

# An evidence-theoretic $k$ -NN rule with parameter optimization<sup>1</sup>

L. M. Zouhal<sup>\*,†</sup> and T. Denœux<sup>\*</sup>

\* Université de Technologie de Compiègne - U.M.R. CNRS 6599 Heudiasyc  
BP 20529 - F-60205 Compiègne cedex - France  
email: [Thierry.Denoeux@hds.utc.fr](mailto:Thierry.Denoeux@hds.utc.fr)

† Centre International des Techniques Informatiques  
Lyonnaise des Eaux

<sup>1</sup>Technical Report Heudiasyc 97/46. To appear in *IEEE Transactions on Systems, Man and Cybernetics*.

## Abstract

This paper presents a learning procedure for optimizing the parameters in the evidence-theoretic  $k$ -nearest neighbor rule, a pattern classification method based on the Dempster-Shafer theory of belief functions. In this approach, each neighbor of a pattern to be classified is considered as an item of evidence supporting certain hypotheses concerning the class membership of that pattern. Based on this evidence, basic belief masses are assigned to each subset of the set of classes. Such masses are obtained for each of the  $k$  nearest neighbors of the pattern under consideration and aggregated using the Dempster's rule of combination. In many situations, this method was found experimentally to yield lower error rates than other methods using the same information. However, the problem of tuning the parameters of the classification rule was so far unresolved. In this paper, we propose to determine optimal or near-optimal parameter values from the data by minimizing an error function. This refinement of the original method is shown experimentally to result in substantial improvement of classification accuracy.

# 1 Introduction

In the classical approach to statistical pattern recognition, the entities to be classified are assumed to be selected by some form of random experiment. The feature vector describing each entity is then a random vector with well-defined – though unknown – probability density function depending on the pattern category. Based on these densities and on the prior probability of each class, posterior probabilities can be defined, and the optimal Bayes decision rule can then theoretically be used for classifying an arbitrary pattern with minimal expected risk. Since the class-conditional densities and prior probabilities are usually unknown, they need to be estimated from the data. A lot of methods have been proposed for building consistent estimators of the posterior probabilities under various assumptions. However, for finite sample size, the resulting estimates generally do not provide a faithful representation of the fundamental uncertainty pertaining to the class of a pattern to be classified. For example, if only a relatively small number of training vectors is available, and a new pattern is encountered that is very dissimilar from all previous patterns, the uncertainty is quite high and this situation of near-ignorance is not reflected by the outputs of a conventional parametric or non parametric statistical classifier, whose principle fundamentally relies on asymptotic assumptions. This problem is particularly acute in situations in which decisions need to be made based on weak information, such as commonly encountered in system diagnosis applications, for example.

As an attempt to provide an answer to the above problem, it was recently suggested to re-formulate the pattern classification problem by considering the following question: Given a training set of finite size containing feature vectors with known (or partly known) classification, and a suitable distance measure, how to characterize the uncertainty pertaining to the class of a new pattern ? In a recent paper [2], an answer to this question was proposed based on the Dempster-Shafer theory of evidence [12]. The approach consists in considering each neighbor of a pattern to be classified as an item of evidence supporting certain hypotheses concerning the class membership of that pattern. Based on this evidence, basic belief masses are assigned to each subset of the set of classes. Such masses are obtained for each of the  $k$  nearest neighbors of the pattern under consideration and aggregated using the Dempster's rule of combination. Given a finite set of actions and losses associated to each action and each class, decisions can then be made by using some generalization of the Bayes decision theory.

In many situations, this method was found experimentally to yield lower error rates than other methods based on the same information. However, the problem of optimizing the parameters involved in the classification rule was so far unresolved. In this paper, we propose to determine optimal or near-optimal parameter values from the data by minimizing a certain error function. This refinement of the original method is shown experimentally to result in substantial improvement of classification accuracy.

The paper is organized as follows. The evidence-theoretic  $k$ -NN rule is first recalled in Section 2. The basic concepts of the Dempster-Shafer theory are assumed to be

known to the reader who is invited to refer to [12] and [14] for detailed presentations, and to [2] for a short introduction. Section 3 describes the learning procedure, as well as an approximation to it allowing to near-optimize the error function very efficiently. Simulation results are then presented in Section 4, and Section 5 concludes the paper.

## 2 The evidence-theoretic $k$ -NN rule

We consider the problem of classifying entities into  $M$  categories or classes. The set of classes is denoted by  $\Omega = \{\omega_1, \dots, \omega_M\}$ . The available information is assumed to consist in a training set  $\mathcal{T} = \{(\mathbf{x}^{(1)}, \omega^{(1)}), \dots, (\mathbf{x}^{(N)}, \omega^{(N)})\}$  of  $N$   $n$ -dimensional patterns  $\mathbf{x}^{(i)}, i = 1, \dots, N$  and their corresponding class labels<sup>1</sup>  $\omega^{(i)}, i = 1, \dots, N$  taking values in  $\Omega$ . The similarity between patterns is assumed to be correctly measured by a certain distance function  $d(\cdot, \cdot)$ .

Let  $\mathbf{x}$  be a new vector to be classified on the basis of the information contained in  $\mathcal{T}$ . Each pair  $(\mathbf{x}^{(i)}, \omega^{(i)})$  constitutes a distinct item of evidence regarding the class membership of  $\mathbf{x}$ . If  $\mathbf{x}$  is “close” to  $\mathbf{x}^{(i)}$  according to the relevant metric  $d$ , then one will be inclined to believe that both vectors belong to the same class. On the contrary, if  $d(\mathbf{x}, \mathbf{x}^{(i)})$  is very large, then the consideration of  $\mathbf{x}^{(i)}$  will leave us in a situation of almost complete ignorance concerning the class of  $\mathbf{x}$ . Consequently, this item of evidence may be postulated to induce a basic belief assignment (BBA)  $m(\cdot|\mathbf{x}^{(i)})$  over  $\Omega$  defined by:

$$m(\{\omega_q\}|\mathbf{x}^{(i)}) = \alpha\phi_q(d^{(i)}) \quad (1)$$

$$m(\Omega|\mathbf{x}^{(i)}) = 1 - \alpha\phi_q(d^{(i)}) \quad (2)$$

$$m(A|\mathbf{x}^{(i)}) = 0, \quad \forall A \in 2^\Omega \setminus \{\Omega, \{\omega_q\}\} \quad (3)$$

where  $d^{(i)} = d(\mathbf{x}, \mathbf{x}^{(i)})$ ,  $\omega_q$  is the class of  $\mathbf{x}^{(i)}$  ( $\omega^{(i)} = \omega_q$ ),  $\alpha$  is a parameter such that  $0 < \alpha < 1$  and  $\phi_q$  is a decreasing function verifying  $\phi_q(0) = 1$  et  $\lim_{d \rightarrow \infty} \phi_q(d) = 0$ . Note that  $m(\cdot|\mathbf{x}^{(i)})$  reduces to the vacuous belief function ( $m(\Omega|\mathbf{x}^{(i)}) = 1$ ) when the distance between  $\mathbf{x}$  and  $\mathbf{x}^{(i)}$  tends to infinity, reflecting a state of total ignorance. When  $d$  denotes the Euclidean distance, a rational choice for  $\phi_q$  was shown in [4] to be:

$$\phi_q(d) = \exp(-\gamma_q d^2) \quad (4)$$

$\gamma_q$  being a positive parameter associated to class  $\omega_q$ .

As a result of considering each training pattern in turn, we obtain  $N$  BBAs that can be combined using the Dempster’s rule of combination to form a resulting BBA  $m$  synthesizing one’s final belief regarding the class of  $\mathbf{x}$ :

$$m = m(\cdot|\mathbf{x}^{(1)}) \oplus \dots \oplus m(\cdot|\mathbf{x}^{(N)}) \quad (5)$$

---

<sup>1</sup>In this paper, we assume for simplicity the class of each training vector to be known with certainty. The more general situation in which the training set is only imperfectly labeled has been introduced in [2]. However, the problem of optimizing the parameters in the general case is not completely solved yet (see Section 5).

Since those training patterns situated far from  $\mathbf{x}$  actually provide very little information, it is sufficient to consider the  $k$  nearest neighbors of  $\mathbf{x}$  in this sum. An alternative definition of  $m$  is therefore:

$$m = m(\cdot|\mathbf{x}^{(i_1)}) \oplus \dots \oplus m(\cdot|\mathbf{x}^{(i_k)}) \quad (6)$$

where  $I_k = \{i_1, \dots, i_k\}$  contains the indices of the  $k$  nearest neighbors of  $\mathbf{x}$  in  $\mathcal{T}$ .

Adopting this latter definition,  $m$  can be shown [2] to have the following expression:

$$m(\{\omega_q\}) = \frac{1}{K} \left( 1 - \prod_{i \in I_{k,q}} (1 - \alpha\phi_q(d^{(i)})) \right) \prod_{r \neq q} \prod_{i \in I_{k,r}} (1 - \alpha\phi_r(d^{(i)})) \quad (7)$$

$$\forall q \in \{1, \dots, M\}$$

$$m(\Omega) = \frac{1}{K} \prod_{r=1}^M \prod_{i \in I_{k,r}} (1 - \alpha\phi_r(d^{(i)})) \quad (8)$$

where  $I_{k,q}$  is the subset of  $I_k$  corresponding to those neighbors of  $\mathbf{x}$  belonging to class  $\omega_q$  and  $K$  is a normalizing factor. Hence, the focal elements of  $m$  are singletons and the whole frame  $\Omega$ . Consequently, the credibility and the plausibility of each class  $\omega_q$  are respectively equal to:

$$\text{bel}(\{\omega_q\}) = m(\{\omega_q\}) \quad (9)$$

$$\text{pl}(\{\omega_q\}) = m(\{\omega_q\}) + m(\Omega) \quad (10)$$

The pignistic probability distribution as defined by Smets [13] is given by:

$$\text{BetP}(\{\omega_q\}) = \sum_{\{A \subseteq \Omega | \omega_q \in A\}} \frac{m(A)}{|A|} = m(\{\omega_q\}) + \frac{m(\Omega)}{M} \quad (11)$$

for  $q = 1, \dots, M$ . Let us now assume that, based on this evidential corpus, a decision has to be made regarding the assignment of  $\mathbf{x}$  to a class, and let us denote by  $\alpha_q$  the action of assigning  $\mathbf{x}$  to class  $\omega_q$ . Let us further assume that the loss incurred in case of a wrong classification is equal to 1, while the loss corresponding to a correct classification is equal to 0. Then, the lower and the upper expected losses [3] associated to action  $\alpha_q$  are respectively equal to:

$$R_*(\alpha_q|\mathbf{x}) = 1 - \text{pl}(\{\omega_q\}) \quad (12)$$

$$R^*(\alpha_q|\mathbf{x}) = 1 - \text{bel}(\{\omega_q\}) \quad (13)$$

The expected loss relative to the pignistic distribution is:

$$R_{bet}(\alpha_q|\mathbf{x}) = 1 - \text{BetP}(\{\omega_q\}) \quad (14)$$

Given the particular form of  $m$ , the three strategies consisting in minimizing  $R_*$ ,  $R^*$  and  $R_{bet}$  lead to the same decision in that case: the pattern is assigned to the class with maximum belief assignment. Other decision strategies including the possibility of pattern rejection as well as the existence of unknown classes are studied in [3, 5].

### 3 Parameter optimization

#### 3.1 The approach

In the above description of the evidence-theoretic  $k$ -NN rule, we left open the question of the choice of parameters  $\alpha$  and  $\gamma = (\gamma_1, \dots, \gamma_q)^t$  appearing in Equations 1 and 4. Whereas the value of  $\alpha$  proves in practice not to be too critical, the tuning of the other parameters was found experimentally to have significant influence on classification accuracy. In [2], it was proposed to set  $\alpha = 0.95$  and  $\gamma_q$  to the inverse of the mean distance between training patterns belonging to class  $\omega_q$ . Although this heuristic yields good results on average, the efficiency of the classification procedure can be improved if these parameters are determined as the values optimizing a performance criterion. Such a criterion can be defined as follows.

Let us consider a training pattern  $\mathbf{x}^{(\ell)}$  belonging to class  $\omega_q$ . The class membership of  $\mathbf{x}$  can be encoded as a vector  $\mathbf{t}^{(\ell)} = (t_1^{(\ell)}, \dots, t_M^{(\ell)})^t$  of  $M$  binary indicator variables  $t_j^{(\ell)}$  defined by  $t_j^{(\ell)} = 1$  if  $j = q$  and  $t_j^{(\ell)} = 0$  otherwise. By considering the  $k$  nearest neighbors of  $\mathbf{x}^{(\ell)}$  in the training set, one obtains a “leave-one-out” BBA  $m^{(\ell)}$  characterizing one’s belief concerning the class of  $\mathbf{x}^{(\ell)}$  if this pattern was to be classified using other training patterns. Based on  $m^{(\ell)}$ , an output vector  $\mathbf{P}^{(\ell)} = (\text{BetP}^{(\ell)}(\{\omega_1\}), \dots, \text{BetP}^{(\ell)}(\{\omega_M\}))^t$  of pignistic probabilities can be computed,  $\text{BetP}^{(\ell)}$  being the pignistic probability distribution associated to  $m^{(\ell)}$ . Ideally, vector  $\mathbf{P}^{(\ell)}$  should as “close” as possible to vector  $\mathbf{t}^{(\ell)}$ , closeness being defined, for example, according to the squared error  $E(\mathbf{x}^{(\ell)})$ :

$$E(\mathbf{x}^{(\ell)}) = (\mathbf{P}^{(\ell)} - \mathbf{t}^{(\ell)})^t (\mathbf{P}^{(\ell)} - \mathbf{t}^{(\ell)}) = \sum_{q=1}^M (P_q^{(\ell)} - t_q^{(\ell)})^2 \quad (15)$$

The mean squared error over the whole training set  $\mathcal{T}$  of size  $N$  is finally equal to:

$$E = \frac{1}{N} \sum_{\ell=1}^N E(\mathbf{x}^{(\ell)}) \quad (16)$$

Function  $E$  can be used as a cost function for tuning the parameter vector  $\gamma$ . The analytical expression for the gradient of  $E(\mathbf{x}^{(\ell)})$  with respect to  $\gamma$  can be calculated, allowing the parameters  $\gamma_q$  to be determined iteratively by a gradient search procedure (see Appendix A). Alternatively, the minimum of function  $E$  can be approximated in one step for large  $N$  using the approach described in the sequel.

#### 3.2 One-step procedure

For an arbitrary training pattern  $\mathbf{x}^{(\ell)}$  and fixed parameters, vector  $\mathbf{P}^{(\ell)}$  can be regarded as a function of two vectors:

1. a vector  $\mathbf{d}^2 = (d^{(i_1)2}, \dots, d^{(i_k)2})^t$  of squared distances between  $\mathbf{x}^{(\ell)}$  and its  $k$  nearest neighbors, and

2. a vector  $\boldsymbol{\omega}$  containing the class labels of these neighbors.

For small  $k$  and large  $N$ ,  $\mathbf{d}^2$  can be assumed to close to zero<sup>2</sup>, allowing each component  $P_q^{(\ell)}$  to be approximated by Taylor series expansion around 0 up to the first order:

$$P_q^{(\ell)}(\mathbf{d}^2) \cong P_q^{(\ell)}(\mathbf{0}) + \nabla_{\mathbf{d}^2} P_q^{(\ell)t} \Big|_{\mathbf{d}^2=\mathbf{0}} \mathbf{d}^2 \quad (17)$$

The first term in this expression can be readily obtained from Equations 7 and 8 by setting  $d^{(i)}$  to 0 for all  $i$ , which leads to:

$$P_q^{(\ell)}(\mathbf{0}) = \frac{(1-\alpha)^k}{K} \left( (1-\alpha)^{-k_q} - 1 + \frac{1}{M} \right) \quad (18)$$

where  $k_q$  is the number of neighbors of  $\mathbf{x}^{(\ell)}$  in class  $\omega_q$  and

$$K = \frac{(1-\alpha)^k}{\sum_{q=1}^M (1-\alpha)^{-k_q} - k + 1} \quad (19)$$

The computation of the first order term:

$$\nabla_{\mathbf{d}^2} P_q^{(\ell)t} \Big|_{\mathbf{d}^2=\mathbf{0}} \mathbf{d}^2 = \sum_{i=1}^k \frac{\partial P_q^{(\ell)}}{\partial d^{(i)2}}(\mathbf{0}) d^{(i)2} \quad (20)$$

is more involved (see appendix B). This term can be shown to be of the form  $\mathbf{A}^{(\ell)}\boldsymbol{\gamma}$ , where  $\mathbf{A}^{(\ell)}$  is a square matrix of size  $M$ . As a consequence, both  $E(\mathbf{x}^{(\ell)})$  and  $E$  can be approximated by quadratic forms of  $\boldsymbol{\gamma}$ , which allows the minimum to be approached directly by solving a system of linear equations.

Figures 1 and 2 show the quality of this approximation in the case of two Gaussian classes with mean vectors  $\boldsymbol{\mu}_1 = (2 \ 0)^t$  and  $\boldsymbol{\mu}_2 = (-2 \ 0)^t$ , respectively, and covariance matrices  $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \mathbf{I}$ . The data set contained 500 samples of each class. Displayed are the mean squared error as a function of  $(\gamma_1, \gamma_2)$  (Figure 1) and its quadratic approximation (Figure 2), for  $k = 5$ . The minima of the two functions differ by less than 0.0007 %, which proves the relevance of the approximation in that case. Note that the quality of the approximation depends on both  $k$  and  $N$ .

## 4 Numerical experiments

The performances of the above methods were compared to those of the voting  $k$ -NN rule with randomly resolved ties, the distance-weighted  $k$ -NN rule [6], the fuzzy  $k$ -NN rule proposed by Keller [9], and the evidence-theoretic rule without parameter optimization [2]. Experiments were carried out on a set of standard artificial and

---

<sup>2</sup>This assumption is justified by the following result [1]: Regarding the training set as a sample drawn from some probability distribution, the  $k$ -th nearest neighbor of  $\mathbf{x}^{(\ell)}$  converges to  $\mathbf{x}^{(\ell)}$  with probability one as the sample size increases with  $k$  fixed.

real-world benchmark classification tasks. The main characteristics of the used data sets are summarized in Table 1.

Data sets  $B_1$  and  $B_2$  were generated using a method proposed in [7]. The data consists in three Gaussian classes in 10 dimensions, which diagonal covariance matrices  $\mathbf{D}_1$ ,  $\mathbf{D}_2$  and  $\mathbf{D}_3$ . The  $i$ -th diagonal element  $D_{qi}$  of  $\mathbf{D}_q$  is defined as a function of two parameters  $a$  and  $b$ :

$$\begin{aligned} D_{1i}(a, b) &= a + (b - a) \frac{i - 1}{n - 1} \\ D_{2i}(a, b) &= a + (b - a) \frac{n - i}{n - 1} \\ D_{3i}(a, b) &= a + (b - a) \min \left( i, \frac{n - i}{n/2} \right) \end{aligned}$$

where  $n = 10$  is the input dimension. The mean vectors and covariance matrices (with  $n = 10$ ) for the three classes were:

$$\begin{array}{ccc} \boldsymbol{\mu}_1 = (0, \dots, 0) & \boldsymbol{\mu}_2 = (1, \dots, 1) & \boldsymbol{\mu}_3 = (1, -1, \dots, 1, -1) \\ \mathbf{D}_1(1, 10) & \mathbf{D}_2(10, 1) & \mathbf{D}_3(1, 10) \end{array}$$

The ionosphere data set (Ion) was collected by a radar system and consists of phased array of 16 high-frequency antennas with a total transmitted power of the order of 6.4 kilowatts [11]. The targets were free electrons in the ionosphere. ‘‘Good’’ radar returns are those showing evidence of some type of structure in the ionosphere. ‘‘Bad’’ returns are those that do not.

The vehicle data set (Veh) was collected from silhouettes by the HIPS (Hierarchical Image Processing System) extension BINATTS. Four model vehicles were used for the experiment: bus, Chevrolet van, Saab 9000 and Opel Manta 400. The data was used to distinguish 3D objects within a 2D silhouette of the objects [11].

The sonar data were used by Gorman and Sejnowski in a study of the classification of sonar signals using a neural network [8]. The task is to discriminate between sonar signals bounced off a metal cylinder and those bounced off a roughly cylindrical rock.

Test error rates are represented as a function of  $k$  in Figures 3 to 12. The results for synthetic data are averages over 6 independent training sets. Test error rates obtained for the values of  $k$  yielding the best classification of training vectors are presented in Table 2.

As can be seen from these results, the evidence-theoretic rule with optimized  $\gamma$  presented in this paper always performed as well or better than the four other rules tested, and significantly improved at the 90 % confidence level over the evidence-theoretic rule with fixed  $\gamma$  on data sets  $B_1$  and  $B_2$  when considering the best results obtained for  $1 \leq k \leq 40$ . However, the most distinctive feature of this rule seems to be its robustness with respect to the number  $k$  of neighbors taken in consideration (Figures 3 to 12). By optimizing  $\gamma$ , the method learns to discard those neighbors whose distance to the pattern under consideration is too high. Practically, this property is of great importance since it relieves the designer of the system from the burden of



searching for the optimal value of  $k$ . When the number of training patterns is large, then the amount of computation may be further reduced by adopting the approximate one-step procedure for optimizing  $\gamma$  which gives reasonably good results for small  $k$  (Figures 3, 5, 7, 9 and 11). However, the use of the exact procedure should be preferred for small and medium-sized training sets.

## 5 Concluding remarks

A technique for optimizing the parameters in the evidence-theoretic  $k$ -NN rule has been presented. The classification rule obtained by this method has proved superior to the voting, distance-weighted and fuzzy rules on a number of benchmark problems. A remarkable property achieved with this approach is the relative insensitivity of the results to the choice of  $k$ .

The method can be generalized in several ways. First of all, one can assume a more general metric than the Euclidean one considered so far, and apply the principles described in this paper to search for the optimal metric [15]. For instance, let  $\Sigma_q$  be a positive definite diagonal matrix with diagonal elements  $\gamma_{q,1}, \dots, \gamma_{q,n}$ . The distance between an input vector  $\mathbf{x}$  and a learning vector  $\mathbf{x}^{(i)}$  belonging to class  $\omega_q$  can be defined as:

$$\begin{aligned} d(\mathbf{x}, \mathbf{x}^{(i)}) &= (\mathbf{x} - \mathbf{x}^{(i)})^t \Sigma_q (\mathbf{x} - \mathbf{x}^{(i)}) \\ &= \sum_{j=1}^n \gamma_{q,j} (x_j - x_j^{(i)})^2 \end{aligned}$$

The parameters  $\gamma_{q,j}$  for  $1 \leq q \leq M$  and  $1 \leq j \leq n$  can then be optimized using exactly the same approach as described in this paper, which may in some cases result in further improvement of classification results. A more general form could even be assumed for  $\Sigma_q$ , with however the risk of a dramatic increase in the number of parameters for large input dimensions.

More fundamentally, the method can also be extended to handle the more general situation in which the class membership of training patterns is itself affected by uncertainty. For example, let us assume that the class of each training pattern  $\mathbf{x}^{(i)}$  is only known to lie in a subset  $A^{(i)}$  of  $\Omega$  (such a situation may typically arise, e.g., in medical diagnosis problems in which some records in a database are related to patients for which only a partial or uncertain diagnosis is available). A natural extension of Equations 1–3 is then:

$$\begin{aligned} m(A^{(i)}|\mathbf{x}^{(i)}) &= \alpha \phi(d^{(i)}) \\ m(\Omega|\mathbf{x}^{(i)}) &= 1 - \alpha \phi(d^{(i)}) \\ m(B|\mathbf{x}^{(i)}) &= 0, \quad \forall B \in 2^\Omega \setminus \{\Omega, A^{(i)}\} \end{aligned}$$

with  $\phi(d) = \exp(-\gamma d^2)$ ,  $\gamma$  being a positive parameter (note that we cannot define a separate parameter for each class in this case, since the class of  $\mathbf{x}^{(i)}$  is only partially known). The BBAs defined in that way correspond to simple belief functions and can

be combined in linear time with respect to the number of classes. For optimizing  $\gamma$ , the error criterion defined in Equation 15 has to be generalized in some way. With the same notations as in Section 3.1, a possible expression for the error concerning pattern  $\mathbf{x}^{(\ell)}$  is:

$$E(\mathbf{x}^{(\ell)}) = \left( \text{BetP}^{(\ell)}(A^{(\ell)}) - 1 \right)^2,$$

which reflects the fact that the pignistic probability of  $\mathbf{x}^{(\ell)}$  belonging to  $A^{(\ell)}$ , given the other training patterns, should be as high as possible. The value of  $\gamma$  minimizing the mean error may then be determined using an iterative search procedure. Experiments with this approach are currently under way and will be reported in future publications.

## A Computation of the derivatives of $E$ w.r.t. $\gamma_q$

Let  $\mathbf{x}^{(\ell)}$  be a training pattern and  $m^{(\ell)}$  the BBA obtained by classifying  $\mathbf{x}^{(\ell)}$  using its  $k$  nearest neighbors in the training set. Function  $m^{(\ell)}$  is computed according to Equations 7 and 8:

$$m^{(\ell)}(\{\omega_q\}) = \frac{1}{K} \left( 1 - \prod_{i \in I_{k,q}^{(\ell)}} (1 - \alpha \phi_q(d^{(\ell,i)})) \right) \prod_{r \neq q} \prod_{i \in I_{k,r}^{(\ell)}} (1 - \alpha \phi_r(d^{(\ell,i)})) \quad (21)$$

$\forall q \in \{1, \dots, M\}$

$$m^{(\ell)}(\Omega) = \frac{1}{K} \prod_{r=1}^M \prod_{i \in I_{k,r}^{(\ell)}} (1 - \alpha \phi_r(d^{(\ell,i)})) \quad (22)$$

where  $I_{k,r}^{(\ell)}$  denotes the set of indices of the  $k$  nearest neighbors of pattern  $\mathbf{x}^{(\ell)}$  in class  $\omega_r$ ,  $d^{(\ell,i)}$  is the distance between  $\mathbf{x}^{(\ell)}$  and  $\mathbf{x}^{(i)}$ , and  $K$  is a normalizing factor. In the following, we shall assume that:

$$\phi_q(d^{(\ell,i)}) = \exp(-\gamma_q d^{(\ell,i)2}) \quad (23)$$

which will simply be denoted by  $\phi_q^{(\ell,i)}$ . To simplify the calculations, we further introduce the *unnormalized* orthogonal sum [13]  $m^{(\ell)'}$  defined as  $m^{(\ell)'}(A) = K m^{(\ell)}(A)$  for all  $A \subseteq \Omega$ . We also denote as  $\bar{m}^{(\ell,i)}$  the unnormalized orthogonal sum of the  $m(\cdot|\mathbf{x}^{(j)})$  for all  $j \in I_k^{(\ell)}$ ,  $j \neq i$ , that is  $m^{(\ell)'} = \bar{m}^{(\ell,i)} \wedge m(\cdot|\mathbf{x}^{(i)})$ , where  $\wedge$  denotes the unnormalized orthogonal sum operator. More precisely, we have:

$$m^{(\ell)'}(\{\omega_q\}) = m^{(\ell)}(\{\omega_q\}|\mathbf{x}^{(i)}) (\bar{m}^{(\ell,i)}(\{\omega_q\}) + \bar{m}^{(\ell,i)}(\Omega)) + m^{(\ell)}(\Omega|\mathbf{x}^{(i)}) \bar{m}^{(\ell,i)}(\{\omega_q\}) \quad \forall q \in \{1, \dots, M\} \quad (24)$$

$$m^{(\ell)' }(\Omega) = \bar{m}^{(\ell,i)}(\Omega) m(\Omega|\mathbf{x}^{(i)}) \quad (25)$$

The error for pattern  $\mathbf{x}^{(\ell)}$  is:

$$E(\mathbf{x}^{(\ell)}) = \sum_{q=1}^M (P_q^{(\ell)} - t_q^{(\ell)})^2 \quad (26)$$

where  $\mathbf{t}^{(\ell)}$  is the class membership vector for pattern  $\mathbf{x}^{(\ell)}$  and  $P_q^{(\ell)}$  is the pignistic probability of class  $\omega_q$  computed from  $m^{(\ell)}$  as  $P_q^{(\ell)} = m^{(\ell)}(\{\omega_q\}) + \frac{m^{(\ell)}(\Omega)}{M}$ .

The derivative of  $E(\mathbf{x}^{(\ell)})$  with respect to each parameter  $\gamma_q$  can be computed as:

$$\frac{\partial E(\mathbf{x}^{(\ell)})}{\partial \gamma_q} = \sum_{i \in I_{k,q}^{(\ell)}} \frac{\partial E(\mathbf{x}^{(\ell)})}{\partial \phi_q^{(\ell,i)}} \frac{\partial \phi_q^{(\ell,i)}}{\partial \gamma_q} \quad (27)$$

with

$$\frac{\partial E(\mathbf{x}^{(\ell)})}{\partial \phi_q^{(\ell,i)}} = \sum_{r=1}^M \frac{\partial E(\mathbf{x}^{(\ell)})}{\partial P_r^{(\ell)}} \frac{\partial P_r^{(\ell)}}{\partial \phi_q^{(\ell,i)}} \quad (28)$$

$$= \sum_{r=1}^M 2(P_r^{(\ell)} - t_r^{(\ell)}) \left[ \frac{\partial m^{(\ell)}(\{\omega_r\})}{\partial \phi_q^{(\ell,i)}} + \frac{1}{M} \frac{\partial m^{(\ell)}(\Omega)}{\partial \phi_q^{(\ell,i)}} \right] \quad (29)$$

and

$$\frac{\partial \phi_q^{(\ell,i)}}{\partial \gamma_q} = -d^{(\ell,i)} 2 \phi_q^{(\ell,i)} \quad (30)$$

The derivatives in Equation 29 can be computed as:

$$\frac{\partial m^{(\ell)}(\{\omega_r\})}{\partial \phi_q^{(\ell,i)}} = \frac{\partial}{\partial \phi_q^{(\ell,i)}} \left( \frac{m^{(\ell)'(\{\omega_r\})}{K} \right) \quad (31)$$

$$= \frac{1}{K^2} \left[ K \frac{\partial m^{(\ell)'(\{\omega_r\})}{\partial \phi_q^{(\ell,i)}} - m^{(\ell)'(\{\omega_r\})} \frac{\partial K}{\partial \phi_q^{(\ell,i)}} \right] \quad (32)$$

$$\frac{\partial m^{(\ell)}(\Omega)}{\partial \phi_q^{(\ell,i)}} = \frac{1}{K^2} \left[ K \frac{\partial m^{(\ell)'(\Omega)}{\partial \phi_q^{(\ell,i)}} - m^{(\ell)'(\Omega)} \frac{\partial K}{\partial \phi_q^{(\ell,i)}} \right] \quad (33)$$

$$\frac{\partial K}{\partial \phi_q^{(\ell,i)}} = \sum_{r=1}^M \frac{\partial m^{(\ell)'(\{\omega_r\})}{\partial \phi_q^{(\ell,i)}} + \frac{\partial m^{(\ell)'(\Omega)}{\partial \phi_q^{(\ell,i)}} \quad (34)$$

Finally:

$$\frac{\partial m^{(\ell)'(\{\omega_r\})}{\partial \phi_q^{(\ell,i)}} = (\bar{m}^{(\ell,i)}(\{\omega_r\}) + \bar{m}^{(\ell,i)}(\Omega)) \alpha \delta_{rq} - \alpha \bar{m}^{(\ell,i)}(\Omega) \quad (35)$$

where  $\delta$  is the Kronecker symbol, and

$$\frac{\partial m^{(\ell)'(\Omega)}{\partial \phi_q^{(\ell,i)}} = -\bar{m}^{(\ell,i)}(\Omega) \quad (36)$$

which completes the calculation of the gradient of  $E(\mathbf{x}^{(\ell)})$  w.r.t.  $\gamma_q$ .

To account for the constraint  $\gamma_q \geq 0$ , we introduce new parameters  $\mu_q$  ( $q = 1, \dots, M$ ) such that:

$$\gamma_q = (\mu_q)^2 \quad (37)$$

and we compute  $\frac{\partial E(\mathbf{x}^{(\ell)})}{\partial \mu_q}$  as:

$$\frac{\partial E(\mathbf{x}^{(\ell)})}{\partial \mu_q} = \frac{\partial E(\mathbf{x}^{(\ell)})}{\partial \gamma_q} \frac{\partial \gamma_q}{\partial \mu_q} = 2\mu_q \frac{\partial E(\mathbf{x}^{(\ell)})}{\partial \gamma_q} \quad (38)$$

## B Linearization

We consider the expansion around  $\mathbf{d}^2 = \mathbf{0}$  of  $P_q^{(\ell)}$  by a Taylor series up to the first order:

$$P_q^{(\ell)}(\mathbf{d}^2) \cong P_q^{(\ell)}(\mathbf{0}) + (\nabla_{\mathbf{d}^2} P_q^{(\ell)}(\mathbf{0}))^t \mathbf{d}^2 \quad (39)$$

where  $P_q^{(\ell)}(\mathbf{0})$  is given by Equations 18. In the following, we shall compute the first order term in the above equation, and deduce from that result a method for determining an approximation to the optimal parameter vector. To simplify the notation, the superscript  $(\ell)$  will be omitted from the following calculations.

As a result of the definition of the pignistic probability, we have:

$$\frac{\partial P_q}{\partial d^{(i)2}} = \frac{\partial m(\{\omega_q\})}{\partial d^{(i)2}} + \frac{1}{M} \frac{\partial m(\Omega)}{\partial d^{(i)2}} \quad (40)$$

The derivatives of  $m(\{\omega_q\})$  and  $m(\Omega)$  can be more conveniently expressed as a function of the unnormalized BBA  $m'$ :

$$\frac{\partial m(\{\omega_q\})}{\partial d^{(i)2}} = \frac{1}{K^2} \left( K \frac{\partial m'(\{\omega_q\})}{\partial d^{(i)2}} - m'(\{\omega_q\}) \frac{\partial K}{\partial d^{(i)2}} \right) \quad (41)$$

$$\frac{\partial m(\Omega)}{\partial d^{(i)2}} = \frac{1}{K^2} \left( K \frac{\partial m'(\Omega)}{\partial d^{(i)2}} - m'(\Omega) \frac{\partial K}{\partial d^{(i)2}} \right) \quad (42)$$

To compute  $\frac{\partial m'(\{\omega_q\})}{\partial d^{(i)2}}$  and  $\frac{\partial m'(\Omega)}{\partial d^{(i)2}}$  we need to distinguish two cases:

**Case 1:**  $i \in I_{k,q}$ . We then have:

$$\begin{aligned} \frac{\partial m'(\{\omega_q\})}{\partial d^{(i)2}} &= -\alpha \gamma_q \exp(-\gamma_q d^{(i)2}) \prod_{j \in I_{k,q}, j \neq i} (1 - \alpha \exp(-\gamma_q d^{(j)2})) \\ &\quad \times \prod_{r \neq q} \prod_{j \in I_{k,r}} (1 - \alpha \exp(-\gamma_q d^{(j)2})) \end{aligned} \quad (43)$$

$$\frac{\partial m'(\Omega)}{\partial d^{(i)2}} = \alpha \gamma_q \exp(-\gamma_q d^{(i)2}) \prod_{r=1}^M \prod_{j \in I_{k,r}, j \neq i} (1 - \alpha \exp(-\gamma_r d^{(j)2})) \quad (44)$$

Setting all distances to 0 in the above equations, we have:

$$\left. \frac{\partial m'(\{\omega_q\})}{\partial d^{(i)2}} \right|_{\mathbf{d}^2=\mathbf{0}} = -\alpha \gamma_q (1 - \alpha)^{k-1} \quad (45)$$

$$\left. \frac{\partial m'(\Omega)}{\partial d^{(i)2}} \right|_{\mathbf{d}^2=\mathbf{0}} = \alpha \gamma_q (1 - \alpha)^{k-1} \quad (46)$$

**Case 2:**  $i \in I_{k,l}, l \neq q$ . We have:

$$\frac{\partial m'(\{\omega_q\})}{\partial d^{(i)2}} = \alpha \gamma_l \exp(-\gamma_l d^{(i)2}) \left( 1 - \prod_{j \in I_{k,q}} (1 - \alpha \exp(-\gamma_q d^{(j)2})) \right)$$

$$\times \prod_{r \neq q} \prod_{j \in I_{k,r}, j \neq i} (1 - \alpha \exp(-\gamma_r d^{(j)2})) \quad (47)$$

$$\frac{\partial m'(\Omega)}{\partial d^{(i)2}} = \alpha \gamma_l \exp(-\gamma_l d^{(i)2}) \prod_{r=1}^M \prod_{j \in I_{k,r}, j \neq i} (1 - \alpha \exp(-\gamma_r d^{(j)2})) \quad (48)$$

Setting the distances to zero in the above equations:

$$\left. \frac{\partial m'(\{\omega_q\})}{\partial d^{(i)2}} \right|_{\mathbf{d}^2 = \mathbf{0}} = \alpha \gamma_l (1 - (1 - \alpha)^{k_q}) (1 - \alpha)^{k - k_q - 1} \quad (49)$$

$$\left. \frac{\partial m'(\Omega)}{\partial d^{(i)2}} \right|_{\mathbf{d}^2 = \mathbf{0}} = \alpha \gamma_l (1 - \alpha)^{k-1} \quad (50)$$

where  $k_q = |I_{k,q}|$ .

The derivatives of  $K$  are simply obtained as follows:

$$\frac{\partial K}{\partial d^{(i)2}} = \sum_{q=1}^M \frac{\partial m'(\{\omega_q\})}{\partial d^{(i)2}} + \frac{\partial m'(\Omega)}{\partial d^{(i)2}} \quad (51)$$

Hence:

$$\left. \frac{\partial K}{\partial d^{(i)2}} \right|_{\mathbf{d}^2 = \mathbf{0}} = \alpha \gamma_q \sum_{r=1, r \neq q}^M (1 - (1 - \alpha)^{k_r}) (1 - \alpha)^{k - k_r - 1} \quad (52)$$

It follows from the preceding calculations that, for  $i \in I_{k,r}$ , the derivatives of  $m'(\{\omega_q\})$ ,  $m'(\Omega)$  and  $K$  for  $\mathbf{d}^2 = \mathbf{0}$  are proportional to  $\gamma_r$ . Since  $m'(\{\omega_q\})$ ,  $m'(\Omega)$  and  $K$  do not themselves depend on  $\gamma$  for  $\mathbf{d}^2 = \mathbf{0}$ , the derivative of  $P_q$  is also proportional to  $\gamma_r$ . Hence, we have:

$$\sum_{i \in I_{k,r}} \left. \frac{\partial P_q}{\partial d^{(i)2}} \right|_{\mathbf{d}^2 = \mathbf{0}} d^{(i)2} = A_{q,r} \gamma_r \quad (53)$$

for all  $r \in \{1, \dots, M\}$ ,  $A_{q,r}$  being some constant not depending on  $\gamma$ . Consequently, we can write:

$$P_q = P_q(\mathbf{0}) + \sum_{r=1}^M A_{q,r} \gamma_r \quad (54)$$

and, expressing this result in matrix form:

$$\mathbf{P} \cong \mathbf{P}(\mathbf{0}) + \mathbf{A} \boldsymbol{\gamma} \quad (55)$$

with  $\mathbf{A} = (A_{i,j})$  is a square matrix of size  $M$ .

The above calculations have been performed for an arbitrary training pattern  $\mathbf{x}$ . Reintroducing the pattern index  $\ell$ , we have:

$$\mathbf{P}^{(\ell)} \cong \mathbf{P}^{(\ell)}(\mathbf{0}) + \mathbf{A}^{(\ell)} \boldsymbol{\gamma} \quad (56)$$

Introducing these terms into the mean squared error, we have:

$$\begin{aligned}
E &= \frac{1}{N} \sum_{\ell=1}^N (\mathbf{P}^{(\ell)} - \mathbf{t}^{(\ell)})^t (\mathbf{P}^{(\ell)} - \mathbf{t}^{(\ell)}) \\
&= \frac{1}{N} \sum_{\ell=1}^N (\mathbf{P}^{(\ell)}(\mathbf{0}) - \mathbf{t}^{(\ell)} + \mathbf{A}^{(\ell)}\boldsymbol{\gamma})^t (\mathbf{P}^{(\ell)}(\mathbf{0}) - \mathbf{t}^{(\ell)} + \mathbf{A}^{(\ell)}\boldsymbol{\gamma}) \tag{57}
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{N} \sum_{\ell=1}^N (\mathbf{P}^{(\ell)}(\mathbf{0}) - \mathbf{t}^{(\ell)})^t (\mathbf{P}^{(\ell)}(\mathbf{0}) - \mathbf{t}^{(\ell)}) + 2\boldsymbol{\gamma}^t \mathbf{A}^{(\ell)t} (\mathbf{P}^{(\ell)}(\mathbf{0}) - \mathbf{t}^{(\ell)}) + \tag{58} \\
&\quad \boldsymbol{\gamma}^t \mathbf{A}^{(\ell)t} \mathbf{A}^{(\ell)} \boldsymbol{\gamma}
\end{aligned}$$

The gradient  $E$  with respect to  $\boldsymbol{\gamma}$  is therefore given by:

$$\nabla_{\boldsymbol{\gamma}} E = \frac{1}{N} \left[ \sum_{\ell=1}^N \mathbf{A}^{(\ell)t} (\mathbf{P}^{(\ell)}(\mathbf{0}) - \mathbf{t}^{(\ell)}) + \sum_{\ell=1}^N \mathbf{A}^{(\ell)t} \mathbf{A}^{(\ell)} \boldsymbol{\gamma} \right] \tag{59}$$

Minimizing  $E$  under the constraint  $\boldsymbol{\gamma} \geq \mathbf{0}$  is a nonnegative least squares problem that may be solved efficiently using, for instance, the algorithm described in [10, page 161].

## References

- [1] T. M. Cover and P. E. Hart. Nearest neighbor pattern classification. *IEEE Transactions on Information Theory*, IT-13(1):21–27, 1967.
- [2] T. Denœux. A  $k$ -nearest neighbor classification rule based on Dempster-Shafer theory. *IEEE Transactions on Systems, Man and Cybernetics*, 25(05):804–813, 1995.
- [3] T. Denœux. Analysis of evidence-theoretic decision rules for pattern classification. *Pattern Recognition*, 30(7):1095–1107, 1997.
- [4] T. Denœux. Application du modèle des croyances transférables en reconnaissance de formes. *Traitement du Signal (In press)*, 1997.
- [5] T. Denœux and G. Govaert. Combined supervised and unsupervised learning for system diagnosis using Dempster-Shafer theory. In P. Borne et al., editor, *CESA '96 IMACS Multiconference. Symposium on Control, Optimization and Supervision*, volume 1, pages 104–109, Lille, July 1996.
- [6] S. A. Dudani. The distance-weighted  $k$ -nearest-neighbor rule. *IEEE Transactions on Systems, Man and Cybernetics*, SMC-6:325–327, 1976.
- [7] J. H. Friedman. Regularized discriminant analysis. *Journal of the American Statistical Association*, 84:165–175, 1989.
- [8] R. P. Gorman and T. J. Sejnowski. Analysis of hidden units in a layered network trained to classify sonar targets. *Neural Networks*, 1:75–89, 1988.
- [9] J. M. Keller, M. R. Gray, and J. A. Givens. A fuzzy  $k$ -NN neighbor algorithm. *IEEE Transactions on Systems, Man and Cybernetics*, SMC-15(4):580–585, 1985.
- [10] C. L. Lawson and R. J. Hanson. *Solving Least Squares Problems*. Prentice-hall, 1974.
- [11] P. M. Murphy and D. W. Aha. *UCI Reposition of machine learning databases [Machine-readable data repository]*. University of California, Department of Information and Computer Science., Irvine, CA, 1994.
- [12] G. Shafer. *A mathematical theory of evidence*. Princeton University Press, Princeton, N.J., 1976.
- [13] P. Smets. The combination of evidence in the Transferable Belief Model. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 12(5):447–458, 1990.
- [14] P. Smets and R. Kennes. The Transferable Belief Model. *Artificial Intelligence*, 66:191–243, 1994.



- [15] L. M. Zouhal. *Contribution à l'application de la théorie des fonctions de croyance en reconnaissance des formes*. PhD thesis, Université de Technologie de Compiègne, 1997.

## **Biography**

### **Lalla Meriem Zouhal**

Lalla Meriem Zouhal received a M.S. in electronics from the Faculté des Sciences Hassan II, Casablanca, Morocco, in 1990, a DEA (Diplôme d'Etudes Approfondies) in System Control from the Université de Technologie de Compiègne, France, in 1992, and a PhD from the same institution in 1997. Her research interests concern pattern classification, Dempster-Shafer theory and Fuzzy logic.

### **Thierry Dencœux**

Thierry Dencœux graduated in 1985 as an engineer from the Ecole Nationale des Ponts et Chaussées in Paris, and earned a PhD from the same institution in 1989. He obtained the "Habilitation à diriger des Recherches" from the Institut National Polytechnique de Lorraine in 1996. From 1989 to 1992, he was employed by the Lyonnaise des Eaux water company, where he was in charge of research projects concerning the application of neural networks to forecasting and diagnosis. Dr. Dencœux joined the Université de Technologie de Compiègne as an assistant professor in 1992. His research interests include artificial neural networks, statistical pattern recognition, uncertainty modeling and data fusion.

## List of Tables

1	Main characteristics of data sets: number of classes ( $M$ ), training set size ( $N$ ), test set size ( $N_t$ ) and input dimension ( $n$ ). . . . .	18
2	Test error rates obtained with the voting, distance-weighted, fuzzy and evidence-theoretic classification rules for the best value of $k$ (in brackets), with 90 % confidence intervals. ETF: evidence-theoretic classifier with fixed $\gamma$ ; ETO: evidence-theoretic classifier with optimized $\gamma$ . . . . .	18

## List of Figures

1	Contour lines of the error function for different values of $\gamma$ , using the gradient-descent method. The optimal value is $\gamma = (0.84, 0.46)^t$ . . . . .	19
2	Contour lines of the error function for different values of $\gamma$ , using the linearization method for pignistic probabilities vectors. The optimal value is $\gamma = (0.76, 0.6)^t$ . . . . .	19
3	Test error rates on data sets $B_1$ as a function of $k$ , for the ETF (x.) and ETO $k$ -NN rules: (:) gradient method and (- -) linearization method. . . . .	20
4	Test error rates on data sets $B_1$ as a function of $k$ , for the voting (-), ETO(:), fuzzy (- -) and distance-weighted (-.) $k$ -NN rules. . . . .	20
5	Test error rates on data sets $B_2$ as a function of $k$ , for the ETF (x.) and ETO $k$ -NN rules: (:) gradient method and (- -) linearization method. . . . .	21
6	Test error rates on data sets $B_2$ as a function of $k$ , for the voting (-), ETO(:), fuzzy (- -) and distance-weighted (-.) $k$ -NN rules. . . . .	21
7	Test error rates on ionosphere data as a function of $k$ , for the ETF (x.) and ETO $k$ -NN rules: (:) gradient method and (- -) linearization method. . . . .	22
8	Test error rates on ionosphere data as a function of $k$ , for the voting (-), ETO(:), fuzzy (- -) and distance-weighted (-.) $k$ -NN rules. . . . .	22
9	Test error rates on vehicle data as a function of $k$ , for the ETF (x.) and ETO $k$ -NN rules: (:) gradient method and (- -) linearization method. . . . .	23
10	Test error rates on vehicle data as a function of $k$ , for the voting (-), ETO(:), fuzzy (- -) and distance-weighted (-.) $k$ -NN rules. . . . .	23
11	Test error rates on sonar data as a function of $k$ , for the ETF (x.) and ETO $k$ -NN rules: (:) gradient method and (- -) linearization method. . . . .	24
12	Test error rates on sonar data as a function of $k$ , for the voting (-), ETO(:), fuzzy (- -) and distance-weighted (-.) $k$ -NN rules. . . . .	24

Table 1: Main characteristics of data sets: number of classes ( $M$ ), training set size ( $N$ ), test set size ( $N_t$ ) and input dimension ( $n$ ).

data set	$M$	$N$	$N_t$	$n$
$B_1$	3	60	1000	10
$B_2$	3	300	1000	10
Ion	2	175	176	34
veh	4	564	282	18
Son	2	104	104	60

Table 2: Test error rates obtained with the voting, distance-weighted, fuzzy and evidence-theoretic classification rules for the best value of  $k$  (in brackets), with 90 % confidence intervals. ETF: evidence-theoretic classifier with fixed  $\gamma$ ; ETO: evidence-theoretic classifier with optimized  $\gamma$ .

data set	voting	ETF	ETO	weighted	fuzzy
$B_1$	$0.41 \pm 0.03$ (8)	$0.37 \pm 0.03$ (13)	$0.32 \pm 0.02$ (27)	$0.37 \pm 0.03$ (21)	$0.40 \pm 0.03$ (5)
$B_2$	$0.31 \pm 0.02$ (14)	$0.29 \pm 0.02$ (13)	$0.24 \pm 0.02$ (17)	$0.28 \pm 0.02$ (15)	$0.29 \pm 0.02$ (16)
Ion	$0.15 \pm 0.04$ (1)	$0.11 \pm 0.04$ (2)	$0.07 \pm 0.03$ (8)	$0.15 \pm 0.04$ (1)	$0.15 \pm 0.04$ (1)
Veh	$0.33 \pm 0.05$ (3)	$0.34 \pm 0.05$ (4)	$0.33 \pm 0.05$ (3)	$0.32 \pm 0.05$ (6)	$0.34 \pm 0.05$ (3)
Son	$0.17 \pm 0.06$ (1)	$0.17 \pm 0.06$ (1)	$0.15 \pm 0.06$ (6)	$0.17 \pm 0.06$ (1)	$0.17 \pm 0.06$ (1)

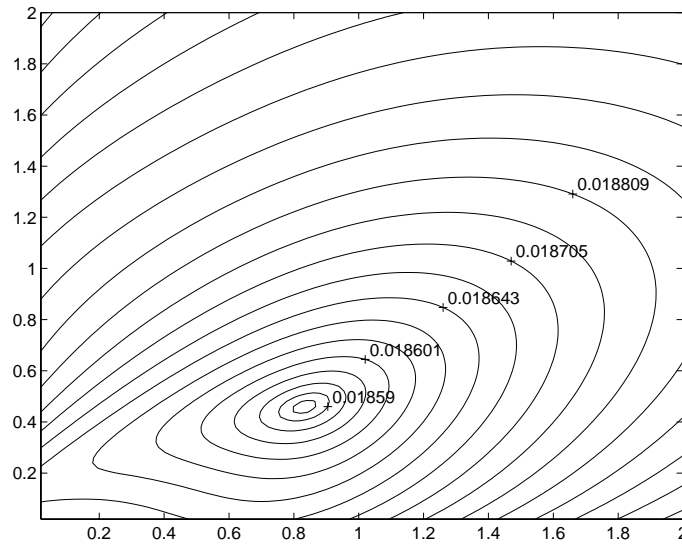


Figure 1: Contour lines of the error function for different values of  $\gamma$ , using the gradient-descent method. The optimal value is  $\gamma = (0.84, 0.46)^t$ .

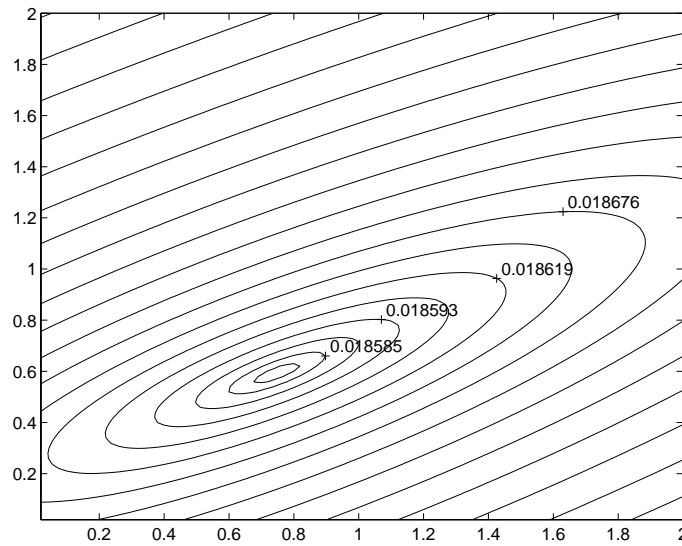


Figure 2: Contour lines of the error function for different values of  $\gamma$ , using the linearization method for pignistic probabilities vectors. The optimal value is  $\gamma = (0.76, 0.6)^t$ .

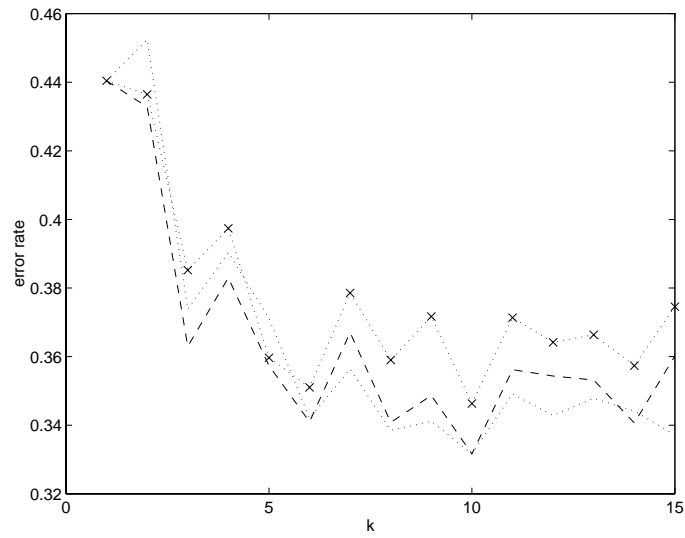


Figure 3: Test error rates on data sets  $B_1$  as a function of  $k$ , for the ETF (x.) and ETO  $k$ -NN rules: (·) gradient method and (- -) linearization method.

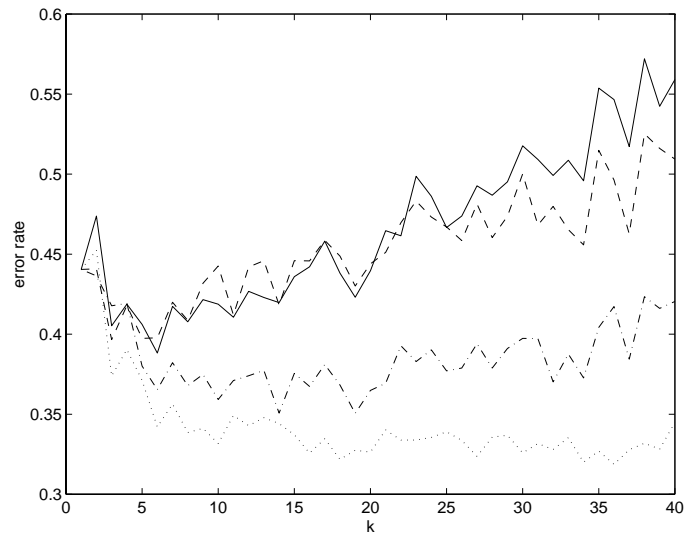


Figure 4: Test error rates on data sets  $B_1$  as a function of  $k$ , for the voting (-), ETO(·), fuzzy (- -) and distance-weighted (-.)  $k$ -NN rules.

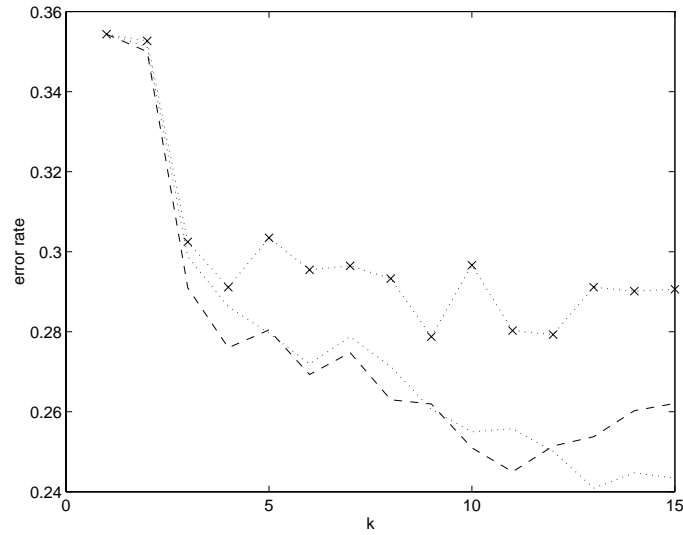


Figure 5: Test error rates on data sets  $B_2$  as a function of  $k$ , for the ETF (x.) and ETO  $k$ -NN rules: (:) gradient method and (- -) linearization method.

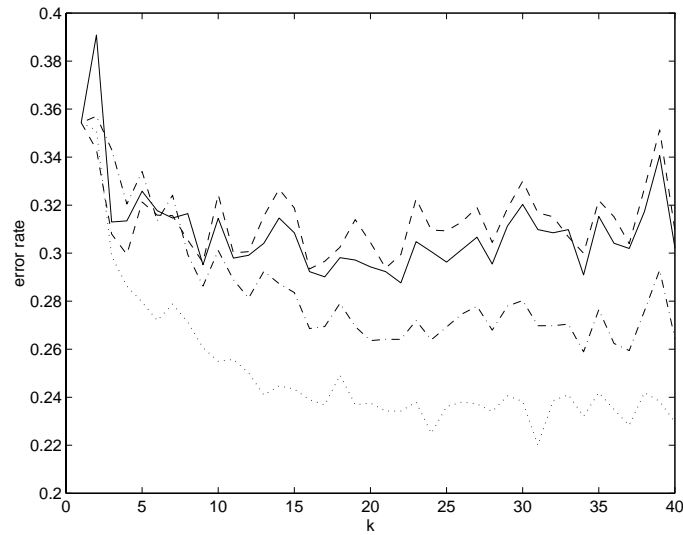


Figure 6: Test error rates on data sets  $B_2$  as a function of  $k$ , for the voting (-), ETO (:), fuzzy (- -) and distance-weighted (-.)  $k$ -NN rules.

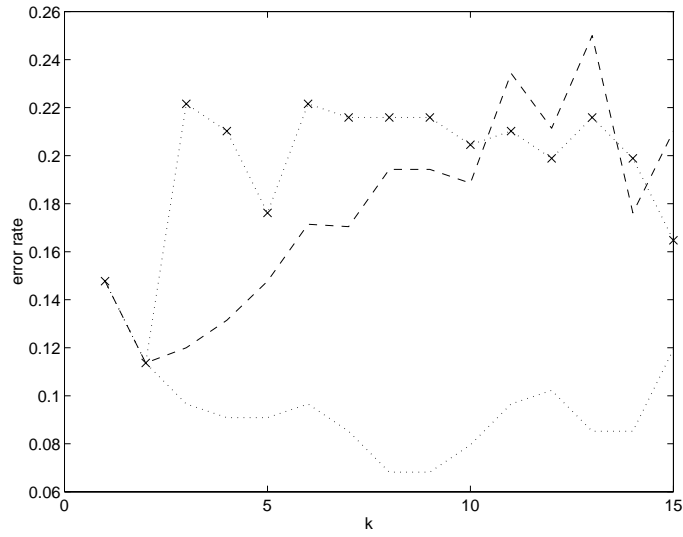


Figure 7: Test error rates on ionosphere data as a function of  $k$ , for the ETF (x.) and ETO  $k$ -NN rules: (:) gradient method and (- -) linearization method.

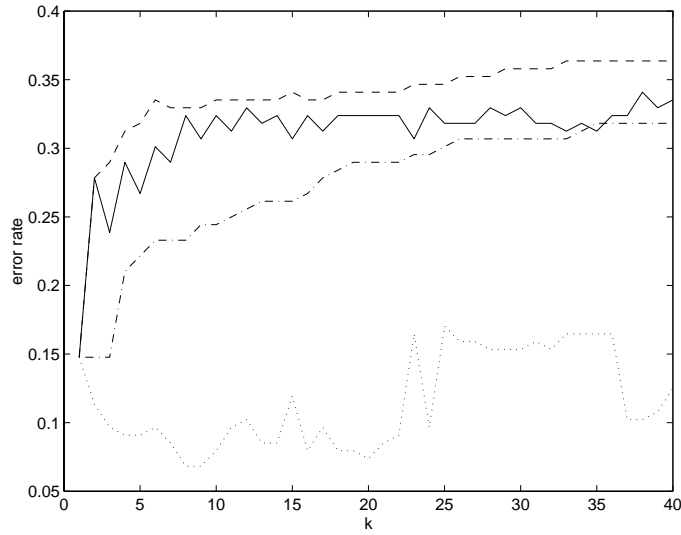


Figure 8: Test error rates on ionosphere data as a function of  $k$ , for the voting (-), ETO(:), fuzzy (- -) and distance-weighted (-.)  $k$ -NN rules.



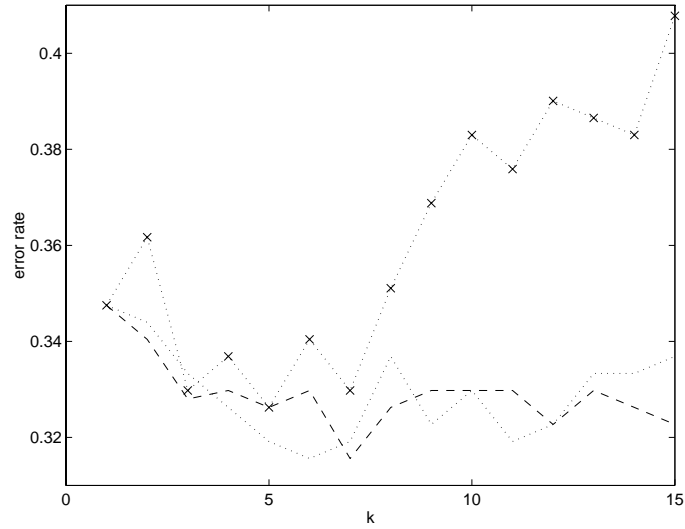


Figure 9: Test error rates on vehicle data as a function of  $k$ , for the ETF (x.) and ETO  $k$ -NN rules: (:) gradient method and (- -) linearization method.

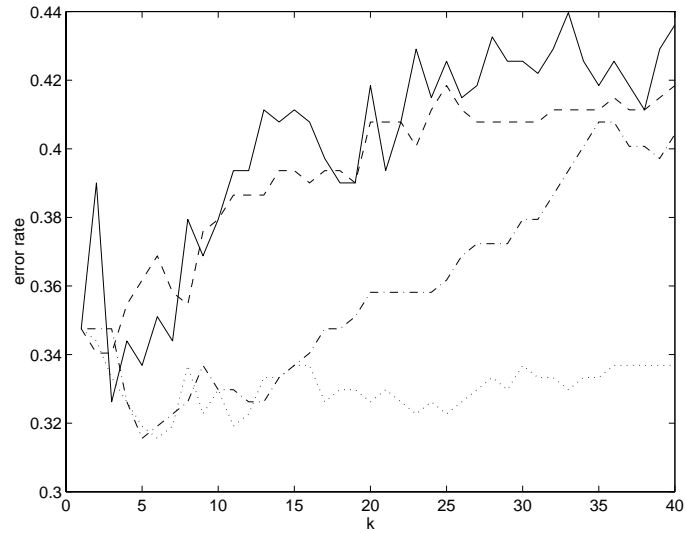


Figure 10: Test error rates on vehicle data as a function of  $k$ , for the voting (-), ETO(:), fuzzy (- -) and distance-weighted (-.)  $k$ -NN rules.

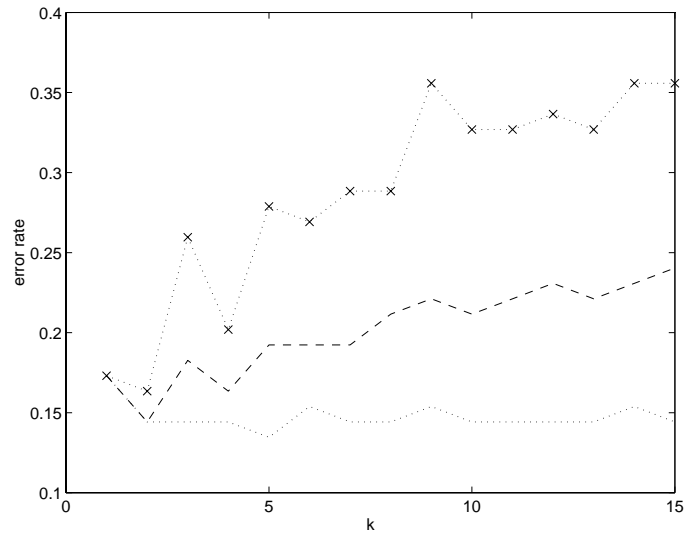


Figure 11: Test error rates on sonar data as a function of  $k$ , for the ETF (x.) and ETO  $k$ -NN rules: (:) gradient method and (-) linearization method.

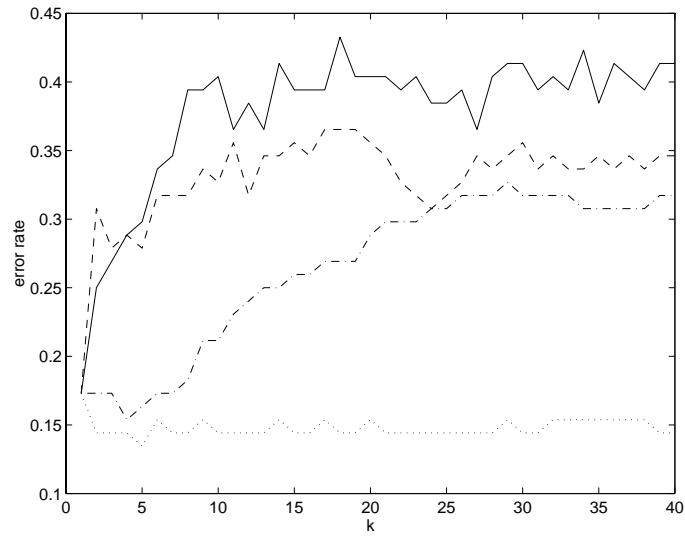


Figure 12: Test error rates on sonar data as a function of  $k$ , for the voting (-), ETO(:), fuzzy (- -) and distance-weighted (-.)  $k$ -NN rules.