



**HAL**  
open science

## An automatic system for the classification of cellular categories in cytological images

Marinette Revenu, Abderrahim Elmoataz, Christine Porquet, Hubert Cardot

► **To cite this version:**

Marinette Revenu, Abderrahim Elmoataz, Christine Porquet, Hubert Cardot. An automatic system for the classification of cellular categories in cytological images. SPIE: Intelligent Robot and Computer Vision XII: Algorithms and Techniques, 1993, Boston, United States. pp.967-970, 10.1117/12.150150 . hal-00975330

**HAL Id: hal-00975330**

**<https://hal.science/hal-00975330>**

Submitted on 8 Apr 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# An automatic system for the classification of cellular categories in cytological images

Marinette Revenu, Abderrahim Elmoataz, Christine Porquet, Hubert Cardot

LAIAC (Laboratoire d'Algorithmique et d'Intelligence Artificielle de Caen)  
ISMRA (Engineering School of Caen)  
6, boulevard du Maréchal Juin  
F-14050 CAEN CEDEX

## ABSTRACT

In this paper, we describe research carried out within the framework of the optimization of an image analyzer dedicated to rapid detection of abnormalities of ploidy in human tumors. The system takes as its input microscopic images of dissociated cells which are to be segmented in order to extract cellular objects, calculate shape and texture measures and identify each category of cell, by means of two classification methods that are compared and discussed: classification based on the Bayes decision rule and classification using neural networks.

## 1. INTRODUCTION

We are concerned with the segmentation of cytological images in order to automatically classify cellular categories and to calculate the DNA rate of epithelial cells only. We are dealing with microscopic images of human esophagus cells, that have been dissociated and tainted according to Feulgen and Rossenbeck methods. These images are provided by the pathology department of the cancer-research centre F. Baclesse of Caen. Calculating the DNA rate gives an idea of the proliferation of cancerous cells.

At present, this cancer-research centre uses two approaches to calculate the DNA rate: flow cytometry and image analysis<sup>1</sup>. In both cases, the biological material is a preparation obtained from thin sections of human esophagus that have been deparaffined and dissociated according to Hedley's method. The cells that can be observed on a coverglass are of several types: epithelial cells, lymphocytes, stromal cells. Fragments and clusters of cells can also be examined. The drawback of the flow cytometry approach stems to the fact that measures are done on all cells, without distinction. On the contrary, the computer system does the same measures by segmenting images and eliminating through sorting lymphocytes, stromal cells as well as debris and clusters. The interest of the latter approach is in getting a measure of the DNA rate more closely related to abnormalities of ploidy in tumors; however, its drawback is linked to an important manual intervention, as the system presently in use necessitates a manual sorting step.

The study we are presenting here comes within the framework of the optimization of the image analyzer and deals with the selection of segmentation, characterization and classification operators.

By using knowledge about cytological images and the specific features of cells, we have devised an analyzer composed of three modules dealing respectively with the segmentation of cellular objects, their characterization and their automatic classification including two techniques: classification based on the Bayes decision rule and classification using neural networks.

In this paper, before describing each of the three modules of our analyzer, we first give an overview of image segmentation issues, independently from any application, in order to show how to select the operators and be directed by a priori knowledge.

## 2. IMAGE SEGMENTATION ISSUES: SELECTION OF OPERATORS

Image segmentation is one of the most active research issues in image processing. Despite much effort, there is no universal segmentation methodology. Each segmentation approach is only efficient for a precise type of image associated to a

specific application. However, by studying the various works described in the literature, we could notice that any segmentation problem can be solved according to three independent steps:

- First step: Pre-segmentation. It consists in obtaining an initial segmentation by extracting region and/or primitives adapted to a given problem. Various segmentation operators based on the boundary or the region approach can be used during this purely numerical step 2,3,4.

- Second step: Characterization. This step is concerned with the evaluation of intrinsic or relational features and measures on the region and/or boundary primitives extracted during the first step. These various attributes will be used in the third step to pass from one segmentation to another through the merging and splitting of primitives. This step must also involve evaluation or control measures of the initial segmentation.

- Third step: Search of the best solution by correcting the results of the initial segmentation and by using operators for merging or splitting region and/or boundary primitives. During this step, which is the most complex of the three, one can consider various merging modes, various scanning modes of the image, various strategies to solve the problem and determine a sequence of operators.

The selection of segmentation operators, the determination of their parameters and the organization of operators into sequences, in order to solve a specific problem are all based on the use of four types of knowledge:

1 - physical knowledge, linked to the formation of images. This type of knowledge is related to the characterization of a class of images (resolution, signal-to-noise ratio, distinctive details...). By taking it into account, one can direct the selection of operators to get an initial segmentation. For instance, if the signal-to-noise ratio is high for a given class of images, edge-detection operators cannot be used. Knowledge about the nature of the noise also enables to direct the selection of noise-reducing operators.

2 - perceptive knowledge, leading to the notion of visual primitives that are a means to understand segmentation mechanisms. This knowledge is related to the nature of primitives (lines, regions, pixels...), their extraction mode, segmentation criteria and rules resulting mainly from perceptive grouping based on the Gestalt theory, according to similitude, proximity, closing, and continuity notions.

3 - semantic knowledge, directly linked to the application: image origin, objective, representation and organization of objects in a class of images.

4 knowledge about segmentation operators, which can be seen as an expertise in image processing. By thoroughly studying these operators, we could bring to the fore their underlying models, parameters, effects on image details, complexity, usage, possible extensions...

In the next section, we are going to illustrate summarily the use of these four types of knowledge in the case of cytological images.

### 3. APPLICATION TO CYTOLOGICAL IMAGES

The images we are processing are 256x256 images coded on 8 bits. They are characterized by cellular objects scattered on a homogeneous background, corresponding either to isolated cells or to clusters of cells, a cluster being composed of several cells overlapping each other. It must also be noticed that, due to the biological preparation mode, debris of cells can also be examined. Besides clusters and debris, the following types of cells must be identified: stromal cells, epithelial cells and lymphocytes. Figure 1 shows an example of cytological image.

Perceptive and semantic knowledge help to describe the segmentation objective (getting precise measures characterizing each type of cell, in order to classify them) and the class of images (images composed of cellular objects on a background). Cellular objects can either be clusters of cells or isolated cells, the latter being distinguished thanks to morphometric features. Each of these objects can be characterized as follows:

- The background is characterized by a homogeneous texture and lighter gray levels than the rest of the image.
- A cluster of cells is has a larger size than the average size of objects and a non-convex form.
- A normal human esophageal epithelial cell shows little texture, a low inner contrast, an oval convex shape and a medium size.
- A lymphocyte is circular, non-textured, darker than epithelial cells and has a convex shape.
- A stromal cell is non-textured, shows an elongated oval shape and a medium size.
- Some debris, cut during the preparation phase, are dissymmetrical in relation to their centre of gravity.



*Fig. 1: an example of cytological image.*

This a priori knowledge directs our plan of treatments, composed of three steps: extraction of cellular objects (and separation of partially touching objects), characterization by shape and texture attributes, and finally classification.

#### 4. THE SEGMENTATION MODULE

The segmentation module deals with the extraction of regions corresponding to cells and the separation of partially touching cells, in order to eliminate noisy regions and cells touching the border of the image. To solve this specific segmentation task, several segmentation plans can be considered: classical plans based either on the boundary or on the region approach, or plans taking into account a priori knowledge related to our application. Let us first present classical segmentation approaches.

##### 4.1 Classical segmentation approaches

As our objective is to get precisely-located boundaries in order to obtain good measures on the corresponding regions, a segmentation based on edge detection seems quite suitable. On the other hand, as we have to detect cells, a region-based segmentation has the advantage of giving closed regions that can be directly associated to cells.

#### 4.1.1. Segmentation based on edge detection

This approach can be summarized as follows:

- application of local operators such as Gradient or Laplacian operators,
- localization of edges and elimination of false edge points,
- closing of boundaries in order to obtain regions.

Although this technique gives precisely-located boundaries, it generates a lot of regions, mainly in the areas of clusters. Moreover, it is unstable, as it depends on too many different parameters: threshold for the elimination of false edge points, threshold for closing boundaries, etc...

#### 4.1.2. Classical segmentation based on region growing

The problem is here to find all the connected components of the image, that verify an homogeneity criterion. The segmentation technique depends on the choice of the criterion and the seeds to start the growing of regions. It gives poor results in border areas, and it must be followed by a merging process, thus increasing the complexity of the approach. It also generates a lot of regions in the areas of clusters.

Both methods are not completely satisfactory and they do not take sufficiently into account the a priori knowledge about our application. That is the reason why we had to consider a specific region/boundary segmentation approach.

#### 4.2. Our region/boundary segmentation approach

As the background of images has a homogeneous texture, one can consider a separation of objects based on the gray level information. The idea is to find the optimal threshold separating objects from the background and to apply an optimal binary thresholding operator. During this step, we use knowledge about texture and gray level intensity.

Then, in order to separate partially touching regions corresponding to clusters, we use knowledge about the shape of cellular objects: the convexity degree of each object is calculated and a region presenting a low convexity degree is re-segmented.

These two steps are now described in detail.

##### 4.2.1. Extraction of regions by optimal binary thresholding

The principle of our binary thresholding is similar to Kolher's method<sup>5</sup>: one has to determine the threshold which detects the maximum of highly contrasted edge points and the minimum of lowly contrasted ones. But we use an original approach to automatically calculate that threshold as follows:

- Image smoothing by using filters that preserve boundary information.
- Calculation of the Gradient image and estimation of the average amplitude of the Gradient  $AM$ .
- Automatic determination of the threshold giving a binary image, the edges of which best match with those of the Gradient image. For a given threshold  $s$ , we evaluate  $NBP(s)$ , the number of edge points, the amplitude of which is higher than  $AM$ . The optimal threshold is the one that maximizes  $NBP(s)$ .

##### 4.2.2. Region labeling and re-segmentation of some regions

The labeling of the image consists in searching all the connected components of the segmented image. At the end of this step, each region is associated to a single label.

Then, several attributes are calculated:

- the bounding rectangle, in order to limit further operations to that area only,
- the surface,

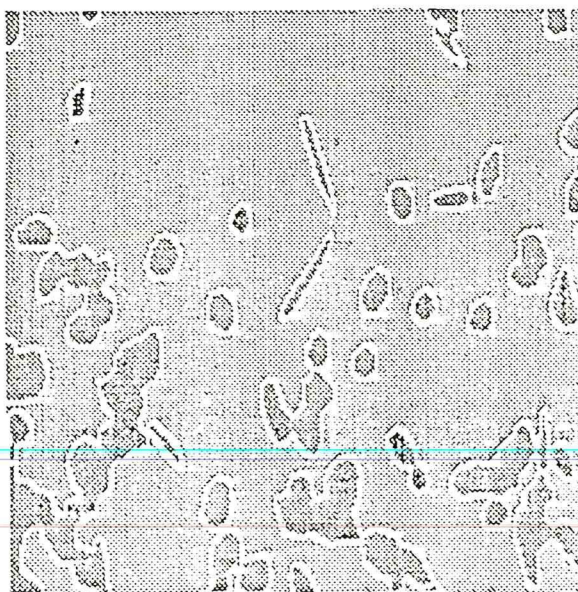
- the convexity degree, which is defined as the surface of the convex envelope of region divided by the surface of the region.

Thanks to this attribute, it is possible to separate regions into two classes: if the convexity degree of the region is lower than a threshold (defined experimentally), then the region is classified as a highly convex object (which corresponds to epithelial, stromal cells and lymphocytes). Otherwise, the region is considered as being either of cluster of cells or a cell which was not correctly segmented. In that case, the region is re-segmented as follows:

- Calculation of the distance of each pixel to the boundary of its region.
- Determination of the maxima of the distance function, on each region and selection of the minimum,  $t_0$ .
- Erosion of size  $t_0$ , followed by a conditional dilatation.

The advantage of this technique is that, only clusters formed by partially touching cells are re-segmented.

After this step, we can eliminate noisy regions, which are characterized by a small size and a low edge-contrast thanks to size and contrast criteria. Results are show on figure 2.



*Fig. 2 : Image resulting of our region/boundary segmentation.*

## 5. THE CHARACTERIZATION MODULE

The next step is concerned with the computation of various shape and texture measures enabling the identification of different types of cells.

In our application, the identification of cells is based on three classes of parameters that must be quantified and according to which characterization operators must be selected: intensity parameters (light, dark), shape parameters (circular, oval, elongated) and texture parameters (fine, coarse). We briefly give some of these parameters:

- **Intensity parameters:** They are calculated either from the gray level histogram or the contrast histogram of a region: average intensity, variance, dissymmetry, entropy, energy.

- **Shape parameters:** The most usual ones are bounding rectangle, surface, perimeter, compactness. Moments give also important information about the shape of the region and are used to calculate elongation and eccentricity parameters that characterize stromal cells. Third moments bring data about the symmetry of regions and are used to distinguish debris. Finally, convexity parameters are used to detect clusters of cells.

- **Texture parameters:** They give an idea of the homogeneity of the chromatin distribution of cellular nuclei, which is an important information because chromatin texture is related to cell activity. We calculate some texture parameters from the co-occurrence matrix, whereas others are determined from the edge image<sup>6</sup>.

Let us now give explanations about the way classification is performed.

## 6. THE CLASSIFICATION MODULE

To identify epithelial cells in cytological images in order to quantify the DNA rate and to establish the ploidy curve, a computer system<sup>7</sup> is presently in use at the pathology department of the cancer-research centre F. Baclesse of Caen. Its purpose is to automate the acquisition and the analysis of images. After a region segmentation step, the pathologist has to select each image zone corresponding to a cell and to give the name of its class, which is stored on disk together with the region feature-vector.

Thanks to this real-size database, we are studying how to make the system completely automatic without any manual intervention from the pathologist. Two approaches have been considered for the elaboration of an automatic classification module: an approach based on the Bayes decision rule and an approach based on artificial neural networks (ANN). This study follows research we did for the development of a handwritten signature verification system<sup>8</sup> and takes up many of its principles. The points in common to both works are the extraction, from an image, of features which are compared to reference data in order to take a decision. A statistical approach and several approaches based on ANN have been simultaneously tested and we could notice that the addition of an ANN module to take the final decision could improve results noticeably.

Any classification system<sup>9, 10</sup> works according to two distinct functioning modes: a training phase in order to find out the elements that will be used during the comparison between reference data and samples, and a classification phase which takes a decision for each presented sample.

We are now going to describe the database, the evaluation criteria used in both classification modules, before presenting each module in detail, the Bayesian module and the ANN one. An analysis of the results of each module should give us some ideas on the specification of a fusion module making the best of each approach.

### 6.1. The training and test databases

To carry out this study, we had several databases containing about 4000 cells descriptions each. Each description consists in the cell class associated to a vector of 40 features related to shape, intensity and texture information. Only epithelial cells must be considered for establishing a diagnostic, but, in order to build the most informative database, pathologists have defined 9 classes. Thanks to this choice, we can take into account the different stages of evolution of epithelial cells within their reproduction cycle (2 classes) and better differentiate elements which do not belong to the epithelial cell classes. In particular, debris and clusters of cells are associated to 3 classes. The name of the 9 classes in the case of normal reference esophagus are: intermediary epithelial cells (1), basal or superficial epithelial cells (2), non-epithelial unidentified cells (3), lymphocytes (4), polynuclears (5), stromal cells (6), clusters (7), cut debris (8) and other debris (9).

These databases represent the expertise acquired by biologists during years, in the analysis of many images. Like all human expertise corresponding to a tedious task, the reproducibility of the classification is not perfect and experts can slightly disagree on some classification decisions. To take this fact into account, we first worked on each file separately. The proportion of the different cellular classes is not the same in all databases and images. As we did not want to favor any class, which would have altered further results, we decided to take P% of cells in each class for the training database, the others cells being used for tests (in practice, about 66%). Another experiment was carried out with different databases for training and test phases. Here, we give our results in the latter case, because they are more significant and can be generalized.

The evaluation of the quality of the classifier is based on three criteria: the global recognition rate, representing the number of correctly identified cells, which should be as high as possible, and the "positive false rate" and the "negative false rate" which should be as low as possible.

- the "positive falses" are, for a given group, the cells wrongly classified by the system, which belong to this group and are not classified in this group.
- the "negative falses" are, for a given group, the wrongly classified cells, which do not belong to this group and are all the same classified in this group.

The confusion matrix contains all the necessary informations to calculate these different rates. It shows precisely in which group a cell from a given group has been classified.

Another analysis of the results can be considered, by grouping together results of several classes. Actually, all the classes that have been chosen to constitute the database are not independent as regards the final result consisting in getting the ploidy curve. In particular, neither the distinction between various types of epithelial cells (classes 1 and 2), nor the distinction between various debris (classes 7, 8 and 9) is significant. What really matters is the elimination of cells that can be classified without doubt as non-epithelial cells. Thus, in the presentation of our results, results for some classes are grouped together and we get a more global quality rate than in previous experiments.

## 6.2. Classification based on the Bayes decision rule

Let  $N_c$  be the number of classes and  $X=(x_1, x_2, \dots, x_{N_p})$  be a feature-vector describing one cell, where  $N_p$  is the number of features. During the training phase, each class  $C_i$  can be represented by a vector  $M_i$  and a covariance matrix  $CoV_i$  whose size is  $N_p \times N_p$ .

Two hypotheses are commonly used:

- in each class, features are distributed according to a normal law,
- features are independent.

In this case, the a priori probability  $P_i(X)$  of a given unknown vector  $X$  is defined as:

$$P_i(X) = \frac{\exp(-D_i(X) / 2)}{2\pi^{N_p/2} \text{Det}(CoV_i)^{1/2}}$$

where  $D_i(X)$  is the Mahalanobis distance between  $X$  and  $M_i$ .

$$D_i(X) = (X - M_i)^T CoV_i^{-1} (X - M_i)$$

If we use the following notation:

$L_i(X) = -\text{Log}(P_i(X))$ , we get :

$$L_i(X) = \frac{1}{2} D_i(X) + \frac{1}{2} N_p \text{Log}(2\pi) + \frac{1}{2} \text{Log}(\text{Det}(CoV_i))$$

or, with our hypotheses:

$$L_i(X) = \frac{1}{2} \sum_{j=1}^{N_p} \frac{(x_j - m_{ij})^2}{\sigma_{ij}} + \frac{1}{2} N_p \text{Log}(2\pi) + \frac{1}{2} \text{Log} \left( \prod_{j=1}^{N_p} \sigma_{ij} \right)$$

where  $i$  denotes a class,  $j$  a feature,  $m_{ij}$  the mean of feature  $j$  in class  $i$  and  $\sigma_{ij}$  the variance of feature  $j$  in class  $i$ .

If  $L_i(X)$  is the minimum of  $L_k(X)$  for  $k=1$  to  $N_c$ , then the unknown vector  $X$  is classified in class  $C_i$ .



## Results:

We present the results we obtained with our system on a database where the number of cells in the training base is 2058 and the number of cells in the test base is 2058. Table 1 shows the distribution of cells in the different classes.

class number	number of cells in the training base	number of cells in the test base
1	210	304
2	716	659
3	195	200
4	330	308
5	11	9
6	5	10
7	35	45
8	176	258
9	380	265

Table 1. The distribution of cells in the databases

The confusion matrix obtained by classification using the Bayes decision rule is shown in Table 2. The global recognition rate is 67%.

class class	1	2	3	4	5	6	7	8	9
1	275	29	0	0	0	0	0	0	0
2	66	456	122	2	5	0	0	8	0
3	1	18	79	47	47	1	0	7	0
4	0	0	0	305	3	0	0	0	0
5	0	0	0	4	5	0	0	0	0
6	0	0	1	0	0	7	0	2	0
7	1	2	9	2	8	0	17	5	1
8	29	53	53	7	9	8	2	88	9
9	1	5	33	10	10	3	1	51	151

Table 2. The confusion matrix in the Bayesian module

### Recognition rates for classes 1 and 2:

As it was already mentioned, results obtained in class 1 and 2 are the most significant. Here, we give rates calculated on each class individually and on class 1 and 2 grouped together.

#### Class 1:

- 275 cells out of 304 have been correctly classified, yielding a recognition rate of 90.46%. The "positive false rate" is 9.54%.

- Out of 1754 cells that do not belong to class 1, 98 were classified in this class, yielding a "negative false rate" of 5.58%.

#### Class 2:

- 456 cells out of 659 have been correctly classified, yielding a recognition rate of 69.20%. The "positive false rate" is 30.80%.

- Out of 1399 cells that do not belong to class 2, 107 were classified in this class, yielding a "negative false rate" of 7.64%.

#### Class 1 + class 2:

- 826 cells out of 963 have been correctly classified, yielding a recognition rate of 85.77%. The "positive false rate" is 14.23%.
- Out of 1095 cells that neither belong to class 1 or 2, 110 were classified either in class 1 or 2, yielding a "negative false rate" of 10.4%.

This approach has the advantage of being simple to implement. On the other hand, classes containing elements, the parameters of which are not all stable, for instance, class 3 of non-epithelial unidentified cells and classes 7, 8, 9 are not well identified.

### 6.3. Classification using neural networks

For the pathologist, classifying a cell is no easy task. His decision cannot always be unambiguous, because cells can take quite different appearances in the course of their evolution. This is one of the reasons why methods based on neural networks seem promising. Their learning and generalizing capabilities ensure that they can cope with the diversity of cells and variations in manual classifications. Our classification module using neural networks is built on the multilayer perceptron model.

#### 6.3.1. Architecture of the multilayer network

The network (fig. 3) has an input layer of 38 units corresponding to the 38 features describing a region of the image, an output layer of 9 units corresponding to the 9 predefined classes and a variable number of hidden layers. The 38 measures that constitute the input of the network are those among the 40 features that are independent of the DNA rate. Cells of one layer are completely connected to the cells of the next layer and each link bears a weight  $W_{ij}$ . The training phase is done using the error Gradient backpropagation algorithm<sup>11</sup>.

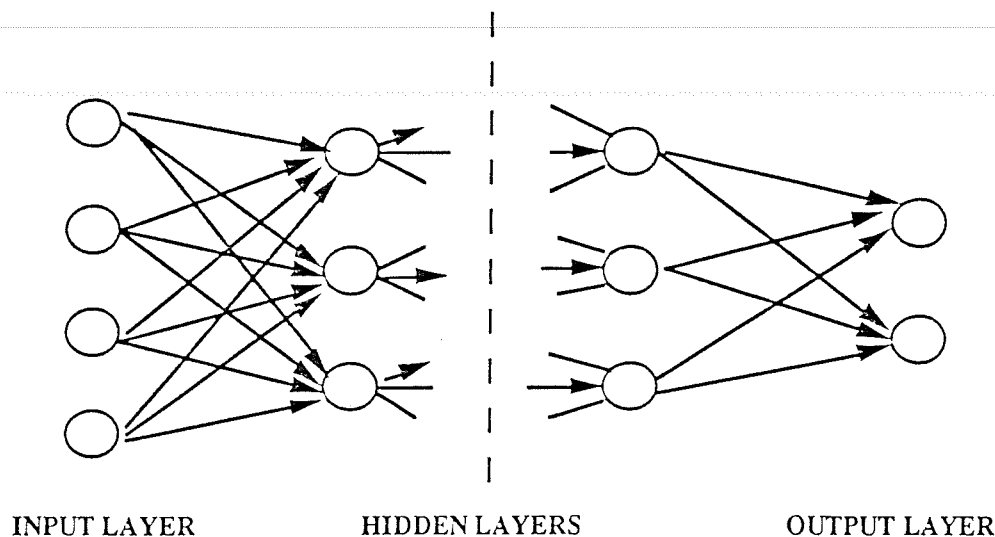


Fig 3 : The multilayer network working on the features of a region.

In multilayer networks, two distinct working phases must be considered: a training phase and a classification phase.

a) The training phase splits up into two steps:

- The propagation of the values of the input layer (vector  $X$ ) to the output layer (vector  $S$ ). The transfer function used in the calculation is a sigmoid.
- The backpropagation of the values of the output layer to the input layer. The goal of this step is to minimize the error between the answer of the network ( $O_i$ ) and the expected answer ( $Y_i$ ). The rules used for the modification of the values of the weights  $W_{ij}$  are based on the derivative of the transfer function.

The adjustment parameters of the network are  $a$ , the adaptation gain in the computation of the new values of  $W_{ij}$ , and  $b_1$  and  $b_2$ , two coefficients of the transfer function and its derivative.

b) The classification phase only consists in presenting as input to the network a feature-vector of a region of the image and in propagating values through the network. The output getting the greatest value gives the class of the element presented as input.

### 6.3.2 Working mode during the training phase

The development of a classifier based on multilayer networks involves several issues that are now discussed: data preparation and selection of a sequence of tests to adjust parameters according to quality criteria mentioned in § 6.1.

a) Data preparation

The input values propagated through the network must be normalized over the interval  $[0..1]$ . This normalization is done in 2 steps:

- a linear transformation based on the minimum and maximum values of each feature, calculated on all the regions of the training base.
- a quantizing of the range of resulting values in the interval  $[0..1]$ , achieved by histogram equalization. Each value  $x$  is replaced by  $T(x)$ , the value of the cumulative distributive function for  $x$ .

Thanks to this transformation, we obtain values for the input feature-vector than are more discriminant for further classification.

b) Choice of the network configuration and adjustment of parameters

The network configuration is determined by the number of hidden layers and the number of cells in each layer. It is impossible to try every configuration but it is generally assumed that 0, 1 or 2 hidden layers are sufficient. On the other hand, the values of parameters  $a$ ,  $b_1$  and  $b_2$ , as well as the number of presentations of samples to the network must be tested for each network configuration. As it is difficult to simultaneously deal with all these choices, because different quality criteria must be optimized, we operated in several steps:

First, the global recognition rate is taken into account to determine an approximate value for all parameters. It yields initial values of 0.07, 0.08 and 0.4 for  $a$ ,  $b_1$  and  $b_2$  respectively, and 300 000 for the number of presentations. We determined the latter value by starting with a number of presentations of 100 000, and by increasing this number progressively, until results deteriorate. Then these values were used to choose the number of hidden layers and cells, by a trial-and-error process.

Secondly, once the configuration and number of presentation are fixed, coefficients  $a$ ,  $b_1$  and  $b_2$  can be adjusted so as to minimize the "positive falses rate" and the "negative falses rate".

Experiments were also carried out to determine whether to take the same sequence of samples for training and classification, or to take different sequences. We present here results obtained in the latter case, because they are more representative of the performances of our neural module.

In order to stabilize the weights  $W_{ij}$  and to get a classification response in accordance to the training base, each sample must be presented to the network several hundreds of times. However, as the training process must not depend on the order in which data are stored in files, samples are chosen randomly. Likewise, weights  $W_{ij}$  are initialized with random values chosen in the interval  $[-1..+1]$ .

## Results:

We present the results we obtained with our system on the database shown in Table 1. The best results were obtained with a single hidden layer of 25 cells and values of 0.05, 0.8 and 0.7 for coefficients  $a$ ,  $b_1$  and  $b_2$  respectively.

The confusion matrix used in the neural moduleneural is shown in Table 3. The global recognition rate is 68,07%.

	1	2	3	4	5	6	7	8	9
1	231	66	0	0	0	0	0	7	0
2	89	460	11	1	0	0	1	96	1
3	63	14	35	43	0	0	2	38	5
4	12	0	0	294	0	0	2	0	0
5	3	0	0	4	2	0	0	0	0
6	2	0	2	0	0	5	0	1	0
7	5	0	0	1	0	0	26	5	8
8	39	7	0	5	0	0	0	119	88
9	16	0	0	10	1	0	0	9	229

Table 3. The confusion matrix in the neural module

### Recognition rates for classes 1 and 2:

#### Class 1:

- 231 cells out of 304 have been correctly classified, yielding a recognition rate of 75.9%. The "positive falses rate" is 24.1%.
- Out of 1754 cells that do not belong to class 1, 229 were classified in this class, yielding a "negative falses rate" of 13.05%.

#### Class 2:

- 460 cells out of 659 have been correctly classified, yielding a recognition rate of 69.8%. The "positive falses rate" is 30.2%.
- Out of 1399 cells that do not belong to class 2, 87 were classified in this class, yielding a "negative falses rate" of 6.21%.

#### Class 1 + class 2:

- 846 cells out of 963 have been correctly classified, yielding a recognition rate of 87.85%. The "positive falses rate" is 12.15%.
- Out of 1095 cells that neither belong to class 1 or 2, 161 were classified either in class 1 or 2, yielding a "negative falses rate" of 14.70%.

The results obtained following this neural approach are quite comparable with those obtained following the Bayesian approach; however the recognition rates are higher.

## 6.4. Comparison of the results and discussion

With the neural approach, the global recognition rate is higher, and it is also the case when grouping together results of classes 1 and 2. However, too many debris are classified in class 1 and 2, which is not the case in the Bayesian module.

We are now studying how to fusion results of both approaches, thanks to a multilayer network without hidden layer, an idea that was already tested in our signature authentication system <sup>12</sup>. 9 inputs of this network would be the 9 outputs of the Bayesian module, and the 9 other inputs would be the 9 outputs of the neural module. The 9 outputs of the network would

give the final decision. Thanks to this scheme, one can determine through automatic learning how each module can play a part in the final decision.

## 7. CONCLUSIONS

In order to develop an automatic system for the classification of esophageal epithelial cells, we described the different image processing operations necessary to calculate features that can be used for classification. Two approaches were tried and compared: a statistical approach based on the Bayes decision rule and a neural approach. Results are promising and we are considering introducing a network dealing with the fusion of the results from each module, in order to get the most out of each approach. We are also carrying out experiments on methods based on decision trees.

## 8. ACKNOWLEDGMENTS

The authors are grateful to Mrs Mandard and Herlin of the pathology department of the cancer-research centre F. Baclesse of Caen, who provided the cytological images and gave us all the explanations about the various types of cells we had to classify. Special thanks also to Mr Masson of the LEI (Electronics and Instrumentation Laboratory).

## 9. REFERENCES

1. E. Masson, P. Herlin, D. Bloyet, F. Duigou, A.M Manadard, "Image analysis optimisation of the routine retrospective DNA ploidy measurement of solid tumors. Part I : Methodological choices", *Acta Stereologica*, 1991.
2. A. Elmoataz, C. Porquet, M. Revenu, " A general-purpose segmentation system using knowledge on images and segmentation operators", *Advances in Intelligent Robotic System SPIE*, vol. 1607, Boston, 1991.
3. M. Coster, J.L Cherman, "Précis d'analyse d'images", *Editions CNRS*, Paris, France, 1985.
4. B. Manfred, I. Werner, " VIPER: a general-purpose digital image processing applied to video microscopy", *Computer Methods and Program in Biomedicine*, vol 26, 1988.
5. R. Kohler, "A segmentation system based on thresholding", *C.G.I.P.*, Vol 15, pp 319-338, 1981.
6. H. Harms, U. Gunzer, H. M. Aus, "Combined local color and texture analysis of stained cell", *C.V.G.I.P.*, Vol 33, pp 364-376, 1986.
7. E. Masson, "Cytométrie automatisée par traitement numérique d'images: application à la caractérisation et à la classification des cellules épithéliales normales et tumorales", Thèse de doctorat, Université de Caen, 1992.
8. H. Cardot, M. Revenu, B. Victorri, M-J. Revillet, "Coopération de réseaux neuronaux pour l'authentification de signatures manuscrites", *Int. Conf. Neuro-Nîmes*, France, 1991.
9. J. Piper, D. Granumd, " On fully automatic feature measurement of banded chromosomes classification", *Cytometry*, 1989.
10. X. Qui, D. Barba, G. Ramstein, "La classification automatique par analyse d'images", *Congrès AFCET-RFIA*, Lyon, France, 1991.
11. D. Rumelhart, G. Hinton, R. Williams, "Learning representations by backpropagating errors", *Nature*, Vol. 323, No. 9, pp. 533-536, 1986.
12. H. Cardot, M. Revenu, B. Victorri, M-J. Revillet, "An Artificial Neural Networks Architecture for Handwritten Signatures Verification", *SPIE Intelligent Information Systems*, Orlando, USA, Avril 1993.