_{\mathbf{w}_k \in W_l} \alpha_k^- d_k \tilde{h}(k, \ell),$$

the derivative of (37) with respect to  $\mathbf{w}_i$  reads:

$$\sum_{||-\ell| \leq \tau} \int_{X_\ell} \frac{4\alpha_i \tilde{E}(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} D_i h(i, c_i, \ell) dp(\mathbf{x}) \quad (46)$$

$$+ \sum_{||-\ell| > \tau} \int_{X_\ell} \frac{-4(1 - d_i / (2\sigma_3^2)) \alpha_i E(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} D_i \tilde{h}(i, \ell) dp(\mathbf{x}) \quad (47)$$

$$+ \sum_{\ell=1}^C \int_{X_\ell} \sum_{|\ell-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_j} \frac{2\alpha_j d_j \tilde{E}(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} \frac{\partial h(j, l, \ell)}{\partial \mathbf{w}_i} d\mathbf{p}(\mathbf{x}) \quad (48)$$

$$+ \sum_{|\ell-c_i| > \tau} \int_{X_\ell} \sum_{|\ell-\ell| > \tau} \sum_{\mathbf{w}_k \in W_k} \frac{2\alpha_k d_k E(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} \frac{\partial \tilde{h}(k, \ell)}{\partial \mathbf{w}_i} d\mathbf{p}(\mathbf{x}). \quad (49)$$

Note that (46) and (47) correspond to updates (29) and (30), respectively. Terms (48) and (49) will vanish as shown in the Appendix.

Similarly, we examine the derivative of (37) with respect to the metric term  $\Omega$ . Denoting  $\Omega(\mathbf{x} - \mathbf{w}_i)(\mathbf{x} - \mathbf{w}_i)^T$  by  $F_i$ , then the derivative reads:

$$\sum_{\ell=1}^C \int_{X_\ell} \frac{4}{(E(\ell) + \tilde{E}(\ell))^2} \left( \tilde{E}(\ell) \sum_{|\ell-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_j} \alpha_j F_j h(j, l, \ell) \right. \quad (50)$$

$$\left. - E(\ell) \sum_{|\ell-\ell| > \tau} \sum_{\mathbf{w}_k \in W_k} (1 - d_k / (2\sigma_3^2)) \alpha_k F_k \tilde{h}(k, \ell) \right) d\mathbf{p}(\mathbf{x}) \quad (51)$$

$$+ \sum_{\ell=1}^C \int_{X_\ell} \frac{2}{(E(\ell) + \tilde{E}(\ell))^2} \left( \tilde{E}(\ell) \sum_{|\ell-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_j} \alpha_j d_j \frac{\partial h(j, l, \ell)}{\partial \Omega} \right. \quad (52)$$

$$\left. - \tilde{E}(\ell) \sum_{|\ell-\ell| > \tau} \sum_{\mathbf{w}_k \in W_k} \alpha_k d_k \frac{\partial \tilde{h}(k, \ell)}{\partial \Omega} \right) d\mathbf{p}(\mathbf{x}). \quad (53)$$

Terms (50) and (51) correspond to the update of  $\Omega$  in (31). As shown in the Appendix, (52) and (53) will vanish. Consequently, with an appropriate choice of learning rate, the update of a-OGMLVQ constitutes a stochastic gradient descent method.

### 3.2. Hyperparameter learning in the proposed algorithm

In this section, we propose to automatically adapt the bandwidth parameters  $\sigma_1$ ,  $\sigma_2$ , and  $\sigma_3$  in the Gaussian weighting functions during training process through minimizing the cost function with respect to the parameters  $\sigma_1$ ,  $\sigma_2$ , and  $\sigma_3$ .

Partial derivative of  $f(\mathbf{x}_i)$  with respect to  $\sigma_1$  is computed as follows:

$$\frac{\partial f(\mathbf{x}_i)}{\partial \sigma_1} = \sum_{j=1}^r \frac{\partial f(\mathbf{x}_i)}{\partial \alpha_j^+} \cdot \frac{\partial \alpha_j^+}{\partial \sigma_1}, \quad (54)$$

where

$$\begin{aligned} \frac{\partial f(\mathbf{x}_i)}{\partial \alpha_j^+} &= \frac{d^A(\mathbf{x}_i, \mathbf{w}_j^+)}{\left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)^2} \\ &\cdot \left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right) \\ &- \frac{d^A(\mathbf{x}_i, \mathbf{w}_j^+)}{\left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)^2} \\ &\cdot \left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) - \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right) \\ &= \frac{2d^A(\mathbf{x}_i, \mathbf{w}_j^+) \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-)}{\left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)^2} \end{aligned}$$

$$= d^A(\mathbf{x}_i, \mathbf{w}_j^+) \cdot \gamma^+, \quad (55)$$

and

$$\frac{\partial \alpha_j^+}{\partial \sigma_1} = \alpha_j^+ \frac{h(c(\mathbf{x}_i), c(\mathbf{w}_j^+))^2}{2\sigma_1^3}. \quad (56)$$

Partial derivatives of  $f(\mathbf{x}_i)$  with respect to  $\sigma_2$  and  $\sigma_3$  read:

$$\frac{\partial f(\mathbf{x}_i)}{\partial \sigma_2} = \sum_{j=1}^r \frac{\partial f(\mathbf{x}_i)}{\partial \alpha_j^-} \cdot \frac{\partial \alpha_j^-}{\partial \sigma_2}, \quad (57)$$

$$\frac{\partial f(\mathbf{x}_i)}{\partial \sigma_3} = \sum_{j=1}^r \frac{\partial f(\mathbf{x}_i)}{\partial \alpha_j^-} \cdot \frac{\partial \alpha_j^-}{\partial \sigma_3}, \quad (58)$$

where

$$\begin{aligned} \frac{\partial f(\mathbf{x}_i)}{\partial \alpha_j^-} &= - \frac{d^A(\mathbf{x}_i, \mathbf{w}_j^-)}{\left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)^2} \\ &\cdot \left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right) \\ &- \frac{d^A(\mathbf{x}_i, \mathbf{w}_j^-)}{\left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)^2} \\ &\cdot \left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) - \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right) \\ &= - \frac{2d^A(\mathbf{x}_i, \mathbf{w}_j^-) \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+)}{\left( \sum_{k=1}^r \alpha_k^+ d^A(\mathbf{x}_i, \mathbf{w}_k^+) + \sum_{k=1}^r \alpha_k^- d^A(\mathbf{x}_i, \mathbf{w}_k^-) \right)^2} \\ &= -d^A(\mathbf{x}_i, \mathbf{w}_j^-) \cdot \gamma^-, \end{aligned} \quad (59)$$

and

$$\frac{\partial \alpha_j^-}{\partial \sigma_2} = \alpha_j^- \frac{(T - h(c(\mathbf{x}_i) - c(\mathbf{w}_j^-)))^2}{2\sigma_2^3}, \quad (60)$$

$$\frac{\partial \alpha_j^-}{\partial \sigma_3} = \alpha_j^- \frac{d^A(\mathbf{x}_i, \mathbf{w}_j^-)}{2\sigma_3^3}. \quad (61)$$

Hence, updating of  $\sigma_1$ ,  $\sigma_2$  and  $\sigma_3$  proceeds as follows:

$$\begin{aligned} \sigma_1 &:= \sigma_1 \\ &- \epsilon \cdot \frac{\gamma^+}{2\sigma_1^3} \cdot \sum_{j=1}^r \alpha_j^+ \cdot d^A(\mathbf{x}_i, \mathbf{w}_j^+) \cdot (h(c(\mathbf{x}_i), c(\mathbf{w}_j^+)))^2, \end{aligned} \quad (62)$$

$$\begin{aligned} \sigma_2 &:= \sigma_2 \\ &+ \epsilon \cdot \frac{\gamma^-}{2\sigma_2^3} \cdot \sum_{j=1}^r \alpha_j^- \cdot d^A(\mathbf{x}_i, \mathbf{w}_j^-) \cdot (T - h(c(\mathbf{x}_i), c(\mathbf{w}_j^-)))^2, \end{aligned} \quad (63)$$

$$\sigma_3 := \sigma_3 + \epsilon \cdot \frac{\gamma^-}{2\sigma_3^3} \cdot \sum_{j=1}^r \alpha_j^- \cdot (d^A(\mathbf{x}_i, \mathbf{w}_j^-))^2. \quad (64)$$

The hyperparameter updation (62)–(64) is similar to that in Robust Soft LVQ algorithm presented in Schneider, Biehl, and Hammer (2010). However the hyperparameters in our approach serve a difference purpose, introducing scales to the ordered label and input spaces.

The overall method is summarized in algorithm 1.

**Algorithm 1** The proposed ordinal regression algorithm based on LVQ

- 1: Initialization : Initialize the prototypes  $\mathbf{w}_j \in \mathbb{R}^d$ . Initialize the metric tensor parameter  $\Omega$  by setting it as identity matrix. Randomly initialize the bandwidth parameters  $\sigma_1, \sigma_2$ , and  $\sigma_3$ .
- 2: **while** a stopping criterion is not reached **do**
- 3: Randomly select a training example  $(\mathbf{x}_i, y_i)$ .
- 4: Determine the correct classes  $C^+$  and incorrect classes  $C^-$  for  $\mathbf{x}_i$  based on (12) and (13), respectively.
- 5: Identify the correct prototype set  $W^+$  and incorrect prototype set  $W^-$  using (14) and (15), respectively.
- 6: Order  $W^+$  and  $W^-$  in increasing order in terms of distance to the input  $\mathbf{x}_i$ , only keep the first  $r$  prototype in the set  $W^+$  and  $W^-$ , where  $r = \min(|W^+|, |W^-|)$ .
- 7: Compute the weights for the correct and incorrect prototypes found in the previous step according to (16) and (17), respectively.
- 8: **for**  $\forall \mathbf{w}_j \in W^+$  **do**
- 9: Update  $\mathbf{w}_j$  according to (29)
- 10: **end for**
- 11: **for**  $\forall \mathbf{w}_j \in W^-$  **do**
- 12: Update  $\mathbf{w}_j$  according to (30)
- 13: **end for**
- 14: Update  $\Omega$  according to (31). Similarly,  $\Omega$  is normalized so that  $\sum_i A_{ii} = 1$  to prevent the algorithm from degeneration.
- 15: Update  $\sigma_1, \sigma_2$ , and  $\sigma_3$  according to (62), (63), and (64), respectively.
- 16: **end while**

#### 4. Experiments

In this section, we study the performance of the proposed ordinal approach. Three performance measures were used:

1. **MZE**: The Mean Zero-one Error (MZE), also known as misclassification rate, is the fraction of incorrect predictions:

$$MZE = \frac{\sum_{i=1}^n I(y_i \neq \hat{y}(\mathbf{x}_i))}{n}, \quad (65)$$

where  $n$  is the number of the test examples.  $I(\cdot)$  is the indicator function whose function value is 1, if the argument  $*$  is true, 0 otherwise. MZE values range from 0 to 1 and reflect a global performance for the classification task (without considering class order). However, this performance measurement alone is not suitable for ordinal regression problems.

2. **MAE**: The Mean Absolute Error (MAE) is the average absolute deviation of the predicted ranks from the true ranks:

$$MAE = \frac{\sum_{i=1}^n |y_i - \hat{y}(\mathbf{x}_i)|}{n}. \quad (66)$$

The MAE values range from 0 to  $C - 1$ . Since MZE does not reflect the category order, MAE is typically used in the ordinal regression literature together with MZE (Sánchez-Monedero et al., 2013). However, this performance measurement is not suitable for imbalanced datasets.

3. **MMAE**: Macroaveraged Mean Absolute Error (MMAE) is specially designed for evaluating ordinal regression problems with imbalanced classes:

$$MMAE = \frac{1}{C} \sum_{j=1}^C \frac{\sum_{y_i=j} |y_i - \hat{y}(\mathbf{x}_i)|}{n_j}, \quad (67)$$

**Table 1**  
Brief description of the ordinal regression benchmark datasets.

Dataset	Dimension	#Training	#Test
Pyrimidines	27	50	24
MachineCPU	6	150	59
Boston	13	300	206
Abalone	8	1000	3 177
Bank	32	3000	5 182
Computer	21	4000	4 182
California	8	5000	15 640
Census	8	6000	16 784

where  $C$  is the number of classes and  $n_j$  is the number of examples of class  $j$ . The MMAE values also range from 0 to  $C - 1$ .

We first compared our approach with  $p$ -OGMLVQ on the ordinal regression benchmark datasets provided by Chu and Keerthi (2005). As a baseline, the performance of GMLVQ was also reported. A brief description of the datasets is given by Table 1. These datasets are not real ordinal datasets but regression ones. These datasets will be turned into ordinal regression through discretization. Following Fouad and Tiño (2012), the continuous targets were discretized into 10 ordinal categories (equal-frequency binning). The input vectors were normalized to zero mean and unit variance. Each dataset was randomly partitioned into training/test splits according to Table 1. The partition was repeated 20 times independently, yielding 20 re-sampled training/test sets.

Prototype initialization and learning rates were set in the same way for all three approaches. Number of prototypes per class was tuned on the training sets, with candidates ranging in  $\{1, 2, 3, 4, 5\}$ , via 5-fold cross validation. For  $p$ -OGMLVQ, the ranking loss threshold  $\tau$  together with the bandwidth parameters  $\sigma$  and  $\sigma_1$  was also tuned through 5-fold cross validation. Candidate values for the hyperparameters are listed as follows:  $\tau \in \{0, 1, 2\}$ ,  $\sigma, \sigma_1 \in \{0.1, 0.3, 0.5, 0.7, 0.8, 0.9, 1, 1.2, 2, 3, 5, 50\}$ . For our approach, we also tuned the ranking loss threshold  $\tau$  via 5-fold cross validation over the same candidate values as in  $p$ -OGMLVQ (the bandwidth parameters in our approach are automatically adapted during training process).

The MZE and MAE results,<sup>4</sup> along with standard deviation across 20 runs, are given by Figs. 3 and 4, respectively. The results show that our  $a$ -OGMLVQ ordinal approach was better than  $p$ -OGMLVQ in terms of MZE and MAE performance measures in most of the datasets, only occasionally our  $a$ -OGMLVQ obtained comparable performance. One possible reason might be that the problem mentioned in Fig. 2 do not happen in that case. Still, our  $a$ -OGMLVQ ordinal approach is consistently better than the baseline GMLVQ approach.

We further compared our ordinal approach to other ordinal approaches: two support vector ordinal approaches proposed in Chu and Keerthi (2005) – support vector ordinal regression with implicit constraints (SVORIM) and support vector ordinal regression with explicit constraints (SVOREX), SVM based reduction framework of ordinal regression to binary classification (RED-SVM) proposed in Li and Lin (2006), and an ordinal approach extended from kernel discriminate analysis proposed in Sun et al. (2010). The MZE and MAE results are reported in Tables 2 and 3, respectively. Our  $a$ -OGMLVQ approach beats the other approaches four times in terms of MZE and twice in terms of MAE. In order to have a general summary of the results, we computed the mean ranking values when comparing the different methods, where rank  $R = 1$  corresponds to the best method, while  $R = 6$  represents for the worse one. The mean rankings  $\bar{R}$  were computed based on the performance on the

<sup>4</sup> MMAE results were not given in this case, since these datasets are balanced.

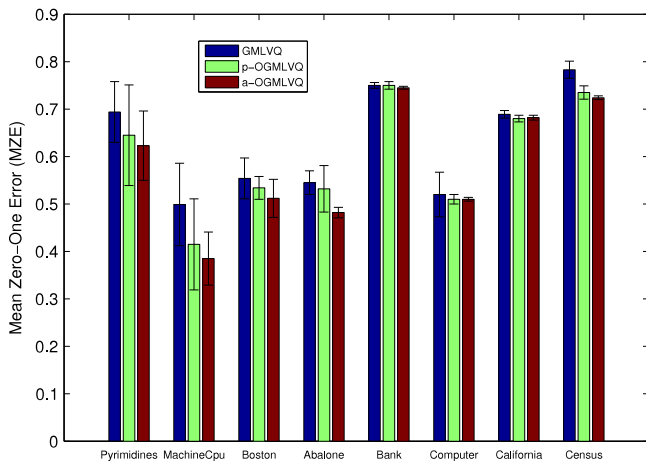


**Table 2**  
Mean zero–one error (MZE) results along with standard deviations ( $\pm$ ) across 20 runs of the ordinal regression benchmark datasets given in Table 1. Our  $\alpha$ -OGMLVQ approach is compared with KDLMR reported in Sun et al. (2010), SVORIM and SVOREX reported in Chu and Keerthi (2005), RED-SVM reported in Li and Lin (2006), and  $p$ -OGMLVQ reported in Fouad and Tiño (2012). The best performance for each dataset has been boldfaced.

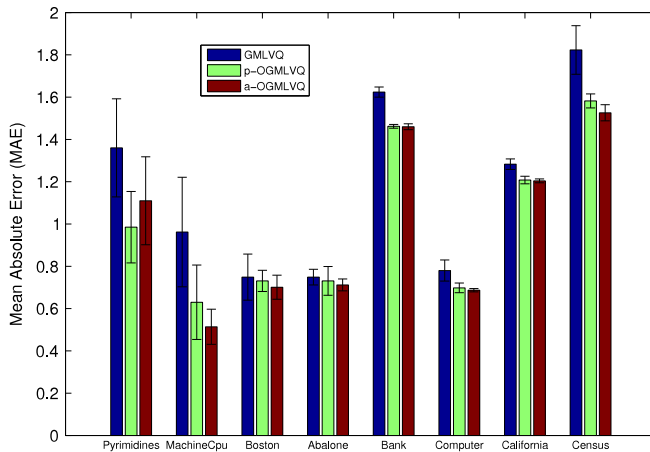
Dataset/Method	KDLOR	SVORIM	SVOREX	RED-SVM	$p$ -OGMLVQ	$\alpha$ -OGMLVQ
Pyrimidines	0.739 $\pm$ 0.055	0.719 $\pm$ 0.066	0.752 $\pm$ 0.063	0.762 $\pm$ 0.021	0.645 $\pm$ 0.106	<b>0.623 <math>\pm</math> 0.073</b>
MachineCPU	0.480 $\pm$ 0.010	0.655 $\pm$ 0.045	0.611 $\pm$ 0.056	0.572 $\pm$ 0.013	0.415 $\pm$ 0.096	<b>0.385 <math>\pm</math> 0.056</b>
Boston	0.560 $\pm$ 0.020	0.561 $\pm$ 0.026	0.569 $\pm$ 0.025	0.541 $\pm$ 0.009	0.534 $\pm$ 0.024	<b>0.512 <math>\pm</math> 0.040</b>
Abalone	0.740 $\pm$ 0.020	0.732 $\pm$ 0.007	0.736 $\pm$ 0.011	0.721 $\pm$ 0.002	0.532 $\pm$ 0.049	<b>0.482 <math>\pm</math> 0.011</b>
Bank	0.745 $\pm$ 0.003	0.751 $\pm$ 0.005	<b>0.744 <math>\pm</math> 0.005</b>	0.751 $\pm$ 0.001	0.750 $\pm$ 0.008	0.745 $\pm$ 0.003
Computer	0.472 $\pm$ 0.020	0.473 $\pm$ 0.005	0.462 $\pm$ 0.005	<b>0.451 <math>\pm</math> 0.002</b>	0.510 $\pm$ 0.010	0.510 $\pm$ 0.004
California	0.643 $\pm$ 0.005	0.639 $\pm$ 0.003	0.640 $\pm$ 0.003	<b>0.613 <math>\pm</math> 0.001</b>	0.680 $\pm$ 0.007	0.682 $\pm$ 0.005
Census	0.711 $\pm$ 0.020	0.705 $\pm$ 0.002	0.699 $\pm$ 0.002	<b>0.688 <math>\pm</math> 0.001</b>	0.735 $\pm$ 0.014	0.724 $\pm$ 0.004

**Table 3**  
Mean absolute error (MAE) results along with standard deviations ( $\pm$ ) across 20 runs of the ordinal regression benchmark datasets given in Table 1. Our  $\alpha$ -OGMLVQ approach is compared with KDLMR reported in Sun et al. (2010), SVORIM and SVOREX reported in Chu and Keerthi (2005), RED-SVM reported in Li and Lin (2006), and  $p$ -OGMLVQ reported in Fouad and Tiño (2012). The best performance for each dataset has been boldfaced.

Dataset	KDLOR	SVOR-IMC	SVOR-EXC	RED-SVM	$p$ -OGMLVQ	$\alpha$ -OGMLVQ
Pyrimidines	1.100 $\pm$ 0.100	1.294 $\pm$ 0.204	1.331 $\pm$ 0.193	1.304 $\pm$ 0.040	<b>0.985 <math>\pm</math> 0.169</b>	1.114 $\pm$ 0.208
MachineCPU	0.690 $\pm$ 0.015	0.990 $\pm$ 0.115	0.986 $\pm$ 0.127	0.842 $\pm$ 0.022	0.630 $\pm$ 0.176	<b>0.514 <math>\pm</math> 0.083</b>
Boston	<b>0.700 <math>\pm</math> 0.035</b>	0.747 $\pm$ 0.049	0.773 $\pm$ 0.049	0.732 $\pm$ 0.013	0.731 $\pm$ 0.050	0.702 $\pm$ 0.057
Abalone	1.400 $\pm$ 0.050	1.361 $\pm$ 0.013	1.391 $\pm$ 0.021	1.383 $\pm$ 0.004	0.731 $\pm$ 0.068	<b>0.712 <math>\pm</math> 0.028</b>
Bank	1.450 $\pm$ 0.020	<b>1.393 <math>\pm</math> 0.011</b>	1.512 $\pm$ 0.017	1.404 $\pm$ 0.002	1.462 $\pm$ 0.009	1.460 $\pm$ 0.014
Computer	0.601 $\pm$ 0.025	0.596 $\pm$ 0.008	0.602 $\pm$ 0.009	<b>0.565 <math>\pm</math> 0.002</b>	0.698 $\pm$ 0.023	0.687 $\pm$ 0.008
California	<b>0.907 <math>\pm</math> 0.004</b>	1.008 $\pm$ 0.005	1.068 $\pm$ 0.005	0.940 $\pm$ 0.001	1.208 $\pm$ 0.018	1.204 $\pm$ 0.009
Census	1.213 $\pm$ 0.003	<b>1.205 <math>\pm</math> 0.007</b>	1.270 $\pm$ 0.007	1.143 $\pm$ 0.002	1.582 $\pm$ 0.018	1.526 $\pm$ 0.038



**Fig. 3.** MZE results for the eight ordinal regression benchmark datasets.



**Fig. 4.** MAE results for the eight ordinal regression benchmark datasets.

test sets. The results are given in Table 4, indicating that our approach is the best in terms of MZE and third in terms of MAE. For

**Table 4**  
Mean ranking for the generalization sets of the ordinal benchmark datasets.

Method	$\bar{R}_{MZE}$	$\bar{R}_{MAE}$
KDLOR	3.94	2.75
SVORIM	4.19	3.25
SVORIEX	3.75	5.00
RED-SVM	3.19	2.88
$p$ -OGMLVQ	3.69	3.88
$\alpha$ -OGMLVQ	2.25	3.25

**Table 5**  
Real ordinal regression datasets used for the experiments ( $n$  is the number of patterns,  $d$  is the number of attributes and  $K$  is the number of classes).

Dataset	$n$	$d$	$K$	Ordered class distribution
Automobile	205	71	6	(3, 22, 67, 54, 32, 27)
Bondrate	57	37	5	(6, 33, 12, 5, 1)
Contact-lenses	24	6	3	(15, 5, 4)
Eucalyptus	736	91	5	(180, 107, 130, 214, 105)
Newthyroid	215	5	3	(30, 150, 35)
Pasture	36	25	3	(12, 12, 12)
Squash-stored	52	51	3	(23, 21, 8)
Squash-unstored	52	52	3	(24, 24, 4)
Tae	151	54	3	(49, 50, 52)
Winequality-red	1599	11	6	(10, 53, 681, 638, 199, 18)

both performance measures the mean rank of  $\alpha$ -OGMLVQ is clearly superior to that of  $p$ -OGMLVQ.

The datasets described in Table 1 are not real ordinal regression tasks but adapted from regression tasks. Thus, we further evaluated our approach on 10 real ordinal regression datasets with imbalanced classes suggested in Sánchez-Monedero et al. (2013), which are briefly described in Table 5. Following the experimental setup in Sánchez-Monedero et al. (2013), 30 different random splits of the datasets were considered, with 75% of the instances for training, the remaining 25% for testing. The data partitions were shared by all the compared methods. The results reported in Tables 6–8 are the average values over 30 runs, together with standard derivations. From the results, we can see that our approach beats the other algorithms once in terms of MZE, twice in terms of MAE, and three times in terms of MMAE. Again, the mean rankings were computed and reported in Table 9. Most

**Table 6**Mean zero-one error (MZE) results along with standard deviations ( $\pm$ ) across 30 runs of the real ordinal datasets. The best performance for each dataset has been boldfaced.

Dataset/Method	KDLOR	SVORIM	SVOREX	RED-SVM	$p$ -OGMLVQ	$\alpha$ -OGMLVQ
Automobile	<b>0.228 <math>\pm</math> 0.058</b>	0.361 $\pm$ 0.076	0.335 $\pm$ 0.068	0.316 $\pm$ 0.055	0.325 $\pm$ 0.099	0.275 $\pm$ 0.055
Bondrate	0.458 $\pm$ 0.087	0.453 $\pm$ 0.092	<b>0.447 <math>\pm</math> 0.096</b>	<b>0.447 <math>\pm</math> 0.073</b>	0.491 $\pm$ 0.090	0.477 $\pm$ 0.112
Contact-lenses	0.411 $\pm$ 0.174	0.367 $\pm$ 0.127	0.350 $\pm$ 0.127	0.300 $\pm$ 0.111	0.233 $\pm$ 0.112	<b>0.222 <math>\pm</math> 0.126</b>
Eucalyptus	0.389 $\pm$ 0.028	0.361 $\pm$ 0.028	0.353 $\pm$ 0.029 $\pm$	<b>0.349 <math>\pm</math> 0.024</b>	0.368 $\pm$ 0.025	0.367 $\pm$ 0.026
Newthyroid	<b>0.028 <math>\pm</math> 0.019</b>	0.031 $\pm$ 0.021	0.033 $\pm$ 0.022	0.031 $\pm$ 0.022	0.055 $\pm$ 0.025	0.043 $\pm$ 0.021
Pasture	0.322 $\pm$ 0.125	0.333 $\pm$ 0.120	0.370 $\pm$ 0.1	0.352 $\pm$ 0.134	<b>0.304 <math>\pm</math> 0.109</b>	0.314 $\pm$ 0.105
Squash-stored	<b>0.297 <math>\pm</math> 0.112</b>	0.361 $\pm$ 0.118	0.372 $\pm$ 0.113	0.336 $\pm$ 0.104	0.394 $\pm$ 0.124	0.389 $\pm$ 0.133
Squash-unstored	<b>0.172 <math>\pm</math> 0.104</b>	0.236 $\pm$ 0.103	0.282 $\pm$ 0.128	0.251 $\pm$ 0.086	0.366 $\pm$ 0.136	0.271 $\pm$ 0.106
Tae	0.445 $\pm$ 0.052	<b>0.410 <math>\pm</math> 0.066</b>	0.419 $\pm$ 0.060	0.478 $\pm$ 0.074	0.440 $\pm$ 0.063	0.428 $\pm$ 0.069
Winequality-red	0.397 $\pm$ 0.017	0.370 $\pm$ 0.022	<b>0.361 <math>\pm</math> 0.022</b>	0.392 $\pm$ 0.022	0.415 $\pm$ 0.016	0.415 $\pm$ 0.024

**Table 7**Mean absolute error (MAE) results along with standard deviations ( $\pm$ ) across 30 runs of the real ordinal datasets. The best performance for each dataset has been boldfaced.

Dataset/Method	KDLOR	SVOR-IMC	SVOR-EXC	RED-SVM	$p$ -OGMLVQ	$\alpha$ -OGMLVQ
Automobile	<b>0.334 <math>\pm</math> 0.076</b>	0.424 $\pm$ 0.090	0.408 $\pm$ 0.089	0.393 $\pm$ 0.073	0.448 $\pm$ 0.122	0.403 $\pm$ 0.104
Bondrate	0.587 $\pm$ 0.107	0.591 $\pm$ 0.102	0.573 $\pm$ 0.121	0.598 $\pm$ 0.008	0.564 $\pm$ 0.121	<b>0.536 <math>\pm</math> 0.137</b>
Contact-lenses	0.539 $\pm$ 0.208	0.506 $\pm$ 0.167	0.489 $\pm$ 0.185	0.378 $\pm$ 0.169	0.338 $\pm$ 0.172	<b>0.327 <math>\pm</math> 0.188</b>
Eucalyptus	0.424 $\pm$ 0.032	0.395 $\pm$ 0.035	0.392 $\pm$ 0.031	<b>0.380 <math>\pm</math> 0.027</b>	0.422 $\pm$ 0.030	0.419 $\pm$ 0.032
Newthyroid	<b>0.028 <math>\pm</math> 0.019</b>	0.031 $\pm$ 0.021	0.033 $\pm$ 0.022	0.032 $\pm$ 0.022	0.055 $\pm$ 0.024	0.043 $\pm$ 0.021
Pasture	0.322 $\pm$ 0.125	0.333 $\pm$ 0.120	0.370 $\pm$ 0.125	0.359 $\pm$ 0.142	<b>0.307 <math>\pm</math> 0.116</b>	0.318 $\pm$ 0.112
Squash-stored	<b>0.308 <math>\pm</math> 0.128</b>	0.372 $\pm$ 0.126	0.382 $\pm$ 0.139	0.346 $\pm$ 0.110	0.415 $\pm$ 0.141	0.410 $\pm$ 0.151
Squash-unstored	<b>0.172 <math>\pm</math> 0.104</b>	0.239 $\pm$ 0.109	0.282 $\pm$ 0.128	0.251 $\pm$ 0.086	0.394 $\pm$ 0.158	0.227 $\pm$ 0.096
Tae	0.473 $\pm$ 0.069	<b>0.461 <math>\pm</math> 0.081</b>	0.485 $\pm$ 0.078	0.515 $\pm$ 0.087	0.551 $\pm$ 0.090	0.536 $\pm$ 0.096
Winequality-red	0.443 $\pm$ 0.019	<b>0.406 <math>\pm</math> 0.022</b>	0.408 $\pm$ 0.023	0.419 $\pm$ 0.021	0.458 $\pm$ 0.020	0.458 $\pm$ 0.022

**Table 8**Mean macro-averaged absolute error (MMAE) results along with standard deviations ( $\pm$ ) across 30 runs of the real ordinal datasets. The best performance for each dataset has been boldfaced.

Dataset/Method	KDLOR	SVOR-IMC	SVOR-EXC	RED-SVM	$p$ -OGMLVQ	$\alpha$ -OGMLVQ
Automobile	<b>0.345 <math>\pm</math> 0.104</b>	0.518 $\pm$ 0.096	0.523 $\pm$ 0.105	0.468 $\pm$ 0.096	0.482 $\pm$ 0.125	0.446 $\pm$ 0.132
Bondrate	1.037 $\pm$ 0.270	1.114 $\pm$ 0.233	1.072 $\pm$ 0.217	1.184 $\pm$ 0.225	0.768 $\pm$ 0.243	<b>0.737 <math>\pm</math> 0.251</b>
Contact-lenses	0.519 $\pm$ 0.280	0.589 $\pm$ 0.259	0.517 $\pm$ 0.303	0.385 $\pm$ 0.198	0.243 $\pm$ 0.166	<b>0.221 <math>\pm</math> 0.182</b>
Eucalyptus	0.426 $\pm$ 0.038	0.420 $\pm$ 0.043	<b>0.411 <math>\pm</math> 0.030</b>	0.414 $\pm$ 0.030	0.450 $\pm$ 0.035	0.477 $\pm$ 0.035
Newthyroid	0.059 $\pm$ 0.040	0.055 $\pm$ 0.042	<b>0.054 <math>\pm</math> 0.042</b>	0.057 $\pm$ 0.049	0.124 $\pm$ 0.056	0.097 $\pm$ 0.051
Pasture	0.322 $\pm$ 0.125	0.333 $\pm$ 0.120	0.370 $\pm$ 0.125	0.359 $\pm$ 0.142	<b>0.307 <math>\pm</math> 0.116</b>	0.318 $\pm$ 0.112
Squash-stored	<b>0.349 <math>\pm</math> 0.156</b>	0.427 $\pm$ 0.148	0.433 $\pm$ 0.172	0.391 $\pm$ 0.149	0.415 $\pm$ 0.171	0.411 $\pm$ 0.181
Squash-unstored	0.309 $\pm$ 0.180	0.367 $\pm$ 0.140	0.426 $\pm$ 0.157	0.348 $\pm$ 0.159	0.488 $\pm$ 0.187	<b>0.228 <math>\pm</math> 0.130</b>
Tae	<b>0.471 <math>\pm</math> 0.070</b>	0.459 $\pm$ 0.081	0.484 $\pm$ 0.079	0.513 $\pm$ 0.086	0.553 $\pm$ 0.091	0.537 $\pm$ 0.096
Winequality-red	1.258 $\pm$ 0.069	1.093 $\pm$ 0.072	1.095 $\pm$ 0.067	<b>1.068 <math>\pm</math> 0.069</b>	1.078 $\pm$ 0.044	<b>1.069 <math>\pm</math> 0.041</b>

**Table 9**

Mean ranking for the generalization sets of the real ordinal datasets.

Method	$\bar{R}_{MZE}$	$\bar{R}_{MAE}$	$\bar{R}_{MMAE}$
KDLOR	3.2000	2.900	3.100
SVORIM	3.150	3.200	3.900
SVORIEX	3.450	3.700	4.100
RED-SVM	3.000	3.300	3.200
$p$ -OGMLVQ	4.550	4.550	3.900
$\alpha$ -OGMLVQ	3.650	3.350	2.800

importantly, our approach has the best mean rank in terms of the most relevant measure MMAE (imbalanced classes). Again, our approach is consistently better than  $p$ -OGMLVQ in all the three performance measures.

Overall, our method retains a competitive edge when compared with the alternative methods considered in this study. Of course, there will be situations where one method is better than the other, but the overall picture emerging from the experiments is that our methods clearly have competitive performance with other methods, while retaining the advantages of prototype based models – e.g. interpretability, natural extension to relevance and metric learning etc. Besides generalization performance, typically there are other criteria to be considered in real world machine learning applications – for example training complexity or model interpretability. In terms of generalization performance, there are many application scenarios where LVQ performs on par with e.g. SVM. Those two approaches are derived in different

frameworks and based on different principles. SVM focuses on boundaries of receptive fields and hence can construct more complex decision boundaries naturally, but at the cost of much more complex, less interpretable models and sufficient data. Often, it is preferable to have simple, interpretable models with generalization performance approaching that of more complex models, while retaining the advantage of constant (as opposed to linear) time classification effort and linear (as opposed to cubic) time training effort. Last but not least, LVQ models are particularly well suited for life-long learning.

## 5. Conclusion

We have proposed a new approach to classify data with ordered categories within the learning vector quantization framework. This approach extends the existing  $p$ -OGMLVQ method to more natural cost function and more intuitive prototype updating rules.

In contrast to  $p$ -OGMLVQ, which pairs the identified correct and incorrect prototypes in distinct couples and then effectively defines a cost for each such couple, our approach defines a single unified cost for all the identified prototypes. Consequently, the updating rule derived from our cost function explicitly reflects the global ordinal relations amongst the prototype classes. Furthermore, we proposed to automatically adapt the bandwidth parameters of the weighting functions during training.

Empirical experiments have been conducted to verify the effectiveness of our ordinal approach. The experimental results show that the performance of our ordinal approach is better than  $p$ -OGMLVQ. We also compared our approach with other ordinal classifiers. The results show that our approach could obtain comparable performance sometimes superior performance, especially in terms of (the most relevant) MMAE performance measure.

### Acknowledgments

Peter Tiño was supported by the EPSRC grant EP/L000296/1. Dr. Fengzhen Tang would like to thank Dr. Bailu Si from Shenyang Institute of Automation, Chinese Academy of Sciences and Dr. Kerstin Bunte from University of Birmingham for the discussion and comments.

### Appendix. Proof of vanishing gradient

In this appendix we will show that the terms (48), (49), (52) and (53) of the gradient of our error function vanish.

Note that derivative of the Heaviside function is the delta function  $\delta$ , which is a symmetric function with  $\delta(x) = 0$  for  $x \neq 0$  and  $\int_{\mathbb{R}} \delta(x) = 1$ . We will first concentrate on terms (48) and (49). The integrand in (48) can be expanded as

$$\sum_{||-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_l} \frac{2\alpha_j d_j \tilde{E}(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} H(g(j, \ell) + f(\ell)) \frac{\partial \mathcal{I}(j, l)}{\partial \mathbf{w}_i} \quad (68)$$

$$+ \sum_{||-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_l} \frac{2\alpha_j d_j \tilde{E}(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} \mathcal{I}(j, l) \delta(g(j, \ell) + f(\ell)) \frac{\partial g(j, \ell)}{\partial \mathbf{w}_i} \quad (69)$$

$$+ \sum_{||-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_l} \frac{2\alpha_j d_j \tilde{E}(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} \mathcal{I}(j, l) \delta(g(j, \ell) + f(\ell)) \frac{\partial f(\ell)}{\partial \mathbf{w}_i}. \quad (70)$$

First, we show that (68) vanishes. Denoting

$$T(j, l, \ell) = \frac{2\alpha_j d_j \tilde{E}(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} \cdot H(g(j, \ell) + f(\ell)) \cdot \delta \left( \sum_{\mathbf{w}_q \in W_l} H(d_q - d_j) - |W_l| \right)$$

by  $T(j, l, \ell)$ , (68) can then be rewritten as

$$\begin{aligned} & \sum_{||-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_l} \left[ \frac{2\alpha_j d_j \tilde{E}(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} \cdot H(g(j, \ell) + f(\ell)) \right. \\ & \left. \delta \left( \sum_{\mathbf{w}_q \in W_{c_i}} H(d_q - d_j) - |W_{c_i}| \right) \right. \\ & \left. \sum_{\mathbf{w}_q \in W_l} \delta(d_q - d_j) \left( \frac{\partial d_q}{\partial \mathbf{w}_i} - \frac{\partial d_j}{\partial \mathbf{w}_i} \right) \right] \\ & = \sum_{||-\ell| \leq \tau} \sum_{\substack{\mathbf{w}_j \in W_l \\ j \neq i}} T(j, l, \ell) \delta(d_i - d_j) 2D_i \\ & + T(i, c_i, \ell) \sum_{||-\ell| \leq \tau} \sum_{\substack{\mathbf{w}_q \in W_l \\ q \neq i}} \delta(d_q - d_i) (-2D_i). \end{aligned}$$

This expression vanishes since  $\delta$  is symmetric and non-vanishing only for  $d_i = d_j$  and  $d_q = d_i$ .

We now show that (69) vanishes. Denoted by  $Q(j, l, \ell)$  the term

$$\frac{2\alpha_j d_j \tilde{E}(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} \cdot \mathcal{I}(j, l) \cdot \delta(g(j, \ell) + f(\ell)),$$

then, (69) yields:

$$\sum_{||-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_l} Q(j, l, \ell) \sum_{|s-\ell| \leq \tau} \sum_{\mathbf{w}_q \in W_s} H(d_q - d_j) \frac{\partial \mathcal{I}(q, s)}{\partial \mathbf{w}_i} \quad (71)$$

$$+ \sum_{||-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_l} Q(j, l, \ell) \times \sum_{|s-\ell| \leq \tau} \sum_{\mathbf{w}_q \in W_s} \mathcal{I}(q, s) \delta(d_q - d_j) \left[ \frac{\partial d_q}{\partial \mathbf{w}_i} - \frac{\partial d_j}{\partial \mathbf{w}_i} \right]. \quad (72)$$

Let

$$A(q, s) = H(d_q - d_j) \cdot \delta \left( \sum_{\mathbf{w}_t \in W_s} H(d_t - d_q) - |W_s| \right),$$

(71) vanishes because

$$\begin{aligned} & \sum_{|s-\ell| \leq \tau} \sum_{\mathbf{w}_q \in W_s} H(d_q - d_j) \frac{\partial \mathcal{I}(q, s)}{\partial \mathbf{w}_i} \\ & = \sum_{|s-\ell| \leq \tau} \sum_{\mathbf{w}_q \in W_s} H(d_q - d_j) \delta \left( \sum_{\mathbf{w}_t \in W_s} H(d_t - d_q) - |W_s| \right) \\ & \quad \sum_{\mathbf{w}_q \in W_s} \delta(d_t - d_q) \left( \frac{\partial d_t}{\partial \mathbf{w}_i} - \frac{\partial d_q}{\partial \mathbf{w}_i} \right) \\ & = \sum_{|s-\ell| \leq \tau} \sum_{\substack{\mathbf{w}_q \in W_s \\ q \neq i}} A(q, s) \delta(d_i - d_q) 2D_i \\ & \quad + A(i, c_i) \sum_{|s-\ell| \leq \tau} \sum_{\substack{\mathbf{w}_t \in W_s \\ t \neq i}} \delta(d_t - d_i) (-2D_i) \end{aligned}$$

vanishes since  $\delta$  is symmetric and non-vanishing only for  $d_i = d_q$  and  $d_t = d_i$ .

(72) will also vanish, since it leads to:

$$\begin{aligned} & \sum_{||-\ell| \leq \tau} \sum_{\substack{\mathbf{w}_j \in W_l \\ j \neq i}} Q(j, l, \ell) \mathcal{I}(q, s) \delta(d_i - d_j) 2D_i \\ & + Q(i, c_i, \ell) \sum_{||-\ell| \leq \tau} \sum_{\substack{\mathbf{w}_q \in W_s \\ q \neq i}} \mathcal{I}(q, s) \delta(d_q - d_i) (-2D_i). \end{aligned}$$

Since both (71) and (72) are vanishing, (69) vanishes.

Finally, we show (70) vanishes. (70) leads to

$$\begin{aligned} & \sum_{||-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_l} Q(j, l, \ell) \left( \frac{\partial \tilde{N}(\ell)}{\partial \mathbf{w}_i} H(N(\ell) - \tilde{N}(\ell)) \right. \\ & \left. + (N(\ell) - \tilde{N}(\ell)) \delta(N(\ell) - \tilde{N}(\ell)) \frac{\partial \tilde{N}(\ell)}{\partial \mathbf{w}_i} \right). \end{aligned}$$

This term vanishes because

$$\frac{\partial \tilde{N}(\ell)}{\partial \mathbf{w}_i} = \sum_{||-\ell| > \tau} \sum_{\mathbf{w}_j \in W_l} \delta(D - d_k) \left( \frac{\partial D}{\partial \mathbf{w}_i} - \frac{\partial d_k}{\partial \mathbf{w}_i} \right) \quad (73)$$

vanishes ( $\delta$  is symmetric and non-vanishing only for  $D = d_k$ ).

In summary, (68), (69), and (70) all vanish and hence the integrand in (48) vanishes as well. We will now show with

analogous arguments that the integrand of (49) vanishes. We have:

$$\sum_{||-\ell|>\tau} \sum_{\mathbf{w}_k \in W_l} \frac{2\alpha_k d_k E(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} H(\tilde{\mathbf{g}}(k, \ell) + \tilde{f}(\ell)) \frac{\partial \tilde{\mathbf{I}}(k, D)}{\partial \mathbf{w}_i} \quad (74)$$

$$+ \sum_{||-\ell|>\tau} \sum_{\mathbf{w}_k \in W_l} \frac{2\alpha_k d_k E(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} \tilde{\mathbf{I}}(k, D) \delta(\tilde{\mathbf{g}}(k, \ell) + \tilde{f}(\ell)) \frac{\partial \tilde{\mathbf{g}}(k, \ell)}{\partial \mathbf{w}_i} \quad (75)$$

$$+ \sum_{||-\ell|>\tau} \sum_{\mathbf{w}_k \in W_l} \frac{2\alpha_k d_k E(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} \tilde{\mathbf{I}}(k, D) \delta(\tilde{\mathbf{g}}(k, \ell) + \tilde{f}(\ell)) \frac{\partial \tilde{f}(\ell)}{\partial \mathbf{w}_i}. \quad (76)$$

As before, we argue that (74), (75), and (76) vanish, separately.

First, (74) vanishes, since

$$\frac{\partial \tilde{\mathbf{I}}(k, D)}{\partial \mathbf{w}_i} = \delta(D - d_k) \left( \frac{\partial D}{\partial \mathbf{w}_i} - \frac{\partial d_k}{\partial \mathbf{w}_i} \right) = 0. \quad (77)$$

This is because  $\delta$  is symmetric and non-vanishing only for  $D = d_k$ . Let us move on to (75). Denoting

$$\frac{2\alpha_k d_k E(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} \tilde{\mathbf{I}}(k, D) \delta(\tilde{\mathbf{g}}(k, \ell) + \tilde{f}(\ell))$$

by  $\tilde{G}(k, \ell)$ , (75) leads to

$$\sum_{||-\ell|>\tau} \sum_{\mathbf{w}_k \in W_l} \tilde{G}(k, \ell) \sum_{||s-\ell|>\tau} \sum_{\mathbf{w}_q \in W_s} H(d_q - d_k) \frac{\partial \tilde{\mathbf{I}}(q, D)}{\partial \mathbf{w}_i} \quad (78)$$

$$+ \sum_{||-\ell|>\tau} \sum_{\mathbf{w}_k \in W_l} \tilde{G}(k, \ell) \times \sum_{||s-\ell|>\tau} \sum_{\mathbf{w}_q \in W_s} \tilde{\mathbf{I}}(q, D) \delta(d_q - d_k) \left( \frac{\partial d_q}{\partial \mathbf{w}_i} - \frac{\partial d_k}{\partial \mathbf{w}_i} \right). \quad (79)$$

Now, (78) vanishes, since we have already shown that  $\frac{\partial \tilde{\mathbf{I}}(q, D)}{\partial \mathbf{w}_i} = 0$ .

The term (79) leads to

$$\sum_{||-\ell|>\tau} \sum_{\substack{\mathbf{w}_k \in W_l \\ k \neq i}} \tilde{G}(k, \ell) \tilde{\mathbf{I}}(i, D) \delta(d_i - d_k) 2D_i + \tilde{G}(i, \ell) \sum_{||s-\ell|>\tau} \sum_{\substack{\mathbf{w}_q \in W_s \\ q \neq i}} \tilde{\mathbf{I}}(q, D) \delta(d_q - d_i) (-2D_i),$$

which vanishes, since  $\delta$  is symmetric and non-vanishing only for  $d_i = d_k$  and  $d_q = d_i$ .

Finally, (76) vanishes as it leads to:

$$\sum_{||-\ell|>\tau} \sum_{\mathbf{w}_k \in W_l} \tilde{G}(k, \ell) \left( -H(\tilde{N}(\ell) - N(\ell)) \frac{\partial \tilde{N}(\ell)}{\partial \mathbf{w}_i} - ((\tilde{N}(\ell) - N(\ell)) \delta(\tilde{N}(\ell) - N(\ell))) \frac{\partial \tilde{N}(\ell)}{\partial \mathbf{w}_i} \right)$$

and this term vanishes, since  $\frac{\partial \tilde{N}(\ell)}{\partial \mathbf{w}_i} = 0$ , as shown in (73).

We now turn our attention to terms (52) and (53):

$$\sum_{||-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_l} \frac{2\alpha_j d_j \tilde{E}(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} H(\mathbf{g}(j, \ell) + f(\ell)) \left( \frac{\partial \mathbf{I}(j, l)}{\partial \Omega} \right) \quad (80)$$

$$+ \sum_{||-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_l} Q(j, l, \ell) \left( \frac{\partial \mathbf{g}(j, \ell)}{\partial \Omega} + \frac{\partial f(\ell)}{\partial \Omega} \right) \quad (81)$$

$$- \sum_{||-\ell|>\tau} \sum_{\mathbf{w}_k \in W_l} \frac{2\alpha_k d_k E(\ell)}{(E(\ell) + \tilde{E}(\ell))^2} H(\tilde{\mathbf{g}}(k, \ell) + \tilde{f}(\ell)) \frac{\partial \tilde{\mathbf{I}}(k, D)}{\partial \Omega} \quad (82)$$

$$- \sum_{||-\ell|>\tau} \sum_{\mathbf{w}_k \in W_l} \tilde{G}(k, \ell) \left( \frac{\partial \tilde{\mathbf{g}}(k, \ell)}{\partial \Omega} + \frac{\partial \tilde{f}(\ell)}{\partial \Omega} \right). \quad (83)$$

The partial derivatives  $\frac{\partial \mathbf{I}(j, l)}{\partial \Omega}$ ,  $\frac{\partial \mathbf{g}(j, \ell)}{\partial \Omega}$ ,  $\frac{\partial f(\ell)}{\partial \Omega}$ ,  $\frac{\partial \tilde{\mathbf{I}}(k, D)}{\partial \Omega}$ ,  $\frac{\partial \tilde{\mathbf{g}}(k, \ell)}{\partial \mathbf{w}_i}$ , and  $\frac{\partial \tilde{f}(\ell)}{\partial \mathbf{w}_i}$  all vanish as

$$\frac{\partial \mathbf{I}(j, l)}{\partial \Omega} = \delta \left( \sum_{\mathbf{w}_q \in W_l} H(d_q - d_j) - |W_l| \right) \cdot \sum_{\mathbf{w}_q \in W_l} \delta(d_q - d_j) \left( \frac{\partial d_q}{\partial \Omega} - \frac{\partial d_j}{\partial \Omega} \right), \quad (84)$$

$$\frac{\partial \tilde{\mathbf{I}}(k, D)}{\partial \Omega} = \delta(D - d_k) \left( \frac{\partial D}{\partial \Omega} - \frac{\partial d_k}{\partial \Omega} \right) \quad (85)$$

and  $\delta$  is symmetric and non-vanishing only for  $d_q = d_j$ , and  $D = d_k$ . We have:

$$\begin{aligned} \frac{\partial \mathbf{g}(j, \ell)}{\partial \Omega} &= \sum_{||s-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_s} \left( \frac{\partial \mathbf{I}(j, l)}{\partial \Omega} H(d_q - d_j) \right. \\ &\quad \left. + \mathbf{I}(q, s) \delta(d_q - d_j) \left( \frac{\partial d_q}{\partial \Omega} - \frac{\partial d_j}{\partial \Omega} \right) \right) \\ &= \sum_{||s-\ell| \leq \tau} \sum_{\mathbf{w}_j \in W_s} \mathbf{I}(q, s) \delta(d_q - d_j) \left( \frac{\partial d_q}{\partial \Omega} - \frac{\partial d_j}{\partial \Omega} \right), \end{aligned}$$

$$\begin{aligned} \frac{\partial \tilde{\mathbf{g}}(k, \ell)}{\partial \Omega} &= \sum_{||s-\ell|>\tau} \sum_{\mathbf{w}_q \in W_s} \left( \frac{\partial \tilde{\mathbf{I}}(k, D)}{\partial \Omega} H(d_q - d_k) \right. \\ &\quad \left. + \tilde{\mathbf{I}}(q, D) \delta(d_q - d_k) \left( \frac{\partial d_q}{\partial \Omega} - \frac{\partial d_k}{\partial \Omega} \right) \right) \\ &= \sum_{||s-\ell|>\tau} \sum_{\mathbf{w}_k \in W_s} \tilde{\mathbf{I}}(q, D) \delta(d_q - d_k) \left( \frac{\partial d_q}{\partial \Omega} - \frac{\partial d_k}{\partial \Omega} \right). \end{aligned}$$

The last two terms vanish since  $\delta$  is symmetric and non-vanishing only for  $d_q = d_j$  and  $d_q = d_k$ .

From

$$\frac{\partial \tilde{N}(\ell)}{\partial \Omega} = \sum_{||s-\ell|>\tau} \sum_{\mathbf{w}_k \in W_s} \frac{\partial \tilde{\mathbf{I}}(k, D)}{\partial \Omega} = 0,$$

it follows that

$$\begin{aligned} \frac{\partial f(\ell)}{\partial \Omega} &= \frac{\partial \tilde{N}(\ell)}{\partial \Omega} H(N(\ell) - \tilde{N}(\ell)) \\ &\quad + (N(\ell) - \tilde{N}(\ell)) \delta(N(\ell) - \tilde{N}(\ell)) \frac{\partial \tilde{N}(\ell)}{\partial \Omega} \\ &= 0, \end{aligned}$$

$$\begin{aligned} \frac{\partial \tilde{f}(\ell)}{\partial \Omega} &= -\frac{\partial \tilde{N}(\ell)}{\partial \Omega} H(\tilde{N}(\ell) - N(\ell)) \\ &\quad - (\tilde{N}(\ell) - N(\ell)) \delta(\tilde{N}(\ell) - N(\ell)) \frac{\partial \tilde{N}(\ell)}{\partial \Omega} \\ &= 0. \end{aligned}$$

## References

- Chu, W., & Ghahramani, Z. (2005). Gaussian process for ordinal regression. *Journal of Machine Learning Research*, 6, 1019–1041.
- Chu, W., & Keerthi, S.S. (2005). New approaches to support vector ordinal regression. In *Proceedings of the 22nd international conference on machine learning (ICML'05)* (pp. 145–152).
- Chu, W., & Keerthi, S. S. (2007). Support vector ordinal regression. *Neural Computation*, 19(3), 792–815.
- Cramer, K., & Singer, Y. (2001). Pranking with ranking. In *Advances in neural information processing systems*, 14 (pp. 641–647). MIT Press.

- Dikkers, H., & Rothkrantz, L. (2012). Support vector machines in ordinal classification: an application to corporate credit scoring. *Neural Network World*, 3(1), 59–70.
- Fouad, S., & Tiño, P. (2012). Adaptive metric learning vector quantization for ordinal classification. *Neural Computation*, 24(11), 2825–2851.
- Frank, E., & Hall, M. (2001). A simple approach to ordinal classification. In *Proceedings of the European conference on machine learning*, vol. 2167 (pp. 145–165).
- Gutiérrez, P. A., & Tiño, P. C. H.-M. (2014). Ordinal regression neural networks based on concentric hyperspheres. *Neural Networks*, 59(0), 51–60.
- Herbrich, R., Graepel, T., & Obermayer, K. (2000). *Large margin rank boundaries for ordinal regression*. (pp. 115–132). Cambridge, MA: MIT Press, Ch. 7.
- Kohonen, T. (1986). *Learning vector quantization for pattern recognition*, Technical report tkkf-a601. Espoo, Finland: Helsinki University of Technology.
- Kohonen, T. (1995). *The handbook of brain theory and neural networks*. (pp. 537–540). MIT Press, Ch. Learning vector quantization.
- Kramer, S., Widmer, G., Pfahringer, B., & de Groeve, M. (2001). Prediction of ordinal classes using regression trees. *Fundamenta Informaticae*, 47(1–2), 1–13.
- Lall, R., Campbell, M. J., Walters, S. J., & Morgan, K. (2012). A review of ordinal regression models applied on health-related quality of life assessments. *Statistic Methods Medical Research*, 11, 49–67.
- Li, L., & Lin, H. (2006). Ordinal regression by extended binary classification. In *Advances in neural information processing system* (pp. 865–872).
- Lin, H.-T., & Li, L. (2012). Reduction from cost-sensitive ordinal ranking to weighted binary classification. *Neural Computation*, 24(5), 1329–1367.
- Sánchez-Monedero, J., Gutiérrez, P. A., Tiño, P., & Hervás-Martínez, C. (2013). Exploitation of pairwise class distances for ordinal classification. *Neural Computation*, 25(9), 2450–2485.
- Sato, A., & Yamada, K. (1996). Generalized learning vector quantization. In D. Touretzky, M. Mozer, & M. Hasselmo (Eds.), *Advances in neural information processing system*, 8 (pp. 423–429). MIT Press.
- Schneider, P., Biehl, P., & Hammer, B. (2009). Adaptive relevance matrices in learning vector quantization. *Neural Computation*, 21(12), 3532–3561.
- Schneider, P., Biehl, M., & Hammer, B. (2010). Hyperparameter learning in probabilistic prototype-based models. *Neurocomputing*, 73(7–9), 1117–1124.
- Seah, C.-W., Tsang, I. W., & Ong, Y.-S. (2012). Transductive ordinal regression. *IEEE Transactions on Neural Networks and Learning Systems*, 23(7), 1074–1086.
- Shashua, A., & Levin, A. (2003). Ranking with large margin principle: Two approaches. In *Advances in neural information processing system*, 15 (pp. 937–944).
- Sun, B.-Y., Li, J., Wu, D. D., Zhang, X.-M., & Li, W. (2010). Kernel discriminant learning for ordinal regression. *IEEE Transaction on Knowledge and Data Engineering*, 22(6), 906–910.
- Tino, P., Schittenkopf, C., & Dorffner, G. (2001). Volatility trading via temporal pattern recognition in quantized financial time series. *Pattern Analysis and Application*, 4(4), 283–299.
- Villmann, B., Hammer (2002). Generalized relevance learning vector quantization. *Neural Networks*, 15(8–9), 1059–1068.
- Yan, H. (2014). Cost-sensitive ordinal regression for fully automatic facial beauty assessment. *Neurocomputing*, 129, 334–342.