Computational intelligence & modeling of crop disease data in Africa

Owomugisha, Godliver

*DOI:*
[10.33612/diss.130773079](10.33612/diss.130773079)

**IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.**

*Document Version*
Publisher's PDF, also known as Version of record

[Link to publication in University of Groningen/UMCG research database](#)

# Chapter 2
## Learning Vector Quantization

**Abstract**

*This chapter serves as a general introduction to the concept of prototype based classification. Under the family of Learning Vector Quantization (LVQ), we introduce different variants we used in our study. In our findings, Generalized Matrix Learning Vector Quantization (GMLVQ) was outstanding in handling both classification and the weighting or selection of features. By name, the classification task of this learning model is based on highly ranked features. In this chapter, we explain the working of this algorithm.*

## 2.1   Introduction

The previous chapter introduced supervised learning, a major branch of machine learning where each example is a pair consisting of an input object and a desired output value. A supervised learning algorithm analyzes the training data and produces an inferred function, which can be used for mapping new examples. An optimal scenario will allow for the algorithm to correctly determine the class labels for unseen instances.

A wide range of supervised learning algorithms are available but the choice of using one algorithm over the other will depend on the problem at hand.

Learning Vector Quantization (LVQ) is a prototype-based supervised classification algorithm and a special case of an artificial neural network, more precisely that applies a winner-take-all (Maass 2000) approach. The classification algorithms was first introduced in (Kohonen 1986) and since then various modifications have been suggested in the literature. LVQ and its variants are a family of classifiers that have found much popularity within the machine learning field for various reasons. For review of relevance learning, LVQ and prototype-based systems in general see, for instance (Biehl et al. 2016).

In this thesis we use a special kind of dataset that is attracting more scientists to investigate crop disease using spectral data. By nature of spectral data, it presents us with high dimensions of the input space and this became our critical consideration. The input feature vectors have very high dimension, the learning problem can be difficult even if the true function only depends on a small number of those features. Hence, choosing a classifier that handles the aspect of feature selection and dimensionality reduction of the input space was of high importance.

Similar to the LVQ family is the $k$-Nearest Neighbour algorithm (Cover and Hart 1967) which skips the training phase and directly uses all the available training data in the classification of a new data point. The algorithm is based on the intuitive "Birds of the same feather flock together" meaning that similar things exist in close proximity. To classify a new data point, $k$ of its closest neighbours are obtained using some measure of closeness, for example, Euclidean distance, and the average class of these $k$ neighbours is awarded to the new data point. However, the algorithm gets significantly slower as the number of training data points $N$ increases since the distance to each point needs to be calculated every time.

Several LVQ variants were proposed e.g LVQ1 and LVQ2.1 by Kohonen aiming at faster convergence and better approximation of Bayesian decision boundaries, respectively. Other variants have been suggested in (Sato and Yamada 1995, Schneider et al. 2007).

In relation to our current study, successful studies were conducted, e.g. (Mwebaze

et al. 2011) developed the corresponding divergence-based LVQ (DLVQ) schemes applied to cassava dataset to identify cassava diseases. A closely related study also in (Mwebaze and Biehl 2016) uses prototype-based classification by exploring the right configuration of type of algorithm and type of features extracted from the leaves to optimally diagnose crop disease in cassava plants. More studies on the use of LVQ schemes are highlighted in different chapters in this thesis.

## 2.2 Learning Vector Quantization and its variants

In this section, we give a description of LVQ variants we investigated on including: Generalized LVQ, Generalized Relevance LVQ and Generalized Matrix LVQ that was very applicable in answering some of our research questions. By definition, we will consider data sets of the form:

$$\{x^{\mu}, y^{\mu}\}_{\mu=1}^{P} \tag{2.1}$$

where $x^{\mu} \in \mathbb{R}^{N}$ are feature vectors and the labels $y^{\mu} \in 1, 2, ...C$ specify their class membership.

The LVQ system is defined by a set of $M$ prototype vectors $w^{j} \in \mathbb{R}^{N}$ which carry labels $c(w^{j}) \in \{1, 2, ...C\}$ such that $W = \{w^{j}, c(w^{j})\}_{j=1}^{M}$.

The system can be set up with one or more prototype vectors per class. Prototype vectors are identified in the feature space and ideally serve as typical representatives of their classes.

A nearest prototype classifier (NPC) assigns a given feature vector $x \in \mathbb{R}^{N}$ to the *closest* prototype with respect to some meaningful distance measure.

Most frequently, standard Euclidean distance $d(w, x)$ is employed. The corresponding NPC assigns x to the class $c(w^{J})$ of the closest prototype with $d(x, w^{J}) \leqslant d(x, w^{j})$ for all $j$.

### 2.2.1 Classical LVQ

This is the a heuristic algorithm introduced in (Kohonen 1986). Prototypes are updated based on how close they are to a presented data point given the class of the prototype and that of the data point. The training process is represented by the following steps:.

1. Randomly select a training sample $(\{x^{\mu}, y^{\mu}\})$

2. Determine the winning prototype $(w^{J})$ with $d(x, w^{J}) \leqslant d(x, w^{j}), j = 1......M$

3. Update $(w^{J})$ according to:

$$w^J \longleftarrow w^J + \eta \cdot (x - w^J), \quad (\text{if } c(w^J) = y).$$
$$w^J \longleftarrow w^J - \eta \cdot (x - w^J), \quad (\text{if } c(w^J) \neq y)$$

Should the same feature vector be observed again, $w^J$ will give rise to a decreased (increased) distance, if the label $c(w^J)$ and $y$ agree (disagree).

### 2.2.2 Generalized LVQ

Generalized LVQ (GLVQ) is one key variant of LVQ that was introduced in (Sato and Yamada 1995). The algorithm uses an objective (cost) function in the training of the LVQ system. The advantage of an objective function based LVQ system is that one can use gradient methods (online or batch) to optimize it. With the training data in form $\{x^\mu, y^\mu\}_{\mu=1}^P$, the cost function is defined by:

$$E(W) = \sum_{\mu=1}^{p} \Phi\left(\frac{d(x^\mu, w^J) - d(x^\mu, w^K)}{d(x^\mu, w^J) + d(x^\mu, w^K)}\right) \tag{2.2}$$

where $w^J$ denotes the closest correct prototype with $c(w^J) = y^\mu$ and $w^K$ is the closest incorrect prototype with $c(w^K) \neq y^\mu$, also termed as winner-takes all rule (Crammer et al. 2002). The logistic sigmoid function $\Phi$ determines the active regions of the algorithm: $\Phi(x) = x$. Training constitutes the minimization of $E(W)$ with respect to the model parameters. The learning algorithm is defined in (Sato and Yamada 1995) in terms stochastic gradient descent.

### 2.2.3 Generalized Matrix LVQ

The cost function in Eq.(2.2) has been extended to other distance metrics. A generalization bound has been derived for GRLVQ which uses an adaptive metric (Hammer and Villmann 2002, Strickert et al. 2001, Villmann et al. 2015). The distance metric is defined as the squared weighted euclidean metric $d^\lambda(x, w) = \sum_i^N \lambda_i(x_i - w_i)^2$ where $\lambda_i \geqslant 0$ and $\sum_i \lambda_i = 1$.

Another powerful extension of the basic LVQ concept is Generalized Matrix LVQ (GMLVQ). The learning algorithm can be seen by the following steps:

1. Randomly select a training sample $(\{x^\mu, y^\mu\})$

2. Determine the closest correct prototype $(w^J)$ with $c(w^J) = y$, and $d^\Lambda(x, w^J) \leqslant d^\Lambda(x, w^j)$ for all $w^j$ with $c(w^j) = y$, and the closest incorrect prototype, $(w^K)$ with $c(w^K) \neq y$, and $d^\Lambda(x, w^K) \leqslant d^\Lambda(x, w^j)$, for all $w^j$ with $c(w^j) \neq y$

**3.** Update (w) according to:

$$w^L \longleftarrow w^J + \epsilon \cdot \Lambda \cdot (x - w^J), \quad (y = c(w^J))$$
$$w^K \longleftarrow w^K - \epsilon \cdot \Lambda \cdot (x - w^K), \quad (y \neq c(w^K))$$

**4.** Update $\Omega$ according to:

$$\Omega \longleftarrow \Omega + \epsilon \cdot \Omega \cdot (x - w^J)(x - w^J)^\top, \quad (\text{if } (y = c(w^J)))$$
$$\Omega \longleftarrow \Omega - \epsilon \cdot \Omega \cdot (x - w^K)(x - w^K)^\top, \quad (\text{if } (y \neq c(w^K)))$$

These processes are followed by a normalization step such that Tr(. . .) = 1 for $\Lambda = \Omega^\top \Omega$. $\epsilon \in [0, 1]$ is the learning rate for the metric parameter and the matrix $\Lambda$ is updated in a such way that the distance $d^\Lambda(x, w^J)$ is decreased in case of a correct classification, while $d^\Lambda(x, w^K)$ increases, if the sample $(x, y)$ is misclassified.

GMLVQ algorithm proposed in (Schneider et al. 2007, Schneider et al. 2009b) employs a full matrix $\Lambda \in \mathbb{R}^{N \times N}$ of relevances that describes the importance of the individual features in the classification task. Here, the distance measure $d^\Lambda(x, w)$ is defined as:

$$d^\Lambda(x, w) = (x - w)^\top \Lambda (x - w) \tag{2.3}$$

where the parameterization $\Lambda = \Omega^\top \Omega$ guarantees that $\Lambda$ is positive semi-definite and that $d^\Lambda(x, w) \geqslant 0$ for arbitrary matrices $\Omega \in \mathbb{R}^{N \times N}$. Now the squared distance reads:

$$d^\Lambda(x, w) = \sum_{ijk} (x_i - w_i) \Omega_{ki} \Omega_{kj} (x_j - w_j). \tag{2.4}$$

To obtain the adaptation formula, we compute the derivatives with respect to w and $\Omega$. The derivative of $d^\Lambda$ with respect to w yields:

$$\nabla_w d^\Lambda = -2\Lambda(x - w) = -2\Omega^\top \Omega(x - w). \tag{2.5}$$

Derivative with respect to single element $\Omega_l m$ gives

$$\frac{\partial d^\Lambda}{\partial \Omega_{lm}} = \sum_j (x_l - w_l) \Omega_{mj} (x_j - w_j) + \sum_i (x_i - w_i) \Omega_{il} (x_m - w_m)$$
$$= (x_l - w_l)[\Omega(x - w)]_m + (x_m - w_m)[\Omega(x - w)]_l \tag{2.6}$$

Thus, we get the update equations:

$$\Delta w_J = \epsilon \cdot \Phi'(\mu(x)) \cdot \mu^+(x) \cdot \Omega^\top \Omega \cdot (x - w_J)$$
$$\Delta w_K = -\epsilon \cdot \Phi'(\mu(x)) \cdot \mu^-(x) \cdot \Omega^\top \Omega \cdot (x - w_K) \tag{2.7}$$

Where

$$\mu^+ = \frac{2d^\Lambda(x - w^K)}{(d^\Lambda(x - w^K) + d^\Lambda(x - w^J))^2}$$

$$\mu^- = \frac{2d^\Lambda(x - w^J)}{(d^\Lambda(x - w^J) + d^\Lambda(x - w^K))^2}$$

These updates correspond to the standard Hebb terms of LVQ, pushing the closest correct prototype ($w_J$) towards the considered data point and the closest wrong prototype ($w_K$) away from the considered data point. For the update of the matrix elements $\Omega_{lm}$ we get

$$\Delta \Omega_{lm} = -\epsilon \cdot \Phi'(\mu(x)) \cdot \left( \mu^+(x) \cdot \left( [\Omega(x - w_J)]_m (x_l - w_{J,l}) + [\Omega(x - w_J)]_l (x_m - w_{J,m}) \right) \right.$$

$$\left. - \mu^-(x) \cdot \left( [\Omega(x - w_K)]_m (x_l - w_{K,l}) + [\Omega(x - w_K)]_l (x_m - w_{K,m}) \right) \right) \tag{2.8}$$

The learning rate for the metric can be chosen independently of the learning rate for the prototypes. After each update, $\Lambda$ is normalised to prevent the algorithm from degeneration. Therefore, we set $\sum_i \Lambda_{ii} = \sum_{ij} \Omega_{ij}^2 = 1$ which fixes the sum of diagonal elements and, here, the sum of eigenvalues. Learning stops after a maximal number of epochs are reached.

At the end of the learning process the algorithm provides a set of prototypes $w_j$, their labels $c(w_j)$, and a task specific discriminative distance $d^\Lambda$. If features have the same magnitude, the diagonal elements $\Lambda_{ii}$ of the dissimilarity matrix can be interpreted as overall relevances of every feature $i$ for the classification. The off-diagonal elements $\Lambda_{ij}$ with $j \neq i$ weigh the pairwise correlations between features $i$ and $j$. High absolute values in the matrix denote highly relevant features, while values near zero can be seen as less important for the classification accuracy.

The above description corresponds to updates based on stochastic gradient descent described in (Schneider et al. 2009a). In our experiments, we employed a batch gradient descent algorithm which uses an automatic step size adaption scheme aiming at minimizing the cost function in Eq (2.2). This form of update is described in (Papari et al. 2011). Specifically, we have used a publicly available implementation of batch GMLVQ, see (Biehl 2017).

# Part I

# Disease Diagnosis with Leaf Images