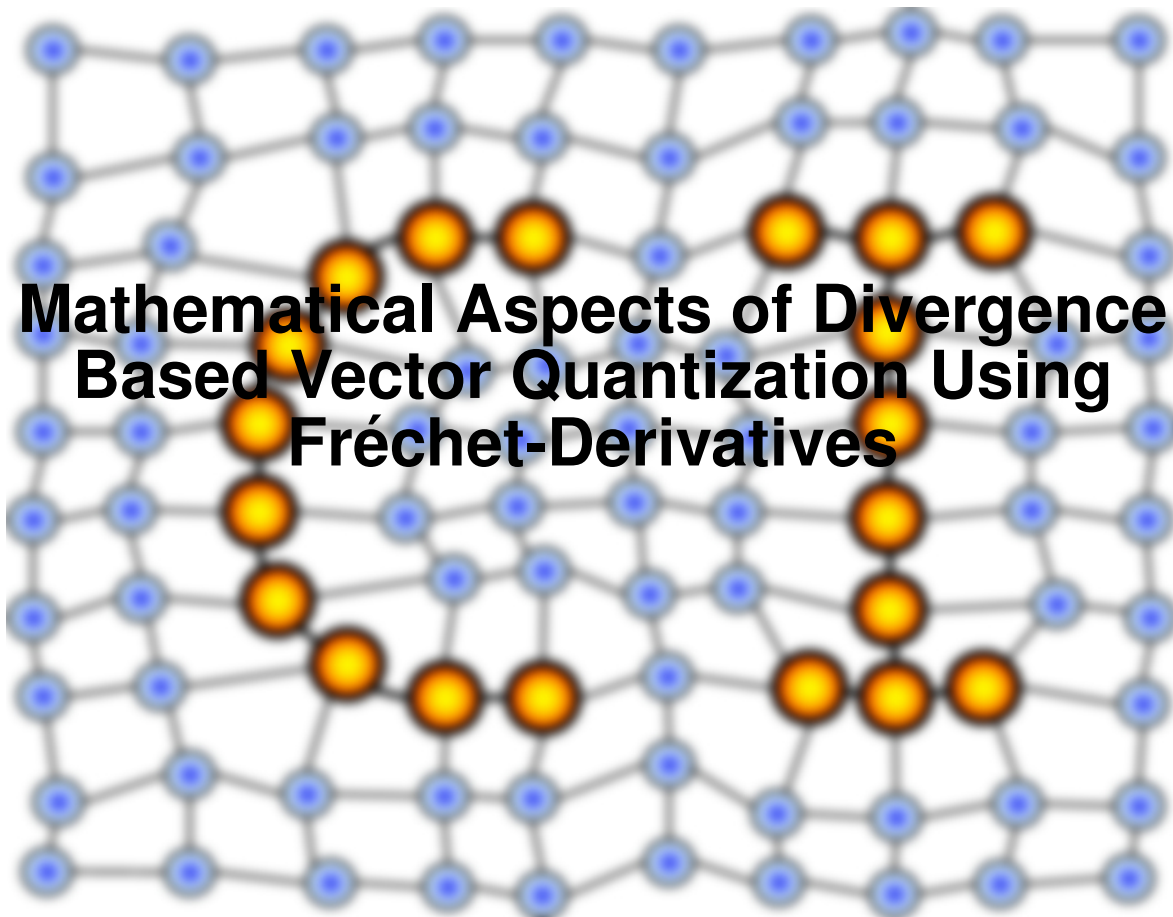


MACHINE LEARNING REPORTS



Mathematical Aspects of Divergence Based Vector Quantization Using Fréchet-Derivatives

Report 03/2009

Submitted: 10.11.2009

Published: 11.11.2009

T. Villmann and S. Haase

University of Applied Sciences Mittweida, Department of MPI

Technikumplatz 17, 09648 Mittweida - Germany

Abstract

Supervised and unsupervised vector quantization methods for classification and clustering traditionally use dissimilarities, frequently taken as Euclidean distances. In this article we investigate the applicability of divergences instead. We deduce the mathematical fundamentals for its utilization in derivative based vector quantization algorithms. It bears on the generalized derivatives known as Fréchet-derivatives. We exemplarily show the application of this methodology for widely applied supervised and unsupervised vector quantization schemes including self-organizing maps, neural gas, and learning vector quantization. Further we show principles for hyperparameter optimization for parametrized divergences in the case of supervised vector quantization to achieve an improved classification accuracy.

1 Introduction

Supervised and unsupervised vector quantization for classification and clustering is strongly associated with the concept of dissimilarity usually judged in terms of distances. The most common choice is the Euclidean metric. Yet, in the last years alternative dissimilarity measures became attractive for advanced data processing. Examples are functional metrics like *Sobolev-distances* or kernel based dissimilarity measures [49],[29]. These metrics take the functional structure of the data into account [28],[36],[39],[46].

Recently, information theory based approaches are proposed considering *divergences* for clustering [2], [22],[30],[16]. For other data processing methods like multi-dimensional scaling (MDS) [27], stochastic neighbor embedding [?], blind source separation [34] or non-negative matrix factorization [6], also divergence based approaches are introduced. In prototype based classification, first approaches utilizing information theoretic approaches were recently proposed [11],[44],[48].

Yet, a systematic analysis of prototype based clustering and classification relying on divergences is not given so far. Further, the respective existing approaches usually are carry out in the so-called batch mode for optimization but *are not available for online learning*. The latter method requires the calculation of the derivatives of the underlying metrics, i.e. divergences here.

In the present contribution we offer a systematic approach for divergence based vector quantization using divergence derivatives. For this purpose, important but general classes of divergences are identified, widely following and extending the scheme introduced by CICHOCKI ET AL. in [7]. The mathematical framework for functional derivatives of divergences is given by the functional-analytic generalization of usual derivatives – the concept of *Fréchet-derivatives* [13],[23].

After characterization of the different classes of divergences and a short introduction of Fréchet-derivatives we apply this framework to the several divergences obtaining generalized derivatives, which can be used for online learning of divergence based methods for supervised and unsupervised vector quantization as well as other gradient based approaches. We explore explicitly for prominent examples the respective derivatives. Thereafter, we exemplarily consider some of the most prominent approaches for unsupervised as well as supervised prototype based vector quantization in the light of divergence based online learning using Fréchet-derivatives. For the latter approaches we additionally provide a gradient learning scheme, called hyperparameter adaptation, for optimization of parameters occurring in parametrized divergences.

2 Characterization of divergences

According to the classification given in CICHOCKI ET AL. [7] one can distinguish at least *three* main classes of divergences, the *Bregman*-divergences, the *Csiszár's* f -divergences and the γ -divergences.

We emphasize at this point that we generally assume that p and ρ are positive measure (densities), not necessarily be normalized. In case of normalized densities we explicitly refer to these as probability densities.

2.1 Bregman divergences

Bregman divergences are defined by generating convex functions Φ in the following way [5]:

Let Φ be a strictly convex real-valued function with the domain \mathcal{L} (the Lebesgue-integrable functions). Further, Φ is assumed to be twice continuously Fréchet-differentiable [23]. A Bregman divergence is defined as $D_{\Phi}^B : \mathcal{L} \times \mathcal{L} \longrightarrow \mathbb{R}^+$ with

$$D_{\Phi}^B(p||\rho) = \Phi(p) - \Phi(\rho) - \frac{\delta\Phi(\rho)}{\delta\rho}(p - \rho) \quad (2.1)$$

whereby $\frac{\delta\Phi(\rho)}{\delta\rho}$ is the Fréchet-derivative of Φ with respect to ρ (see sec. 3.1.1).

Examples

1. *generalized Kullback-Leibler-divergence* for non-normalized p and ρ [7]:

$$D_{GKL}(p||\rho) = \int p(\mathbf{x}) \log\left(\frac{p(\mathbf{x})}{\rho(\mathbf{x})}\right) d\mathbf{x} - \int p(\mathbf{x}) - \rho(\mathbf{x}) d\mathbf{x} \quad (2.2)$$

with the generating function

$$\Phi(f) = \int f \cdot \log f - f d\mathbf{x}.$$

If p and ρ are normalized densities (probability densities) $D_{GKL}(p||\rho)$ is reduced to the usual Kullback-Leibler-divergence [26],[24]:

$$D_{KL}(p||\rho) = \int p(\mathbf{x}) \log\left(\frac{p(\mathbf{x})}{\rho(\mathbf{x})}\right) d\mathbf{x}, \quad (2.3)$$

which is related to the *Shannon-entropy* [42]

$$H_S(p) = - \int p(\mathbf{x}) \log(p(\mathbf{x})) d\mathbf{x} \quad (2.4)$$

via

$$D_{KL}(p||\rho) = V_S(p, \rho) - H_S(p)$$

where

$$V_S(p, \rho) = - \int p(\mathbf{x}) \log(\rho(\mathbf{x})) d\mathbf{x}$$

is *Shannon's cross entropy*.

2. *Itakura-Saito*-divergence [21]:

$$D_{IS}(p||\rho) = \int \left[\frac{p}{\rho} - \log \left(\frac{p}{\rho} \right) - 1 \right] d\mathbf{x} \quad (2.5)$$

with the generating function

$$\Phi(f) = - \int \log(f) d\mathbf{x} .$$

If we assume that p and ρ are general densities then an important subset of Bregman divergences belong to the class of β -divergences [10], which are defined, following CICHOCKI ET AL., as

$$D_{\beta}(p||\rho) = \int p \cdot \frac{p^{\beta-1} - \rho^{\beta-1}}{\beta - 1} d\mathbf{x} - \int \frac{p^{\beta} - \rho^{\beta}}{\beta} d\mathbf{x} \quad (2.6)$$

with $\beta \neq 1$ and $\beta \neq 0$. In the limit $\beta \rightarrow 1$ the divergence $D_{\beta}(p, \rho)$ becomes the generalized Kullback-Leibler-divergence (2.2)¹. The limit $\beta \rightarrow 0$ gives the Itakura-Saito-divergence (2.5). Further, β -divergences are related to the density power divergences \widehat{D}_{β} introduced in [3] by

$$\widehat{D}_{\beta}(p||\rho) = \frac{1}{(1 + \beta)} D_{\beta}(p||\rho) .$$

2.2 Csiszár's f -divergences

Given a *convex* function $f : [0, \infty) \rightarrow \mathbb{R}$ with $f(1) = 0$ (without loss of generality). The f -divergences D_f for general densities p and ρ are given by

$$D_f(p||\rho) = \int \rho(\mathbf{x}) \cdot f \left(\frac{p(\mathbf{x})}{\rho(\mathbf{x})} \right) d\mathbf{x} \quad (2.7)$$

with the definitions $0 \cdot f \left(\frac{0}{0} \right) = 0$, $0 \cdot f \left(\frac{a}{0} \right) = \lim_{x \rightarrow 0} x \cdot f \left(\frac{a}{x} \right) = \lim_{u \rightarrow \infty} a \cdot \frac{f(u)}{u}$ [9],[31],[43]. It corresponds to a *generalized* f -entropy [7] of the form

$$H_f(p) = - \int f(p) d\mathbf{x} . \quad (2.8)$$

Yet, CICHOCKI ET AL. also suggested a generalization, we refer as *generalized* f -divergence [7]. In that divergence, f has not longer to be convex. It is proposed to be

$$D_f^G(p||\rho) = c_f \int p - \rho d\mathbf{x} + \int \rho(\mathbf{x}) \cdot f \left(\frac{p(\mathbf{x})}{\rho(\mathbf{x})} \right) d\mathbf{x} \quad (2.9)$$

with $c_f = f'(1) \neq 0$. Hence, in case of p and ρ being probability densities, the first term vanishes such that the usual form of f -divergences is obtained but without the

¹We remark here that the relations $\frac{p^{\gamma} - \rho^{\gamma}}{\gamma} \xrightarrow{\gamma \rightarrow 0} \log \frac{p}{\rho}$ and $\frac{p^{\gamma} - 1}{\gamma} \xrightarrow{\gamma \rightarrow 0} \log p$ hold.

convexity assumption on f . Thus, as a famous example the *Hellinger divergence* [31],[43]:

$$D_H(p||\rho) = \int (\sqrt{p} - \sqrt{\rho})^2 dx \quad (2.10)$$

with the generating function $f(u) = (\sqrt{u} - 1)^2$ with $u = \frac{p}{\rho}$. We remark, $D_H(p||\rho)$ is not a f -divergence for general densities p and ρ because f is not convex in that case, whereas for probability densities it is a f -divergence according to the *Cichocki-f-divergence* properties [7].

As the β -divergences in case of Bregman divergences, one can identify here an important subset of the f -divergences, the so-called α -divergences according to the definition given in [7]:

$$D_\alpha(p||\rho) = \frac{1}{\alpha(\alpha-1)} \int [p^\alpha \rho^{1-\alpha} - \alpha \cdot p + (\alpha-1)\rho] dx \quad (2.11)$$

with the generating f -function

$$f(u) = u \frac{(u^{\alpha-1} - 1)}{\alpha^2 - \alpha} + \frac{1-u}{\alpha}$$

and $u = \frac{p}{\rho}$. In the limit $\alpha \rightarrow 1$ the generalized Kullback-Leibler-divergence D_{GKL} (2.2) is obtained. Further, in [7] is stated that the β -divergences can be generated from the α -divergences applying the non-linear transforms

$$p \rightarrow p^{\beta+2} \text{ and } \rho \rightarrow \rho^{\beta+2} .$$

The α -divergences are closely related to the generalized *Rényi-divergences* [1],[7]:

$$D_\alpha^{GR}(p||\rho) = \frac{1}{\alpha-1} \log \left(\int [p^\alpha \rho^{1-\alpha} - \alpha \cdot p + (\alpha-1)\rho + 1] dx \right) \quad (2.12)$$

for non-normalized ρ and p , whereas for normalized densities the usual *Rényi-divergence* [37],[38]

$$D_\alpha^R(p||\rho) = \frac{1}{\alpha-1} \log \left(\int p^\alpha \rho^{1-\alpha} dx \right) \quad (2.13)$$

is obtained². The divergence $D_\alpha^R(p||\rho)$ is based on the *Rényi-entropy*

$$H_\alpha(p) = -\frac{1}{\alpha-1} \log \left(\int p^\alpha dx \right) . \quad (2.14)$$

²Notify that a careful transformation of the parameter α is required for exact transformations between both divergences. For details see [1] p.84ff and [7] p.104.

2.3 γ -divergences

A class of very robust divergences with respect to outliers has been proposed by FUJISAWA&EGUCHI [14]³. It is called γ -divergences defined for un-normalized ρ and p as

$$D_\gamma(p||\rho) = \log \left[\frac{\left(\int p^{\gamma+1} d\mathbf{x} \right)^{\frac{1}{\gamma(\gamma+1)}} \cdot \left(\int \rho^{\gamma+1} d\mathbf{x} \right)^{\frac{1}{\gamma+1}}}{\left(\int p \cdot \rho^\gamma d\mathbf{x} \right)^{\frac{1}{\gamma}}} \right] \quad (2.15)$$

$$\begin{aligned} &= \frac{1}{\gamma+1} \log \left[\left(\int p^{\gamma+1} d\mathbf{x} \right)^{\frac{1}{\gamma}} \cdot \left(\int \rho^{\gamma+1} d\mathbf{x} \right) \right] \\ &\quad - \log \left[\left(\int p \cdot \rho^\gamma d\mathbf{x} \right)^{\frac{1}{\gamma}} \right]. \end{aligned} \quad (2.16)$$

In the limit $\gamma \rightarrow 0$ $D_\gamma(\rho||p)$ becomes the usual Kullback-Leibler-divergence (2.3) $D_{KL}(\hat{\rho}||\hat{p})$ with normalized densities

$$\hat{\rho} = \frac{\rho}{\int \rho d\mathbf{x}} \text{ and } \hat{p} = \frac{p}{\int p d\mathbf{x}}.$$

For $\gamma = 1$ the γ -divergence becomes the *Cauchy-Schwarz-divergence*

$$D_{CS}(p||\rho) = \frac{1}{2} \log \left(\int \rho^2(\mathbf{x}) d\mathbf{x} \cdot \int p^2(\mathbf{x}) d\mathbf{x} \right) - \log(V(p, \rho)) \quad (2.17)$$

with

$$V(p, \rho) = \int p(\mathbf{x}) \cdot \rho(\mathbf{x}) d\mathbf{x} \quad (2.18)$$

being the *cross correlation potential*. The *Cauchy-Schwarz-divergence* $D_{CS}(p||\rho)$ was introduced by J. PRINCIPE considering the Cauchy-Schwarz-inequality for norms [35].

3 Derivatives of divergences – a functional analytic approach

In this section we provide the mathematical formalism of generalized derivatives for functionals p and ρ . It is known as *Fréchet-derivatives* or *functional derivatives*. First, we briefly reconsider the theory of functional derivatives including Fréchet- and Gâteaux-derivatives and its relation to directional derivatives. Thereafter we investigate the above divergence classes within this framework. In particular, we explain their Fréchet-derivatives.

³The divergence $D_\gamma(p||\rho)$ is proposed to be robust for $\gamma \in [0, 1]$ with existence of $D_{\gamma=0}$ in the limit $\gamma \rightarrow 0$. A detailed analysis of robustness is given in [14].

3.1 Functional derivatives

3.1.1 Fréchet-derivatives and Fréchet-derivatives of a functional

Suppose, X and Y are *Banach spaces*, $U \subset X$ is open and $F : X \rightarrow Y$. F is called *Fréchet differentiable* at $x \in X$, if there exists a *bounded linear operator* $\frac{F[x]}{\delta x} : X \rightarrow Y$, such that for $h \in X$ the limit

$$\lim_{h \rightarrow 0} \frac{\left\| F(u+h) - F(u) - \frac{F[u]}{\delta u} [h] \right\|_Y}{\|h\|_X} = 0.$$

This general definition can be focussed for functional mapping: Let L be a functional mapping from a linear, functional *Banach-space* B to \mathbb{R} . Further, let B be equipped with a norm $\|\cdot\|$, and $f, h \in B$ are two functionals. The *Fréchet-derivative* of L at point f is formally defined as

$$\lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} (L[f + \varepsilon h] - L[f]) =: \frac{L[f]}{\delta f} [h].$$

The existence and continuity of the limes is equivalent to the existence and continuity of the derivative. For a detailed introduction we refer to [23].

Yet, we recall two main properties of the Fréchet-derivative for functionals, which are important in studying divergences: First, if L is linear then

$$L[f + \varepsilon h] - L[f] = \varepsilon L[h]$$

and, hence, $\frac{L[f]}{\delta f} [h] = L$. Further, an analogon of the chain rule known from usual differential calculus can be stated:

Let $F : \mathbb{R} \rightarrow \mathbb{R}$ be a continuously-differentiable mapping. We consider the functional

$$L[f] = \int F(f(x)) dx$$

Then the Fréchet-derivative is obtained as $\frac{L[f]}{\delta f} = F'(f)$, which can be seen from

$$\begin{aligned} \frac{1}{\varepsilon} (L[f + \varepsilon h] - L[f]) &= \frac{1}{\varepsilon} \int F(f(x) + \varepsilon h(x)) - F(f(x)) dx \\ &= \frac{1}{\varepsilon} \int F'(f(x)) \cdot \varepsilon h(x) + \mathcal{O}(\varepsilon^2 h(x)^2) dx \\ &\rightarrow_{\varepsilon \rightarrow 0} \int F'(f(x)) \cdot h(x) dx \end{aligned}$$

and utilization of the linear property of the integral operator.

Last but not least we state the following important remark: The Fréchet derivative in finite-dimensional spaces is the usual derivative. In particular, it is represented in coordinates by the *Jacobi matrix*. Thus, the Fréchet derivative is a generalization of the *directional derivatives*.

3.1.2 Gâteaux-derivatives

Gâteaux-derivatives are also a generalization of the concept of directional derivatives and can be seen in the middle between Fréchet-derivatives and usual derivatives.

Suppose, X and Y are locally convex topological vector spaces (for example, *Banach spaces*), $U \subset X$ is open and $F : X \rightarrow Y$. The Gâteaux-differential of F at $u \in U$ is *in the direction* $v \in X$ defined as

$$dF(u; v) = \lim_{\tau \rightarrow 0} \frac{F(u + \tau \cdot v) - F(u)}{\tau} = T_u(v) ,$$

if the limit exists, and the operator $T_u(v) : X \rightarrow Y$ is bounded, $\tau \in \mathbb{R}$. The value

$$T_u(v) = \left. \frac{d}{d\tau} F(u + \tau \cdot v) \right|_{\tau=0}$$

is denoted as *Gâteaux-derivative* at u , if the limit exists for all $v \in X$, and one says that F is Gâteaux differentiable at u .

If F is Fréchet differentiable, then it is also Gâteaux differentiable, and its Fréchet and Gâteaux derivatives are identical $\frac{F[u]}{\delta u} [v] = T_u(v)$ and, hence, $T_u(v)$ is linear. The converse is clearly not true, since the Gâteaux derivative may fail to be linear or continuous.⁴

3.2 Fréchet-derivatives for the different divergence classes

We are now ready to investigate functional derivatives of divergences. In particular we focus on Fréchet-derivatives.

3.2.1 Bregman-divergences

We investigate the Fréchet-derivative for the Bregman-divergences (2.1) and formally obtain

$$\frac{D_{\Phi}^B(p||\rho)}{\delta \rho} = \frac{\Phi(p)}{\delta \rho} - \frac{\Phi(\rho)}{\delta \rho} - \frac{\delta \left[\frac{\delta \Phi(\rho)}{\delta \rho} (p - \rho) \right]}{\delta \rho} \quad (3.1)$$

with

$$\frac{\delta \left[\frac{\delta \Phi(\rho)}{\delta \rho} (p - \rho) \right]}{\delta \rho} = \frac{\delta^2 [\Phi(\rho)]}{\delta \rho^2} (p - \rho) - \frac{\delta \Phi(\rho)}{\delta \rho} .$$

In case of the generalized Kullback-Leibler-divergence (2.2) this reads as

$$\frac{D_{GKL}(p||\rho)}{\delta \rho} = -\frac{p}{\rho} + 1 \quad (3.2)$$

⁴In fact, it is even possible for the Gâteaux derivative to be linear and continuous but for the Fréchet derivative to fail to exist.

whereas for the usual Kullback-Leibler-divergence (2.3)

$$\frac{D_{GKL}(p||\rho)}{\delta\rho} = -\frac{p}{\rho} \quad (3.3)$$

is obtained.

Further for the subset of β -divergences (2.6) we have

$$\frac{D_\beta(p||\rho)}{\delta\rho} = -p \cdot \rho^{\beta-2} + \rho^{\beta-1} . \quad (3.4)$$

3.2.2 f -divergences

For f -divergences (2.7) the Fréchet-derivative is

$$\begin{aligned} \frac{D_f(p||\rho)}{\delta\rho} &= f\left(\frac{p(\mathbf{x})}{\rho(\mathbf{x})}\right) + \rho(\mathbf{x}) \frac{\partial f(u)}{\partial u} \frac{u}{\delta\rho} \\ &= f\left(\frac{p(\mathbf{x})}{\rho(\mathbf{x})}\right) + \rho(\mathbf{x}) \frac{\partial f(u)}{\partial u} \cdot \frac{-p}{\rho^2} \end{aligned} \quad (3.5)$$

with $u = \frac{p}{\rho}$. As a famous example we get for the Hellinger divergence (2.10)

$$\frac{D_H(p||\rho)}{\delta\rho} = 1 - \sqrt{\frac{p}{\rho}} . \quad (3.6)$$

The subset of α -divergences (2.11) can be handled by

$$\frac{D_\alpha(p||\rho)}{\delta\rho} = -\frac{1}{\alpha} (p^\alpha \rho^{-\alpha} - 1) . \quad (3.7)$$

The generalized Rényi-divergences (2.12) are treated according to

$$\begin{aligned} \frac{D_\alpha^{GR}(p||\rho)}{\delta\rho} &= -\frac{p^\alpha \rho^{-\alpha} - 1}{\int [p^\alpha \rho^{1-\alpha} - \alpha \cdot p + (\alpha - 1) \rho + 1] d\mathbf{x}} \\ &= \frac{\alpha}{\int [p^\alpha \rho^{1-\alpha} - \alpha \cdot p + (\alpha - 1) \rho + 1] d\mathbf{x}} \frac{D_\alpha(p||\rho)}{\delta\rho} , \end{aligned} \quad (3.8)$$

which is reduced to

$$\frac{D_\alpha^R(p||\rho)}{\delta\rho} = -\frac{p^\alpha \rho^{-\alpha}}{\int p^\alpha \rho^{1-\alpha} d\mathbf{x}} \quad (3.9)$$

in case of the usual Rényi-divergences (2.13).

3.2.3 γ -divergences

For the γ -divergences we rewrite (2.15) in the form

$$D_\gamma(p||\rho) = \frac{1}{\gamma + 1} \ln F_1 - \ln F_2$$

with $F_1 = (\int p^{\gamma+1} d\mathbf{x})^{\frac{1}{\gamma}} \cdot (\int \rho^{\gamma+1} d\mathbf{x})$ and $F_2 = (\int p \cdot \rho^\gamma d\mathbf{x})^{\frac{1}{\gamma}}$. Then we get

$$\frac{D_\gamma(p||\rho)}{\delta\rho} = \frac{1}{\gamma+1} \frac{1}{F_1} \frac{F_1}{\delta\rho} - \frac{1}{F_2} \frac{F_2}{\delta\rho}$$

with

$$\begin{aligned} \frac{F_1}{\delta\rho} &= \left(\int p^{\gamma+1} d\mathbf{x} \right)^{\frac{1}{\gamma}} \frac{\left(\int \rho^{\gamma+1} d\mathbf{x} \right)}{\delta\rho} \\ &= \left(\int p^{\gamma+1} d\mathbf{x} \right)^{\frac{1}{\gamma}} (\gamma+1) \rho^\gamma \end{aligned}$$

and

$$\begin{aligned} \frac{F_2}{\delta\rho} &= \frac{1}{\gamma} \left(\int p \cdot \rho^\gamma d\mathbf{x} \right)^{\frac{1}{\gamma}-1} \frac{p \cdot \rho^\gamma}{\delta\rho} \\ &= \left(\int p \cdot \rho^\gamma d\mathbf{x} \right)^{\frac{1}{\gamma}-1} p \rho^{\gamma-1}, \end{aligned}$$

such that $\frac{D_\gamma(p||\rho)}{\delta\rho}$ finally yields

$$\frac{D_\gamma(p||\rho)}{\delta\rho} = \frac{\rho^\gamma}{\left(\int \rho^{\gamma+1} d\mathbf{x} \right)} - \frac{p \rho^{\gamma-1}}{\left(\int p \cdot \rho^\gamma d\mathbf{x} \right)}. \quad (3.10)$$

Again as an important special case with $\gamma = 1$, the Fréchet-derivative of the Cauchy-Schwarz-divergence (2.17) is derived as

$$\frac{D_{CS}(p||\rho)}{\delta\rho} = \frac{\rho}{\left(\int \rho^2 d\mathbf{x} \right)} - \frac{p}{V(p, \rho)}. \quad (3.11)$$

4 Divergence based Vector Quantization using Fréchet-derivatives

Supervised and unsupervised vector quantization frequently are described in terms of dissimilarities or distances. Suppose, data are given as data vectors $\mathbf{v} \in \mathbb{R}^n$. We further assume that the vectors are discrete representations of continuous positive valued functions $p(x)$ with $v_i = p(x_i)$, $i = 1 \dots n$.

We now focus on prototype based vector quantization, i.e. data processing (clustering or classification) is realized using prototypes $\mathbf{w} \in \mathbb{R}^n$ as representatives, whereby the dissimilarity between data points as well as between data and prototypes are determined by dissimilarity measures ξ (not necessarily fulfilling triangle inequality or symmetry restrictions).

Frequently, such algorithms optimize a somewhat cost function E depending on the dissimilarity between the data points and the prototypes, i.e. usually one has $E = E(\xi(\mathbf{v}_i, \mathbf{w}_k))$ and $i = 1 \dots N$ the number of data and $k = 1 \dots C$ the number

of prototypes. This cost function can be a variant of the usual classification error in supervised learning or modified mean squared error of the dissimilarities $\xi(\mathbf{v}_i, \mathbf{w}_k)$.

If $E = E(\xi(\mathbf{v}_i, \mathbf{w}_k))$ is differentiable with respect to ξ , and ξ differentiable with respect to the prototype \mathbf{w} , then a stochastic gradient minimization is a widely used optimization scheme for E . This methodology implies the calculation of the derivatives $\frac{\partial \xi}{\partial \mathbf{w}_k}$, which has now to be considered in the light of the above functional analytic investigations for divergence measures.

If we identify the prototypes as discrete realizations of a function $\rho(x)$ and further require that p and ρ are positive functions (measures), the dissimilarity measure ξ can be chosen as a discrete variant of a divergence. The derivative $\frac{\partial \xi}{\partial \mathbf{w}}$ has to be replaced in this scenario by the Fréchet-derivative $\frac{\xi}{\delta \rho}$ in the continuous case, which reduces to usual derivatives in the discrete case (see remark in sec. 3.1.1). This is formally achieved by replacing p and ρ by their vectorial counterparts \mathbf{v} and \mathbf{w} in the formulae of the divergences provided in sec. 3.2 and further translating integrals into sums.

In the following we give prominent examples of unsupervised and supervised vector quantization, which can be optimized by gradient methods using the above introduced frame work.

4.1 Unsupervised Vector Quantization

4.1.1 Basic Vector Quantization

Unsupervised vector quantization is a class of algorithm for distributing prototypes $W = \{\mathbf{w}_k\}_A$, $\mathbf{w}_k \in \mathbb{R}^n$ such that data points $\mathbf{v} \in V \subseteq \mathbb{R}^n$ are faithfully represented in terms of a dissimilarity measure ξ . Thereby, $C = \text{card}(A)$ is the cardinality of the index set A . More formally, the data point \mathbf{v} is represented by this prototype $\mathbf{w}_{s(\mathbf{v})}$ minimizing the dissimilarity $\xi(\mathbf{v}, \mathbf{w}_k)$, i.e.

$$\mathbf{v} \mapsto s(\mathbf{v}) = \underset{k \in A}{\operatorname{argmin}} \xi(\mathbf{v}, \mathbf{w}_k) . \quad (4.1)$$

The aim of the algorithm is to distribute the prototypes in such a way that the quantization error

$$E_{VQ} = \frac{1}{2} \int P(\mathbf{v}) \xi(\mathbf{v}, \mathbf{w}_{s(\mathbf{v})}) d\mathbf{v} \quad (4.2)$$

is minimized. In its simplest form basic vector quantization (VQ) leads to a (stochastic) gradient descent on E_{VQ} with

$$\Delta \mathbf{w}_{s(\mathbf{v})} = -\varepsilon \cdot \frac{\partial \xi(\mathbf{v}, \mathbf{w}_{s(\mathbf{v})})}{\partial \mathbf{w}_{s(\mathbf{v})}} \quad (4.3)$$

for prototype update of the winning prototype $\mathbf{w}_{s(\mathbf{v})}$ according to (4.1), also known as the online variant of LBG-algorithm (C -means) [32],[52]. Here, ε is a small positive value called learning rate. As we see, the update (4.3) take into account the

derivative of the dissimilarity measure ξ with respect to the prototype. Beside the common choice of ξ being the squared Euclidean distance, the choice is given to the user with the restriction of differentiability. Hence, we are here allowed to apply divergences using its derivatives in the sense of Fréchet-derivatives.

4.1.2 Self-Organizing Maps and Neural Gas

There exist several variants of the basic vector quantization scheme to avoid local minima or to realize a projective mapping. For example, the latter can be obtained introducing an topological structure in A , usually a regular grid structure. The resulting vector quantization scheme is the Self-Organizing *Map* (SOM) introduced by T. KOHONEN [25]. The respective cost function (in the variant of Heskes, [17]) is

$$E_{\text{SOM}} = \frac{1}{2K(\sigma)} \int P(\mathbf{v}) \sum_{\mathbf{r} \in A} \delta_{\mathbf{r}}^{s(\mathbf{v})} \sum_{\mathbf{r}' \in A} h_{\sigma}^{\text{SOM}}(\mathbf{r}, \mathbf{r}') \xi(\mathbf{v}, \mathbf{w}_{\mathbf{r}}) d\mathbf{v} \quad (4.4)$$

with the so-called neighborhood function

$$h_{\sigma}^{\text{SOM}}(\mathbf{r}, \mathbf{r}') = \exp\left(-\frac{\|\mathbf{r} - \mathbf{r}'\|_A}{2\sigma^2}\right)$$

and $\|\mathbf{r} - \mathbf{r}'\|_A$ is the distance in A according to the topological structure. $K(\sigma)$ is a normalization constant depending on the neighborhood range σ . For this SOM, the mapping rule (4.1) is modified to

$$\mathbf{v} \mapsto s(\mathbf{v}) = \operatorname{argmin}_{\mathbf{r} \in A} \sum_{\mathbf{r}' \in A} h_{\sigma}^{\text{SOM}}(\mathbf{r}, \mathbf{r}') \cdot \xi(\mathbf{v}, \mathbf{w}_{\mathbf{r}}) \quad (4.5)$$

which yields in the limit $\sigma \rightarrow 0$ the original mapping (4.1). The prototype update for *all* prototypes then is given as [17]:

$$\Delta \mathbf{w}_{\mathbf{r}} = -\varepsilon h_{\sigma}^{\text{SOM}}(\mathbf{r}, s(\mathbf{v})) \frac{\partial \xi(\mathbf{v}, \mathbf{w}_{\mathbf{r}})}{\partial \mathbf{w}_{\mathbf{r}}}. \quad (4.6)$$

As above, the utilization of divergence based update is straightforward also for SOM.

If the aspect of projective mapping can be ignored while keeping the neighborhood cooperativeness aspect to avoid local minima in vector quantization, then the Neural Gas algorithm (NG) is an alternative to SOM presented by T. MARTINETZ [33]. The cost function of NG to be minimized writes as

$$E_{\text{NG}} = \frac{1}{2C(\sigma)} \sum_{j \in A} \int P(\mathbf{v}) h_{\sigma}^{\text{NG}}(\mathbf{v}, W, j) \xi(\mathbf{v}, \mathbf{w}_j) d\mathbf{v} \quad (4.7)$$

with

$$h_{\sigma}^{\text{NG}}(\mathbf{v}, W, i) = \exp\left(-\frac{k_i(\mathbf{v}, W)}{\sigma}\right), \quad (4.8)$$

with the rank function

$$k_i(\mathbf{v}, W) = \sum_j \theta(\xi(\mathbf{v}, \mathbf{w}_i) - \xi(\mathbf{v}, \mathbf{w}_j)) . \quad (4.9)$$

The mapping is realized as in basic VQ (4.1) and the prototype update for all prototypes is similar to that of SOM

$$\Delta \mathbf{w}_i = -\varepsilon h_\sigma^{NG}(\mathbf{v}, W, i) \frac{\partial \xi(\mathbf{v}, \mathbf{w}_i)}{\partial \mathbf{w}_i} . \quad (4.10)$$

Again, the incorporation of divergences is obvious also for NG.

4.1.3 Further vector quantization approaches

There exist a long list of other vector quantization approaches, like kernelized SOMs [18],[20],[19], Generative Topographic Mapping (GTM) [4], Soft Topographic Mapping [15] etc. to name just a few. Most of them utilize the Euclidean metric and the respective derivatives for adaptation. Thus, the idea of divergence based processing can be transferred to these in a similar manner.

Yet, recently, a somewhat reverse SOM has been proposed for embedding data into an embedding space \mathcal{S} , called *Exploration Machine (XOM)* [51]. This XOM can be seen as a projective structure preserving mapping of the input data into the embedding space and shows, therefore, similarities to MDS. In the XOM approach the data points $\mathbf{v}_k \in V \subseteq \mathbb{R}^n$, $k = 1, \dots, N$ are uniquely associated with prototypes $\mathbf{w}_k \in \mathcal{S}$ in the embedding space \mathcal{S} and $W = \{\mathbf{w}_k\}_{k=1}^N$. The dissimilarity $\xi_{\mathcal{S}}$ in the embedding space usually is chosen to be the quadratic Euclidean metric. Further, a *hypothesis* about the topological structure of the data \mathbf{v}_k to be embedded is formulated for the embedding space \mathcal{S} by defining a probability distribution $P_{\mathcal{S}}(s)$ for so-called *sampling vectors* $s \in \mathcal{S}$. The cost function of XOM is defined as

$$E_{\text{XOM}} = \frac{1}{2K(\sigma)} \int_{\mathcal{S}} P_{\mathcal{S}}(s) \sum_{k=1}^N \delta_k^{k^*(s)} \sum_{j=1}^N h_\sigma^{\text{XOM}}(\mathbf{v}_k, \mathbf{v}_j) \cdot \xi_{\mathcal{S}}(s, \mathbf{w}_j) ds \quad (4.11)$$

with the mapping rule

$$k^*(s) = \operatorname{argmin}_{i=1, \dots, N} \sum_{j=1}^N h_\sigma^{\text{XOM}}(\mathbf{v}_i, \mathbf{v}_j) \cdot \xi_{\mathcal{S}}(s, \mathbf{w}_j) . \quad (4.12)$$

As in usual SOMs the neighborhood cooperativeness is given by a Gaussian

$$h_\sigma^{\text{XOM}}(\mathbf{v}_k, \mathbf{v}_j) = \exp\left(-\frac{\xi_V(\mathbf{v}_k, \mathbf{v}_j)}{2\sigma^2}\right)$$

with the data dissimilarity $\xi_V(\mathbf{v}_k, \mathbf{v}_j)$ defined as Euclidean distance in the original XOM. The update of the prototypes in the embedding space is obtained in complete

analogy to SOM as

$$\Delta \mathbf{w}_i = -\varepsilon h_{\sigma}^{XOM}(\mathbf{v}_i, \mathbf{v}_{k^*(s)}) \frac{\partial \xi_S(\mathbf{s}, \mathbf{w}_i)}{\partial \mathbf{w}_i}. \quad (4.13)$$

As one can see, we can apply divergences to both ξ_V and ξ_S . In case of the latter one, the prototype update (4.13) has to be changed accordingly using the respective Fréchet-derivatives.

4.2 Learning Vector Quantization

Learning Vector Quantization (LVQ) is the supervised counterpart of basic VQ. Now the data $\mathbf{v} \in V \subseteq \mathbb{R}^n$ to be learned are equipped with class information \mathbf{c}_v . Suppose, we have K classes, we define $\mathbf{c}_v \in [0, 1]^K$. If $\sum_{k=1}^K c_k = 1$ the labeling is *probabilistic* and *possibilistic* otherwise. In case of a probabilistic labeling with $\mathbf{c}_v \in \{0, 1\}^K$ the labeling is called *crisp*.

We now briefly explore, how divergences can be used also for supervised learning. Again we start with the widely applied basic LVQ-approaches and outline afterwards the procedure for some more sophisticated methods without any claim of completeness.

4.2.1 Basic LVQ algorithms

The basic LVQ-schemes are invented by T. KOHONEN [25]. For standard LVQ a crisp data labeling is assumed. Further, the prototypes \mathbf{w}_j with labels y_j correspond to the K classes in such a way that at least one prototype is assigned to each class. For simplicity, we take exactly one prototype for each class now. The task is to distribute the prototypes in such a manner that the classification error is reduced. The respective algorithms LVQ1...LVQ3 are heuristically motivated.

As in the unsupervised vector quantization, the similarity between data and prototypes for LVQs are judged by a dissimilarity measure $\xi(\mathbf{v}, \mathbf{w}_j)$. Beside some small modifications the basic LVQ-schemes LVQ1...LVQ3 mainly consist in determination of the most proximate prototype(s) $\mathbf{w}_{s(\mathbf{v})}$ for given \mathbf{v} according to the mapping rule (4.1) and subsequent adaptation. Depending on the agreement of \mathbf{c}_v and $y_{s(\mathbf{v})}$ the adaptation of the prototype(s) takes place according to

$$\Delta \mathbf{w}_{s(\mathbf{v})} = \alpha \cdot \varepsilon \cdot \frac{\partial \xi(\mathbf{v}, \mathbf{w}_{s(\mathbf{v})})}{\partial \mathbf{w}_{s(\mathbf{v})}} \quad (4.14)$$

and $\alpha = 1$ iff $\mathbf{c}_v = y_{s(\mathbf{v})}$, and $\alpha = -1$ otherwise.

A popular generalization of these standard algorithms is the *generalized LVQ* (GLVQ) introduced by SATO&YAMADA [40]. In GLVQ the classification error is replaced by a dissimilarity based cost function, which is, of course, closely related to the classification error but not identical.

For a given data point \mathbf{v} with class label $c_{\mathbf{v}}$ the two best matching prototypes with respect to the data metric ξ , usually the quadratic Euclidian, are determined: $\mathbf{w}_{s^+(\mathbf{v})}$ has minimum distance $\xi^+ = \xi(\mathbf{v}, \mathbf{w}_{s^+(\mathbf{v})})$ under the constraint that the class labels are identically: $y_{s^+(\mathbf{v})} = c_{\mathbf{v}}$. The other best prototype $\mathbf{w}_{s^-(\mathbf{v})}$ has minimum distance $\xi^- = \xi(\mathbf{v}, \mathbf{w}_{s^-(\mathbf{v})})$ supposing the class labels are different: $y_{s^-(\mathbf{v})} = c_{\mathbf{v}}$. Then the classifier function $\mu(\mathbf{v})$ is defined as

$$\mu(\mathbf{v}) = \frac{\xi^+ - \xi^-}{\xi^+ + \xi^-} \quad (4.15)$$

being negative in case of a correct classification. The value $\xi^+ - \xi^-$ yields the hypothesis margin of the classifier [8]. Then the *generalized* LVQ (GLVQ) is derived as gradient descent on the cost function

$$E_{\text{GLVQ}} = \sum_{\mathbf{v}} \mu(\mathbf{v}) \quad (4.16)$$

with respect to the prototypes. In each learning step, for a given data point, both $\mathbf{w}_{s^+(\mathbf{v})}$ and $\mathbf{w}_{s^-(\mathbf{v})}$ are adapted in parallel. Taking the derivatives $\frac{\partial E_{\text{GLVQ}}}{\partial \mathbf{w}_{s^+(\mathbf{v})}}$ and $\frac{\partial E_{\text{GLVQ}}}{\partial \mathbf{w}_{s^-(\mathbf{v})}}$ we get for the updates

$$\Delta \mathbf{w}_{s^+(\mathbf{v})} = \epsilon^+ \cdot \theta^+ \cdot \frac{\partial \xi(\mathbf{v}, \mathbf{w}_{s^+(\mathbf{v})})}{\partial \mathbf{w}_{s^+(\mathbf{v})}} \quad (4.17)$$

and

$$\Delta \mathbf{w}_{s^-(\mathbf{v})} = -\epsilon^- \cdot \theta^- \cdot \frac{\partial \xi(\mathbf{v}, \mathbf{w}_{s^-(\mathbf{v})})}{\partial \mathbf{w}_{s^-(\mathbf{v})}} \quad (4.18)$$

with the scaling factors

$$\theta^+ = \frac{2 \cdot \xi^-}{(\xi^+ + \xi^-)^2} \text{ and } \theta^- = \frac{2 \cdot \xi^+}{(\xi^+ + \xi^-)^2} \cdot \quad (4.19)$$

The values ϵ^+ and $\epsilon^- \in (0, 1)$ are the learning rates.

Obviously, the distance measure ξ could be replaced for all these LVQ schemes by one of the introduced divergences. This offers a new possibility for information theoretic learning in classification schemes, which differs from the previous approaches significantly. These earlier approaches stress the information optimum class representation whereas here the expected information loss in terms of the applied divergence measure is optimized [45],[44],[48].

4.2.2 Advanced Learning Vector Quantization

Apart from the basic LVQ schemes, many more sophisticated prototype based learning schemes are proposed for classification learning. Here we will only restrict ourself to such approaches which can deal with probabilistic or possibilistic labeled training data (uncertain decisions), and which are additionally related to the basic

unsupervised and supervised vector quantization algorithms mentioned in this paper so far.

In particular, we focus on the Fuzzy-labeled SOM (FLSOM) and the very similar Fuzzy-labeled NG (FLNG) [50],[47]. Both approaches extend the cost function of its unsupervised counterpart in the following shorthand manner

$$E_{\text{FLSOM/FLNG}} = (1 - \beta) E_{\text{SOM/NG}} + \beta E_{\text{FL}}$$

where E_{FL} measures the classification accuracy . The factor $\beta \in [0, 1]$ is a factor balancing unsupervised and supervised learning. The classification accuracy term E_{FL} is defined as

$$E_{\text{FL}} = \frac{1}{2} \int P(\mathbf{v}) \sum_{\mathbf{r}} g_{\gamma}(\mathbf{v}, \mathbf{w}_{\mathbf{r}}) \psi(\mathbf{c}_{\mathbf{v}}, \mathbf{y}_{\mathbf{r}}) d\mathbf{v} \quad (4.20)$$

where $g_{\gamma}(\mathbf{v}, \mathbf{w}_{\mathbf{r}})$ is a Gaussian kernel describing a neighborhood range in the data space

$$g_{\gamma}(\mathbf{v}, \mathbf{w}_{\mathbf{r}}) = \exp\left(-\frac{\xi(\mathbf{v}, \mathbf{w}_{\mathbf{r}})}{2\gamma^2}\right) . \quad (4.21)$$

using the dissimilarity $\xi(\mathbf{v}, \mathbf{w}_{\mathbf{r}})$ in the data space. $\psi(\mathbf{c}_{\mathbf{v}}, \mathbf{y}_{\mathbf{r}})$ judges the dissimilarities between label vectors of data and prototypes. $\psi(\mathbf{c}_{\mathbf{v}}, \mathbf{y}_{\mathbf{r}})$ is originally suggested to be the quadratic Euclidean distance.

Note that E_{FL} depends on the dissimilarity in the data space $\xi(\mathbf{v}, \mathbf{w}_{\mathbf{r}})$ via $g_{\gamma}(\mathbf{v}, \mathbf{w}_{\mathbf{r}})$. Hence, prototype adaptation in FLSOM/FLNG is influenced by the classification accuracy:

$$\frac{\partial E_{\text{FLSOM/NG}}}{\partial \mathbf{w}_{\mathbf{r}}} = \frac{\partial E_{\text{SOM/NG}}}{\partial \mathbf{w}_{\mathbf{r}}} + \frac{\partial E_{\text{FL}}}{\partial \mathbf{w}_{\mathbf{r}}} \quad (4.22)$$

which yields

$$\begin{aligned} \Delta \mathbf{w}_{\mathbf{r}} &= -\epsilon(1 - \beta) \cdot h_{\sigma}^{\text{SOM/NG}}(\mathbf{r}, s(\mathbf{v})) \frac{\partial \xi(\mathbf{v}, \mathbf{w}_{\mathbf{r}})}{\partial \mathbf{w}_{\mathbf{r}}} \\ &+ \epsilon\beta \frac{1}{4\gamma^2} \cdot g_{\gamma}(\mathbf{v}, \mathbf{w}_{\mathbf{r}}) \frac{\partial \xi(\mathbf{v}, \mathbf{w}_{\mathbf{r}})}{\partial \mathbf{w}_{\mathbf{r}}} \psi(\mathbf{c}_{\mathbf{v}}, \mathbf{y}_{\mathbf{r}}) . \end{aligned} \quad (4.23)$$

The label adaptation is only influenced by the second part E_{FL} . The derivative $\frac{\partial E_{\text{FL}}}{\partial \mathbf{y}_{\mathbf{r}}}$ yields

$$\Delta \mathbf{y}_{\mathbf{r}} = \epsilon_l \beta \cdot g_{\gamma}(\mathbf{v}, \mathbf{w}_{\mathbf{r}}) \frac{\partial \psi(\mathbf{c}_{\mathbf{v}}, \mathbf{y}_{\mathbf{r}})}{\partial \mathbf{y}_{\mathbf{r}}} \quad (4.24)$$

with learning rate $\epsilon_l > 0$ [50],[47]. This label learning leads to a weighted average $\mathbf{y}_{\mathbf{r}}$ of the fuzzy labels $\mathbf{c}_{\mathbf{v}}$ of those data \mathbf{v} , which are close to the associated prototypes according to $\xi(\mathbf{v}, \mathbf{w}_{\mathbf{r}})$.

It should be noted at this point that a similar approach can easily be installed also for XOM in an analog manner yielding FLXOM.

Clearly, beside the possibility of choosing a divergence measure for $\xi(\mathbf{v}, \mathbf{w}_r)$ as in the unsupervised case, there is no contradiction to do so also for the label dissimilarity $\psi(\mathbf{c}_v, \mathbf{y}_r)$ in these FL-methods. As before, the simple plugin of the respective discrete divergence variants and their Fréchet-derivatives modifies the algorithms such that a semi-supervised learning can be proceeded relying on divergences for both variants.

4.2.3 Hyperparameter learning for α -, β -, and γ -divergences

Considering the parametrized divergence families of γ -, α -, and β -divergences, one could further think about the optimal choice of the so-called hyperparameters γ , α , and β as suggested in a similar manner for other parametrized LVQ-algorithms [41]. In case of supervised learning schemes for classification based on differentiable cost functions, the optimization can be handled as an object of a gradient descent based adaptation procedure. Thus, the parameter is optimized in dependence of the classification task at hand.

Suppose, the classification accuracy for a certain approach is given as

$$E = E(\xi_\eta, W)$$

depending on a *parametrized divergence* ξ_η with parameter η . If E and ξ_η are both differentiable with respect to η according to

$$\frac{\partial E(\xi_\eta, W)}{\partial \eta} = \frac{\partial E}{\partial \xi_\eta} \cdot \frac{\partial \xi_\eta}{\partial \eta},$$

a gradient based optimization is derived by

$$\Delta \eta = -\varepsilon \frac{\partial E(\xi_\eta, W)}{\partial \eta} = -\varepsilon \frac{\partial E}{\partial \xi_\eta} \cdot \frac{\partial \xi_\eta}{\partial \eta}$$

depending on the derivative $\frac{\partial \xi_\eta}{\partial \eta}$ for a certain choice of the divergence ξ_η .

We assume in the following that the the (positive) measures p and ρ are continuously differentiable. Than, considering derivatives of parametrized divergences $\frac{\partial \xi_\eta}{\partial \eta}$ with respect to the parameter η , it is allowed to interchange integration and differentiation, if the resulting integral exists [12]. Hence, we can differentiate parametrized divergences with respect to their hyperparameter in that case. For the several α -, β -, and γ -divergences characterized in sec. 2 we obtain after some elementary calculations (see Appendix):

- β -divergence $D_\beta(p||\rho)$ from (2.6)

$$\begin{aligned} \frac{\partial D_\beta(p||\rho)}{\partial \beta} &= \frac{1}{\beta - 1} \int p \left(p^{\beta-1} \ln p - \rho^{\beta-1} \ln \rho - \frac{(p^{\beta-1} - \rho^{\beta-1})}{(\beta - 1)} \right) dx \\ &\quad - \int (p^\beta \ln p - \rho^\beta \ln \rho) \frac{1}{\beta} - \frac{1}{\beta^2} (p^\beta - \rho^\beta) dx \end{aligned}$$

- α -divergence $D_\alpha(p||\rho)$ from (2.11)

$$\begin{aligned} \frac{\partial D_\alpha(p||\rho)}{\partial \alpha} &= -\frac{(2\alpha-1)}{\alpha^2(\alpha-1)^2} \int [p^\alpha \rho^{1-\alpha} - \alpha \cdot p + (\alpha-1)\rho] d\mathbf{x} \\ &\quad + \frac{1}{\alpha(\alpha-1)} \int p^\alpha \rho^{1-\alpha} (\ln p - \ln \rho) - p + \rho d\mathbf{x} \end{aligned}$$

- generalized Rényi-divergence $D_\alpha^{GR}(p||\rho)$ from (2.12)

$$\begin{aligned} \frac{\partial D_\alpha^{GR}(p||\rho)}{\partial \alpha} &= -\frac{1}{(\alpha-1)^2} \log \left(\int [p^\alpha \rho^{1-\alpha} - \alpha \cdot p + (\alpha-1)\rho + 1] d\mathbf{x} \right) \\ &\quad + \frac{1}{\alpha-1} \frac{\int p^\alpha \rho^{1-\alpha} (\ln p - \ln \rho) - p + \rho d\mathbf{x}}{\int [p^\alpha \rho^{1-\alpha} - \alpha \cdot p + (\alpha-1)\rho + 1] d\mathbf{x}} \end{aligned}$$

- Rényi-divergence $D_\alpha^R(p||\rho)$ from (2.13)

$$\begin{aligned} \frac{\partial D_\alpha^R(p||\rho)}{\partial \alpha} &= -\frac{1}{(\alpha-1)^2} \log \left(\int p^\alpha \rho^{1-\alpha} d\mathbf{x} \right) \\ &\quad + \frac{1}{\alpha-1} \frac{\int p^\alpha \rho^{1-\alpha} (\ln p - \ln \rho) d\mathbf{x}}{\int p^\alpha \rho^{1-\alpha} d\mathbf{x}} \end{aligned}$$

- γ -divergence $D_\gamma(p||\rho)$ from (2.15)

$$\begin{aligned} \frac{\partial D_\gamma(p||\rho)}{\partial \gamma} &= -\frac{(2\gamma+1)}{\gamma^2(\gamma+1)^2} \ln \left(\int p^{\gamma+1} d\mathbf{x} \right) + \frac{\int p^{\gamma+1} \ln p d\mathbf{x}}{(\gamma+1)\gamma \int p^{\gamma+1} d\mathbf{x}} \\ &\quad - \frac{1}{(\gamma+1)^2} \ln \left(\int \rho^{\gamma+1} d\mathbf{x} \right) + \frac{\int \rho^{\gamma+1} \ln \rho d\mathbf{x}}{(\gamma+1) \int \rho^{\gamma+1} d\mathbf{x}} \\ &\quad + \frac{1}{\gamma^2} \ln \left(\int p \cdot \rho^\gamma d\mathbf{x} \right) - \frac{\int p \rho^\gamma \ln \rho d\mathbf{x}}{\gamma \int p \cdot \rho^\gamma d\mathbf{x}} \end{aligned}$$

5 Conclusion

In this article we provide the mathematical foundation for divergence based supervised and unsupervised vector quantization bearing on the derivatives of the applied divergences. For this purpose, we first characterized the main sub-classes of divergences, Bregman-, α -, β -, γ -, and f -divergences following [7]. The mathematical framework of Fréchet-derivatives is then used to calculate the functional divergence derivatives.

We exemplarily show the utilization of this methodology for famous examples of supervised and unsupervised vector quantization including SOM, NG, and GLVQ. In particular, we explained that the divergences can be taken as suitable dissimilarity measures for data, which leads to the usage of the respective Fréchet-derivatives in the online learning schemes. Further, we declare, how a parameter adaptation could be integrated in supervised learning to achieve improved classification results in case of the parametrized α -, β -, and γ -divergences.

6 Appendix – Calculation of the derivatives of the parametrized divergences with respect to the hyperparameters

We assume for the differentiation of the divergences with respect to their hyperparameters that the (positive) measures p and ρ are continuously differentiable. Then, considering derivatives of divergences, it is allowed to interchange integration and differentiation, if the resulting integral exists [12].

6.1 β -divergence

The β -divergence is according to (2.6)

$$\begin{aligned} D_\beta(p||\rho) &= \int p \cdot \frac{p^{\beta-1} - \rho^{\beta-1}}{\beta - 1} dx - \int \frac{p^\beta - \rho^\beta}{\beta} dx \\ &= I_1(\beta) - I_2(\beta) \end{aligned}$$

We treat both integrals independently.

$$\begin{aligned} \frac{\partial I_1(\beta)}{\partial \beta} &= \int \frac{\partial \left[p \cdot \frac{p^{\beta-1} - \rho^{\beta-1}}{\beta-1} \right]}{\partial \beta} dx \\ &= \int p \left(\frac{\partial [p^{\beta-1} - \rho^{\beta-1}]}{\partial \beta} \frac{1}{\beta-1} - \frac{(p^{\beta-1} - \rho^{\beta-1})}{(\beta-1)^2} \right) dx \\ &= \frac{1}{\beta-1} \int p \left(p^{\beta-1} \ln p - \rho^{\beta-1} \ln \rho - \frac{(p^{\beta-1} - \rho^{\beta-1})}{(\beta-1)} \right) dx \end{aligned}$$

$$\begin{aligned} \frac{\partial I_2(\beta)}{\partial \beta} &= \int \frac{\partial \left[\frac{p^\beta - \rho^\beta}{\beta} \right]}{\partial \beta} dx \\ &= \int \frac{\partial [p^\beta - \rho^\beta]}{\partial \beta} \frac{1}{\beta} - \frac{1}{\beta^2} (p^\beta - \rho^\beta) dx \\ &= \int (p^\beta \ln p - \rho^\beta \ln \rho) \frac{1}{\beta} - \frac{1}{\beta^2} (p^\beta - \rho^\beta) dx \end{aligned}$$

Thus

$$\begin{aligned} \frac{\partial D_\beta(p||\rho)}{\partial \beta} &= \frac{1}{\beta-1} \int p \left(p^{\beta-1} \ln p - \rho^{\beta-1} \ln \rho - \frac{(p^{\beta-1} - \rho^{\beta-1})}{(\beta-1)} \right) dx \\ &\quad - \int (p^\beta \ln p - \rho^\beta \ln \rho) \frac{1}{\beta} - \frac{1}{\beta^2} (p^\beta - \rho^\beta) dx \end{aligned}$$

if the integral exists for an appropriate choice of β .

6.2 α -divergences

We consider the α -divergence (2.11)

$$\begin{aligned} D_\alpha(p||\rho) &= \frac{1}{\alpha(\alpha-1)} \int [p^\alpha \rho^{1-\alpha} - \alpha \cdot p + (\alpha-1)\rho] dx. \\ &= \frac{1}{\alpha(\alpha-1)} I(\alpha) \end{aligned}$$

We have

$$\begin{aligned} \frac{\partial D_\alpha(p||\rho)}{\partial \alpha} &= \frac{\partial \left[\frac{1}{\alpha(\alpha-1)} \right]}{\partial \alpha} I(\alpha) + \frac{1}{\alpha(\alpha-1)} \frac{\partial I(\alpha)}{\partial \alpha} \\ &= -\frac{(2\alpha-1)}{\alpha^2(\alpha-1)^2} I(\alpha) + \frac{1}{\alpha(\alpha-1)} \frac{\partial I(\alpha)}{\partial \alpha} \end{aligned}$$

The derivative $\frac{\partial I(\alpha)}{\partial \alpha}$ yields

$$\begin{aligned} \frac{\partial I(\alpha)}{\partial \alpha} &= \int \frac{\partial [p^\alpha \rho^{1-\alpha} - \alpha \cdot p + (\alpha-1)\rho]}{\partial \alpha} dx \\ &= \int p^\alpha \rho^{1-\alpha} (\ln p - \ln \rho) - p + \rho dx \end{aligned}$$

and, finally we get

$$\begin{aligned} \frac{\partial D_\alpha(p||\rho)}{\partial \alpha} &= -\frac{(2\alpha-1)}{\alpha^2(\alpha-1)^2} \int [p^\alpha \rho^{1-\alpha} - \alpha \cdot p + (\alpha-1)\rho] dx \\ &\quad + \frac{1}{\alpha(\alpha-1)} \int p^\alpha \rho^{1-\alpha} (\ln p - \ln \rho) - p + \rho dx \end{aligned}$$

6.3 Rényi-divergences

Considering the *generalized* Rényi-divergence $D_\alpha^{GR}(p||\rho)$ from (2.12)

$$\begin{aligned} D_\alpha^{GR}(p||\rho) &= \frac{1}{\alpha-1} \log \left(\int [p^\alpha \rho^{1-\alpha} - \alpha \cdot p + (\alpha-1)\rho + 1] dx \right) \\ &= \frac{1}{\alpha-1} \log I(\alpha) \end{aligned}$$

we get:

$$\frac{\partial D_\alpha^{GR}(p||\rho)}{\partial \alpha} = -\frac{1}{(\alpha-1)^2} \log I(\alpha) + \frac{1}{\alpha-1} \frac{1}{I(\alpha)} \frac{\partial I(\alpha)}{\partial \alpha}$$

with

$$\begin{aligned} \frac{\partial I(\alpha)}{\partial \alpha} &= \int \frac{\partial [p^\alpha \rho^{1-\alpha} - \alpha \cdot p + (\alpha-1)\rho + 1]}{\partial \alpha} dx \\ &= \int p^\alpha \rho^{1-\alpha} (\ln p - \ln \rho) - p + \rho dx \end{aligned}$$

Summarizing the differentiation yields

$$\begin{aligned} \frac{\partial D_\alpha^{GR}(p||\rho)}{\partial \alpha} &= -\frac{1}{(\alpha-1)^2} \log \left(\int [p^\alpha \rho^{1-\alpha} - \alpha \cdot p + (\alpha-1)\rho + 1] d\mathbf{x} \right) \\ &\quad + \frac{1}{\alpha-1} \frac{\int p^\alpha \rho^{1-\alpha} (\ln p - \ln \rho) - p + \rho d\mathbf{x}}{\int [p^\alpha \rho^{1-\alpha} - \alpha \cdot p + (\alpha-1)\rho + 1] d\mathbf{x}} \end{aligned}$$

We now turn to the usual Rényi-divergence $D_\alpha^R(p||\rho)$ from (2.13)

$$D_\alpha^{GR}(p||\rho) = \frac{1}{\alpha-1} \log \left(\int p^\alpha \rho^{1-\alpha} d\mathbf{x} \right)$$

We analogously achieve

$$\begin{aligned} \frac{\partial D_\alpha^{GR}(p||\rho)}{\partial \alpha} &= -\frac{1}{(\alpha-1)^2} \log \left(\int p^\alpha \rho^{1-\alpha} d\mathbf{x} \right) \\ &\quad + \frac{1}{\alpha-1} \frac{\int p^\alpha \rho^{1-\alpha} (\ln p - \ln \rho) d\mathbf{x}}{\int p^\alpha \rho^{1-\alpha} d\mathbf{x}} \end{aligned}$$

6.4 γ -divergences

The remaining divergences are the γ -divergences (2.15):

$$\begin{aligned} D_\gamma(p||\rho) &= \frac{1}{\gamma+1} \ln \left[\left(\int p^{\gamma+1} d\mathbf{x} \right)^{\frac{1}{\gamma}} \cdot \left(\int \rho^{\gamma+1} d\mathbf{x} \right) \right] - \ln \left[\left(\int p \cdot \rho^\gamma d\mathbf{x} \right)^{\frac{1}{\gamma}} \right] \\ &= \frac{1}{\gamma+1} \ln \left[\left(\int p^{\gamma+1} d\mathbf{x} \right)^{\frac{1}{\gamma}} \right] + \frac{1}{\gamma+1} \ln \left[\left(\int \rho^{\gamma+1} d\mathbf{x} \right) \right] - \ln \left[\left(\int p \cdot \rho^\gamma d\mathbf{x} \right)^{\frac{1}{\gamma}} \right] \\ &= \frac{1}{(\gamma+1)\gamma} \ln I_1(\gamma) + \frac{1}{\gamma+1} \ln I_2(\gamma) - \frac{1}{\gamma} \ln I_3(\gamma) \end{aligned}$$

The derivative is obtained according to

$$\begin{aligned} \frac{\partial D_\gamma(p||\rho)}{\partial \gamma} &= -\frac{(2\gamma+1)}{\gamma^2(\gamma+1)^2} \ln I_1(\gamma) + \frac{1}{(\gamma+1)\gamma I_1(\gamma)} \frac{\partial I_1(\gamma)}{\partial \gamma} \\ &\quad - \frac{1}{(\gamma+1)^2} \ln I_2(\gamma) + \frac{1}{(\gamma+1)I_2(\gamma)} \frac{\partial I_2(\gamma)}{\partial \gamma} \\ &\quad + \frac{1}{\gamma^2} \ln I_3(\gamma) - \frac{1}{\gamma I_3(\gamma)} \frac{\partial I_3(\gamma)}{\partial \gamma} \end{aligned}$$

Next, we calculate the derivatives $\frac{\partial I_1(\gamma)}{\partial \gamma}$, $\frac{\partial I_2(\gamma)}{\partial \gamma}$ and $\frac{\partial I_3(\gamma)}{\partial \gamma}$:

$$\begin{aligned} \frac{\partial I_1(\gamma)}{\partial \gamma} &= \int \frac{\partial (p^{\gamma+1})}{\partial \gamma} d\mathbf{x} \\ &= \int p^{\gamma+1} \ln p d\mathbf{x} \end{aligned}$$

$$\begin{aligned} \frac{\partial I_2(\gamma)}{\partial \gamma} &= \int \frac{\partial (\rho^{\gamma+1})}{\partial \gamma} d\mathbf{x} \\ &= \int \rho^{\gamma+1} \ln \rho d\mathbf{x} \end{aligned}$$

$$\begin{aligned}\frac{\partial I_3(\gamma)}{\partial \gamma} &= \int \frac{\partial (p \cdot \rho^\gamma)}{\partial \gamma} d\mathbf{x} \\ &= \int p \rho^\gamma \ln \rho d\mathbf{x}\end{aligned}$$

Collecting all intermediate results we finally have

$$\begin{aligned}\frac{\partial D_\gamma(p||\rho)}{\partial \gamma} &= -\frac{(2\gamma + 1)}{\gamma^2(\gamma + 1)^2} \ln \left(\int p^{\gamma+1} d\mathbf{x} \right) + \frac{\int p^{\gamma+1} \ln p d\mathbf{x}}{(\gamma + 1) \gamma \int p^{\gamma+1} d\mathbf{x}} \\ &\quad - \frac{1}{(\gamma + 1)^2} \ln \left(\int \rho^{\gamma+1} d\mathbf{x} \right) + \frac{\int \rho^{\gamma+1} \ln \rho d\mathbf{x}}{(\gamma + 1) \int \rho^{\gamma+1} d\mathbf{x}} \\ &\quad + \frac{1}{\gamma^2} \ln \left(\int p \cdot \rho^\gamma d\mathbf{x} \right) - \frac{\int p \rho^\gamma \ln \rho d\mathbf{x}}{\gamma \int p \cdot \rho^\gamma d\mathbf{x}}.\end{aligned}$$

References

- [1] S.-I. Amari. *Differential-Geometrical Methods in Statistics*. Springer, 1985.
- [2] A. Banerjee, S. Merugu, I. Dhillon, and J. Ghosh. Clustering with bregman divergences. *Journal of Machine Learning Research*, 6:1705–1749, 2005.
- [3] A. Basu, I. Harris, N. Hjort, and M. Jones. Robust and efficient estimation by minimising a density power divergence. *Biometrika*, 85(3):549–559, 1998.
- [4] C. M. Bishop, M. Svensén, and C. K. I. Williams. GTM: The generative topographic mapping. *Neural Computation*, 10:215–234, 1998.
- [5] L. Bregman. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics*, 7(3):200–217, 1967.
- [6] A. Cichocki, H. Lee, Y.-D. Kim, and S. Choi. Non-negative matrix factorization with α -divergence. *Pattern Recognition Letters*, 2008.
- [7] A. Cichocki, R. Zdunek, A. Phan, and S.-I. Amari. *Nonnegative Matrix and Tensor Factorizations*. Wiley, Chichester, 2009.
- [8] K. Crammer, R. Gilad-Bachrach, A. Navot, and A. Tishby. Margin analysis of the LVQ algorithm. In *Proc. NIPS 2002*, <http://www-2.cs.cmu.edu/Groups/NIPS/NIPS2002/NIPS2002preproceedings/index.html>, 2002.
- [9] I. Csiszár. Information-type measures of differences of probability distributions and indirect observations. *Studia Sci. Math. Hungaria*, 2:299–318, 1967.

- [10] S. Eguchi and Y. Kano. Robustifying maximum likelihood estimation. Technical Report 802, Tokyo-Institute of Statistical Mathematics, Tokyo, 2001.
- [11] D. Erdogmus. *Information theoretic Learning: Renyi's entropy and its application to adaptive systems training*. PhD thesis, University of Florida, 2002.
- [12] G. Fichtenholz. *Differential- und Integralrechnung*, volume II. Deutscher Verlag der Wissenschaften, Berlin, 9th ed. edition, 1964.
- [13] B. A. Frigyik, S. Srivastava, and M. Gupta. An introduction to functional derivatives. Technical Report UWEETR-2008-0001, Dept of Electrical Engineering, University of Washington, 2008.
- [14] H. Fujisawa and S. Eguchi. Robust parameter estimation with a small bias against heavy contamination. *Journal of Multivariate Analysis*, 99:2053–2081, 2008.
- [15] T. Graepel, M. Burger, and K. Obermayer. Self-organizing maps: generalizations and new optimization techniques. *Neurocomputing*, 21(1–3):173–90, 1998.
- [16] A. Hegde, D. Erdogmus, T. Lehn-Schioler, Y. Rao, and J. Principe. Vector quantization by density matching in the minimum Kullback-Leibler-divergence sense. In *Proc. of the International Joint Conference on Artificial Neural Networks (IJCNN) - Budapest*, pages 105–109, IEEE Press, 2004.
- [17] T. Heskes. Energy functions for self-organizing maps. In E. Oja and S. Kaski, editors, *Kohonen Maps*, pages 303–316. Elsevier, Amsterdam, 1999.
- [18] M. M. V. Hulle. *Faithful Representations and Topographic Maps*. Wiley Series and Adaptive Learning Systems for Signal Processing, Communications, and Control. Wiley and Sons, New York, 2000.
- [19] M. M. V. Hulle. Joint entropy maximization in kernel-based topographic maps. *Neural Computation*, 14(8):1887–1906, 2002.
- [20] M. M. V. Hulle. Kernel-based topographic map formation achieved with an information theoretic approach. *Neural Networks*, 15:1029–1039, 2002.
- [21] F. Itakura and S. Saito. Analysis synthesis telephony based on the maximum likelihood method. In *Proc. of the 6th. International Congress on Acoustics*, volume C, pages 17–20. Tokyo, 1968.
- [22] E. Jang, C. Fyfe, and H. Ko. Bregman divergences and the self organising map. In C. Fyfe, D. Kim, S.-Y. Lee, and H. Yin, editors, *Intelligent Data Engineering*

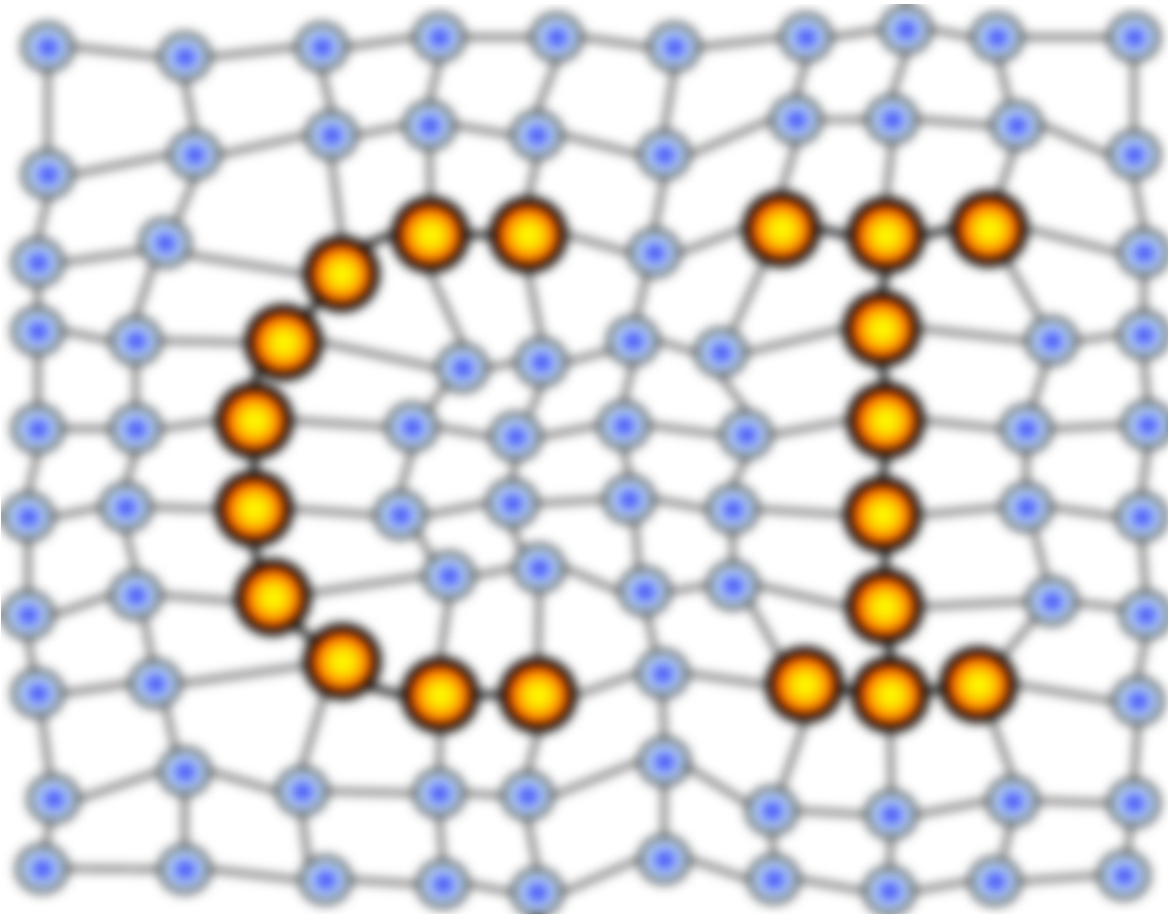
and Automated Learning IDEAL 2008, LNCS 5326, pages 452–458. Springer, 2008.

- [23] I. Kantorowitsch and G. Akilow. *Funktionalanalysis in normierten Räumen*. Akademie-Verlag, Berlin, 2nd, revised edition, 1978.
- [24] J. Kapur. *Measures of Information and their Application*. Wiley, New Delhi, 1994.
- [25] T. Kohonen. *Self-Organizing Maps*, volume 30 of *Springer Series in Information Sciences*. Springer, Berlin, Heidelberg, 1995. (Second Extended Edition 1997).
- [26] S. Kullback and R. Leibler. On information and sufficiency. *Annals of Mathematical Statistics*, 22:79–86, 1951.
- [27] P. Lai and C. Fyfe. Bregman divergences and multi-dimensional scaling. In M. Köppen, N. Kasabov, and G. Coghill, editors, *Proceedings of the International Conference on Information Processing 2008 (ICONIP)*, volume Part II of LNCS 5507, pages 935–942. Springer, 2009.
- [28] J. Lee and M. Verleysen. Generalization of the l_p norm for time series and its application to self-organizing maps. In M. Cottrell, editor, *Proc. of Workshop on Self-Organizing Maps (WSOM) 2005*, pages 733–740, Paris, Sorbonne, 2005.
- [29] J. Lee and M. Verleysen. *Nonlinear Dimensionality Reduction*. Information Sciences and Statistics. Springer Science+Business Media, New York, 2007.
- [30] T. Lehn-Schiøler, A. Hegde, D. Erdogmus, and J. Principe. Vector quantization using information theoretic concepts. *Natural Computing*, 4(1):39–51, 2005.
- [31] F. Liese and I. Vajda. On divergences and informations in statistics and information theory. *IEEE Transaction on Information Theory*, 52(10):4394–4412, 2005.
- [32] Y. Linde, A. Buzo, and R. Gray. An algorithm for vector quantizer design. *IEEE Transactions on Communications*, 28:84–95, 1980.
- [33] T. M. Martinetz, S. G. Berkovich, and K. J. Schulten. 'Neural-gas' network for vector quantization and its application to time-series prediction. *IEEE Trans. on Neural Networks*, 4(4):558–569, 1993.
- [34] M. Mihoko and S. Eguchi. Robust blind source separation by beta divergence. *Neural Computation*, 14:1859–1886, 2002.
- [35] J. C. Principe, J. F. III, and D. Xu. Information theoretic learning. In S. Haykin, editor, *Unsupervised Adaptive Filtering*. Wiley, New York, NY, 2000.

- [36] J. Ramsay and B. Silverman. *Functional Data Analysis*. Springer Science+Media, New York, 2nd edition, 2006.
- [37] A. Renyi. On measures of entropy and information. In *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability*. University of California Press, 1961.
- [38] A. Renyi. *Probability Theory*. North-Holland Publishing Company, Amsterdam, 1970.
- [39] F. Rossi, N. Delannay, B. Conan-Gueza, and M. Verleysen. Representation of functional data in neural networks. *Neurocomputing*, 64:183–210, 2005.
- [40] A. Sato and K. Yamada. Generalized learning vector quantization. In D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo, editors, *Advances in Neural Information Processing Systems 8. Proceedings of the 1995 Conference*, pages 423–9. MIT Press, Cambridge, MA, USA, 1996.
- [41] P. Schneider, M. Biehl, and B. Hammer. Hyperparameter learning in robust soft LVQ. In M. Verleysen, editor, *Proceedings of the European Symposium on Artificial Neural Networks ESANN*, pages 517–522. d-side publications, 2009.
- [42] C. Shannon. A mathematical theory of communication. *Bell System Technical Journal*, 27:379–432, 1948.
- [43] I. Taneja and P. Kumar. Relative information of type s , Csiszár’s f -divergence, and information inequalities. *Information Sciences*, 166:105–125, 2004.
- [44] K. Torkkola. Feature extraction by non-parametric mutual information maximization. *Journal of Machine Learning Research*, 3:1415–1438, 2003.
- [45] K. Torkkola and W. Campbell. Mutual information in learning feature transformations. In *Proc. Of International Conference on Machine Learning ICML’2000*, Stanford, CA, 2000.
- [46] T. Villmann. Sobolev metrics for learning of functional data - mathematical and theoretical aspects. *Machine Learning Reports*, 1(MLR-03-2007):1–15, 2007. ISSN:1865-3960, http://www.uni-leipzig.de/compint/mlr/mlr_01_2007.pdf.
- [47] T. Villmann, B. Hammer, F.-M. Schleif, T. Geweniger, and W. Herrmann. Fuzzy classification by fuzzy labeled neural gas. *Neural Networks*, 19:772–779, 2006.
- [48] T. Villmann, B. Hammer, F.-M. Schleif, W. Hermann, and M. Cottrell. Fuzzy classification using information theoretic learning vector quantization. *Neurocomputing*, 71:3070–3076, 2008.

- [49] T. Villmann and F.-M. Schleif. Functional vector quantization by neural maps. In *Proceedings of WHISPERS 2009*, page in press, 2009.
- [50] T. Villmann, F.-M. Schleif, M. Kostrzewa, A. Walch, and B. Hammer. Classification of mass-spectrometric data in clinical proteomics using learning vector quantization methods. *Briefings in Bioinformatics*, 9(2):129–143, 2008.
- [51] A. Wismüller. The exploration machine – a novel method for data visualization. In J. Principe and R. Miikkulainen, editors, *Advances in Self-Organizing Maps – Proceedings of the 7th International Workshop WSOM 2009, St. Augustine, FL, USA*, LNCS5629, pages 344–352, Berlin, 2009. Springer.
- [52] P. L. Zador. Asymptotic quantization error of continuous signals and the quantization dimension. *IEEE Transaction on Information Theory*, (28):149–159, 1982.

MACHINE LEARNING REPORTS



Impressum

Machine Learning Reports

ISSN: 1865-3960

▽ Publisher/Editors

Prof. Dr. rer. nat. Thomas Villmann & Dr. rer. nat. Frank-Michael Schleif
University of Applied Sciences Mittweida
Technikumplatz 17, 09648 Mittweida - Germany●
<http://www.uni-leipzig.de/compint>

▽ Copyright & Licence

Copyright of the articles remains to the authors. Requests regarding the content of the articles should be addressed to the authors. All article are reviewed by at least two researchers in the respective field.

▽ Acknowledgments

We would like to thank the reviewers for their time and patience.