

Repairs to GLVQ: A New Family of Competitive Learning Schemes

Nicolaos B. Karayiannis, *Member, IEEE*, James C. Bezdek, *Fellow, IEEE*, Nikhil R. Pal, *Member, IEEE*,
Richard J. Hathaway, *Member, IEEE*, and Pin-I Pai

Abstract—First, we identify an algorithmic defect of the generalized learning vector quantization (GLVQ) scheme that causes it to behave erratically for a certain scaling of the input data. We show that GLVQ can behave incorrectly because its learning rates are reciprocally dependent on the sum of squares of distances from an input vector to the node weight vectors. Finally, we propose a new family of models—the GLVQ-F family—that remedies the problem. We derive competitive learning algorithms for each member of the GLVQ-F model and prove that they are invariant to all scalings of the data. We show that GLVQ-F offers a wide range of learning models since it reduces to LVQ as its weighting exponent (a parameter of the algorithm) approaches one from above. As this parameter increases, GLVQ-F then transitions to a model in which either all nodes may be excited according to their (inverse) distances from an input or in which the winner is excited while losers are penalized. And as this parameter increases without limit, GLVQ-F updates all nodes equally. We illustrate the failure of GLVQ and success of GLVQ-F with the IRIS data.

I. INTRODUCTION

PROTOTYPE generating clustering algorithms attempt to organize unlabeled feature vectors into natural groups and represent them compactly with one or more prototypes (codevectors, paradigms, templates, vector quantizers, etc.) for each cluster. Treatments of many classical approaches to this problem are given in the texts by Kohonen [1], Bezdek [2], Duda and Hart [3], Tou and Gonzalez [4], Hartigan [5], and Jain and Dubes [6]. Kohonen's work on learning vector quantization (LVQ) for unlabeled data has become particularly timely in recent years [7].

Let $X = \{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^p$ denote the unlabeled data, and c denote the number of nodes in the competitive layer. The LVQ model is shown in Fig. 1. The input layer is connected directly to the competition or output layer. The i th node in the output layer is associated with a weight (or prototype) vector v_i . The p components $\{v_{ij}\}$ of v_i are often regarded as weights or connection strengths of the edges that connect the p inputs to node i .

Manuscript received February 3, 1995; revised September 11, 1995. This work was supported in part by the Faculty Research Committee of Georgia Southern University and the Georgia Southern Foundation.

N. B. Karayiannis and P.-I. Pai are with the Department of Electrical and Computer Engineering, University of Houston, Houston, TX 77204 USA.

J. C. Bezdek is with the Department of Computer Science, University of West Florida, Pensacola, FL 32514 USA.

N. R. Pal is with the Machine Intelligence Unit, Indian Statistical Institute, Calcutta 700035, India.

R. J. Hathaway is with the Mathematics and Computer Science Department, Georgia Southern University, Statesboro, GA 30460-3093 USA.

Publisher Item Identifier S 1045-9225(96)06501-5.

The prototypes $V = (v_1, v_2, \dots, v_c)$, $v_i \in \mathbb{R}^p$ for $1 \leq i \leq c$, are the (unknown) vector quantizers we seek. In this context *learning* refers to finding values for the $\{v_{ij}\}$. When an input vector x is submitted to this network, distances are computed between x and each v_k . The output nodes "compete," a (minimum distance) winner node, say v_i , is found, and it is then updated using one of several update rules. Fig. 2 gives a brief specification of the LVQ algorithm (some writers call it simple competitive learning) that is used in Section IV.

The update scheme in (2), shown in Fig. 2, has the simple geometric interpretation shown in Fig. 3. The winning prototype $v_{i,t-1}$ is simply rotated toward the current data point x_k by moving along the vector $(x_k - v_{i,t-1})$ which connects it to x_k .

The amount by which $v_{i,t-1}$ is shifted to arrive at $v_{i,t}$ depends on the value of the *learning rate* parameter $\alpha_{i,t-1}$, which varies from zero to one. There is no update if $\alpha_{i,t-1} = 0$, and when $\alpha_{i,t-1} = 1$, $v_{i,t}$ becomes x_k ($v_{i,t}$ is just a convex combination of x_k and $v_{i,t-1}$). This process continues until termination via LVQ4 (or by iterate limit T), at which time the terminal prototypes can be used, for example, to find the nearest prototype hard c -partition of X via LVQ5.

LVQ ordinarily uses the Euclidean distance in (1), shown in Fig. 2. This choice corresponds roughly to the update rule shown in (2), since $\nabla_v(\|x - v\|^2) = -2(x - v)$. The origin of this learning rule comes about by assuming that each $x \in \mathbb{R}^p$ is distributed according to an unknown time-invariant probability density function $f(x)$. LVQ's objective is to find a set of v_i 's such that the expected value of the square of the discretization error

$$E(\|x - v_i\|^2) = \iint_{\mathbb{R}^p} \dots \int \|x - v_i\|^2 f(x) dx \quad (3)$$

is minimized. In this expression v_i is the winning prototype for each x , and will of course vary as x ranges over \mathbb{R}^p . A sample function of the optimization problem defined by minimization of E over \mathbb{R}^p is $e = \|x - v_i\|^2$. An optimal set of v_i 's can be approximated by applying local gradient descent to $e = \|x - v_i\|^2$ for every data point x drawn from f . Kohonen has shown [8] that under certain assumptions on the learning rates, steepest descent optimization of E is also possible, and both strategies (stochastic approximation and steepest descent) lead to update rule (2).

LVQ attempts to minimize an objective function that places all of its emphasis on the winning prototype for each data point. This is reflected in (2), which alters only the winner

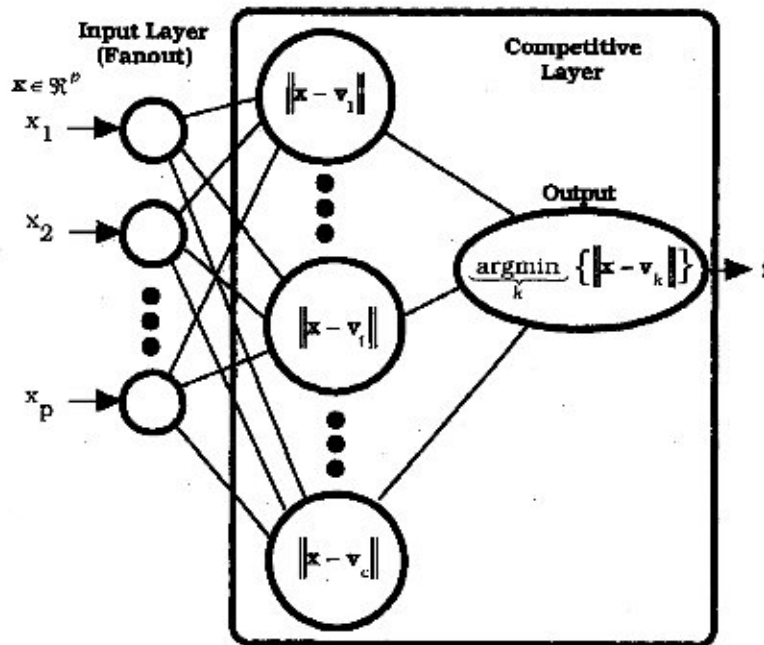


Fig. 1. The LVQ competitive learning network.

LVQ1. Given unlabeled data set $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathcal{R}^p$. Pick c , $T = \text{maximum number of updating steps}$, and $\varepsilon > 0$.

LVQ2. Initialize $\mathbf{V}_0 = (\mathbf{v}_{1,0}, \dots, \mathbf{v}_{c,0}) \in \mathcal{R}^{cp}$; Initialize learning rate $\alpha_0 \in (1,0)$

For $t = 1, 2, \dots, T$:

For $k = 1, 2, \dots, n$:

a. Find $i = \underset{r}{\operatorname{arg\,min}} \{ \|\mathbf{x}_k - \mathbf{v}_r\| \}$. (1)

b. Update the winner by

$$\mathbf{v}_{i,t} = \mathbf{v}_{i,t-1} + \alpha_t (\mathbf{x}_k - \mathbf{v}_{i,t-1}). \quad (2)$$

Next k

LVQ3. Compute

$$E_t = E_t = \|\mathbf{V}_t - \mathbf{V}_{t-1}\| = \sum_{k=1}^n \sum_{r=1}^c |\mathbf{v}_{rk,t} - \mathbf{v}_{rk,t-1}|.$$

LVQ4. If $E_t \leq \varepsilon$ stop;

Else $\alpha_t \leftarrow \alpha_t (1 - t/T)$;

Next t .

Labeling Phase (optional)

LVQ5. Find $i = \underset{s}{\operatorname{arg\,min}} \{ \|\mathbf{x}_k - \mathbf{v}_{s,t}\| \}$, and mark

\mathbf{x}_k in X with label i .

Fig. 2. The (Unlabeled Data) LVQ Algorithm. LVQ1-4 are the learning phase.

for each \mathbf{x} submitted. This ignores global information about the geometric structure of the data that is represented in

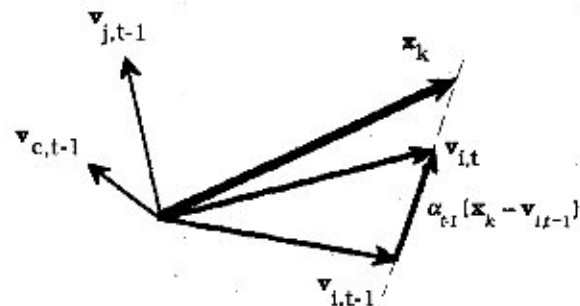


Fig. 3. Updating the winning LVQ prototype.

the remaining $(c - 1)$ distances from \mathbf{x} to the nonwinner prototypes. In this sense LVQ updating is somewhat like using the sup norm, $\|\mathbf{x} - \mathbf{v}\|_\infty = \max_{j=1}^p |x_j - v_j|$, to measure distance in \mathcal{R}^p . The extreme value in the absolute difference between pairs of coordinates dominates all others, and the distance produced simply ignores them. In the same way, LVQ updating is a very harsh local strategy that ignores global relationships between the input datum and nonwinner prototypes.

Using an inner product norm on \mathcal{R}^p ameliorates the harshness of the sup norm by counting contributions from each pair of coordinate differences in the overall distance calculation. In the same way, we think that nonwinner nodes in an LVQ network should be allowed to influence the update of the winner, and perhaps, be updated themselves. This presents two questions. First, which nodes should be accounted for (*what is the update neighborhood*)? Second, how much influence should each nonwinner node exert (*what is the learning rate distribution*)? Modifications of LVQ are usually motivated by a desire to answer one or both of these questions. Variations of LVQ that update all c quantizers simultaneously during each

updating epoch include: the soft competition scheme (SCS) of Yair *et al.* [9]; the fuzzy learning vector quantization (FLVQ) method of Tsao *et al.* [10]; and the generalized learning vector quantization (GLVQ) algorithm of Pal *et al.* [11].

The objectives of this paper are to demonstrate, analyze and remedy a defect of GLVQ that occurs for a particular scaling of the input data. Independently and at the roughly the same time as this study, Gonzalez *et al.* performed a somewhat different analysis and presented further examples of the erratic behavior of GLVQ [12]. Their results are different than the ones given in this article in several respects. First, the data used in their examples that showed sensitivity to data scaling were randomly generated subsets of \mathbb{R}^2 . They also demonstrated computationally that GLVQ could exhibit descent, ascent and even limit cycle behavior with respect to its loss function, and that the algorithm was sensitive to the number (c) of prototypes being sought. On the other hand, no replacement for or correction to GLVQ as given here was offered by Gonzalez *et al.* Nonetheless, their paper is well worth reading as a companion paper to this one, and readers are encouraged to study both.

Section II reviews GLVQ, demonstrates the problem with a numerical example and identifies the cause of the problem. Section III defines, illustrates, and analyzes a new infinite family of competitive learning models that overcomes the defect. We prove that algorithms used to optimize members of this new family are invariant to positive scalings of the data; and give a rough analysis of the limiting behavior of these algorithms as one of their parameters approaches its limits from above and below. Section IV contains our conclusions and some ideas for further research.

II. GENERALIZED LVQ

Let $\mathbf{x} \in \mathbb{R}^p$ be a stochastic input vector distributed according to a time invariant probability distribution $f(\mathbf{x})$, and let i be the best matching node as in (2). Let $L_{\mathbf{x}}$ be a loss function which measures the locally weighted mismatch (error) of \mathbf{x} with respect to the winner

$$L_{\mathbf{x}} = L(\mathbf{x}; \mathbf{v}_1, \dots, \mathbf{v}_c) = \sum_{r=1}^c g_r \|\mathbf{x} - \mathbf{v}_r\|^2, \quad \text{where} \quad (4a)$$

$$g_i = \begin{cases} 1 & \text{if } i = \arg \min_r \{\|\mathbf{x} - \mathbf{v}_r\|\} \\ \frac{1}{\sum_{j=1}^c \|\mathbf{x} - \mathbf{v}_j\|^2} & \text{otherwise} \end{cases} \quad (4b)$$

Notice especially that (4b) specifies only two values for the (c) weights $\{g_r\}$; the weight $g_i = 1$ for the winner, and a different but equal weight g_r for all ($c-1$) nonwinners. When $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n, \dots\}$ is a set of samples from $f(\mathbf{x})$ drawn at times $t = 1, 2, \dots, n, \dots$, the objective of GLVQ is to find a set of c \mathbf{v}_r 's, say $\mathbf{V} = \{\mathbf{v}_r\}$, such that the locally weighted error functional $L_{\mathbf{x}}$ defined with respect to the winner \mathbf{v}_i is minimized over \mathbf{X} . In other words, we seek a solution for the

optimization problem

$$\underbrace{\text{minimize}}_{\mathbf{V} \in \mathbb{R}^{cp}} \left\{ L(\mathbf{V}) = \iint_{\mathbb{R}^p} \dots \int \sum_{r=1}^c g_r \|\mathbf{x} - \mathbf{v}_r\|^2 f(\mathbf{x}) d\mathbf{x} \right\}. \quad (5)$$

Pal *et al.* derived update rules for solving (5) based on minimization of $L_{\mathbf{x}_k}$ via the method of steepest descent [11]. These rules are

$$\mathbf{v}_{i,t} = \mathbf{v}_{i,t-1} + \alpha_{t-1} \underbrace{\frac{D^2 - D + \|\mathbf{x}_k - \mathbf{v}_{i,t-1}\|^2}{D^2}}_{\beta_{ik,t-1}} (\mathbf{x}_k - \mathbf{v}_{i,t-1}) \quad \text{for the winner node } i, \text{ and} \quad (6a)$$

$$\mathbf{v}_{j,t} = \mathbf{v}_{j,t-1} - \alpha_{t-1} \underbrace{\frac{\|\mathbf{x}_k - \mathbf{v}_{j,t-1}\|^2}{D^2}}_{\gamma_{jk,t-1}} (\mathbf{x}_k - \mathbf{v}_{j,t-1}) \quad \text{for } j = 1, 2, \dots, c, j \neq i \quad (6b)$$

where \mathbf{x}_k is the current input vector, and

$$D = \sum_{r=1}^c \|\mathbf{x}_k - \mathbf{v}_{r,t-1}\|^2.$$

The GLVQ algorithm is defined by replacing update (2) in LVQ by (6). To avoid possible oscillations of the solution, the recommendation in [11] for learning rate α_t was that it should satisfy two conditions: as $t \rightarrow \infty$; $\alpha_t \rightarrow 0$ and $\sum \alpha_t = \infty$. For example, the learning rate $\alpha_t \leftarrow \alpha_0(1 - t/T)$ at LVQ4 is fine—note that it depends only on t and T , the iteration counter and limit, respectively. The overall learning rates in (6) thus become products of two factors; for the winner, $\alpha_{t-1} \cdot \beta_{ik,t-1}$; and for the nonwinners, $\alpha_{t-1} \cdot \gamma_{jk,t-1}$, where β and γ are as shown in (6a) and (6b). It is important to observe that β and γ both contain D^2 in their denominators.

From (6b) it follows that when the match to the winner is perfect ($\|\mathbf{x} - \mathbf{v}_{i,t-1}\|^2 = 0$), nonwinner nodes are not updated and GLVQ reduces to LVQ. In [11] it was also asserted that as the match between \mathbf{x} and the winner node \mathbf{v}_i decreased, the updating impact on nonwinner nodes increased. In view of (4b), this impact was always thought to produce nonwinner updates that were less than the update applied to the winner. However, such is not always the case. Karayiannis and Pai [13], [14] discovered that using GLVQ on the IRIS data after scaling it with $(\mathbf{x}_k \leftarrow \mathbf{x}_k/10)$ produced far different results than those reported in [11].

To appreciate this, we again use Anderson's IRIS data as an experimental data set [15]. IRIS contains 50 (physically labeled) vectors in \mathbb{R}^3 for each of $c = 3$ classes of IRIS subspecies, and it is important for this example to know that the minimum and maximum feature values in IRIS are 0.1 and 7.7. IRIS has been used in many papers to illustrate various clustering (unsupervised) and classifier (supervised) designs. Typical (resubstitution) error rates for supervised designs are zero-five mistakes; and for unsupervised designs, around 16 "mistakes" on classifiers that are subsequently designed with clustering outputs (usually prototypes).

TABLE 1
GLVQ ON IRIS AND IRIS/10 FOR $T = 500$ AND $\alpha_0 = 0.6$

	Initial prototypes				Terminal V on IRIS/10				Physical Class Means/10			
v_1	0.430	0.200	0.100	0.010	0.588	0.320	0.348	0.113	0.500	0.342	0.146	0.024
v_2	0.610	0.320	0.395	0.130	0.586	0.321	0.349	0.113	0.593	0.277	0.426	0.132
v_3	0.790	0.440	0.690	0.250	0.586	0.321	0.347	0.112	0.658	0.297	0.555	0.202
	Initial prototypes				Terminal V on IRIS				Physical Class Means			
v_1	4.300	2.000	1.000	0.100	5.007	3.423	1.472	0.250	5.006	3.428	1.462	0.246
v_2	6.100	3.200	3.950	1.300	5.881	2.742	4.387	1.433	5.936	2.770	4.260	1.326
v_3	7.900	4.400	6.900	2.500	6.847	3.078	5.708	2.058	6.588	2.974	5.552	2.026

Nearest prototype classifier design based on (any) unsupervised prototype generating algorithm is done as follows. First, the prototypes are found by exercising the training phase on all of the unlabeled data until termination. At this point, the terminal prototypes $V = \{v_k; 1 \leq k \leq c\}$ have numerical (not physical) labels. Once the prototypes are relabeled (if needed) using a relabeling algorithm such as the one discussed in [11], they are used to define a crisp nearest prototype (1-NP) classifier. This classifier is simple to implement, needing, besides V , only a definition of how to interpret (compute) the concept of nearest. The usual choice is to measure distances using a norm on \mathbb{R}^p ; and among the norms, the most common choice is an inner product norm. This is the classifier we will use to compare sets of prototypes.

1) *Crisp Nearest Prototype (1-NP) Inner Product Classifier:* Given prototypes $V = \{v_k; 1 \leq k \leq c\}$ and $z \in \mathbb{R}^p$

Decide $z \in i$ (i.e., label z as class i)

$$\Leftrightarrow \|z - v_i\|_A < \|z - v_j\|_A; 1 < j \leq c, j \neq i. \quad (7)$$

In (7) A is any positive definite $p \times p$ weight matrix—it renders the norm in (7) an inner product norm, $\|z - v_i\|_A = \sqrt{(z - v_i)^T A (z - v_i)}$. Equation (7) defines a crisp classifier that has piecewise-linear decision boundaries, even though its parameters (V) may come from a fuzzy or probabilistic algorithm. Our calculations use $A = I$, the $p \times p$ identity matrix, which gives the Euclidean norm. Ties in (7) are resolved arbitrarily. Note that (7) is the same operation as LVQ5.

How do we use 1-NP machines to compare unsupervised learning algorithms? The method employed here is to derive prototypes V from labeled data without using the labels (that is, we pretend there are no labels) during the training phase. Then, the data are submitted a 1-NP design based on V and crisp labels are assigned to each datum via (7). Finally, we count the number of errors by comparing the computed labels to the target (given) labels. Error counts are conveniently tabulated using the confusion matrix C that can be constructed during this process. The observed resubstitution error rate (in percent) is

$$E(X, V) = 100 * \left(1 - \left(\frac{\# \text{ right}}{\# \text{ tried}} \right) \right) = 100 * (1 - \langle \text{tr}(C) \rangle) \quad (8)$$

where C is the $(c \times c)$ confusion matrix

$$C = [c_{ij}] = [\# \text{ labeled class } j \text{ but were really class } i]. \quad (9)$$

In the examples below where $X = \text{IRIS}$, errors made by (7) using any V are aggregated by resubmitting IRIS to the nearest prototype rule and comparing the label assigned by (7) to the physical label of each vector in the data.

Good initialization of competitive learning schemes that find vector quantizers from unlabeled data is very important [11]. A typical choice for the initial prototypes is to make random draws from the uniform distribution on $[0, 1]$ for each coordinate of each prototype. This might work well when the data have values in this range, and random draws can always be scaled to the general range of the data. Another popular method is to simply draw (c) distinct vectors from data set X , and use them as the initial prototypes.

Random initialization or c draws from X do not generally allocate prototypes uniformly over the input space of the data. An initialization method that guarantees this follows. For data set $X = \{x_1, \dots, x_n\} \subset \mathbb{R}^p$, let data point k and initial prototype i be $x_k = \{x_{1k}, x_{2k}, \dots, x_{pk}\}^T$ and $v_i = \{v_{1i}, v_{2i}, \dots, v_{pi}\}^T$, respectively. Compute the feature ranges

Minimum of feature j :

$$m_j = \min_k \{x_{jk}\}; j = 1, 2, \dots, p; \quad \text{and} \quad (10)$$

Maximum of feature j :

$$M_j = \max_k \{x_{jk}\}; j = 1, 2, \dots, p. \quad (11)$$

With these, compute the j th component of the i th initial prototype as

$$v_{ji} = m_j + (i - 1) \left(\frac{M_j - m_j}{c - 1} \right) \quad i = 1, 2, \dots, c; \quad j = 1, 2, \dots, p. \quad (12)$$

Formula (12) disperses initial prototype values uniformly along each feature range $[m_j, M_j]$. For example, $v_1 = m = [m_1, m_2, \dots, m_p]^T$, $v_c = M = [M_1, M_2, \dots, M_p]^T$, and so on.

Table I shows the results of $T = 500$ iterations of GLVQ using the initialization in (12), with an initial learning rate $\alpha_0 = 0.6$ on IRIS and IRIS/10, where the latter data set is produced from IRIS by the scaling ($x_k \leftarrow x_k/10$). You can see that the prototypes produced by GLVQ on IRIS are very good estimates of the physically labeled subclass mean vectors, while the estimates produced by the same algorithm on IRIS/10 are: 1) nearly equal to each other; and 2) not close to any subclass centroid. This will have an obvious (disastrous!)

TABLE II
ERROR RATES AND CONFUSION MATRICES USING GLVQ ON IRIS AND IRIS/10

	IRIS/10	IRIS
Confusion Matrix	$\begin{pmatrix} 50 & 0 & 0 \\ 1 & 1 & 48 \\ 0 & 1 & 49 \end{pmatrix}$	$\begin{pmatrix} 50 & 0 & 0 \\ 0 & 47 & 3 \\ 0 & 14 & 36 \end{pmatrix}$
Error rate	$\epsilon(\text{IRIS}/10, \mathbf{V}) = 33.3\%$	$\epsilon(\text{IRIS}, \mathbf{V}) = 11.3\%$

effect on any nearest* prototype classifier designed with the terminal quantizers for IRIS/10 shown in Table I.

Table II shows the confusion matrices and error rates produced by the nearest prototype rule (7) after the centroids obtained by GLVQ in these two cases were relabeled as needed by the relabeling algorithm in [11]. Simply dividing each datum by 10 triples the error rate!

Table II shows that GLVQ is very sensitive to simple scaling of the data. To see why, let $\Delta \mathbf{v}_{r,t} = \mathbf{v}_{r,t} - \mathbf{v}_{r,t-1}$ for $r = 1$ to c , and recall that $\gamma_{ik,t-1} = \|\mathbf{x}_k - \mathbf{v}_{i,t-1}\|^2/D^2$ in (6b). With these conventions equations (6a) and (6b) become

$$\Delta \mathbf{v}_{i,t} = \alpha_{i-1} \left[\gamma_{ik,t-1} - 1 - \frac{1}{D} \right] (\mathbf{x}_k - \mathbf{v}_{i,t-1})$$

for the winner node i , and (6a')

$$\Delta \mathbf{v}_{r,t} = \alpha_{r-1} \cdot \gamma_{ik,t-1} (\mathbf{x}_k - \mathbf{v}_{r,t-1})$$

for $r = 1, 2, \dots, c, r \neq i$. (6b')

It is easy to see that if

$$D = \sum_{r=1}^c \|\mathbf{x}_k - \mathbf{v}_{r,t-1}\|^2 < 1$$

$$\left[\gamma_{ik,t-1} + 1 - \frac{1}{D} \right] < \gamma_{ik,t-1}.$$

Since α_{r-1} is the same for winners and nonwinners, the (percent) change in the winner node vector may be *less than* changes that are made to the other $(c-1)$ node vectors in this case. Since the minimum and maximum values in IRIS/10 are 0.01 and 0.77, this has apparently happened during the application of GLVQ to IRIS/10. More generally, there may be a scaling of *any* data set which produces feature vector magnitudes that make this happen. In this case, GLVQ may terminate at a very undesirable or misleading set of vector quantizers of cluster structure in X because the nonwinner nodes will be pulled to the data more strongly than the winner node is. This results in all node vectors migrating to the same point in \mathbb{R}^p as they did for IRIS/10 in Table I.

III. A FUZZY MODIFICATION OF GLVQ

The failure of GLVQ on IRIS/10 leads back to the objective function shown in (4). Substitution of

$$D = \sum_{r=1}^c \|\mathbf{x}_k - \mathbf{v}_{r,t-1}\|^2$$

into (4b) shows the problem just exposed

$$g_i = \left\{ \begin{array}{ll} 1 & \text{if } i = \arg \min_r \{ \|\mathbf{x}_k - \mathbf{v}_r\| \} \\ \frac{1}{D}, & \text{otherwise} \end{array} \right\}, \quad 1 \leq i \leq c. \quad (4b')$$

In this form, it is easy to see that when $D < 1$ GLVQ behaves *exactly opposite* to what was desired.¹ Thus, we are led to propose a new model and algorithm that is related to GLVQ in that it is defined by the same optimization problem—viz. (5). The modification arises (as do others proposed by various authors) by changing the definition of the weights $\{g_r\}$ at (4b).

Here are the properties we want the $\{g_r\}$ to have: 1) the magnitude of g_r should be *inversely proportional* to $\|\mathbf{x}_k - \mathbf{v}_r\|$; 2) each g_r should be in $[0, 1]$; and 3) the sum of the $\{g_r\}$ should be one. There are many choices for the $\{g_r\}$ that satisfy 1)–3). For example, Yair *et al.* [9] proposed in their SCS scheme learning rates that use posterior probabilities (which satisfy these constraints) as one of their learning rate factors.

Many authors, beginning with Huntsberger and Ajji-marangsee [16], have felt that a natural choice for g_r which satisfies 1)–3) is (some function of) the membership value obtained from the fuzzy c-means (FCM) formula [2] for any $m > 1$. The basic formula for the memberships of an input \mathbf{x} over the c nodes is (assuming $\|\mathbf{x} - \mathbf{v}_j\| > 0 \forall j$ and $m > 1$)

$$w_r = \left(\sum_{j=1}^c \left(\frac{\|\mathbf{x} - \mathbf{v}_r\|^{2/(m-1)}}{\|\mathbf{x} - \mathbf{v}_j\|^{2/(m-1)}} \right) \right)^{-1} \quad r = 1, 2, \dots, c. \quad (13)$$

Various investigators have used *functions* of this formula in different ways: e.g., the PLVQ (or 1/KCN) [10], GFCM [17] and UFCL [18] algorithms are all based on (13) in one way or another. The use of (13) is intuitively appealing because FCM is a successful and robust batch clustering algorithm, and we may hope that part of its utility can be transferred to competitive learning schemes via the formula itself. Using (13) as we have here is difficult to justify mathematically (except as it does satisfy 1)–3) above), so this modification is best viewed as a heuristically based idea. Substituting the weights in (13) in (4a) yields

$$L_{\mathbf{x}} = L(\mathbf{x}; \mathbf{v}_1, \dots, \mathbf{v}_c) = \sum_{r=1}^c w_r \|\mathbf{x} - \mathbf{v}_r\|^2$$

$$= \sum_{r=1}^c \left(\sum_{j=1}^c \left(\frac{\|\mathbf{x} - \mathbf{v}_r\|^{2/(m-1)}}{\|\mathbf{x} - \mathbf{v}_j\|^{2/(m-1)}} \right) \right)^{-2} \|\mathbf{x} - \mathbf{v}_r\|^2. \quad (14)$$

Assuming that $\|\mathbf{x} - \mathbf{v}_j\| > 0 \forall j$ and $m > 1$, the gradient of $L_{\mathbf{x}}$ with respect to \mathbf{v}_k is (see Appendix A)

$$\nabla_{\mathbf{v}_k} L = \nabla_{\mathbf{v}_k} \left(\sum_{r=1}^c w_r \|\mathbf{x} - \mathbf{v}_r\|^2 \right)$$

$$= \left(\frac{-2}{m-1} \right) (\mathbf{x} - \mathbf{v}_k)$$

$$\cdot \left[(m-2)w_k + \left(\sum_{r=1}^c \left(\frac{\|\mathbf{x} - \mathbf{v}_k\|^{2/(m-1)}}{\|\mathbf{x} - \mathbf{v}_r\|^{2/(m-1)}} \right)^{2-m} \right) w_k^2 \right]. \quad (15a)$$

And for the important special case $m = 2$, (15a) reduces to

$$\nabla_{\mathbf{v}_k} L = -2w_k^2 (\mathbf{x} - \mathbf{v}_k). \quad (15b)$$

¹In fact, there are values of D for which $[\gamma_{ik,t-1} + 1 - 1/D] < 0$. And for $D \gg 1$, GLVQ reverts to LVQ.

TABLE III
GLVQ-F ($m = 2$) ON IRIS AND IRIS/10 FOR $T = 500$ AND $\alpha_0 = 0.6$

	Initial prototypes/10				Terminal V on IRIS/10				Physical Class Means/10			
v_1	0.430	0.200	0.100	0.010	0.500	0.341	0.148	0.025	0.500	0.342	0.146	0.024
v_2	0.610	0.320	0.395	0.130	0.588	0.276	0.435	0.139	0.593	0.277	0.426	0.132
v_3	0.790	0.440	0.690	0.250	0.676	0.305	0.563	0.205	0.658	0.297	0.555	0.202
	Initial Prototypes				Terminal V on IRIS				Physical Class Means			
v_1	4.300	2.000	1.000	0.100	5.002	3.412	1.483	0.254	5.006	3.428	1.462	0.246
v_2	6.100	3.200	3.950	1.300	5.884	2.762	4.358	1.393	5.936	2.770	4.260	1.326
v_3	7.900	4.400	6.900	2.500	6.769	3.054	5.637	2.057	6.588	2.974	5.552	2.026

TABLE IV
ERROR RATES AND CONFUSION MATRICES USING GLVQ-F ($m = 2$) ON IRIS AND IRIS/10

		IRIS/10	IRIS
GLVQ-F	Confusion Matrix	$\begin{pmatrix} 50 & 0 & 0 \\ 0 & 47 & 3 \\ 0 & 13 & 37 \end{pmatrix}$	$\begin{pmatrix} 50 & 0 & 0 \\ 0 & 47 & 3 \\ 0 & 13 & 37 \end{pmatrix}$
	Error rate	$E(\text{IRIS/10}, V) = 10.7\%$	$E(\text{IRIS}, V) = 10.7\%$

Using (15a), the learning rule obtained by applying gradient descent to L_x at iterate t for input vector x becomes, under the assumptions just stated

$$v_{j,t} = v_{j,t-1} + \left(\frac{2c\alpha_{t-1}}{m-1} \right) \cdot \left[(m-2)u_j + \left(\sum_{r=1}^c \left(\frac{\|x - v_{j,t-1}\|^{2/m-1}}{\|x - v_{r,t-1}\|^{2/m-1}} \right)^{3-m} \right) u_j^2 \right] \cdot (x - v_{j,t-1}); \quad j = 1, 2, \dots, c. \quad (16a)$$

For $m = 2$ gradient (15b) yields the much simpler form

$$v_{j,t} = v_{j,t-1} + \alpha_{t-1} (2cu_j^2)(x - v_{j,t-1}); \quad j = 1, 2, \dots, c. \quad (16b)$$

Since $0 < u_r < 1 = \sum_r u_r$ and $u_i < u_j \Leftrightarrow \|x - v_i\| > \|x - v_j\|$, the overall learning rate $\alpha_{t-1} (2cu_j^2)$ for node r (at $m = 2$) is inversely proportional to its distance from the data point—that is, the closer the node vector, the larger the update. This update rule does not have the problem of GLVQ that was illustrated in Section III. Using (16a) instead of (2) in our specification of LVQ results in an infinite family (one for each m) of competitive learning schemes we shall call GLVQ-F. (Perhaps GLVQ1 is a more natural name, but we wish to avoid possible confusion with Kohonen's LVQ1, which is a learning algorithm that uses labeled data.) We mention that the factor $(2c)$ in the learning rate $\alpha_{t-1} (2cu_j^2)$ is unimportant. The calculations discussed below used the simpler and equivalent expression $\alpha_{t-1} u_j^2$ as the overall learning rate for node r at iterate t in (16b). The overall learning rates in (16a) at any value of m do not satisfy the desirable properties stipulated for the weights $\{g_r\}$; rather, these properties are built into the rates as part of one factor.

We remark that learning rule (16b) was proposed explicitly by Park and Dagher in [17], and implicitly by Chung and Lee in [18] (they used any $m > 1$ but for $m \neq 2$ they did not arrive at (15a)). However, neither of these authors justified the use of (16b) via the well-defined optimization problem at (5); and neither used (16b) as the update rule the way we do. Instead, they mistakenly identified (16b) as part of a scheme to

optimize the k th term of the fuzzy c -means objective function, i.e.,

$$\min_{(U_i, V)} \left\{ J_{m,k}(U_k, V; X) = \sum_{i=1}^c (u_{ik})^m \|x_k - v_i\|_A^2; m > 1 \right\}.$$

A generalized form of (16b), different from (16a), used in the UFCL method proposed in [18] is very unstable to values of the FGM parameter (m) for m less or greater than $m = 2$. We shall return to this point later.

To see that GLVQ-F overcomes the problem of scaling that GLVQ suffers from, consider Tables III and IV, which contain the results of applying GLVQ-F (with $m = 2$) to IRIS and IRIS/10 using exactly the same initializations as in Table I.

The terminal prototypes found by GLVQ-F on both IRIS and IRIS/10 agree well with the physical class means. Differences in the second and third prototypes with the second and third class means indicate that the physical labels of the second and third subsets of IRIS do not correspond exactly with the (geometric) structure of the data that represent the flowers. This is a problem of good data representation, not bad terminal prototypes.

The fact that GLVQ-F makes 16 classification errors on IRIS and IRIS/10 instead of 17 is not important, because changes in T and α_0 can affect its outputs. (Applying the 1-NP rule to IRIS using the physical class means for V results in 11 errors.) What is important is that these results show that GLVQ-F does not suffer from the scaling problem that plagues GLVQ.

Table III suggests that GLVQ-F is affected very little by the scale of the data. We can show (the proof is in Appendix B) that GLVQ-F is scale invariant in the sense of the following proposition, which assures us that GLVQ-F cannot suffer from the scaling problem that GLVQ does.

1) *Proposition:* Let $X = \{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^p$ and let $V_0 = (v_{1,0}, v_{2,0}, \dots, v_{c,0}); v_{i,0} \in \mathbb{R}^p$, be a set of initial

prototypes. Let τ be a fixed positive number and define sets of scaled data and initial prototypes by

$$\begin{aligned} Y &= \{y_1, y_2, \dots, y_n\} = \tau X \{\tau x_1, \tau x_2, \dots, \tau x_n\} \quad \text{and} \\ W_0 &= (w_{1,0}, w_{2,0}, \dots, w_{c,0}) \\ &= \tau V_0 = (\tau v_{1,0}, \tau v_{2,0}, \dots, \tau v_{c,0}). \end{aligned}$$

Then applying GLVQ-F to X , initialized by V_0 , is equivalent to applying GLVQ-F to Y , initialized by W_0 , in the sense that

$$w_{j,t} = \tau v_{j,t}; \quad j = 1, \dots, c \quad \text{and} \quad t = 0, 1, \dots \quad (17)$$

We remark that whenever $\alpha_0, \{\alpha_k\}$ and T are the same for the X and Y problems in this proposition, and when the termination criterion for the two problems satisfies $\epsilon_Y = \tau \epsilon_X$, that the terminal prototypes satisfy $W = \tau V$.

Since GLVQ-F is well defined for any m in $(1, \infty)$, users must choose a value for this parameter before the algorithm can be implemented. Computational experiments we have done suggest that $m = 2$ may be the most reliable value. To investigate this further, we can take limits of the gradient in (15a) as m approaches its upper and lower limits. Here are the results. First (see Appendix C), we have that

$$\lim_{m \rightarrow \infty} \{v_{j,t}\} = v_{j,t-1} + \left(\frac{2\alpha_{t-1}}{c} \right) (x - v_{j,t-1}) \quad (18)$$

$$j = 1, 2, \dots, c.$$

This establishes that all c prototypes are adjusted equally in the sense that the *learning rate* is the same for all c weight vectors in (16) as parameter m grows without bound. At the other extreme, when m approaches one from above, we find that (cf. Appendix D)

$$\lim_{m \rightarrow 1^+} \{v_{j,t}\} = \begin{cases} v_{i,t-1} + (2\alpha_{t-1})(x - v_{i,t-1}), & i = \arg \min_k \{\|x - v_{k,t-1}\|\} \\ 0, & j \neq i \end{cases} \quad (19)$$

The right side of (19) is, up to the scale factor 2 (which can be dropped without loss), identical to (2). Thus, GLVQ-F reduces to LVQ when m approaches 1 from above.

In some cases we examined, gradient (15a) becomes $(-2)(x - v_i)$ (positive number) for the winner prototype but $(-2)(x - v_k)$ (negative number) for the loser prototypes. This corresponds to rewarding the losers (excitation) and penalizing the winners (inhibition). Combining this observation with (18) and (19) suggests that the parameter m controls the overall learning rates applied to the updates in an interesting way. Evidently, m can be used to transition the model from strict excitation of the winner (m very close to one) to mixed excitation-inhibition (for m close to one) to unequal excitation of all nodes (as m gets larger), to equal excitation of all nodes (as m goes to infinity). Thus, our belief is that the GLVQ-F model can represent a wide range of biologically inspired learning paradigms through adjustments to a simple weighting parameter. At present, however, the exact mechanism for the transition from one state to another is unknown to us (and may be a data-dependent event). This aspect of the model deserves a careful study before further assertions about it can be made with much confidence.

IV. CONCLUSIONS

An experimental investigation with the well known IRIS data using the GLVQ algorithm given in [11] shows that GLVQ produces good results on IRIS, but terrible results on IRIS/10. Combining our results with those reported by Gonzalez *et al.* [12] leads to an inescapable conclusion: GLVQ is best regarded as a faltered step on the way to better competitive learning models.

An infinite family of competitive learning schemes called GLVQ-F are given that overcomes the data scaling problem. We have shown that GLVQ-F is invariant to scaling. More importantly, GLVQ-F offers users a wide range of neural-like responses. We have shown that GLVQ-F reduces to LVQ as its weighting exponent approaches one from above. As m increases, it then transitions to a model in which either all nodes may be excited inversely proportionally to their distances from an input or in which the winner is penalized while losers are excited. And as m increases without limit, GLVQ-F updates all nodes equally. This behavior is probably related to epochs of ascent, descent and limit cycling of the optimization being performed, and needs to be carefully studied before stronger conclusions are warranted. A small by-product of this investigation is a very good initialization strategy for iterative algorithms that produce estimates of prototypes from unlabeled data.

A general characteristic of LVQ-like schemes is their sequential nature, which essentially updates one or more quantizers after each look at a datum (this includes the GLVQ-F model). There are batch versions of these schemes that attempt to overcome this limitation—most of them simply average the updates and apply them after a pass through X . This stabilizes the movement of the prototypes somewhat, but it is still the case that each term in the update equation is determined entirely by local information (at a single data point). We believe that batch models such as hard and fuzzy c -means [2] and descending FLVQ [10], whose algorithms attempt to minimize a global function of all n data points, have a better chance of producing useful vector quantizers in most if not all application domains. We plan an investigation that shows this to be the case.

An important point brought up by Gonzalez *et al.* [12] that was not studied here deserves special mention. It is that the quality of the prototypes may importantly depend on c , the number of them. If c is low, the scheme we have outlined here, as well as that proposed by Yair *et al.* [9] may insulate the model from local instability by distributing the updates across all c prototypes. But since the overall learning rates in these schemes satisfy the constraints $0 < \sum_{k=1}^c \alpha_{k,t} < 1$, if the update for any prototype (the winner) is large, the constraints force the other updates to become minute. Thus, as c increases (say, e.g., as is the case in image compression, where $c = 256$ is not uncommon), it may well be that the original LVQ (winner take all) strategy or some less stringent variation than GLVQ-F becomes more and more reliable (and attractive). This important possibility will be the subject of a future investigation.

APPENDIX A
COMPUTATION OF (15a)

See (15a) shown at the bottom of the page.

APPENDIX B
PROOF OF (17)

1) *Proposition:* Let $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subset \mathbb{R}^p$ and let $\mathbf{V}_0 = (\mathbf{v}_{1,0}, \mathbf{v}_{2,0}, \dots, \mathbf{v}_{c,0}), \mathbf{v}_{i,0} \in \mathbb{R}^p$, be a set of initial prototypes. Let τ be a fixed positive number and define sets of scaled data and initial prototypes by $Y = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n\} = \tau X = \{\tau \mathbf{x}_1, \tau \mathbf{x}_2, \dots, \tau \mathbf{x}_n\}$ and $\mathbf{W}_0 = (\mathbf{w}_{1,0}, \mathbf{w}_{2,0}, \dots, \mathbf{w}_{c,0}) = \tau \mathbf{V}_0 = (\tau \mathbf{v}_{1,0}, \tau \mathbf{v}_{2,0}, \dots, \tau \mathbf{v}_{c,0})$. Then applying GLVQ-F to X , initialized by \mathbf{V}_0 , is equivalent to applying GLVQ-F to Y , initialized by \mathbf{W}_0 , in the sense that

$$\mathbf{w}_{j,t} - \tau \mathbf{v}_{j,t}; \quad j = 1, \dots, c \text{ and } t = 0, \dots \quad (17)$$

Proof: We use induction. Result (17) holds for $t = 0$ by hypothesis. We now assume it holds for arbitrary $t - 1$, and

show it must hold for t . By (16a), for $j = 1, \dots, c$

$$\begin{aligned} \mathbf{w}_{j,t} &= \mathbf{w}_{j,t-1} + \left(\frac{2\alpha_t u_j}{m-1} \right) \left\{ u_j \sum_{k=1}^c \left(\frac{\|\mathbf{y} - \mathbf{w}_j\|^{2/(m-1)}}{\|\mathbf{y} - \mathbf{w}_k\|^{2/(m-1)}} \right)^{2-m} \right. \\ &\quad \left. + (m-2) \right\} (\mathbf{y} - \mathbf{w}_{j,t-1}) \\ &= \tau \mathbf{v}_{j,t-1} + \left(\frac{2\alpha_t u_j}{m-1} \right) \\ &\quad \cdot \left\{ u_j \sum_{k=1}^c \left(\frac{\|\tau \mathbf{x} - \tau \mathbf{v}_j\|^{2/(m-1)}}{\|\tau \mathbf{x} - \tau \mathbf{v}_k\|^{2/(m-1)}} \right)^{2-m} + (m-2) \right\} \\ &\quad \cdot (\tau \mathbf{x} - \tau \mathbf{v}_{j,t-1}) \\ &= \tau \left[\mathbf{v}_{j,t-1} + \left(\frac{2\alpha_t u_j}{m-1} \right) \right. \\ &\quad \cdot \left\{ u_j \sum_{k=1}^c \left(\frac{\|\mathbf{x} - \mathbf{v}_j\|^{2/(m-1)}}{\|\mathbf{x} - \mathbf{v}_k\|^{2/(m-1)}} \right)^{2-m} + (m-2) \right\} \\ &\quad \left. \cdot (\mathbf{x} - \mathbf{v}_{j,t-1}) \right] = \tau \mathbf{v}_{j,t}. \end{aligned}$$

$$\begin{aligned} \nabla_{\mathbf{v}_k} L &= \nabla_{\mathbf{v}_k} \left(\sum_{r=1}^c u_r \|\mathbf{x} - \mathbf{v}_r\|^2 \right) = \nabla_{\mathbf{v}_k} \left(\sum_r \frac{\|\mathbf{x} - \mathbf{v}_r\|^2}{\sum_j \frac{\|\mathbf{x} - \mathbf{v}_j\|^{2/(m-1)}}{\|\mathbf{x} - \mathbf{v}_j\|^{2/(m-1)}}} \right) = \nabla_{\mathbf{v}_k} \left(\sum_r \frac{(\|\mathbf{x} - \mathbf{v}_r\|^2)^{(m-2)/(m-1)}}{\sum_j (\|\mathbf{x} - \mathbf{v}_j\|^2)^{-1/(m-1)}} \right) \\ &= \nabla_{\mathbf{v}_k} \left(\frac{\sum_r (\|\mathbf{x} - \mathbf{v}_r\|^2)^{(m-2)/(m-1)}}{\sum_j (\|\mathbf{x} - \mathbf{v}_j\|^2)^{-1/(m-1)}} \right) = \left[\frac{\left(\frac{m-2}{m-1} \right) \|\mathbf{x} - \mathbf{v}_k\|^{-1/(m-1)} (-2)(\mathbf{x} - \mathbf{v}_k) \sum_j (\|\mathbf{x} - \mathbf{v}_j\|^2)^{-1/(m-1)}}{\left(\sum_j (\|\mathbf{x} - \mathbf{v}_j\|^2)^{-1/(m-1)} \right)^2} \right. \\ &\quad \left. - \frac{\sum_r (\|\mathbf{x} - \mathbf{v}_r\|^2)^{(m-2)/(m-1)} \left(\frac{-1}{m-1} \right) (\|\mathbf{x} - \mathbf{v}_k\|^2)^{-m/(m-1)} (-2)(\mathbf{x} - \mathbf{v}_k)}{\left(\sum_j (\|\mathbf{x} - \mathbf{v}_j\|^2)^{-1/(m-1)} \right)^2} \right] \\ &= \left(\frac{-2}{m-1} \right) (\mathbf{x} - \mathbf{v}_k) \left[\frac{(m-2)(\|\mathbf{x} - \mathbf{v}_k\|^{-2/(m-1)})}{\sum_j (\|\mathbf{x} - \mathbf{v}_j\|^{-2/(m-1)})} \right. \\ &\quad \left. + \frac{\sum_r (\|\mathbf{x} - \mathbf{v}_r\|^{2/(m-1)})^{m-2} (\|\mathbf{x} - \mathbf{v}_k\|^{2/(m-1)})^{2-m} (\|\mathbf{x} - \mathbf{v}_k\|^{-2/(m-1)})^2}{\left(\sum_j (\|\mathbf{x} - \mathbf{v}_j\|^{-2/(m-1)}) \right)^2} \right] \\ &= \left(\frac{2}{m-1} \right) (\mathbf{x} - \mathbf{v}_k) \left[(m-2) u_k + \left(\sum_r \left(\frac{\|\mathbf{x} - \mathbf{v}_k\|^{2/(m-1)}}{\|\mathbf{x} - \mathbf{v}_r\|^{2/(m-1)}} \right)^{2-m} \right) u_k^2 \right]. \end{aligned} \quad (15a)$$

To complete the proof, observe that the membership values $\{u_r\}$ for X and Y are identical, because they involve ratios of norm values, so that the scaling factor τ cancels out, just as in

$$\begin{aligned} \frac{\|y - w_r\|^{2/(m-1)}}{\|y - w_k\|^{2/(m-1)}} &= \frac{\|\tau x - v_r\|^{2/(m-1)}}{\|\tau x - v_k\|^{2/(m-1)}} \\ &= \frac{\tau^{2/(m-1)} \|x - v_r\|^{2/(m-1)}}{\tau^{2/(m-1)} \|x - v_k\|^{2/(m-1)}} \\ &= \frac{\|x - v_r\|^{2/(m-1)}}{\|x - v_k\|^{2/(m-1)}}. \end{aligned}$$

APPENDIX C COMPUTATION OF (18)

Rewrite the gradient equation (15a)

$$\begin{aligned} \nabla_{v_k} L &= \left(\frac{2}{m-1} \right) (x - v_k) \\ &\quad \cdot \left[(m-2)u_k - \left(\sum_r \left(\frac{\|x - v_k\|^{2/(m-1)}}{\|x - v_r\|^{2/(m-1)}} \right)^{2-m} \right) u_k^2 \right] \end{aligned} \quad (15a)$$

$$\begin{aligned} &= (-2)(x - v_k) \\ &\quad \cdot \left[u_k + \frac{\left(\sum_r \left(\frac{\|x - v_k\|^{2/(m-1)}}{\|x - v_r\|^{2/(m-1)}} \right)^{4-2m/(m-1)} - \frac{1}{u_k} \right) u_k^2}{m-1} \right] \\ &= (-2)(x - v_k)[u_k - F(k, m)]. \end{aligned} \quad (15a')$$

Evaluating the limit in (15a') as $1 \leftarrow m$ or $m \rightarrow \infty$ is complicated by the fact that u_k depends on m in a fairly messy way. To simplify the argument in the following analysis, it is assumed that $\|x - v_r\| > 0$ for $r = 1, \dots, c$ with no ties in datum to prototype distances. We first consider the case $m \rightarrow \infty$, where it is well known that

$$\lim_{m \rightarrow \infty} \{u_k\} = 1/c$$

for all k [2]. The numerator of $F(k, m)$ is easily seen to approach the number

$$\left(\sum_r \left(\frac{\|x - v_k\|}{\|x - v_r\|} \right)^{-2} \cdot c \right) \left(\frac{1}{c} \right)^2$$

as $m \rightarrow \infty$, while the denominator $(m-1) \rightarrow \infty$. This shows that

$$\lim_{m \rightarrow \infty} \{F(k, m)\} = 0$$

for all k , which, combined with

$$\lim_{m \rightarrow \infty} \{u_k\} = 1/c$$

for all k , establishes that all c prototypes are adjusted equally by (16a) in the limit, i.e.,

$$\lim_{m \rightarrow \infty} \{v_{j,t}\} = v_{j,t-1} + \left(\frac{2\alpha_{t-1}}{c} \right) (x - v_{j,t-1}) \quad j = 1, 2, \dots, c. \quad (18)$$

APPENDIX D COMPUTATION OF (19)

We consider evaluating (15a') as $1 \leftarrow m$. Let

$$\delta_{k,r} = \left(\frac{\|x - v_k\|}{\|x - v_r\|} \right)^2 \quad \text{and} \quad \Delta_k = \max\{\delta_{k,1}, \dots, \delta_{k,c}\}.$$

Note that either $\Delta_k = 1$ or $\Delta_k > 1$. Now, $F(k, m)$ can be written as

$$\begin{aligned} F(k, m) &= \frac{\left(\frac{\sum_r \delta_{k,r}^{-1} \delta_{k,r}^{1/(m-1)}}{\left(\sum_r \delta_{k,r}^{1/(m-1)} \right)^2} \right) - \left(\frac{1}{\sum_r \delta_{k,r}^{1/(m-1)}} \right)}{m-1} \\ &= \frac{\left(\frac{1}{m-1} \right) * \left(\sum_r (\delta_{k,r}^{-1} - 1) \delta_{k,r}^{1/(m-1)} \right)}{\left(\sum_r \delta_{k,r}^{1/(m-1)} \right) * \left(\sum_r \delta_{k,r}^{1/(m-1)} \right)}. \end{aligned} \quad (D.1)$$

We now argue that the limit of (D.1) is zero for both possible cases (that is, for the winner and nonwinner nodes). If $\Delta_k = 1$, then the denominator of (D.1) equals 1 as $1 \leftarrow m$, while the summation in the numerator has nonzero terms only for $\delta_{k,r} < 1$. An application of L'Hopital's rule in this case will easily establish that the numerator goes to zero. On the other hand, if $\Delta_k > 1$, then it is easy to show that

$$\begin{aligned} \frac{\left(\frac{1}{m-1} \right)}{\left(\sum_r \delta_{k,r}^{1/(m-1)} \right)} &\rightarrow 0 \quad \text{and} \\ \frac{\left(\sum_r (\delta_{k,r}^{-1} - 1) \delta_{k,r}^{1/(m-1)} \right)}{\left(\sum_r \delta_{k,r}^{1/(m-1)} \right)} &\rightarrow \Delta_k - 1. \end{aligned}$$

So in both possible cases, we have

$$\lim_{1 \leftarrow m} \{F(k, m)\} = 0.$$

It is well known that

$$\lim_{1 \leftarrow m} \{u_k\} = 1$$

for the winning index and 0 otherwise [2]. This means that in the limit $1 \leftarrow m$, GLVQ- β becomes LVQ; that is (16a) does not update loser prototypes, and updates the winner prototype v_i by

$$\lim_{1 \leftarrow m} \{v_{i,t}\} = v_{i,t-1} + (2\alpha_{t-1})(x - v_{i,t-1}). \quad (19)$$

REFERENCES

- [1] T. Kohonen, *Self-Organization and Associative Memory*, 3rd ed. Berlin: Springer-Verlag, 1989.
- [2] J. Bezdek, *Pattern Recognition with Fuzzy Objective Function Algorithms*. New York: Plenum, 1981.
- [3] R. Duda and P. Hart, *Pattern Classification and Scene Analysis*. New York: Wiley, 1973.
- [4] J. Tou and R. Gonzalez, *Pattern Recognition Principles*. Reading, MA: Addison Wesley, 1974.
- [5] J. Hartigan, *Clustering Algorithms*. New York: Wiley, 1975.
- [6] A. Jain and R. Dubes, *Algorithms that Cluster Data*. Englewood Cliffs, NJ: Prentice-Hall, 1988.
- [7] Y.-H. Pao, *Adaptive Pattern Recognition and Neural Networks*. Reading, MA: Addison-Wesley, 1989.
- [8] T. Kohonen, "Self-organizing maps: Optimization approach," in *Artificial Neural Networks*, T. Kohonen, K. Makisara, O. Simula, and J. Kangas, Eds. New York: Elsevier, 1991, pp. 981-990.
- [9] B. Yalc, K. Zeger, and A. Gersho, "Competitive learning and soft competition for vector quantizer design," *IEEE Trans. Signal Processing*, vol. 40, pp. 94-109, 1992.
- [10] E. C. K. Tsao, J. C. Bezdek, and N. R. Pal, "Fuzzy Kohonen clustering networks," *Pattern Recognition*, vol. 27, no. 5, pp. 757-764, 1994.
- [11] N. R. Pal, J. C. Bezdek, and E. C. K. Tsao, "Generalized clustering networks and Kohonen's self-organizing scheme," *IEEE Trans. Neural Networks*, vol. 4, pp. 549-558, 1993.
- [12] A. I. Gonzalez, M. Graua, and A. D'Anjou, "An analysis of the GLVQ algorithm," *IEEE Trans. Neural Networks*, vol. 6, pp. 1012-1016, 1995.
- [13] N. B. Karayiannis and P.-I. Pai, "A fuzzy algorithm for learning vector quantization," in *Proc. 1994 IEEE Int. Conf. Syst., Man, Cybern.*, Piscataway, NJ, 1994, pp. 126-131.
- [14] ———, "A family of fuzzy algorithms for learning vector quantization," in *Intelligent Engineering Systems through Artificial Neural Networks*, H. Dagli, B. R. Fernandez, J. Ghosh, and R. T. Kumara, Eds., vol. 4. New York: Amer. Soc. Mech. Eng., 1994, pp. 219-224.
- [15] E. Anderson, "The IRIS's of the Gaspe peninsula," *Bull. Amer. Iris Soc.*, vol. 59, pp. 2-5, 1935.
- [16] T. Hunsberger and P. Ajjimarangsee, "Parallel self-organizing feature maps for unsupervised pattern recognition," *Int. J. General Syst.*, vol. 16, pp. 357-372, 1989.
- [17] D. C. Park and I. Dagher, "Gradient-based fuzzy c-means (GBFCM) algorithm," in *Proc. IEEE Int. Conf. Neural Networks*, vol. 3, pp. 1626-1631, 1994.
- [18] F. L. Chung and T. Lee, "Fuzzy competitive learning," *Neural Networks*, vol. 7, no. 3, pp. 539-551, 1992.

Nicolaos B. Karayiannis (S'85-M'91), for a photograph and biography, see p. 426 of the March 1996 issue of this TRANSACTIONS.



James C. Bezdek (M'80-SM'90-F'92) received the Ph.D. degree from Cornell University, Ithaca, NY, in 1973.

He is a Professor of Computer Science at the University of West Florida, Pensacola. His interests include pattern recognition, fishing, computational neural networks, skiing, image processing, blues music, medical computing, and motorcycles.

Dr. Bezdek is Editor-in-Chief of the IEEE TRANSACTIONS ON FUZZY SYSTEMS.



Nikhil R. Pal (M'91) received the M.Tech. and Ph.D. degrees in computer science from the Indian Statistical Institute, Calcutta, in 1984 and 1991, respectively.

He is an Associate Professor in the Machine Intelligence Unit of Physics. He was with the Hindustan Motors Ltd., W. B., from 1984 to 1985 and with the Duplop India Ltd., W. B., from 1985 to 1987. In 1987 he joined the Computer Science Unit of the Indian Statistical Institute, Calcutta. During August 1991 to February 1993 and July 1994 to December

1994 he visited the University of West Florida, Pensacola. He was a Guest Lecturer at the University of Calcutta. His research interests include image processing, pattern recognition, fuzzy sets and systems, uncertainty measures, genetic algorithms, and neural networks.

Dr. Pal is an Associate Editor of the IEEE TRANSACTIONS ON FUZZY SYSTEMS and the *International Journal of Approximate Reasoning*.



Richard J. Hathaway (M'93) received the B.S. degree in applied mathematics from the University of Georgia, Athens, in 1979 and the Ph.D. degree in mathematical sciences from Rice University, Houston, TX, in 1983.

He is a Professor in the Mathematics and Computer Science Department of Georgia Southern University, Statesboro. His research interests include pattern recognition, statistical computing, and nonlinear optimization.

Pin-I Pai, for a photograph and biography, see this issue, p. 1211.