

INSTITUTO DE COMPUTAÇÃO  
UNIVERSIDADE ESTADUAL DE CAMPINAS

**Supervised Pattern Classification  
based on Optimum-Path Forest**

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

Technical Report - IC-08-20 - Relatório Técnico

September - 2008 - Setembro

The contents of this report are the sole responsibility of the authors.  
O conteúdo do presente relatório é de única responsabilidade dos autores.

# Supervised Pattern Classification based on Optimum-Path Forest

João P. Papa

Alexandre X. Falcão

Celso T. N. Suzuki

Institute of Computing, State University of Campinas, Brazil  
Email: afalcao@ic.unicamp.br, {papa.joaopaulo,celso.suzuki}@gmail.com

## Abstract

We present an approach for supervised classification, which interprets a training set as a complete graph, identifies prototypes in all classes, and computes an optimum-path forest rooted at them. The class of a sample in a tree is assumed to be the same of its root. A test sample is classified by identifying which tree would contain it. We show how to improve performance from the errors on an evaluation set, without increasing the training set. The advantages over others approaches are demonstrated using several experiments.

## 1 Introduction

Patterns are usually represented by feature vectors (set of measures or observations) obtained from samples of a dataset [1]. Two fundamental problems in pattern recognition are: (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). 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 (semi-supervised learning [2–4]). 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) non-separable classes with arbitrary shapes. Any reasonable approach should handle (a) and (b), being (c) the most interesting challenge. An artificial neural network with multi-layer perceptrons (ANN-MLP), for example, can address (a) and (b), but not (c) [5]. As an unstable classifier, collections of ANN-MLP [6] can improve its performance up to some unknown limit of classifiers [7]. Support vector machine (SVM) has been proposed to overcome the problem, by assuming linearly separable classes in a higher-dimensional feature space [8]. 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 multi-class problem [9]. Tang and Mazzoni [10] proposed a method to reduce the number of support vectors in the multi-class 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. [11] 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 non-support vectors. However, in all SVM approaches, the assumption of separability may also not be valid in any space of finite dimension [12].

We propose a supervised classifier based on *optimum-path forest* (OPF), which is fast, simple, multi-class, parameter independent, does not make any assumption about shape, and can handle some degree of separability between classes. The training set is thought of as a complete graph, whose the nodes are the samples and the arcs link all pairs of nodes. The arcs are weighted by the distance 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  $t$ . By estimating prototypes as the closest samples from distinct classes, our approach can handle all three cases with the maximum arc-weight function. In the case of overlapping between classes, these prototypes will be class defenders in the overlapped regions of the feature space (Figure 1c).

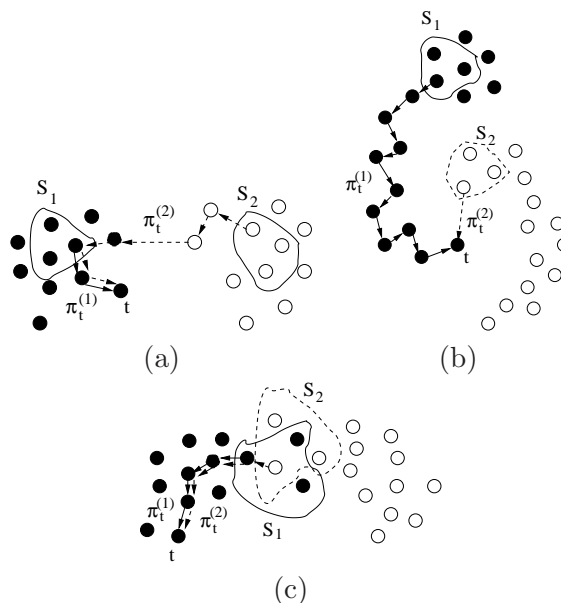


Figure 1: Examples of 2D feature spaces using two classes: (a) linearly separable, (b) piecewise linearly separable, and (c) non-separable classes with arbitrary shapes. Prototypes can be identified in each class, forming the sets  $S_1$  and  $S_2$ . Every sample  $t$  can be connected to a prototype 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.

The optimum paths from the prototypes to other samples are computed by the image foresting transform (IFT) — a tool for the design of image processing operators based on connectivity [13] — which is extended here from the image domain to the feature space. 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. We can exploit other connectivity functions, but this work presents only the results for the maximum arc-weight function. 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. Note the difference between the proposed method with the maximum arc-weight function and the nearest neighbor approach [14]. A

test/training sample may be assigned to a given class, even when its closest labeled sample is from another class (Figure 1b).

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 [15]. Its classification uses the direct distance between samples rather than the connectivity values based on sequences. The dataset partition obtained by the minimization of the maximum arc-weight function in the feature space is equivalent to the image segmentation by IFT-watershed transform from labeled markers [16, 17]. In our case, the markers are the prototypes and we have a special way to estimate them. Similar important relations can be obtained with other image operators, such as relative-fuzzy connected segmentation [18–21]

There are several graph-based approaches for pattern classification [22–25]. Most of them is unsupervised and the idea of modeling supervised classification as an optimum-path forest problem seems to be novel. We have also proposed an unsupervised approach based on optimum-path forest [26, 27]. Besides the learning strategy, the present work differs from that in the graph model, connectivity function, and estimation of prototypes. Previous versions of this work have also been published [28–30]. 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 can not 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 evaluation 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 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. 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 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 paper describes the OPF classifier in Section 2, presents a general learning algorithm in Section 3, which follows the same aforementioned strategy for all classifiers, shows results that compare the OPF classifier with SVM [8], ANN-MLP [5] and  $k$ -NN [15] in Section 4, and states conclusions in Section 5.

## 2 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 (3D pixels), and contours in this paper. 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 pseudo-test 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  $i$ ,  $i = 1, 2, \dots, c$ , 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  $d(s, t)$  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 [31, 32]. 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 4 use shape [33], texture [30] and color [34] descriptors based on this definition.

Our problem consists of projecting a classifier which can predict the correct label  $\lambda(s)$  of any sample  $s \in Z_3$ . We propose a classifier which creates a discrete optimal partition of the feature space such that any sample  $s \in Z_3$  can be classified according to this partition. This partition is an optimum-path forest (OPF) computed on  $Z_1$  by the image foresting transform (IFT) algorithm [13].

Let  $(Z_1, A)$  be a complete graph whose nodes are the training samples and any pair of samples defines an arc in  $A = Z_1 \times Z_1$  (Figure 2a). 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  $\pi_t = \langle s_1, s_2, \dots, t \rangle$  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  $(s, t)$ .

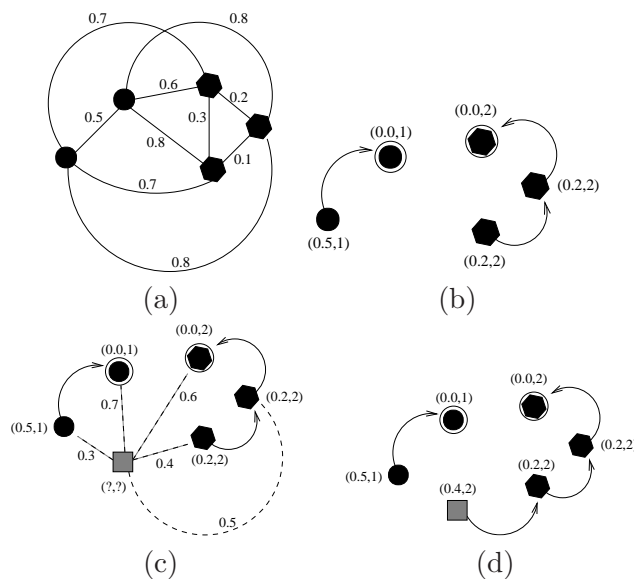


Figure 2: (a) Complete weighted graph for a simple training set. (b) Resulting optimum-path forest for  $f_{\max}$  and two given prototypes (circled nodes). The entries  $(x, y)$  over the nodes are, respectively, the cost and the label of the samples. The directed arcs indicate the predecessor nodes 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.

The OPF algorithm may be used with any *smooth* connectivity function which can group samples with similar properties [13]. 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

$$(a) f(\pi_s) \leq f(\pi_t),$$

(b)  $\pi_s$  is optimum,

(c) for any optimum path  $\tau_s$ ,  $f(\tau_s \cdot \langle s, t \rangle) = f(\pi_t)$ .

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

$$\begin{aligned} f_{\max}(\langle s \rangle) &= \begin{cases} 0 & \text{if } s \in S, \\ +\infty & \text{otherwise} \end{cases} \\ f_{\max}(\pi_s \cdot \langle s, t \rangle) &= \max\{f_{\max}(\pi_s), d(s, t)\} \end{aligned} \quad (1)$$

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 OPF algorithm minimizes  $f_{\max}$ , by storing the minimum costs in a map  $C$ ,

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

and assigning one optimum path  $P^*(t)$  from  $S$  to every sample  $t \in Z_1$ . Its result is 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 Figure 2b). 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  $L(t) \leftarrow \lambda(R(t))$ .

### 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.      $C(s) \leftarrow 0$ ,  $P(s) \leftarrow \text{nil}$ ,  $L(s) \leftarrow \lambda(s)$ , 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  $Q$ .

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 non-decreasing order of 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.

## 2.1 Training

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) [35] and

optimum-path tree for  $f_{\max}$  [21, 36]. The training essentially consists of finding  $S^*$  and an OPF classifier rooted at  $S^*$ .

By computing an 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 (Figure 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.

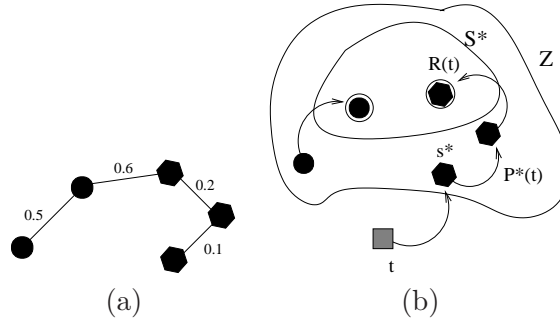


Figure 3: (a) MST of the graph shown in Figure 2a where the optimum prototypes share the arc of weight 0.6. (b) The classification of a test sample (gray square)  $t$  as in Figure 2c assigns the optimum path  $P^*(t)$  from  $R(t) \in S^*$  to  $t$  passing through  $s^*$ .

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  $Z_1$  (Figure 2b). 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 above aims to block these passages, reducing the chances of samples in any given class be reached by optimum paths from prototypes of other classes.

## 2.2 Classification

For any sample  $t \in Z_3$ , we consider all arcs connecting  $t$  with samples  $s \in Z_1$ , as though  $t$  were part of the training graph (Figure 2c). 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  $R(t) \in S^*$  (Figure 3b). This path can be identified incrementally, by evaluating the optimum cost  $C(t)$  as

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

Let the node  $s^* \in Z_1$  be the one that satisfies Equation 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  $t$  (Figure 2d). 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  $Z_3$ .



### 3 Learning from errors on the evaluation set

This section shows how to use a third evaluation set  $Z_2$  to improve the composition of samples in  $Z_1$  without increasing its size. We present a general learning algorithm that interchanges misclassified samples of  $Z_2$  with randomly selected samples of  $Z_1$  (under certain constraints) during a few iterations. This procedure assumes that the most informative samples can be obtained from the errors. The accuracy values along the iterations define a *learning curve*. If the misclassified samples are informative, then the accuracy of the classifier, retrained with the next set  $Z_1$ , should increase on the next set  $Z_2$ . At the end, we select the most accurate instance of each classifier to test it on  $Z_3$ . By comparing the accuracies of the classifiers on  $Z_3$ , before and after the learning process, we can evaluate their learning capacity from the errors.

The accuracy  $Acc$  is measured by taking into account that the classes may have different sizes in  $Z_3$  (similar definition is applied for  $Z_2$ ). 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  $NZ_3(i)$ ,  $i = 1, 2, \dots, c$ , be the number of samples in  $Z_3$  from each class  $i$ . We define

$$e_{i,1} = \frac{FP(i)}{|Z_3| - |NZ_3(i)|} \quad \text{and} \quad e_{i,2} = \frac{FN(i)}{|NZ_3(i)|}, \quad i = 1, \dots, c \quad (4)$$

where  $FP(i)$  and  $FN(i)$  are the false positives and false negatives, respectively. That is,  $FP(i)$  is the number of samples from other classes that were classified as being from the class  $i$  in  $Z_3$ , and  $FN(i)$  is the number of samples from the class  $i$  that were incorrectly classified as being from other classes in  $Z_3$ . 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 accuracy  $Acc$  of the classification is written as

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

Algorithm 2 outputs a *learning curve* and the instance of classifier with highest accuracy on  $Z_2$  after  $T$  iterations. It is used for OPF, SVM, ANN-MLP and  $k$ -NN, by changing Lines 3 and 16 – 17.

#### 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  $\mathcal{L}$  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. For each iteration  $I = 1, 2, \dots, T$ , do
2.      $LM \leftarrow \emptyset$
3.     Train OPF/SVM/ANN-MLP/ $k$ -NN with  $Z_1$ .
4.     For each class  $i = 1, 2, \dots, c$ , do
5.          $FP(i) \leftarrow 0$  and  $FN(i) \leftarrow 0$ .
6.     For each sample  $t \in Z_2$ , do
7.         Use the classifier obtained in Line 3 to classify  $t$
8.         with a label  $L(t)$ .
9.         If  $L(t) \neq \lambda(t)$ , then
10.              $FP(L(t)) \leftarrow FP(L(t)) + 1$ .
11.              $FN(\lambda(t)) \leftarrow FN(\lambda(t)) + 1$ .
12.              $LM \leftarrow LM \cup t$ .
13.     Compute  $Acc$  by Equation 6 and save the current
14.     instance of the classifier.
15.      $\mathcal{L}(I) \leftarrow Acc$



16.  $\left[ \begin{array}{l} \text{While } LM \neq \emptyset \\ \left[ \begin{array}{l} LM \leftarrow LM \setminus t \\ \text{Replace } t \text{ by a randomly selected sample of the} \\ \text{same class in } Z_1, \text{ under some constraints.} \end{array} \right. \end{array} \right.$
- 17.
- 18.
- 19.
20. Select the instance of the classifier with highest accuracy  $\mathcal{L}(I)$ .

In OPF, Line 3 is implemented by computing  $S^* \subset Z_1$  as described in Section 2.1 and the predecessor map  $P$ , label map  $L$  and cost map  $C$  by Algorithm 1. The classification is done by setting  $L(t) \leftarrow L(s^*)$ , where  $s^* \in Z_1$  is the sample that satisfies Equation 3. The constraints in Lines 18 – 19 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 [37] with Radial Basis Function (RBF) kernel, parameter optimization and the one-versus-one strategy for the multi-class problem to implement Line 3.

We use the Fast Artificial Neural Network Library (FANN) [38] 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 [39]. In Line 3, the ANN-MLP is trained by back propagation. There is no constraint in Lines 18 – 19. However, we keep the weights of the neurons as initial setting for training in the next iteration. For  $k$ -NN, training in Line 3 involves the computation of the value of  $k$  which provides the highest accuracy on  $Z_1$  according to the Leave-One-Out approach [40]. Lines 18 – 19 are implemented without constraints.

Lines 4 – 5 initialize the false positive and false negative arrays. The classification of each sample is performed in Lines 6 – 12, updating the false positive and false negative arrays. Misclassified samples are stored in the list  $LM$  (Line 12). Line 12 saves the current instance of the classifier and computes the accuracy at iteration  $I$ , which is stored in the learning curve  $\mathcal{L}$  (Lines 13 – 14). The inner loop in Lines 16 – 19 changes the misclassified samples of  $Z_2$  by randomly selected samples of  $Z_1$ , under the aforementioned 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  $Z_1$ .

Dataset Code	Dataset Name	objects number	classes number
B <sub>1</sub>	MPEG-7	1400	70
B <sub>2</sub>	COREL	1607	49
B <sub>3</sub>	Brodatz	208	13
B <sub>4</sub>	WBC	699	2
B <sub>5</sub>	IS	2310	7
B <sub>6</sub>	LR	5000	20
B <sub>7</sub>	Brain	1578	2
B <sub>8</sub>	Cone-torus	400	3
B <sub>9</sub>	Saturn	200	2
B <sub>10</sub>	Petals	100	4
B <sub>11</sub>	Boat	100	3

Table 1: Description of the datasets.

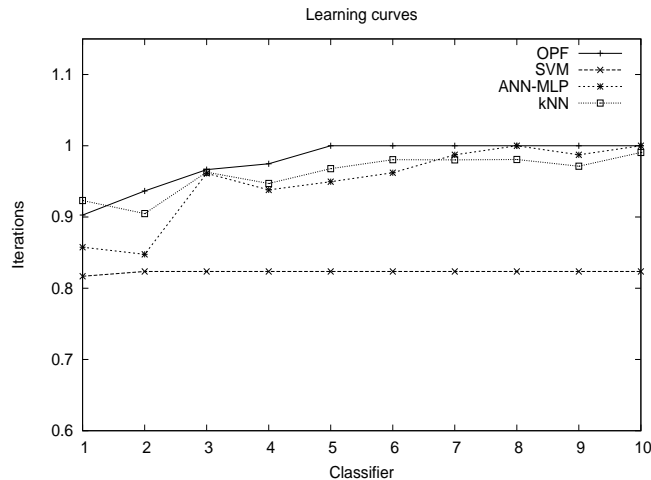


Figure 4: Learning curve of each classifier for the dataset  $B_8$  using descriptor  $D_{10}$  (Tables 1 and 2).

## 4 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 1 presents the 11 datasets used in the experiments, with diverse types of samples. The dataset MPEG-7 [41] uses shape images (Figure 5), COREL [42] uses color images (Figure 6), presented here in gray-scale, and Brodatz [43] uses texture images (Figure 7). These datasets allow to evaluate the performance of the classifiers using shape, color and texture descriptors, respectively. The remaining datasets already provide their feature vectors: WBC - Wisconsin Breast Cancer, IS - Image Segmentation, and LR - Letter Recognition [44]; Brain [45]; and Cone-torus, Saturn, Petals, and Boat [46]. 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 (Figure 8).

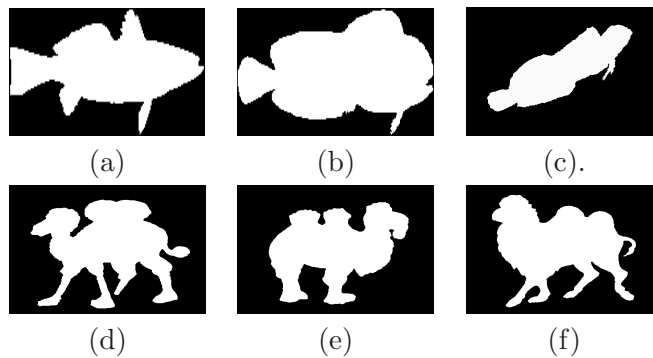


Figure 5: Examples of the MPEG-7 shapes from the classes (a)-(c) fish and (d)-(f) camel.

Table 2 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) [47], Moment Invariants (MI) [48]

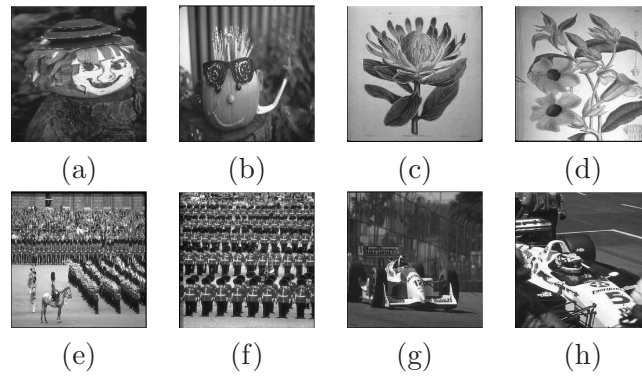


Figure 6: Examples of the Corel images from the classes (a)-(b) pumpkin, (c)-(d) flowers, (e)-(f) British army, and (g)-(h) race cars.

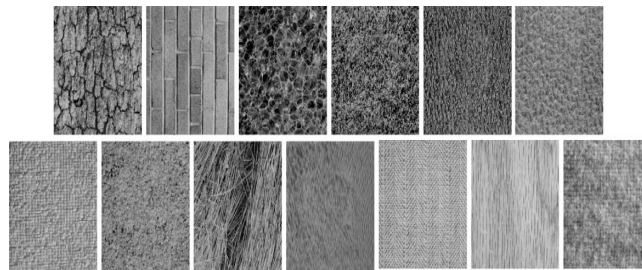


Figure 7: Texture images from the Brodatz dataset. Each image, from left to right and from top to bottom, represents a class: bark, brick, bubbles, grass, leather, pigskin, raffia, sand, straw, water, weave, wood, and wool.

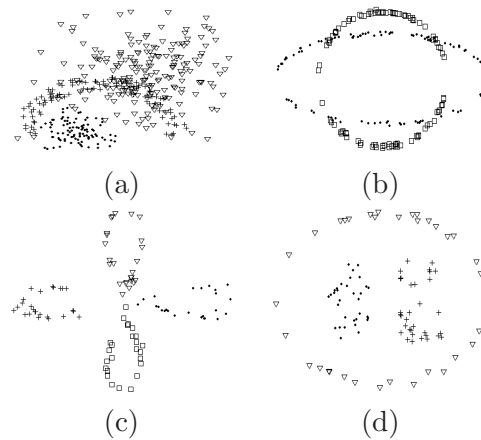


Figure 8: Datasets of 2D points: (a) Cone-torus, (b) Saturn, (c) Petals, and (d) Boat.

and multiscale fractal dimensions (MSF) [49] 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 [33]. They use L2 metric and optimal correspondence subsequence (OCS) [31], 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. Descriptor  $D_6$  classifies pixels into border/interior regions and computes color histograms for each region [34]. It uses as distance function

Descriptor Code	Feature extraction algorithm	Distance function
D <sub>1</sub>	FC	L2
D <sub>2</sub>	MI	L2
D <sub>3</sub>	MSF	L2
D <sub>4</sub>	BAS	L2
D <sub>5</sub>	BAS	OCS
D <sub>6</sub>	BIC	dLog
D <sub>7</sub>	CHIST	L1
D <sub>8</sub>	TEX	RIM
D <sub>9</sub>	OWN	L2
D <sub>10</sub>	XY	L2

Table 2: Descriptors used in the experiments.

the L1 metric between the logarithm of the histograms (dLog). Color images are also represented by color histograms (CHIST) [50] and compared with L1 metric in the descriptor D<sub>7</sub>. Descriptor D<sub>8</sub> uses steerable pyramid decomposition to create texture features (TEX), which are compared by a rotation-invariant texture matching (RIM) [30]. Descriptor D<sub>9</sub> represents all feature vectors (OWN) already available in the datasets from B<sub>4</sub> to B<sub>7</sub> and D<sub>10</sub> represents the 2D-point (XY) features of the datasets from B<sub>8</sub> to B<sub>11</sub> (Table 1). Their distance function is L2. Finally, the combinations between datasets and descriptors are summarized in Table 3.

Dataset Code	Descriptor Code
B <sub>1</sub>	D <sub>1</sub> ,D <sub>2</sub> ,D <sub>3</sub> ,D <sub>4</sub> ,D <sub>5</sub>
B <sub>2</sub>	D <sub>6</sub> ,D <sub>7</sub>
B <sub>3</sub>	D <sub>8</sub>
B <sub>4</sub>	D <sub>9</sub>
B <sub>5</sub>	D <sub>9</sub>
B <sub>6</sub>	D <sub>9</sub>
B <sub>7</sub>	D <sub>9</sub>
B <sub>8</sub>	D <sub>10</sub>
B <sub>9</sub>	D <sub>10</sub>
B <sub>10</sub>	D <sub>10</sub>
B <sub>11</sub>	D <sub>10</sub>

Table 3: Datasets and the respective descriptors used in the experiments.

Some classifiers assume the Euclidean feature space implicitly, as the ANN-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

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

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

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

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

Table 4: Accuracy results  $x \pm y(z)$  on  $Z_3$  without using  $Z_2$ :  $x$  - mean accuracy,  $y$  - its standard deviation and  $z$  - mean kappa. The best accuracies are indicated in bold and the best value of  $k$  is shown in brackets for the  $k$ -NN.

This can be observed in Tables 4 and 5 for dataset  $B_1$  (MPEG-7) with  $D_4$  (BAS with L2) and  $D_5$  (BAS with OCS). Therefore,  $K'$  was 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 3. 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 and standard deviation of the accuracy values and mean value of kappa [51]. Section 4.1 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 4.2. The average computational time of each classifier for training and classification is divided by the number of samples and reported in Section 4.3.

#### 4.1 Accuracy results on $Z_3$ without using $Z_2$

The results in Table 4 are presented as  $x \pm y(z)[k]$ , where  $x$ ,  $y$ ,  $z$  and  $k$  are the mean accuracy, its standard deviation, mean kappa coefficient [51] 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. 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. However, OPF and SVM presented equivalent overall performances, being one better than the other depending on the case. The differences in accuracy along the lines of  $B_1 (D_4)$  and  $B_1 (D_5)$  indicate the importance of

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

Table 5: 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. The best accuracies are indicated in bold and the best value of  $k$  is shown in brackets for the  $k$ -NN.

the specific distance, except for the ANN-MLP which does not use distance function. The results along the five first lines also indicate the importance of the descriptor for dataset  $B_1$ . It is clear that  $D_5$  is the best descriptor for  $B_1$ .

## 4.2 Accuracy results on $Z_3$ with learning on $Z_2$

In order to evaluate the ability of each classifier in learning from the errors in  $Z_2$  without increasing the size of  $Z_1$ , we executed Algorithm 2 for  $T = 3$  iterations. The results are presented in Table 5.

Figure 9 shows the plot of the mean accuracy values of Table 5 for each classifier. Observe that in most cases the general learning algorithm improved the performance of the classifiers with respect to the results in Table 4.

## 4.3 Efficiency results

Table 6 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 2) 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 [37], 443 times faster than the fast ANN-MLP [38], and 1.3 times faster than our implementation of a  $k$ -NN classifier.

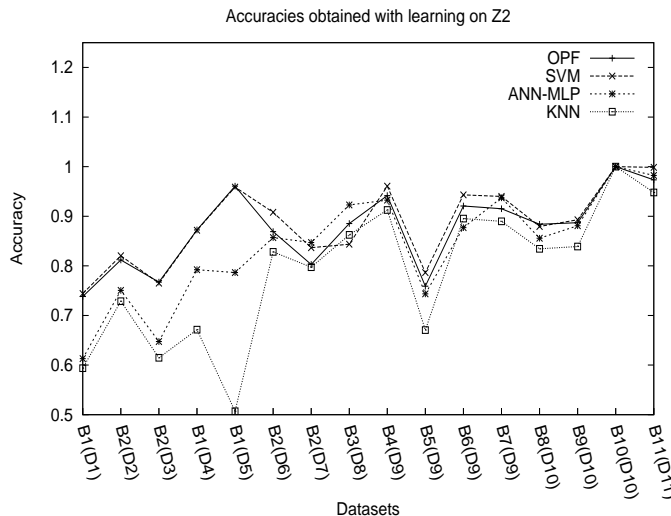


Figure 9: Accuracies of OPF, SVM, ANN-MLP and  $k$ -NN on  $Z_3$  with learning on  $Z_2$ .

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 second per sample and SVM spent 1.304. For 100,000 samples, this represents 8.67 minutes using OPF and 36.22 hours using SVM. In the case of 2D/3D images, for example, the number of pixels 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, due to their variations in noise, inhomogeneity, and protocols.

## Acknowledgments

The authors thank Paulo Miranda (IC-UNICAMP) and Nelson Mascarenhas (UFSCar) for their previous contributions and the financial support from CNPq and FAPESP.

## 5 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 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 theoretical advantages. It is fast, simple, multi-class, parameter independent, does not make any assumption about the shape of the classes, and can handle some degree of separability between classes.

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



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

Table 6: Mean execution times in seconds for training and classification divided by the number of samples. The best times are indicated in bold.

other algorithms to estimate prototypes and to learn from the errors in the evaluation set. We are currently investigating these aspects and the use of Genetic Programming [52] (GP) for arc-weight estimation in OPF. We can combine the distances from shape, color and texture descriptors by GP, for example, and use the result as arc weight [53, 54].

## References

- [1] R.O. Duda, P.E. Hart, and D.G. Stork. *Pattern Classification*. Wiley-Interscience, 2 edition, 2000.
- [2] A. Blum and T. Mitchell. Combining labeled and unlabeled data with co-training. In *Proceedings of the 9th Annual Conference on Computational Learning Theory*, pages 92–100, New York, NY, USA, 1998. ACM Press.
- [3] T. Joachims. Transductive inference for text classification using support vector machines. In *Proceedings of the 16th International Conference on Machine Learning*, pages 200–209, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc.
- [4] X. Zhu. Semi-supervised learning literature survey. Technical Report 1530, Computer Sciences, University of Wisconsin-Madison, 2006.
- [5] S. Haykin. *Neural networks: a comprehensive foundation*. Prentice Hall, 1994.
- [6] L. I. Kuncheva. *Combining Pattern Classifiers: Methods and Algorithms*. Wiley-Interscience, 2004.
- [7] L. Reyzin and R. E. Schapire. How boosting the margin can also boost classifier complexity. In *Proc. of the 23th Intl. Conf. on Machine learning*, pages 753–760, New York, NY, USA, 2006. ACM Press.

- [8] B.E. Boser, I.M. Guyon, and V.N. Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the 5th Workshop on Computational Learning Theory*, pages 144–152, New York, NY, USA, 1992. ACM Press.
- [9] K. Duan and S. S. Keerthi. Which is the best multiclass svm method? an empirical study. *Multiple Classifier Systems*, pages 278–285, 2005.
- [10] B. Tang and D. Mazzone. Multiclass reduced-set support vector machines. In *Proceedings of the 23th International Conference on Machine learning*, pages 921–928, New York, NY, USA, 2006. ACM Press.
- [11] N. Panda, E. Y. Chang, and G. Wu. Concept boundary detection for speeding up svms. In *Proc. of the 23th Intl. Conference on Machine learning*, pages 681–688, New York, NY, USA, 2006. ACM Press.
- [12] R. Collobert and S. Bengio. Links between perceptrons, MLPs and SVMs. In *Proceedings of the 21th International Conference on Machine learning*, page 23, New York, NY, USA, 2004. ACM Press.
- [13] A.X. Falcão, J. Stolfi, and R.A. Lotufo. The image foresting transform: Theory, algorithms, and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(1):19–29, Jan 2004.
- [14] T. Cover and P. Hart. Nearest neighbor pattern classification. *IEEE Transactions on Information Theory*, 13(1):21–27, 1967.
- [15] K. Fukunaga and P. M. Narendra. A branch and bound algorithms for computing k-nearest neighbors. *IEEE Transactions on Computers*, 24(7):750–753, 1975.
- [16] R.A. Lotufo and A.X. Falcão. The ordered queue and the optimality of the watershed approaches. In *Mathematical Morphology and its Applications to Image and Signal Processing (ISMM'00)*, volume 18, pages 341–350. Kluwer, Jun 2000.
- [17] R. Audigier and R.A. Lotufo. Watershed by image foresting transform, tie-zone, and theoretical relationship with other watershed definitions. In *Mathematical Morphology and its Applications to Signal and Image Processing (ISMM)*, pages 277–288, Rio de Janeiro, RJ, Oct 2007. MCT/INPE.
- [18] P. K. Saha and J. K. Udupa. Relative fuzzy connectedness among multiple objects: Theory, algorithms, and applications in image segmentation. *Computer Vision and Image Understanding*, 82(1):42–56, 2001.
- [19] R. Audigier and R.A. Lotufo. Seed-relative segmentation robustness of watershed and fuzzy connectedness approaches. In *XX Brazilian Symposium on Computer Graphics and Image Processing (SIBGRAPI)*, pages 61–68, Belo Horizonte, MG, Oct 2007. IEEE CPS.
- [20] G.T. Herman and B.M. Carvalho. Multiseeded segmentation using fuzzy connectedness. *IEEE Trans on Pattern Analysis and Machine Intelligence*, 23:460–474, May 2001.
- [21] A. Rocha P.A.V. Miranda, A.X. Falcão and F.P.G. Bergo. Object delineation by  $\kappa$ -connected components. *EURASIP Journal on Advances in Signal Processing*, 2008. to appear.
- [22] L. J. Hubert. Some applications of graph theory to clustering. *Psychometrika*, 39(3):283–309, 1974.
- [23] C.T. Zahn. Graph-theoretical methods for detecting and describing gestalt clusters. *IEEE Transactions on Computers*, C-20(1):68–86, Jan. 1971.

- [24] A. K. Jain and R. C. Dubes. *Algorithms for Clustering Data*. Prentice-Hall, Inc., 1988.
- [25] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905, Aug 2000.
- [26] L. M. Rocha, A. X. Falcão, and L. G. P. Meloni. A robust extension of the mean shift algorithm using optimum path forest. In *8th Intl. Workshop on Combinatorial Image Analysis*, pages 29–38, Buffalo-NY, USA, 2008. RPS (ISBN 978-981-08-0228-8).
- [27] 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. In *The Fifth IEEE International Symposium on Biomedical Imaging: From Nano to Macro (ISBI)*, pages 428–431, 2008.
- [28] J.P. Papa, A.X. Falcão, C.T.N. Suzuki, and N.D.A. Mascarenhas. A discrete approach for supervised pattern recognition. In *12th International Workshop on Combinatorial Image Analysis*, volume 4958, pages 136–147. LNCS Springer Berlin/Heidelberg, 2008.
- [29] 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. In *Mathematical Morphology and its Applications to Image and Signal Processing (ISMM'07)*, pages 337–348. MCT/INPE, 2007.
- [30] 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 Journal on Advances in Signal Processing*, 2008:1–16, 2008.
- [31] Y. P. Wang and T. Pavlidis. Optimal correspondence of string subsequences. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 12(11):1080–1087, 1990.
- [32] Y. Rubner, C. Tomasi, and L. J. Guibas. A metric for distributions with applications to image databases. In *Proc. of the 6th Intl. Conference on Computer Vision*, page 59. IEEE Computer Society, 1998.
- [33] N. Arica and F. T. Y. Vural. BAS: A Perceptual Shape Descriptor Based on the Beam Angle Statistics. *Pattern Recognition Letters*, 24(9-10):1627–1639, June 2003.
- [34] R. O. Stehling, M. A. Nascimento, and A. X. Falcão. A compact and efficient image retrieval approach based on border/interior pixel classification. In *Proc. of the 11th Intl. Conf. on Inf. and Knowledge Management*, pages 102–109. ACM Press, 2002.
- [35] T. Cormen, C. Leiserson, and R. Rivest. *Introduction to Algorithms*. MIT, 1990.
- [36] C. Allène, J. Y. Audibert, M. Couprie, J. Cousty, and R. Keriven. Some links between min-cuts, optimal spanning forests and watersheds. In *Mathematical Morphology and its Applications to Image and Signal Processing*, pages 253–264. MCT/INPE, 2007.
- [37] 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>.
- [38] S. Nissen. *Implementation of a Fast Artificial Neural Network Library (FANN)*, 2003. Department of Computer Science University of Copenhagen (DIKU). Software available at <http://leenissen.dk/fann/>.
- [39] S. C. Huang and Y. F. Huang. Bounds on the number of hidden neurons in multilayer perceptrons. *IEEE Transactions on Neural Networks*, 2(1):47–55, 1991.

- [40] R. Kohavi. A study of cross-validation and bootstrap for accuracy estimation and model selection. In *IJCAI*, pages 1137–1145, 1995.
- [41] MPEG-7. The generic multimedia content description standard, part 1. *IEEE MultiMedia*, 09(2):78–87, 2002.
- [42] Corel Corporation. Corel stock photo images, -. <http://www.corel.com>.
- [43] P. Brodatz. *Textures: A Photographic Album for Artists and Designers*. Dover, New York, 1966.
- [44] D.J. Newman A. Asuncion. UCI machine learning repository, 2007.
- [45] 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 Transactions on Medical Imaging*, 17(3):463–468, 1998.
- [46] L. Kuncheva. Artificial data. *School of Informatics, University of Wales, Bangor*, 1996. <http://www.informatics.bangor.ac.uk/~kuncheva>.
- [47] E. Persoon and K. Fu. Shape Discrimination Using Fourier Descriptors. *IEEE Transactions on Systems, Man, and Cybernetics*, 7(3):170–178, 1977.
- [48] M.K. Hu. Visual Pattern Recognition by Moment Invariants. *IRE Transactions on Information Theory*, 8(2):179–187, 1962.
- [49] R. Torres, A. X. Falcão, and L.F. Costa. A graph-based approach for multiscale shape analysis. *Pattern Recognition*, 37(6):1163–1174, 2004.
- [50] M. Swain and D. Ballard. Color Indexing. *International Journal of Computer Vision*, 7(1):11–32, 1991.
- [51] J. Cohen. A coefficient of agreement for nominal scales. In *Educational and Psychological Measurement*, volume 20, pages 37–46, 1960.
- [52] J. R. Koza. *Genetic Programming: On the Programming of Computers by Means of Natural Selection (Complex Adaptive Systems)*. The MIT Press, 1992.
- [53] R.S. Torres, A.X. Falcão, M.A. Gonçalves, B. Zhang, W. Fan, E. A. Fox, and P. Calado. A new framework to combine descriptors for content-based image retrieval. In *Proc. of 14th ACM Conf. on Inf. and Knowledge Management*, pages 335–336, Bremen, Germany, Nov 2005.
- [54] R.S. Torres, A.X. Falcão, M.A. Goncalves, J.P. Papa, B. Zhang, W. Fan, and E.A. Fox. A genetic programming framework for content-based image retrieval. *Pattern Recognition*, 2008. (accepted).