# Supervised Pattern Classification Based on Optimum-Path Forest

**J. P. Papa, A. X. Falcão, C. T. N. Suzuki**

Institute of Computing, University of Campinas, Campinas, SP, Brazil

**ABSTRACT:** We present a supervised classification method which represents each class by one or more optimum-path trees rooted at some key samples, called *prototypes*. The training samples are nodes of a complete graph, whose arcs are weighted by the distances between the feature vectors of their nodes. Prototypes are identified in all classes and the minimization of a *connectivity function* by dynamic programming assigns to each training sample a minimum-cost path from its most strongly connected prototype. This competition among prototypes partitions the graph into an optimum-path forest rooted at them. The class of the samples in an optimum-path tree is assumed to be the same of its root. A test sample is classified similarly, by identifying which tree would contain it, if the sample were part of the training set. By choice of the graph model and connectivity function, one can devise other optimum-path forest classifiers. We present one of them, which is fast, simple, multiclass, parameter independent, does not make any assumption about the shapes of the classes, and can handle some degree of overlapping between classes. We also propose a general algorithm to learn from errors on an evaluation set without increasing the training set, and show the advantages of our method with respect to SVM, ANN-MLP, and *k*-NN classifiers in several experiments with datasets of various types. © 2009 Wiley Periodicals, Inc. Int J Imaging Syst Technol, 19, 120–131, 2009; Published online in Wiley InterScience (www.interscience.wiley.com). DOI 10.1002/ima.20188

**Key words:** supervised learning; image foresting transform; pattern recognition; image analysis; graph-search algorithms

## I. INTRODUCTION

Patterns are usually represented by feature vectors (set of measures or observations) obtained from samples of a dataset (Duda et al., 2000). Two fundamental problems in pattern recognition are as follows: (i) the identification of natural groups (clustering) composed by samples with similar patterns and (ii) the classification of each sample in one of $c$ possible classes (labels). The dataset is usually divided in two parts, a training set and a test set, being the first used to project the classifier and the second used for validation, by measuring its classification errors (accuracy). This process must be also repeated several times with randomly selected training and test samples to achieve a conclusion about the statistics of its accuracy (robustness and precision). While problem (i) has no prior information about the labels of the samples, the training in problem (ii) can count with unlabeled samples (unsupervised learning), labeled samples (supervised learning) or part of the samples labeled and the other part unlabeled (semisupervised learning) (Blum and Mitchell, 1998; Joachims, 1999; Zhu, 2006). Our focus is on the supervised learning approaches.

Figure 1 illustrates three typical cases in 2D feature spaces using two classes: (a) linearly separable, (b) piecewise linearly separable, and (c) nonseparable classes with arbitrary shapes. Any reasonable approach should handle (a) and (b), being (c) the most interesting challenge. An artificial neural network with multilayer perceptrons (ANN-MLP), for example, can address (a) and (b), but not (c) (Haykin, 1994). As an unstable classifier, collections of ANN-MLP (Kuncheva, 2004) can improve its performance up to some unknown limit of classifiers (Reyzin and Schapire, 2006). Support vector machines (SVMs) have been proposed to overcome the problem, by assuming linearly separable classes in a higher-dimensional feature space (Boser et al., 1992). Its computational cost rapidly increases with the training set size and the number of support vectors. As a binary classifier, multiple SVMs are required to solve a multiclass problem (Duan and Keerthi, 2005). Tang and Mazzoni (2006) proposed a method to reduce the number of support vectors in the multiclass problem. Their approach suffers from slow convergence and high computational cost, because they first minimize the number of support vectors in several binary SVMs, and then share these vectors among the machines. Panda et al. (2006) presented a method to reduce the training set size before computing the SVM algorithm. Their approach aims to identify and remove samples likely related to nonsupport vectors. However, in all SVM approaches, the assumption of separability may also not be valid in any space of finite dimension (Collobert and Bengio, 2004).

We propose a supervised classifier based on *optimum-path forest* (OPF), which is fast, simple, multiclass, parameter independent, does not make any assumption about the shapes of the classes, and can handle some degree of overlapping between classes. The training set is thought of as a complete graph, whose nodes are the samples and arcs link all pairs of nodes. The arcs are weighted by the
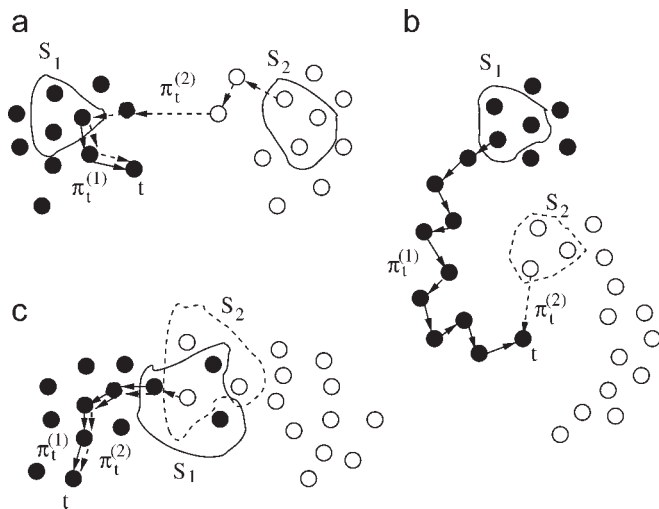
---

**Figure 1.** Examples of 2D feature spaces using two classes: (a) linearly separable, (b) piecewise linearly separable, and (c) nonseparable classes with arbitrary shapes. Prototypes can be identified in each class, forming the sets $S_1$ and $S_2$. Every sample $t$ can be connected to a prototype in $S_i$, $i = 1, 2$, by a sequence $\pi_t^{(i)}$ of distinct samples. The classification is done based on optimal connections to the prototypes.

distances between the feature vectors of their corresponding nodes. Any sequence of distinct samples forms a path connecting the terminal nodes and a *connectivity function* assigns a cost to that path (e.g., the maximum arc-weight along it). The idea is to identify prototypes in each class such that every sample is assigned to the class of its most strongly connected prototype. That is, the one which offers to it a minimum-cost path, considering all possible paths from the prototypes. Figure 1 shows two sets of prototypes, $S_1$ and $S_2$, in classes 1 and 2. The connection from $S_i$ to a sample $t$ is represented by a path $\pi_t^{(i)}$ with terminus $t$ and root in some prototype of $S_i$, $i = 1, 2$. In all cases, the optimum path (to which the maximum arc-weight is minimum) comes from a prototype of the same class of $t$. Our approach can handle all three cases with the maximum arc-weight function and prototypes estimated as the closest samples from distinct classes. In the case of overlapping between classes, these prototypes work as class defenders in the overlapped regions of the feature space (Fig. 1c).

The classifier is an optimum-path forest rooted at the prototypes. That is, each training sample belongs to one optimum-path tree rooted at its most strongly connected prototype. The classification of a test sample evaluates the optimum paths from the prototypes to this sample incrementally, as though it were part of the forest, and assigns to it the label of the most strongly connected root. Note the difference between the proposed method with the maximum arc-weight function and the nearest neighbor approach (Cover and Hart, 1967). A test/training sample may be assigned to a given class, even when its closest labeled sample is from another class (Fig. 1b).

The optimum paths from the prototypes to the other samples are computed by the algorithm of the image foresting transform (IFT)—a tool for the design of image processing operators based on connectivity (Falcão et al., 2004)—which is extended here from the image domain to the feature space. The IFT algorithm is essentially Dijkstra's algorithm (Cormen et al., 1990) modified for multiple sources and more general path-value functions (Falcão et al., 2004).

It first identifies the minima (maxima) of the path-value function as source nodes and then propagates optimum paths from those sources in a nondecreasing (nonincreasing) order of optimum-path values, partitioning the graph into an optimum-path forest rooted at the source nodes. It is a dynamic programming strategy in which, by choice of the path-value function (Eq. 1), we force the prototypes to be the roots of the forest.

The dataset partition by the proposed classifier in the feature space is equivalent to an image segmentation by the IFT-watershed transform from labeled markers (Lotufo and Falcão, 2000; Audigier and Lotufo, 2007b) in the image domain. Similar important relations can be obtained with other image operators, such as relative-fuzzy connected segmentation (Herman and Carvalho, 2001; Saha and Udupa, 2001; Audigier and Lotufo, 2007a; Miranda et al., 2008). In our case, the markers are the prototypes and we have a special way to estimate them. Figure 2 helps to understand this comparison and why the proposed method works in the feature space, when prototypes are estimated as the closest samples from distinct classes. Figure 2a shows an image with one internal marker (white) and one external marker (black) for an object of interest. The pixels are the nodes of a graph whose arcs link the 8-neighbors of each pixel. The arc weights are dissimilarity values between pixels, computed based on their image properties. The dissimilarity function between pixels plays the same role of the distance function between samples and distinct classes are represented by object and background. The connectivity function is the maximum arc-weight along the path. Figure 2b gives an idea of the arc weights by displaying the complement of a gradient-like image, which is created by assigning to each pixel the maximum among the arc weights between it and its eight neighbors. By selecting markers around the weaker parts (lower arc weights) of the boundary (Fig. 2a), we force the minimum-cost paths from internal and external markers to meet first at the weaker parts of the object's boundary, blocking these passages for paths from the other side. Therefore, possible paths from one side to the other will have costs higher than paths from the same side with respect to each marker. The optimum-path propagation from both markers describes an ordered region growing (flooding) process where the wavefronts from each marker meet at the object's boundary (Fig. 2c). The object is defined by the optimum-path forest rooted at the pixels of the internal marker. In the case of multiple internal and external markers, the object is composed of multiple internal forests. Three frames of this process are presented in Figures 2d–2f. Note that internal (external) pixels, which are only reachable by high-cost paths, are initially surrounded by optimum paths from the internal (external) marker and finally conquered by this marker. We can also exploit other connectivity functions, but this work presents only the results for the maximum arc-weight function.

Supervised classification based on prototypes is not new. For example, methods such as the $k$-nearest neighbors ($k$-NN) use all training samples as prototypes (Fukunaga and Narendra, 1975). Its classification relies on the direct distance between samples. As far as we know, our approach is the first to consider optimum-path forests rooted at automatically selected prototypes in the feature space. Besides, by changing the graph model and path-value function, one can derive other types of optimum-path forest classifiers, such as the unsupervised learning approach proposed in (Cappabianco et al., 2008; Rocha et al., 2008), which also relies on a different strategy to estimate prototypes. Most approaches for pattern classification based on graphs and/or paths in graphs are either unsupervised (Zahn, 1971; Hubert, 1974; Jain and Dubes, 1988;
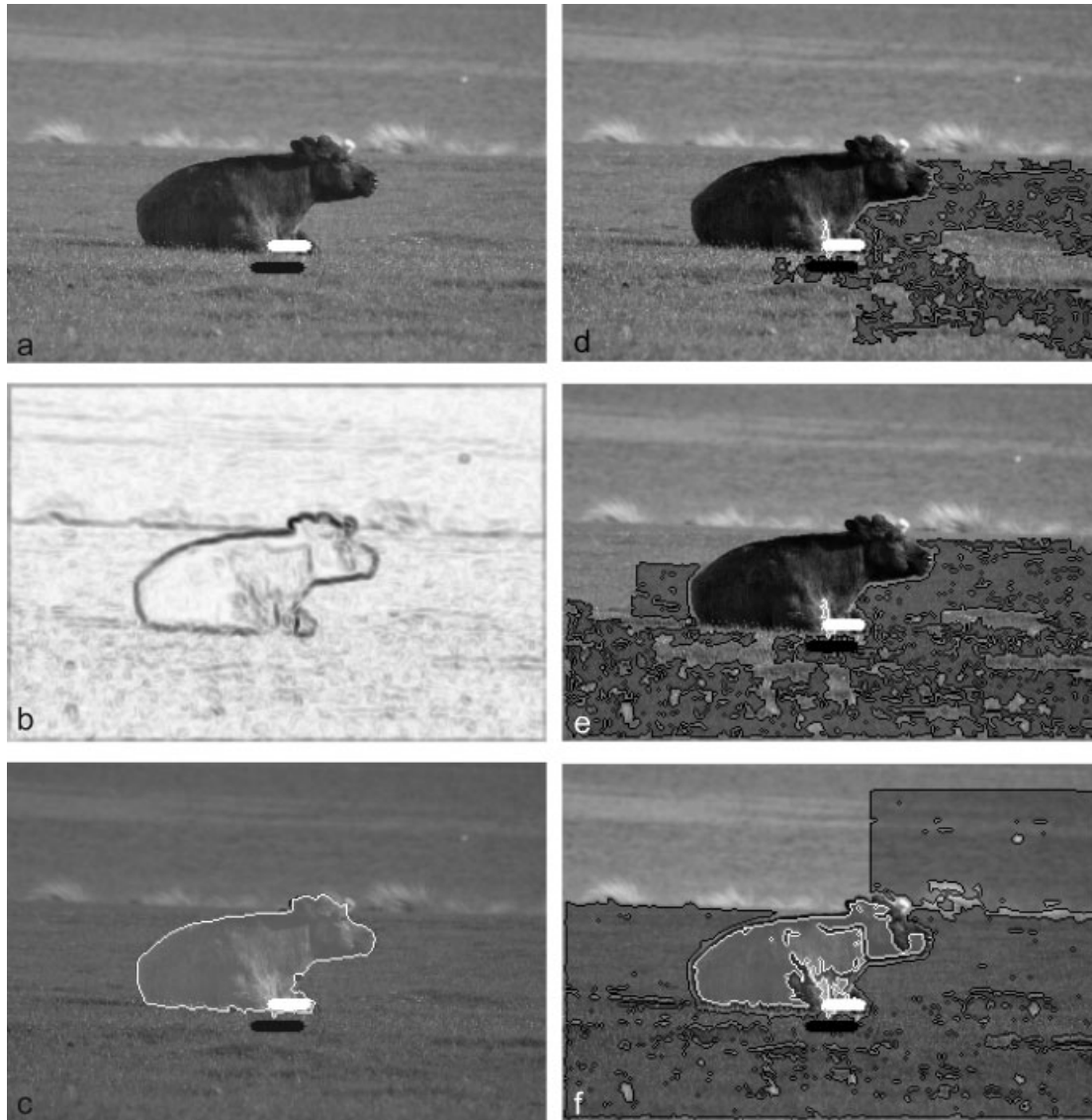
**Figure 2.** IFT-watershed segmentation. (a) Image with internal (white) and external (black) markers. (b) The complement of a gradient-like image which gives an idea of the arc weights. The markers are selected around the weaker parts of the boundary (brighter values in b). (c) The result of segmentation and (d–f) three frames of the IFT flooding process that leads to (c).

Shi and Malik, 2000) or semisupervised (Kulis et al., 2005; Zhou et al., 2005; Callut et al., 2008; Kumar and Kummamuru, 2008). The proposed method can be easily extended to semisupervised classification, given that the optimum-path forest can include unlabeled nonprototype samples. Previous versions of it have also been published (Papa et al., 2007, 2008a,b; Montoya-Zegarra et al., 2008; Spadotto et al., 2008). We have simplified the learning procedure with better results, corrected some mistakes, improved explanations and added several experiments using more datasets, baseline classifiers, and image descriptors based on texture, shape, and color.

Other contribution of this work concerns learning algorithms, which can teach a classifier from its errors on a third evaluation set without increasing the size of the training set. As the samples in the test set cannot be seen during the project, the evaluation set is necessary for this purpose. The basic idea is to randomly interchange samples of the training set with misclassified samples of the evalua-

tion set, retrain the classifier, and evaluate it again, repeating this procedure during a few iterations. The effectiveness is measured by comparing the results on the unseen test set before and after the learning algorithm. It is expected an improvement in performance for any stable classifier.

The learning with fixed training set size is usually required in large datasets with thousands/millions of samples (e.g., pixels/voxels in 2D/3D images). It also stems from applications where the classifier is part of an expert system, which performs a laborious data analysis (sometimes inviable for human beings) and emits its opinion to a human expert. The human expert may agree or not based on other evidences, but the feedback about the classification errors is important to improve performance in a future analysis. The diagnosis of parasites from microscopy images of biological slides is an example (Falcão et al., 2008). The human visual inspection is very difficult and error prone in several situations due to the amount of impurities and small sizes of some parasites (e.g., protozoa in
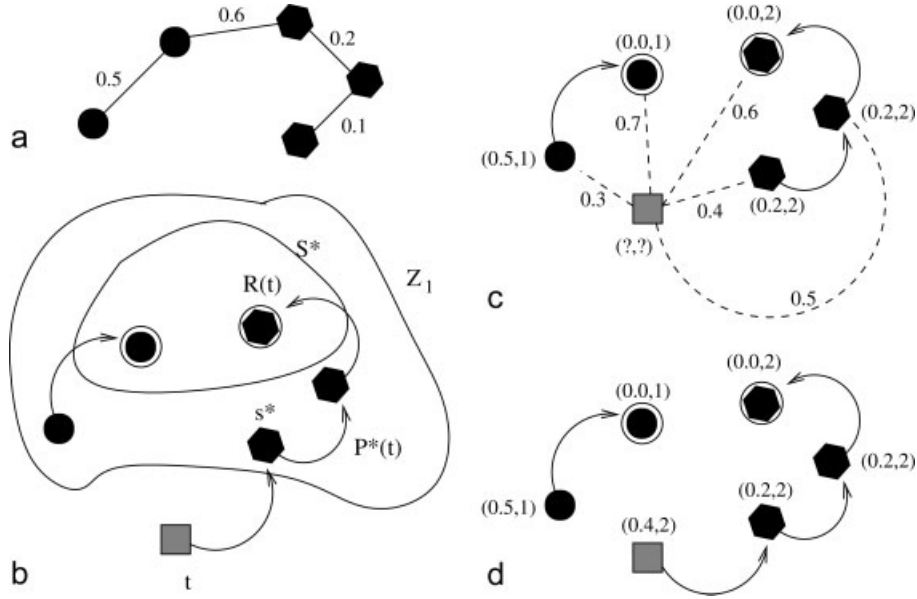
**Figure 3.** a) Complete weighted graph for a simple training set. (b) Resulting optimum-path forest for $f_{max}$ and two given prototypes (circled nodes). The entries (x,y) over the nodes are, respectively, the cost and the label of the samples. The directed arcs indicate the predecessor nodes in the optimum path. (c) Test sample (gray square) and its connections (dashed lines) with the training nodes. (d) The optimum path from the most strongly connected prototype, its label 2, and classification cost 0.4 are assigned to the test sample. The test sample is classified in the class hexagon, although its nearest training sample is from the class circle.

samples of feces). We aim to improve the performance of the expert system along time of use and we are taking into account the fact that computers have a limited storage and processing capacity for the training set.

This article describes the supervised OPF classifier in Section II, presents a general learning algorithm in Section III, which follows the same aforementioned strategy for all classifiers, shows results that compare the OPF classifier with SVM (Boser et al., 1992), ANN-MLP (Haykin, 1994) and $k$-NN (Fukunaga and Narendra, 1975) in Section IV, and states conclusions in Section V.

## II. OPTIMUM-PATH FOREST CLASSIFIER

Let $Z_1$, $Z_2$, and $Z_3$ be training, evaluation, and test sets with $|Z_1|$, $|Z_2|$, and $|Z_3|$ samples of a given dataset. We use samples as points, images, voxels, and contours in this article. As already explained, this division of the dataset is necessary to validate the classifier and evaluate its learning capacity from the errors. $Z_1$ is used to project the classifier and $Z_3$ is used to measure its accuracy, being the labels of $Z_3$ kept unseen during the project. A pseudotest on $Z_2$ is used to teach the classifier by randomly interchanging samples of $Z_1$ with misclassified samples of $Z_2$. After learning, it is expected an improvement in accuracy on $Z_3$.

Let $\lambda(s)$ be the function that assigns the correct label $i$, $i = 1, 2, \ldots, c$, of class $i$ to any sample $s \in Z_1 \cup Z_2 \cup Z_3$, $S \subset Z_1$ be a set of prototypes from all classes, and $v$ be an algorithm which extracts $n$ features (color, shape, texture properties) from any sample $s \in Z_1 \cup Z_2 \cup Z_3$ and returns a vector $\vec{v}(s)$. The distance $d(s,t) \geq 0$ between two samples, $s$ and $t$, is the one between their feature vectors $\vec{v}(s)$ and $\vec{v}(t)$. One can use any distance function suitable for the extracted features. The most common is the Euclidean norm $\|\vec{v}(t) - \vec{v}(s)\|$, but some image features require special distance

algorithms (Wang and Pavlidis, 1990). A pair $(v, d)$ then describes how the samples of a dataset are distributed in the feature space. Therefore, we call $(v, d)$ a *descriptor* and the experiments in Section IV use shape (Arica and Vural, 2003), texture (Montoya-Zegarra et al., 2008), and color (Stehling et al., 2002) descriptors based on this definition.

Our problem consists of projecting a classifier which can predict the correct label $\lambda(s)$ of any sample $s \in Z_3$. Training consists of finding a special set $S^* \subset Z_1$ of prototypes and a discrete optimal partition of $Z_1$ in the feature space (i.e., an optimum-path forest rooted in $S^*$). The classification of a sample $s \in Z_3$ (or $s \in Z_2$) is done by evaluating the optimum paths incrementally, as though it were part of the forest, and assigning to it the label of the most strongly connected prototype.

**A. Training.** Let $(Z_1, A)$ be a complete graph whose nodes are the training samples and any pair of samples defines an arc in $A = Z_1 \times Z_1$ (Fig. 3a). The arcs do not need to be stored and so the graph does not need to be explicitly represented. A path is a sequence of distinct samples $\pi_t = \langle s_1, S_2, \ldots, t \rangle$ with terminus at a sample $t$. A path is said *trivial* if $\pi_t = \langle t \rangle$. We assign to each path $\pi_t$ a cost $f(\pi_t)$ given by a connectivity function $f$. A path $\pi_t$ is said optimum if $f(\pi_t) \leq f(\tau_t)$ for any other path $\tau_t$. We also denote by $\pi_s \cdot \langle s,t \rangle$ the concatenation of a path $\pi_s$ and an arc $(s,t)$.

We will address the connectivity function $f_{max}$.

$$f_{max}(\langle s \rangle) = \begin{cases} 0 & \text{if } s \in S, \\ +\infty & \text{otherwise} \end{cases}$$

$$f_{max}(\pi_s \cdot \langle s, t \rangle) = \max\{f_{max}(\pi_s), d(s, t)\} \quad (1)$$

such that $f_{max}(\pi_s \cdot \langle s,t \rangle)$ computes the maximum distance between adjacent samples along the path $\pi_s \cdot \langle s,t \rangle$. The minimization of $f_{max}$ assigns to every sample $t \in Z_1$ an optimum path $P^*(t)$ from the set

$S \subset Z_1$ of prototypes, whose minimum cost $C(t)$ is

$$C(t) = \min_{\forall \pi_t \in (Z_1, A)} \{f_{\max}(\pi_t)\}. \qquad (2)$$

The minimization of $f_{\max}$ is computed by Algorithm 1, called OPF algorithm, which is an extension of the general image foresting transform (IFT) algorithm (Falcão et al., 2004) from the image domain to the feature space, here specialized for $f_{\max}$. As explained in Section I, this process assigns one optimum path from $S$ to each training sample $t$ in a nondecreasing order of minimum cost, such that the graph is partitioned into an optimum-path forest $P$ (a function with no cycles which assigns to each $t \in Z_1 \setminus S$ its predecessor $P(t)$ in $P^*(t)$ or a marker $nil$ when $t \in S$, as shown in Fig. 3b). The root $R(t) \in S$ of $P^*(t)$ can be obtained from $P(t)$ by following the predecessors backwards along the path, but its label is propagated during the algorithm by setting $L(t) \leftarrow \lambda(R(t))$.

Lines 1–3 initialize maps and insert prototypes in $Q$. The main loop computes an optimum path from $S$ to every sample $s$ in a nondecreasing order of minimum cost (Lines 4–11). At each iteration, a path of minimum cost $C(s)$ is obtained in $P$ when we remove its last node $s$ from $Q$ (Line 5). Ties are broken in $Q$ using first-in-first-out policy. That is, when two optimum paths reach an ambiguous sample $s$ with the same minimum cost, $s$ is assigned to the first path that reached it. Note that $C(t) > C(s)$ in Line 6 is false when $t$ has been removed from $Q$ and, therefore, $C(t) \neq +\infty$ in Line 9 is true only when $t \in Q$. Lines 8–11 evaluate if the path that reaches an adjacent node $t$ through $s$ is cheaper than the current path with terminus $t$ and update the position of $t$ in $Q$, $C(t)$, $L(t)$, and $P(t)$ accordingly.

**Algorithm 1**—OPF Algorithm

**Input:** A training set $Z_1$, $\lambda$-labeled prototypes $S \subset Z_1$ and the pair $(v,d)$ for feature vector and distance computations.
**Output:** Optimum-path forest P, cost map C and label map L.
**Auxiliary:** Priority queue Q and cost variable cst.

1. For each $s \in Z_1 \setminus S$, set $C(s) \leftarrow +\infty$.
2. For each $s \in S$, do
3.    $C(s) \leftarrow; 0$, $P(s) \leftarrow nil$, $L(s) \leftarrow \lambda(s)$ and insert s in Q.
4. While Q is not empty, do
5.    Remove from Q a sample s such that C(s) is minimum.
6.    For each $t \in Z_1$ such that $t \neq s$ and $C(t) > C(s)$, do
7.       Compute $cst \leftarrow \max\{C(s), d(s,t)\}$.
8.       If $cst < C(t)$, then
9.          If $C(t) \neq +\infty$, then remove t from Q.
10.          $P(t) \leftarrow s$, $L(t) \leftarrow L(s)$, and $C(t) \leftarrow cst$.
11.          Insert t in Q.

One can use other *smooth* connectivity functions, as long as they group samples with similar properties (Falcão et al., 2004). A function $f$ is smooth in $(Z_1, A)$ when for any sample $t \in Z_1$, there exists an optimum path $\pi_t$ which either is trivial or has the form $\pi_s \cdot \langle s,t \rangle$, where

1. $f(\pi_s) \leq f(\pi_t)$,
2. $\pi_s$ is optimum,
3. for any optimum path $\tau_s, f(\tau_s \cdot \langle s,t \rangle) = f(\pi_t)$.

We say that $S^*$ is an optimum set of prototypes when Algorithm 1 minimizes the classification errors in $Z_1$. $S^*$ can be found by exploiting the theoretical relation between minimum-spanning tree (MST) (Cormen et al., 1990) and optimum-path tree for $f_{\max}$ (Allène et al., 2007; Miranda et al., 2008).

By computing a MST in the complete graph $(Z_1, A)$, we obtain a connected acyclic graph whose nodes are all samples of $Z_1$ and the arcs are undirected and weighted by the distances $d$ between adjacent samples (Fig. 3a). The spanning tree is optimum in the sense that the sum of its arc weights is minimum as compared to any other spanning tree in the complete graph. In the MST, every pair of samples is connected by a single path which is optimum according to $f_{\max}$. That is, the minimum-spanning tree contains one optimum-path tree for any selected root node.

The optimum prototypes are the closest elements of the MST with different labels in $Z_1$. By removing the arcs between different classes, their adjacent samples become prototypes in $S^*$ and Algorithm 1 can compute an optimum-path forest in $Z_1$ (Fig. 3b). Note that, a given class may be represented by multiple prototypes (i.e., optimum-path trees) and there must exist at least one prototype per class.

It is not difficult to see that the optimum paths between classes tend to pass through the same removed arcs of the minimum-spanning tree. The choice of prototypes as described earlier aims to block these passages, reducing the chances of samples in any given class be reached by optimum paths from prototypes of other classes.

**B. Classification.** For any sample $t \in Z_3$, we consider all arcs connecting $t$ with samples $s \in Z_3$, as though $t$ were part of the training graph (Fig. 3c). Considering all possible paths from $S^*$ to $t$, we find the optimum path $P^*(t)$ from $S^*$ and label $t$ with the class $\lambda(R(t))$ of its most strongly connected prototype $R(t) \in S^*$ (Fig. 3b). This path can be identified incrementally, by evaluating the optimum cost $C(t)$ as

$$C(t) = \min\{\max\{C(s), d(s,t)\}\}, \quad \forall s \in Z_1. \qquad (3)$$

Let the node $s^* \in Z_1$ be the one that satisfies Eq. (3) (i.e., the predecessor $P(t)$ in the optimum path $P^*(t)$). Given that $L(s^*) = \lambda(R(t))$, the classification simply assigns $L(s^*)$ as the class of $t$ (Fig. 3d). An error occurs when $L(s^*) \neq \lambda(t)$.

Similar procedure is applied for samples in the evaluation set $Z_2$. In this case, however, we would like to use misclassified samples of $Z_2$ to learn the distribution of the classes in the feature space and improve the classification performance on $Z_3$.

## III. LEARNING FROM ERRORS ON THE EVALUATION SET

There are many situations that limit the size of $Z_1$: large datasets, limited computational resources, and high computational time as required by some approaches. Mainly in applications with large datasets, it would be interesting to select for $Z_1$ the most informative samples, such that the accuracy of the classifier is little affected by this size limitation. It is also important to show that a classifier can improve its performance along time of use, when we are able to teach it from its errors. This section presents a general learning algorithm which uses a third evaluation set $Z_2$ to improve the composition of samples in $Z_1$ without increasing its size.

From an initial choice of $Z_1$ and $Z_2$, the algorithm projects an instance $I$ of a given classifier from $Z_1$ and evaluates it on $Z_2$. The

misclassified samples of $Z_2$ are randomly selected and replaced by samples of $Z_1$ (under certain constraints). This procedure assumes that the most informative samples can be obtained from the errors. The new sets $Z_1$ and $Z_2$ are then used to repeat the process during a few iterations $T$. The instance of classifier with highest accuracy is selected along the iterations. The accuracy values $L(I)$ obtained for each instance $I$ form a *learning curve*, whose nondecreasing monotonic behavior indicates a positive learning rate for the classifier. Afterwards, by comparing the accuracies of the classifier on $Z_3$, before and after the learning process, we can evaluate its learning capacity from the errors.

The accuracies $L(I)$, $I = 1,2\dots,T$, are measured by taking into account that the classes may have different sizes in $Z_2$ (similar definition is applied for $Z_3$). If there are two classes, for example, with very different sizes and a classifier always assigns the label of the largest class, its accuracy will fall drastically due to the high error rate on the smallest class.

Let $NZ_2(i)$, $i = 1,2\dots,c$, be the number of samples in $Z_2$ from each class $i$. We define

$$e_{i,1} = \frac{\text{FP}(i)}{|Z_2| - |NZ_2(i)|} \text{ and } e_{i,2} = \frac{\text{FN}(i)}{|NZ_2(i)|}, i = 1, \dots, c \quad (4)$$

where $\text{FP}(i)$ and $\text{FN}(i)$ are the false positives and false negatives, respectively. That is, $\text{FP}(i)$ is the number of samples from other classes that were classified as being from the class $i$ in $Z_2$, and $\text{FN}(i)$ is the number of samples from the class $i$ that were incorrectly classified as being from other classes in $Z_2$. The errors $e_{i,1}$ and $e_{i,2}$ are used to define

$$E(i) = e_{i,1} + e_{i,2}, \quad (5)$$

where $E(i)$ is the partial sum error of class $i$. Finally, the accuracies $L(I)$, $I = 1,2\dots,T$, are written as

$$L(I) = \frac{2c - \sum_{i=1}^{c} E(i)}{2c} = 1 - \frac{\sum_{i=1}^{c} E(i)}{2c}. \quad (6)$$

Algorithm 2 presents this learning procedure which has been used for OPF, SVM, ANN-MLP, and $k$-NN, by changing Lines 4 and 19–20.

In OPF, Line 4 is implemented by computing $S^* \subset Z_1$ as described in Section II.A and the predecessor map $P$, label map $L$, and cost map $C$ by Algorithm 2. The classification is done by setting $L(t)L(s^*)$, where $s^* \in Z_1$ is the sample that satisfies Eq. (3). The constraints in Lines 19–20 refer to keep the prototypes out of the sample interchanging process between $Z_1$ and $Z_2$. We do the same with the support vectors in SVM. However, they may be selected for interchanging in future iterations if they are no longer prototypes or support vectors. For SVM, we use the latest version of the LibSVM package (Chang and Lin, 2001) with radial basis function (RBF) kernel, parameter optimization and the one-versus-one strategy for the multiclass problem to implement Line 4.

We use the fast artificial neural network library (FANN) [56] to implement the ANN-MLP. The network configuration is $x : y : z$, where $x = n$ (number of features), $y = |Z_1| - 1$, and $z = c$ (number of classes) are the number of neurons in the input, hidden, and output layers, respectively (Huang and Huang, 1991). In Line 4, the ANN-MLP is trained by back propagation. There is no constraint in

Lines 19–20. However, we keep the weights of the neurons as initial setting for training in the next iteration. For $k$-NN, training in Line 4 involves the computation of the value of $k$ which provides the highest accuracy on $Z_1$ according to the Leave-One-Out approach (Kohavi, 2005). Lines 19–20 are implemented without constraints.

Lines 5–6 initialize the false positive and false negative arrays for accuracy computation. The classification of each sample is performed in Lines 7–13, updating the false positive and false negative arrays. Misclassified samples are stored in the list $LM$ (Line 13). Line 14 computes the accuracy $L(I)$ and Lines 15–16 save the best instance of classifier so far. The inner loop in Lines 17–20 changes the misclassified samples of $Z_2$ by randomly selected samples of $Z_1$, under the aforementioned constraints.

**Algorithm 2**—General Learning Algorithm

**Input:** Training and evaluation sets, $Z_1$ and $Z_2$, labeled by $\lambda$, number $T$ of iterations, and the pair $(v,d)$ for feature vector and distance computations.
**Output:** Learning curve $L$ and the OPF/SVM/ANN-MLP/$k$-NN classifier with highest accuracy.
**Auxiliary:** Arrays FP and FN of sizes $c$ for false positives and false negatives and list $LM$ of misclassified samples.

1.  Set MaxAcc = $-1$.
2.  For each iteration $I = 1,2,\dots T$, do
3.      $LM = \emptyset$
4.      Train OPF/SVM/ANN-MLP/k-NN with Z1.
5.      For each class $i = 1,2,\dots c$, do
6.          FP($i$) = 0 and FN($i$) = 0.
7.      For each sample $t \in Z_2$, do
8.          Use the classifier obtained in Line 4 to classify
9.          $t$ with a label $L(t)$.
10.         If $L(t) \neq \lambda(t)$, then
11.             FP($L(t)$) = FP($L(t)$) + 1.
12.             FN($\lambda(t)$) = FN($\lambda(t)$) + 1.
13.             $LM = LM \cup t$.
14.     Compute accuracy $L(I)$ by Eq. (6).
15.     If $L(I) > MaxAcc$ then save the current instance
16.     of the classifier and set $MaxAcc = L(I)$.
17.     While $LM \neq \emptyset$;
18.         $LM = LM \setminus t$
19.         Replace t by a randomly selected sample of the
20.         of the same class in $Z_1$, under some constraints.

Figure 4 illustrates the learning curve of each classifier for the same dataset and descriptor. Oscillations indicate instability of the classifier (e.g., ANN – MLP) or presence of outliers. The monotonic behavior of the OPF's learning curve is usually observed. Nevertheless, the choice of the classifier instance with highest accuracy aims to avoid outliers in $Z_1$.

## IV. EVALUATION

This section presents the datasets, descriptors, and experiments that compare OPF with SVM, ANN-MLP, and $k$-NN in accuracy and efficiency (computational time).

Table I presents the 11 datasets used in the experiments, with diverse types of samples. The dataset MPEG-7 (2002) uses shape images (Fig. 5), COREL (Corel, 2007) uses color images, and Brodatz (1966) uses texture images (Fig. 6). These datasets allow to evaluate the performance of the classifiers using shape, color, and
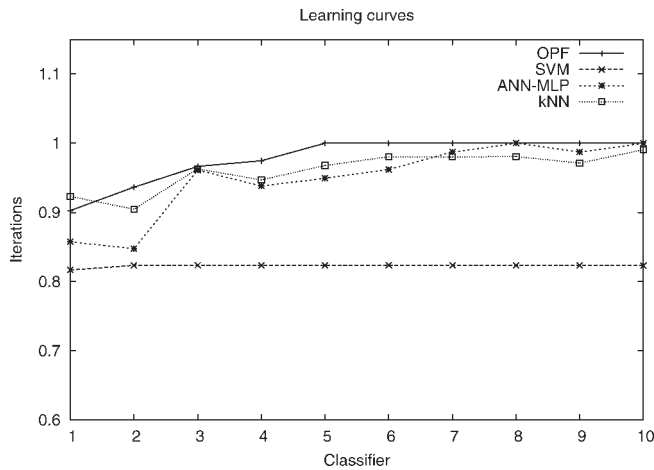
**Figure 4.** Learning curve of each classifier for the dataset $B_8$ using descriptor $D_{10}$ (Tables I and II).

texture descriptors, respectively. The remaining datasets already provide their feature vectors: WBC—Wisconsin Breast Cancer, IS—Image Segmentation, and LR—Letter Recognition (Asuncion and Newman, 2007); Brain (Collins, 1998); and Cone-torus, Saturn, Petals, and Boat (Kuncheva, 1996). The dataset Brain uses voxels as samples from gray and white matter in magnetic resonance images of brain phantoms, with various levels of noise and inhomogeneity that produce outliers. The features are the minimum, maximum, and intensity within a small 3D neighborhood of each voxel. The last four datasets use the $(x, y)$ coordinates of 2D points as features (Fig. 7).

Table II shows 10 different possibilities of combining feature extraction $v$ and distance function $d$ to form descriptors $(v,d)$. Some descriptors were designed for shape ($D_1$–$D_5$), color ($D_6$–$D_7$), and texture ($D_8$) images. Descriptors $D_1$, $D_2$, and $D_3$ use the Fourier coefficients (FC) (Persoon and Fu, 1977), moment invariants (MI) (Hu, 1962), and multiscale fractal dimensions (MSF) (Torres et al., 2004) as shape features, respectively, and Euclidean norm (L2) as distance function. Descriptors $D_4$ and $D_5$ compute three statistical measures, called bean angle statistics (BAS), for each sample on a contour (Arica and Vural, 2003). They use L2 metric and optimal correspondence subsequence (OCS) (Wang and Pavlidis, 1990), respectively, for comparison between feature vectors, illustrating the importance of special distance functions such as OCS. The comparison among descriptors from $D_1$ to $D_5$ using a same classifier illustrates their ability in representing the shapes of a given dataset.

**Table I.** Description of the datasets.

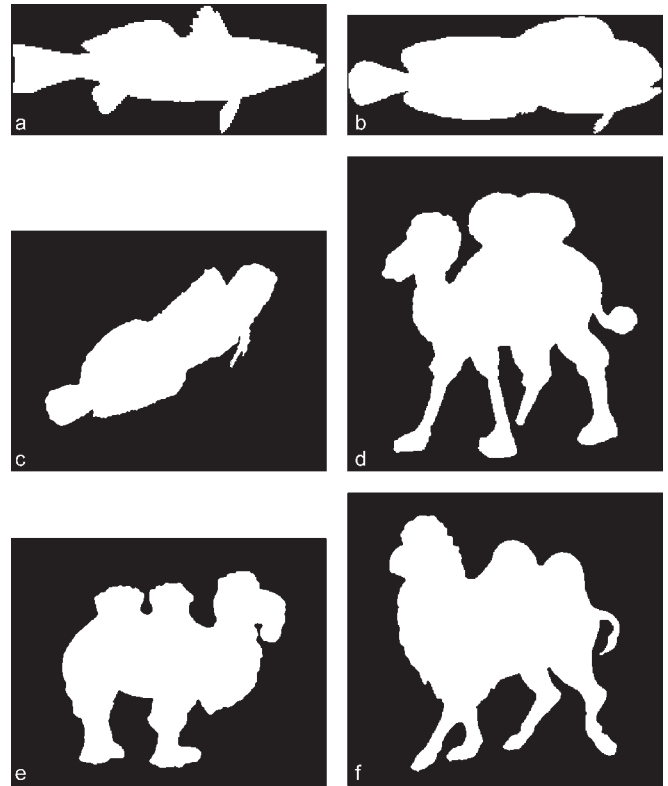| Dataset Code | Dataset Name | Objects Number | Classes Number |
|---|---|---|---|
| $B_1$ | MPEG-7 | 1,400 | 70 |
| $B_2$ | COREL | 1,607 | 49 |
| $B_3$ | Brodatz | 208 | 13 |
| $B_4$ | WBC | 699 | 2 |
| $B_5$ | IS | 2,310 | 7 |
| $B_6$ | LR | 5,000 | 20 |
| $B_7$ | Brain | 1,578 | 2 |
| $B_8$ | Cone-torus | 400 | 3 |
| $B_9$ | Saturn | 200 | 2 |
| $B_{10}$ | Petals | 100 | 4 |
| $B_{11}$ | Boat | 100 | 3 |



**Figure 5.** Examples of the MPEG-7 shapes from the classes (a)–(c) fish and (d)–(f) camel.

Descriptor $D_6$ classifies pixels into border/interior regions and computes color histograms for each region (Stehling et al., 2002). It uses as distance function the L1 metric between the logarithm of the histograms (dLog). Color images are also represented by color histograms (CHIST) (Swain and Ballard, 1991) and compared with L1 metric in the descriptor $D_7$. Descriptor $D_8$ uses steerable pyramid decomposition to create texture features (TEX), which are compared by a rotation-invariant texture matching (RIM) (Montoya-Zegarra et al., 2008). Descriptor $D_9$ represents all feature vectors (OWN) already available in the datasets from $B_4$ to $B_7$ and $D_{10}$ represents the 2D-point (XY) features of the datasets from $B_8$ to $B_{11}$ (Table I). Their distance function is L2. Finally, the combinations between datasets and descriptors are summarized in Table III.

MLP, and others have the distance function embedded in the model, as in the radial basis function (RBF) of the SVM. When the distance function is L2, we use as RBF for SVM

$$K(s,t) = \exp^{-\gamma \|(\vec{v}(s) - \vec{v}(t))\|^2}, \tag{7}$$

where $s$ and $t$ are two samples (one is support vector) and $\vec{v}(s)$ and $\vec{v}(t)$ are their feature vectors. The constant $\gamma$ is found by parameter optimization. In the case of special distances $d$, we have observed a considerable improvement in the SVM's performance when we replace its RBF by

$$K'(s,t) = \exp^{-\gamma d^2(s,t)}. \tag{8}$$

This can be observed in Tables IV and V for dataset $B_1$ (MPEG-7) with $D_4$ (BAS with L2) and $D_5$ (BAS with OCS). Therefore, $K'$ was
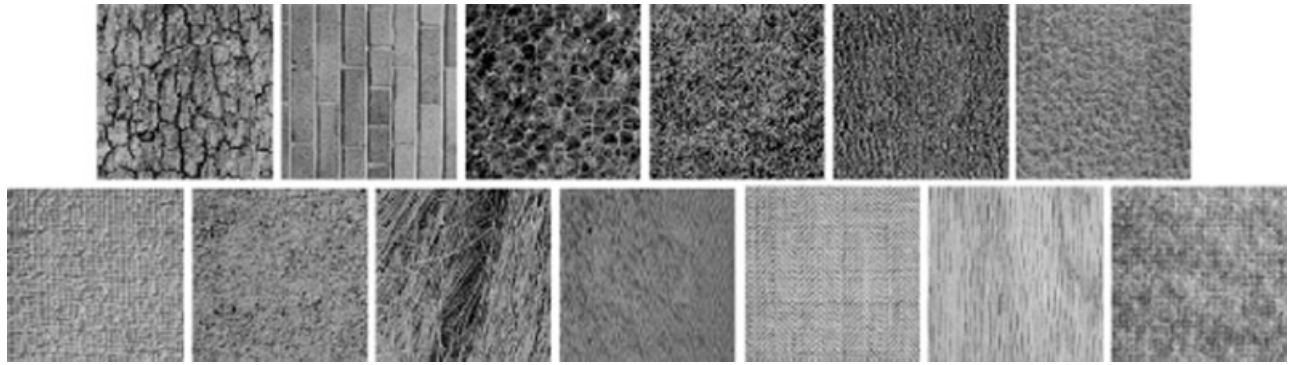
**Figure 6.** Texture images from the Brodatz dataset. Each image, from left to right and from top to bottom, represents a class: bark, brick, bubbles, grass, leather, pigskin, raffia, sand, straw, water, weave, wood, and wool.

used in all experiments involving SVM and special distance functions $d$, and $K$ was used for L2.

The experiments evaluate the accuracy on $Z_3$ and the computational time of each classifier, OPF, SVM, ANN-MLP, and $k$-NN, for each pair dataset and descriptor presented in Table III. In all experiments, the datasets were divided into three parts: a training set $Z_1$ with 30% of the samples, an evaluation set $Z_2$ with 20% of the samples, and a test set $Z_3$ with 50% of the samples. These samples were randomly selected and each experiment was repeated 10 times with different sets $Z_1$, $Z_2$, and $Z_3$ to compute mean (robustness) and standard deviation (precision) of the accuracy values and mean value of kappa (Cohen, 1960). Section IV.A presents the accuracy results of training on $Z_1$ and testing on $Z_3$. The accuracy results of training on $Z_1$, with learning from the errors in $Z_2$, and testing on $Z_3$ are presented in Section IV.B. The average computational time of each classifier for training and classification is divided by the number of samples and reported in Section IV.C.

## A. Accuracy Results on $Z_3$ Without Using $Z_2$.

The results in Table IV are presented as $x \pm y(z)[k]$, where $x$, $y$, $z$, and $k$ are the mean accuracy, its standard deviation, mean kappa coefficient

(Cohen, 1960), and the best value of $k$ obtained for $k$-NN, respectively. Values of kappa below 0.80 indicate the difficulty in classifying some datasets using the respective descriptors. Good descriptors tend to better separate the classes in the feature space, reducing overlap and so facilitating the classification. The results in $B_1$ (MPEG-7), for example, indicate that $D_5$ outperforms the remaining descriptors. Besides, $D_4$ and $D_5$ differ only in the distance function and the results indicate that OCS (Wang and Pavlidis, 1990) outperforms the Euclidean metric. Similarly, one may conclude that $D_6$ (BIC in Stehling et al., 2002) outperforms $D_7$, (color histogram in Swain and Ballard, 1991) in $B_2$ (COREL). Irrespective of that, we are comparing the relative performance of the classifiers.

Most accuracies of OPF and SVM were clearly higher than those of ANN-MLP and $k$-NN. OPF and SVM presented equivalent overall performances, being one better than the other depending on the case. Considering only the cases where the best $k$ is 1 in $k$-NN, we can observe that the criterion of OPF to assign the label of the most strongly connected root to a sample is really more accurate than the label of the closest sample. The instability of ANN-MLP is reflected by the standard deviations, which are about 10 times higher than the standard deviations obtained by the other classifiers. Because of the overlapping between classes, the accuracies of the classifiers in $B_8$ and $B_9$ are lower than their accuracies in $B_{10}$ and $B_{11}$. Because of the quality of the descriptors, similar observation explains the increasing order of accuracy in $B_1$ with $D_1$, $D_3$, $D_2$, $D_4$, and $D_5$.

## B. Accuracy Results on $Z_3$ with Learning on $Z_2$.

To evaluate the ability of each classifier in learning from the errors in $Z_2$ without



**Figure 7.** Datasets of 2D points: (a) cone-torus, (b) saturn, (c) petals, and (d) boat.

**Table II.** Descriptors used in the experiments.

| Descriptor Code | Feature Extraction Algorithm | Distance Function |
|---|---|---|
| $D_1$ | FC | L2 |
| $D_2$ | MI | L2 |
| $D_3$ | MSF | L2 |
| $D_4$ | BAS | L2 |
| $D_5$ | BAS | OCS |
| $D_6$ | BIC | dLog |
| $D_7$ | CHIST | L1 |
| $D_8$ | TEX | RIM |
| $D_9$ | OWN | L2 |
| $D_{10}$ | XY | L2 |

**Table III.** Datasets and the respective descriptors used in the experiments.

| Dataset Code | Descriptor Code |
|---|---|
| $B_1$ | $D_1, D_2, D_3, D_4, D_5$ |
| $B_2$ | $D_6, D_7$ |
| $B_3$ | $D_8$ |
| $B_4$ | $D_9$ |
| $B_5$ | $D_9$ |
| $B_6$ | $D_9$ |
| $B_7$ | $D_9$ |
| $B_8$ | $D_{10}$ |
| $B_9$ | $D_{10}$ |
| $B_{10}$ | $D_{10}$ |
| $B_{11}$ | $D_{10}$ |

increasing the size of $Z_1$, we executed Algorithm 2 for $T = 3$ iterations. The results are presented in Table V.

We can observe that the conclusions drawn from Table V remain the same with respect to the overall performance of the classifiers. In most cases, the general learning algorithm improved the performance of the classifiers with respect to their results in Table IV.

**C. Efficiency Results.** Table VI shows the mean execution time in seconds divided by the number of samples that each classifier takes for training and classification (without learning on $Z_2$) and for each dataset and descriptor.

Note that OPF is extremely fast, except when it uses descriptor $D_5$ (Table II) because of the respective distance function computation. Similar effect can be observed in SVM and $k$-NN. Given that ANN-MLP does not use distance functions, it is free of this problem. On average, the results indicate that our most recent implementation of OPF was about 72 times faster than the latest implementation of SVM (Chang and Lin, 2001), 443 times faster than the fast ANN-MLP (Nissen, 2003), and 1.3 times faster than our implementation of a $k$-NN classifier.

The importance of speed in pattern recognition seems to not have caught much attention. Most of the computational time is spent in training, which is done only once in many applications.

However, take the first case of $B_1$ and $D_1$, for example, where OPF spent 0.0052 s per sample and SVM spent 1.304. For 100,000 samples, this represents 8.67 min using OPF and 36.22 h using SVM. In the case of 2D/3D images, for example, the number of pixels/voxels ranges from thousands to millions, and the time for training becomes a burden. In medical imaging, it is very likely that a new training has to be done for every 3D image, because of their variations in noise, inhomogeneity, and protocols.

## V. CONCLUSIONS

We presented a discrete approach for supervised classification (OPF), which computes an optimum-path forest on a training set and classifies samples with the label of their most strongly connected root in the forest. We also proposed a general learning algorithm, which usually improves performance of the classifiers without increasing the training set. The source code of the supervised OPF is available at http://www.ic.unicamp.br/~afalcao/libopf.

We compared OPF with SVM, ANN-MLP, and $k$-NN using several datasets and descriptors. These experiments involved datasets with shape, color, and texture properties, and datasets commonly used by the machine learning community. The advantage of OPF over the others in computational time is notorious and impressive, which is crucial in the case of large datasets. It can be more or less accurate than SVM, depending on the case, but its accuracy is usually superior to those of ANN-MLP and $k$-NN. OPF also presents some interesting properties. It is fast, simple, multiclass, parameter independent, does not make any assumption about the shape of the classes, and can handle some degree of overlapping between classes.

The OPF classifiers are being successfully used in some real applications: the supervised approach is being used for oropharyngeal dysphagia identification (Spadotto et al., 2008), laryngeal pathology detection (Papa et al., 2008b), and diagnosis of parasites from optical microscopy images (Falcão et al., 2008), and the unsupervised approach is being used for the separation of gray-matter and white-matter in magnetic resonance images of the brain (Cappabianco et al., 2008). In the first three applications, the supervised OPF outperforms SVM in accuracy and efficiency. In all cases, there is no human interaction, however, we also intend to evaluate

**Table IV.** Accuracy results $x \pm y(z)$ on $Z_3$ without using $Z_2$: $x$—mean accuracy, $y$—its standard deviation, and $z$—mean kappa.

| Dataset (Descriptor) | Classifiers | | | |
|---|---|---|---|---|
| | OPF | SVM | ANN-MLP | $k$-NN |
| $B_1(D_1)$ | **71.71 ± 0.01 (0.49)** | 70.07 ± 0.01 (0.40) | 57.28 ± 0.44 (0.14) | 59.38 ± 0.01 (0.17) [1] |
| $B_1(D_2)$ | 79.48 ± 0.01 (0.59) | **82.15 ± 0.01 (0.64)** | 71.48 ± 0.26 (0.46) | 72.04 ± 0.01 (0.64) [1] |
| $B_1(D_3)$ | **75.95 ± 0.01 (0.51)** | 74.49 ± 0.01 (0.50) | 62.98 ± 0.39 (0.25) | 60.16 ± 0.01 (0.19) [1] |
| $B_1(D_4)$ | **87.37 ± 0.01 (0.74)** | 87.05 ± 0.01 (0.75) | 77.99 ± 0.34 (0.57) | 66.55 ± 0.01 (0.67) [1] |
| $B_1(D_5)$ | **95.72 ± 0.01 (0.89)** | 94.92 ± 0.01 (0.88) | 76.29 ± 0.04 (0.55) | 50.70 ± 0.01 (0.31) [1] |
| $B_2(D_6)$ | 86.74 ± 0.01 (0.75) | **90.65 ± 0.01 (0.83)** | 83.07 ± 0.10 (0.64) | 82.83 ± 0.01 (0.70) [1] |
| $B_2(D_7)$ | 80.25 ± 0.01 (0.63) | **83.37 ± 0.01 (0.70)** | 80.07 ± 0.10 (0.61) | 78.03 ± 0.01 (0.61) [1] |
| $B_3(D_8)$ | **88.85 ± 0.02 (0.77)** | 84.27 ± 0.01 (0.68) | 86.97 ± 0.21 (0.73) | 84.52 ± 0.01 (0.80) [1] |
| $B_4(D_9)$ | 93.87 ± 0.01 (0.88) | **95.46 ± 0.01 (0.90)** | 92.83 ± 0.20 (0.86) | 91.85 ± 0.01 (0.81) [3] |
| $B_5(D_9)$ | **79.37 ± 0.01 (0.68)** | 78.35 ± 0.01 (0.59) | 73.35 ± 0.10 (0.68) | 65.89 ± 0.01 (0.41) [2] |
| $B_6(D_9)$ | 90.35 ± 0.01 (0.80) | **93.35 ± 0.01 (0.90)** | 84.72 ± 0.10 (0.73) | 87.20 ± 0.01 (0.79) [2] |
| $B_7(D_9)$ | 90.53 ± 0.01 (0.81) | **93.86 ± 0.01 (0.88)** | 92.94 ± 0.09 (0.85) | 86.39 ± 0.01 (0.73) [1] |
| $B_8(D_{10})$ | **87.29 ± 0.01 (0.71)** | 85.54 ± 0.02 (0.71) | 85.33 ± 0.02 (0.69) | 81.34 ± 0.01 (0.65) [7] |
| $B_9(D_{10})$ | **88.10 ± 0.03 (0.76)** | 86.90 ± 0.05 (0.73) | 83.60 ± 0.54 (0.67) | 81.90 ± 0.02 (0.62) [1] |
| $B_{10}(D_{10})$ | **1.0 ± 0.0 (1.0)** | **1.0 ± 0.0 (1.0)** | **1.0 ± 0.0 (1.0)** | **1.0 ± 0.0 (1.0) [21]** |
| $B_{11}(D_{10})$ | 96.76 ± 0.01 (0.93) | **99.55 ± 0.01 (0.99)** | 97.20 ± 0.36 (0.94) | 93.19 ± 0.01 (0.89) [1] |

The best accuracies are indicated in bold and the best value of $k$ is shown in brackets for the $k$-NN.

**Table V.** Accuracy results $x \pm y(z)$ on $Z_3$ with learning on $Z_2$: $x$—mean accuracy, $y$—its standard deviation, and $z$—mean kappa.

| | Classifiers | | | |
|---|---|---|---|---|
| Dataset (Descriptor) | OPF | SVM | ANN-MLP | $k$-NN |
| $B_1(D_1)$ | **75.94 ± 0.01 (0.51)** | 74.42 ± 0.01 (0.48) | 61.39 ± 0.06 (0.22) | 59.38 ± 0.01 (0.17) [1] |
| $B_1(D_2)$ | 81.20 ± 0.01 (0.60) | **82.03 ± 0.01 (0.64)** | 75.06 ± 0.28 (0.50) | 72.88 ± 0.01 (0.44) [3] |
| $B_1(D_3)$ | **76.72 ± 0.01 (0.53)** | 76.52 ± 0.01 (0.53) | 64.76 ± 0.18 (0.29) | 61.48 ± 0.01 (0.26) [3] |
| $B_1(D_4)$ | **87.65 ± 0.01 (0.75)** | 87.18 ± 0.01 (0.73) | 79.23 ± 0.22 (0.58) | 67.14 ± 0.01 (0.68) [1] |
| $B_1(D_5)$ | **96.08 ± 0.01 (0.90)** | 95.87 ± 0.01 (0.90) | 78.65 ± 0.04 (0.57) | 50.70 ± 0.01 (0.31) [1] |
| $B_2(D_6)$ | 86.90 ± 0.01 (0.76) | **90.80 ± 0.02 (0.84)** | 85.70 ± 0.06 (0.75) | 82.83 ± 0.01 (0.73) [1] |
| $B_2(D_7)$ | 80.29 ± 0.01 (0.63) | 83.66 ± 0.01 (0.71) | **84.70 ± 0.08 (0.79)** | 79.73 ± 0.01 (0.61) [1] |
| $B_3(D_8)$ | 88.54 ± 0.02 (0.77) | 84.37 ± 0.01 (0.68) | **92.29 ± 0.16 (0.84)** | 86.26 ± 0.02 (0.70) [1] |
| $B_4(D_9)$ | 94.17 ± 0.01 (0.89) | **96.07 ± 0.01 (0.91)** | 93.26 ± 0.19 (0.87) | 91.26 ± 0.01 (0.81) [9] |
| $B_5(D_9)$ | **79.90 ± 0.01 (0.70)** | 78.65 ± 0.01 (0.57) | 74.35 ± 0.10 (0.70) | 67.06 ± 0.01 (0.24) [2] |
| $B_6(D_9)$ | 92.07 ± 0.01 (0.84) | **94.31 ± 0.01 (0.93)** | 87.72 ± 0.10 (0.79) | 89.54 ± 0.01 (0.81) [9] |
| $B_7(D_9)$ | 91.53 ± 0.01 (0.83) | **94.00 ± 0.01 (0.88)** | 93.73 ± 0.06 (0.87) | 88.99 ± 0.01 (0.76) [1] |
| $B_8(D_{10})$ | **88.38 ± 0.01 (0.72)** | 87.95 ± 0.03 (0.74) | 85.58 ± 0.01 (0.70) | 83.43 ± 0.01 (0.68) [11] |
| $B_9(D_{10})$ | **89.30 ± 0.02 (0.78)** | **89.30 ± 0.03 (0.78)** | 88.10 ± 0.43 (0.76) | 83.90 ± 0.02 (0.63) [1] |
| $B_{10}(D_{10})$ | **1.0 ± 0.0 (1.0)** | **1.0 ± 0.0 (1.0)** | **1.0 ± 0.0 (1.0)** | **1.0 ± 0.0 (1.0)** [5] |
| $B_{11}(D_{10})$ | 97.35 ± 0.02 (0.94) | **99.85 ± 0.01 (0.99)** | 98.23 ± 0.26 (0.96) | 94.81 ± 0.03 (0.86) [1] |

The best accuracies are indicated in bold and the best value of $k$ is shown in brackets for the $k$-NN.

the supervised OPF for interactive segmentation of brain tissues, where the user selects training markers. In this case the method becomes similar to an IFT-watershed approach, except for the fact that it works in the feature space with no spatial connectivity constraint, which is important for tissues with disconnected voxels.

Applications with large datasets definitely favor OPF with respect to SVM. We must say that, as a discrete approach, the performance of OPF may be reduced for small training sets, if the number of samples are not enough to represent the classes. In SVM, this may also be a problem, but as it estimates a decision hyperplane, it has a chance to divide the feature space with separation between classes. Too much overlapping between classes may also represent an advantage for SVM with respect to OPF, because its transformation to a higher-dimensional space may separate the classes, solving the problem.

The OPF classifier is an important contribution for pattern recognition and related fields, which also opens new research problems. One can investigate the optimum-path forest classification using incomplete graphs (e.g., graphs where the arcs are between $k$-nearest neighbors), different connectivity functions, and other algorithms to estimate prototypes and to learn from the errors in the evaluation set. The use of genetic programming (Koza, 1992) (GP) for arc-weight estimation in OPF is also an alternative to deal with class overlapping, by combining the distances from multiple descriptors in a nonlinear way (Torres et al., 2008).

## ACKNOWLEDGMENTS

**Table VI.** Mean execution times in seconds for training and classification divided by the number of samples.

| | Classifiers | | | |
|---|---|---|---|---|
| Dataset (Descriptor) | OPF | SVM | ANN-MLP | $k$-NN |
| $B_1(D_1)$ | **0.0052** | 1.304 | 0.8973 | 0.0450 |
| $B_1(D_2)$ | **0.0010** | 0.0599 | 0.3453 | 0.0034 |
| $B_1(D_3)$ | **0.0012** | 0.0742 | 0.3603 | 0.0038 |
| $B_1(D_4)$ | **0.0019** | 0.2859 | 0.5043 | 0.0126 |
| $B_1(D_5)$ | 5.8400 | 7.3926 | **0.5031** | 5.8432 |
| $B_2(D_6)$ | **0.0185** | 0.1813 | 0.2553 | 0.0199 |
| $B_2(D_7)$ | **0.0010** | 0.0997 | 0.2491 | 0.0254 |
| $B_3(D_8)$ | **0.2163** | 0.2333 | 0.2891 | 0.2164 |
| $B_4(D_9)$ | **0.0019** | 0.0091 | 0.01505 | 0.0042 |
| $B_5(D_9)$ | 0.0018 | 0.0693 | 0.400 | **0.0010** |
| $B_6(D_9)$ | **0.0012** | 0.0320 | 0.0800 | 0.0017 |
| $B_7(D_9)$ | **0.0011** | 0.1662 | 3.5625 | 0.01960 |
| $B_8(D_{10})$ | **0.0010** | 0.1281 | 1.8540 | 0.0011 |
| $B_9(D_{10})$ | 0.0011 | 0.1445 | 0.3828 | **0.0002** |
| $B_{10}(D_{10})$ | 0.0018 | 0.0078 | 0.0020 | **0.0003** |
| $B_{11}(D_{10})$ | **0.0019** | 0.0594 | 0.0039 | **0.0019** |

The best times are indicated in bold.

## REFERENCES

C. Allène, J.Y. Audibert, M. Couprie, J. Cousty, and R. Keriven, "Some links between min-cuts, optimal spanning forests and watersheds," Mathematical Morphology and its Applications to Image and Signal Processing, MCT/INPE, 2007, pp. 253–264.

N. Arica and F.T.Y. Vural, BAS: A perceptual shape descriptor based on the beam angle statistics. Pattern Recognit Lett 24 (2003), 1627–1639.

A. Asuncion and D.J. Newman, UCI machine learning repository, University of California, Irvine, CA, 2007.

R. Audigier and R.A. Lotufo, "Seed-relative segmentation robustness of watershed and fuzzy connectedness approaches," XX Brazilian Symposium on Computer Graphics and Image Processing (SIBGRAPI), IEEE CPS, Belo Horizonte, MG, 2007a, pp. 61–68.

R. Audigier and R.A. Lotufo, "Watershed by image foresting transform, tie-zone, and theoretical relationship with other watershed definitions," Mathematical Morphology and its Applications to Signal and Image Processing (ISMM), MCT/INPE, Rio de Janeiro, RJ, 2007b, pp. 277–288.

A. Blum and T. Mitchell, "Combining labeled and unlabeled data with co-training," Proceedings of the 9th Annual Conference on Computational Learning Theory, ACM Press, New York, NY, 1998, pp. 92–100.

B.E. Boser, I.M. Guyon, and V.N. Vapnik, "A training algorithm for optimal margin classifiers," Proceedings of the 5th Workshop on Computational Learning Theory, ACM Press, New York, NY, 1992, pp. 144–152.

P. Brodatz, Textures: a photographic album for artists and designers, Dover, New York, 1966.

J. Callut, K. Fançoisse, and M. Saerens, "Semi-supervised classification in graphs using bounded random walks," Proceedings of the 17th Annual Machine Learning Conference of Belgium and the Netherlands (Benelearn), 2008, pp. 67–68.

F.A.M. Cappabianco, A.X. Falcão, and L.M. Rocha, "Clustering by optimum path forest and its application to automatic GM/WM classification in MR-T1 images of the brain," The Fifth IEEE International Symposium on Biomedical Imaging: From Nano to Macro (ISBI), 2008, pp. 428–431.

C.C. Chang and C.J. Lin, LIBSVM: A library for support vector machines, 2001, Software available at url: http://www.csie.ntu.edu.tw/~cjlin/libsvm.

J. Cohen, A coefficient of agreement for nominal scales. Educ Psychol Meas 20 (1960), 37–46.

D. Collins, A. Zijdenbos, V. Kollokian, J. Sled, N. Kabani, C. Holmes, and A. Evans, Design and construction of a realistic digital brain phantom, IEEE Trans Med Imaging 17 (1998), 463–468.

R. Collobert and S. Bengio, "Links between perceptrons, mlps and svms," Proceedings of the 21th International Conference on Machine Learning, ACM Press, New York, NY, 2004, p. 23.

Corel Corporation, Corel stock photo images, 2007, available at: http://www.corel.com.

T. Cormen, C. Leiserson, and R. Rivest, Introduction to algorithms, MIT, Cambridge, England, 1990.

T. Cover and P. Hart, Nearest neighbor pattern classification, IEEE Trans Inf Theory 13 (1967), 21–27.

K. Duan and S.S. Keerthi, "Which is the best multiclass SVM method? An empirical study," Multiple Classifier Systems, 2005, pp. 278–285.

R.O. Duda, P.E. Hart, and D.G. Stork, Pattern classification, 2nd ed., Wiley-Interscience, New York, 2000.

A.X. Falcão, J. Stolfi, and R.A. Lotufo, The image foresting transform: Theory, algorithms, and applications, IEEE Trans Pattern Anal Mach Intell 26 (2004), 19–29.

A.X. Falcão, C.T.N. Suzuki, J.F. Gomes, J.P. Papa, L.C.S. Dias, and S.H. Shimizu, A system for diagnosing intestinal parasites by computerized image analysis, PCTWO/2008/064442, 2008.

K. Fukunaga and P.M. Narendra, A branch and bound algorithms for computing k-nearest neighbors, IEEE Trans Comput 24 (1975), 750–753.

S. Haykin, Neural networks: A comprehensive foundation, Prentice Hall, Upper Saddle River, NJ, 1994.

G.T. Herman and B.M. Carvalho, Multiseeded segmentation using fuzzy connectedness, IEEE Trans Pattern Anal Mach Intell 23 (2001), 460–474.

M.K. Hu, Visual pattern recognition by moment invariants, IRE Trans Inf Theory 8 (1962), 179–187.

S.C. Huang and Y.F. Huang, Bounds on the number of hidden neurons in multilayer perceptrons, IEEE Trans Neural Networks 2 (1991), 47–55.

L.J. Hubert, Some applications of graph theory to clustering, Psychometrika 39 (1974), 283–309.

A.K. Jain and R.C. Dubes, Algorithms for clustering data, Prentice-Hall, Upper Saddle River, NJ, 1988.

T. Joachims, "Transductive inference for text classification using support vector machines," Proceedings of the 16th International Conference on Machine Learning, Morgan Kaufmann Publishers, San Francisco, CA, 1999, pp. 200–209.

L. Kuncheva, Artificial data, School of Informatics, University of Wales, Bangor, 1996, available at: http://www.informatics.bangor.ac.uk/kuncheva.

L.I. Kuncheva, Combining pattern classifiers: Methods and algorithms, Wiley-Interscience, New York, 2004.

B. Kulis, S. Basu, I. Dhillon, and R. Mooney, "Semi-supervised graph clustering: A kernel approach," ICML '05: Proceedings of the 22nd International Conference on Machine Learning, ACM, New York, NY, 2005, pp. 457–464.

N. Kumar and K. Kummamuru, Semisupervised clustering with metric learning using relative comparisons, IEEE Trans Knowledge Data Eng 20 (2008), 496–503.

R. Kohavi, "A study of cross-validation and bootstrap for accuracy estimation and model selection," IJCAI, 1995, pp. 1137–1145.

J.R. Koza, Genetic programming: On the programming of computers by means of natural selection (complex adaptive systems), MIT, Cambridge, MA, 1992.

R.A. Lotufo and A.X. Falcão, "The ordered queue and the optimality of the watershed approaches," Mathematical Morphology and its Applications to Image and Signal Processing (ISMM'00), Kluwer, 2000, pp. 341–350.

MPEG-7. MPEG-7: The generic multimedia content description standard, Part 1. IEEE MultiMedia 09 (2002), 78–87.

P.A.V. Miranda, A.X. Falcão, A. Rocha, and F.P.G. Bergo, Object delineation by $\kappa$-connected components, EURASIP J Adv Signal Process 2008 (2008), 1–15.

J.A. Montoya-Zegarra, J.P. Papa, N.J. Leite, R.S. Torres, and A.X. Falcão, Learning how to extract rotation-invariant and scale-invariant features from texture images, EURASIP J Adv Signal Process 2008 (2008), 1–16.

S. Nissen, Implementation of a fast artificial neural network library (FANN), Department of Computer Science University of Copenhagen (DIKU), 2003, Software available at http://leenissen.dk/fann/.

N. Panda, E.Y. Chang, and G. Wu, "Concept boundary detection for speeding up svms," Proceedings of the 23th international conference on machine learning, ACM Press, New York, NY, 2006, pp. 681–688.

J.P. Papa, A.X. Falcão, P.A.V. Miranda, C.T.N. Suzuki, and N.D.A. Mascarenhas, "Design of robust pattern classifiers based on optimum-path forests," Mathematical Morphology and its Applications to Image and Signal Processing (ISMM'07), MCT/INPE, 2007, pp. 337–348.

J.P. Papa, A.X. Falcão, C.T.N. Suzuki, and N.D.A. Mascarenhas, "A discrete approach for supervised pattern recognition," 12th International Workshop on Combinatorial Image Analysis, LNCS, Springer, Berlin/Heidelberg, 2008a, pp. 136–147.

J.P. Papa, A.A. Spadotto, A.X. Falcão, and J.C. Pereira, "Optimum path forest classifier applied to laryngeal pathology detection," Proceedings of the 15th International Conference on Systems, Signals, and Image Processing, IEEE, Bratislava, Slovakia, 2008b, pp. 249–252.

E. Persoon and K. Fu, Shape discrimination using Fourier descriptors, IEEE Trans Syst, Man, Cybernet 7 (1977), 170–178.

L. Reyzin and R.E. Schapire, "How boosting the margin can also boost classifier complexity," Proceedings of the 23th International Conference on Machine Learning, ACM Press, New York, NY, 2006, pp. 753–760.

L.M. Rocha, A.X. Falcão, and L.G.P. Meloni, "A robust extension of the mean shift algorithm using optimum path forest," 8th International Workshop on Combinatorial Image Analysis, Buffalo, NY, 2008, pp. 29–38.

P.K. Saha and J.K. Udupa, Relative fuzzy connectedness among multiple objects: Theory, algorithms, and applications in image segmentation, Comput Vis Image Understand 82 (2001), 42–56.

R.O. Stehling, M.A. Nascimento, and A.X. Falcão, "A compact and efficient image retrieval approach based on border/interior pixel classification," Proceedings of the 11th International Conference on Information and Knowledge Management, ACM Press, New York, NY, 2002, pp. 102–109.

J. Shi and J. Malik, Normalized cuts and image segmentation, IEEE Trans Pattern Anal Mach Intell 22 (2000), 888–905.

A.A. Spadotto, J.C. Pereira, R.C. Guido, J.P. Papa, A.X. Falcão, A.R. Gatto, P.C. Cola, and A.O. Schelp, "Oropharyngeal dysphagia identification using wavelets and optimum path forest," Proceedings of the 3rd IEEE International Symposium on Communications, Control and Signal Processing, St. Julians, Malta, Greece, 2008, pp. 735–740.

M. Swain and D. Ballard, Color indexing, Int J Comput Vis 7 (1991), 11–32.

B. Tang and D. Mazzoni, "Multiclass reduced-set support vector machines," Proceedings of the 23th International Conference on Machine Learning, ACM Press, New York, NY, 2006, pp. 921–928.

R. Torres, A.X. Falcão, and L.F. Costa, A graph-based approach for multiscale shape analysis, Pattern Recognit 37 (2004), 1163–1174.

R.S. Torres, A.X. Falcão, M.A. Gonçalves, J.P. Papa, B. Zhang, W. Fan, and E.A. Fox, A genetic programming framework for content-based image retrieval, Pattern Recognit 42 (2009), 283–292.

Y.P. Wang and T. Pavlidis, Optimal correspondence of string subsequences, IEEE Trans Pattern Anal Mach Intell 12 (1990), 1080–1087.

C.T. Zahn, Graph-theoretical methods for detecting and describing gestalt clusters, IEEE Trans Comput C-20 (1971), 68–86.

D. Zhou, B. Schölkopf, and T. Hofmann, Semi-supervised learning on directed graphs, Adv Neural Inf Process Syst 17 (2005), 1633–1640.

X. Zhu, Semi-supervised learning literature survey, Technical report 1530, Computer Sciences, University of Wisconsin-Madison, 2006.