



UNIVERSIDAD NACIONAL DE COLOMBIA

Classification of squamous cell cervical cytology

Luz Helena Camargo Casallas

Universidad Nacional de Colombia
Faculty of Medicine - Engineering Faculty
Bogotá D.C. Colombia
2012

Classification of squamous cell cervical cytology

Luz Helena Camargo Casallas

Thesis submitted in partial fulfillment of the requirements for the degree of:
Magister en Ingeniería Biomédica

Advisor:
Ph.D. Eduardo Romero Castro

Research Line:
Interpretation of Medical Images
Grupo de Investigación:
BioIngenium

Universidad Nacional de Colombia
Departamento de Imágenes Diagnósticas
Bogotá D.C., Colombia
2012

A mi hijo Jhon Sebastián.

*A mis padres Hector y Helena por su apoyo
y dedicación.*

*A mis hermanos Esperanza, Constanza y
Yesid, por alentarme a no desfallecer en mis
sueños.*

Acknowledgments

Agradezco al profesor Eduardo Romero Director del Grupo de Investigación BioIngenium por su alto nivel de compromiso, apoyo incondicional y paciencia durante el proceso de realización de la tesis. Al Profesor Fabio González por su apoyo.

A Gloria Díaz, Angel Cruz, Fabián Narváez por sus valiosos y oportunos aportes. A Nancy Landi?ez, Fabio Martínez, Gerardo Tibamoso, David Romo, Juan Carlos Caicedo, Lorenza Henao quienes siempre me brindaron una voz de aliento y a todos los integrantes del Grupo de Investigación BioIngenium, quienes me apoyaron para cumplir mis metas.

Abstract

Cervical cancer occurs significantly in women in developing countries every day and produces a high number of casualties, with a large economic and social cost. The World Health Organization, in the fight against cervical cancer, promotes early detection screening programs by different detection techniques such as conventional cytology (Pap), cytology liquid medium (CML), DNA test Human Papillomavirus (HPV), staining with dilute acetic acid and Lugol's iodine solution. Conventional cytology is the most used technique, being widely accepted, inexpensive, and with quality control mechanisms. The test has shown a sensitivity of 38 % to 84 % and a specificity of 90 % in multiple studies and has been considered as the choice test for screening [14].

The cervical cancer is not a public health problems in developed contries since more than three decades, among others because of implementation of other tests such as the CML which has increased the sensitivity to a figures that vary between 76 % and 99 %. This test in particular produces a thin monolayer of cells that are examined. In our contries this technique is really far from being applied because of its high cost. In consequence, the conventional cytology has remained in practice as the only possible examination of the cervix pathology.

In this technique, a sample of cells from the transformation zone of the cervix is taken, using a brush or wooden spatula, spread onto a slide and fixed with a preservative solution. This sample is then sent to a laboratory for staining and microscopic examination to determine whether cells are normal or not. This task requires time and expertise for the diagnosis. Attempting to alleaviate the work burden from the number of examinations in clinical routine scenario, some researchers have proposed the development of computational tools to detect and classify the cells of the transformation cervix zone. In the present work the transformation zone is firstly characterized using color and texture descriptors defined in the MPEG-7 standard, and the tissue descriptors are used as the input to a bank of binary classifiers, obtaining a precision of 90 % and a sensitivity of 83 %. Unlike traditional approaches that extract cell features from previously segmented cells, the present strategy is completely independent of the particular shape.

Yet most works in the domain report higher precision rates, the images used in these works for training and evaluation are really different from what is obtained in the cytology laboratories in Colombia. Overall, most of these methods are applied to monolayer techniques and therefore the recognition rates are better from what we found in the present investigacion. However, the main aim of the present work was thus to develop a strategy applicable to our real conditions as a pre-screening method, case in which the method should be robust to many random factors that contaminate the image capture. A segmentation strategy is very easily misleded by all these factor so that our method should use characteristics independently of the segmentation quality, while the reading time is minimized, as well as the intra-observer variability, facilitating thereby real application of such screening tools.

Keywords: Cervical Cancer, Papanicolaou, MPEG-7, Color Layout Descriptor, Scalable Color Descriptor, Edge Histogram Descriptor

Content

1. Introduction	2
2. Description and classification system of cervical cytology cells	5
2.1. Squamous epithelium	5
2.2. Papanicolaou staining	6
2.3. Cellular elements	7
2.3.1. Superficial cells	7
2.3.2. Intermediate Cells	7
2.3.3. Parabasal cells	8
2.3.4. Basal cells	8
2.4. Cellular alterations	9
2.4.1. Benign cellular alterations	9
2.4.2. Malignant cellular alterations	10
2.5. Diagnosis classification systems	10
2.5.1. Histological classification: WHO-NIC	11
2.5.2. Cytological classification: Bethesda system	11
2.6. Computer-assisted diagnosis for the analysis of cervical cytology	14
2.6.1. Image Preprocessing	15
2.6.2. Classification of cervical cells	17
3. Classification of cervical smear images using MPEG-7 descriptors	20
3.1. Methods	20
3.1.1. Classification based on global MPEG-7 descriptors	20
3.1.2. Classification models	23
3.1.3. Database	24
3.1.4. Experimental setup	26
3.2. Results	27
3.3. Discussion	31
4. Conclusions	33
A. Clasificación de células escamosas en citología cervical usando descriptores MPEG-7	34

B. Pap smear cell image classification using global MPEG-7 descriptors	36
References	44

List of figures

2-1. Histological location of the type of cells observed in cytology [56]	5
2-2. Superficial squamous cervical cells. Pap staining.	7
2-3. Intermediate squamous cervical cells. Pap staining.	8
2-4. Parabasal cells. Pap staining.	8
2-5. Cervical cells: Note the different sources of noise. A. Presence of blood cells. B. Inflammatory cells. C.Overlapping cells.	14
2-6. Computer-assisted diagnosis for detection and classification of cervical cell, including segmentation, feature extraction and classification	15
2-7. RGB space transformation to HSI. Component of A. Original Image, B. Sa- turation, C. Intensity, D. Hue.	16
2-8. Preprocessing A. Gaussian filter, B. Median filter , D. Midpoint filter.	16
3-1. Overview of proposed method.	21
3-2. Extraction process descriptor of spatial distribution of color	22
3-3. Extraction process descriptor of color distribution	22
3-4. The edge histogram descriptor identifies five types of edges and stores the values in an array 1D	23
3-5. Cells from the database B_1	25
3-6. Cells from the database B_2	26
3-7. Error classification for the B_1 and B.2 databases, in the two class problem.	27
3-8. Error classification for the B_1 database, together with the two-class problem. The panels are as follows: Color Layout (upper panel), Scalable Color (middle panel) and Edge Histogram descriptor (lower panel). Left plots correspond to the SVM polynomial function, whilst SVM radial-basis function kernel is depicted in the right pannel. The blue color represents the parameters with the best performance for each case.	28
3-9. Classification error reported in the B_1 and B.2 database, for the seven Class.	29
3-10. Error classification for the B_1 database, together with the seven class pro- blem. Left plots correspond to the SVM polynomial function, whilst SVM Radial-Basis function kernel is depicted on the right panels as follows: Co- lor layout (upper panel), Scalable Color (middle panel) and Edge Histogram descriptor (lower panel). Clear blue color represents parameters with the best performance for each case.	30

3-11. Histogram for isolated and complete images	32
B-1. Error for different methods in B.1	41
B-2. Error for different methods in B.2	42

List of tables

2-1. Classification systems. Modified from OMS 2007	11
2-2. Morphological characteristics of abnormal cells	14
2-3. The numbers of cases are maintained in each of the methods. (%) indicates the percentage of accuracy of detection of each method	19
3-1. Database	25

1. Introduction

Worldwide, the cervical cancer was the second most common type of cancer among women in 2002. In that year, 493,000 new cases and 274,000 deaths were reported [21]. These data increased in 2008, with nearly 529,000 new cases and 275,000 deaths reported [22]. About 80 % of these deaths occurred in developing countries [14, 21, 22, 63]. an important figure since this problem is still a health public problem in our countries. In Colombia, cervical cancer has a significant impact, in 2002 with 6,815 new cases [63] and between 2000 and 2009, nearly 1,800 deaths reported annually [27], making this disease a main cause of death in women between 35 and 64 until 2005 [63].

The main cause of cervical cancer is the infection with one or more oncogenic strains of human papillomavirus (HPV) [14, 53]. Most HPV infections resolve spontaneously, but those cases in which that infection persist may lead to the premalignant tumor, and if not treated, to cancer [51]. Previous lesions caused by HPV can usually take 10 to 20 years to become invasive cancer. Cervical cancer is preventable, by early detection and treatment of precancerous lesions, but most women in Colombia do not have access to effective screening programmes, and they are often detected when it is too late [14].

Cervical smear screening is the most popular method used to detect abnormal cells in the cervix. This method was proposed by Papanicolaou in 1940 [59], a reason why the method is called the Papanicolaou test, Pap smear or Pap test. This technique was devised to detect morphological changes, i.e., precancerous and cancerous cells. The method consists in taking a sample from the transformation zone of the cervix using a small spatula or brush. The sample is then smeared onto a glass slide, and fixed with a solution to preserve cells. The slides are sent to a laboratory where they are stained and examined by a cytotechnologist using a microscope [14]. The method is considered as effective, simple, fast and economical for detecting early cervical cancer [15]. It has been shown that performing cervical cytology reduces the probability of dying [16], up to a 60 % [48].

Cytological examination allows to identify abnormalities in squamous and glandular cells, an activity that requires time and expertise. In many developing countries, where a cytotechno-

logist must examine at 4X or 10X objectives, spending in average 6 to 10 minutes per slide and about 300 slides per day [12], the reading control is a fundamental health problem. This control is performed by a pathologist that performs a second read on every slide reported as abnormal and/or positive and a randomly selected sample of ten percent (10%) of the whole set reported as negative [11]. Such amount of work and factors like cell overlapping, uneven illumination, poor contrast and presence of blood and microorganisms result in examinations that are quite prone to mistakes. In developed countries where the price of such technicians is very high, several automatic algorithms has been integrated in the screening procedures from the seventies. A first obstacle to apply such solutions in Colombia is that the kind of examination is very different, i.e., these software solutions were developed to deal with screen monolayer preparations while in Colombia examinations are performed using the conventional cytology. It is then crucial to develop automated Pap test that can help our cytologists to screening samples.

Currently there are different solutions in the market, such as the AutoPap Primary Screening System (TriPath Imaging), ThinPrep Imaging System (Cytoc) and Path (Molecular Diagnostics); the first system complies with FDA regulations [24]. However, these tools require liquid-based cytology, that is to say, they use monolayer preparations in which cells are well separated and therefore they can be easily segmented. Straight application of these techniques is probably impossible in our countries because of the high cost of these monolayer preparatinos. Lately, there has been then an increasing investigation to automatically detect cell abnormalities in cervix samples. Nevertheless, these works depend on the quality of the cell segmentation, and again an impossible task in microscopical samples processed as conventional cytology since therein the cell boundary is hidden by inhomogeneities, other cells, some microorganisms or bacterias.

A number of automated methods to segment cell nuclei have been so far proposed [3, 4, 9, 30, 64]. However, segmentation is easily affected by noise and uneven illumination or by blurry contours. The aim of the present project has been to set global features that capture main characteristics with no need of segmenting the image objects. Global descriptors have been previously used in other problems [1, 19, 31, 76], but to the best of our knowledge none of them has focused on cervical images. In this work, we have proposed a global representation of the visual content by using some MPEG-7 descriptors, followed by a classification stage using either a K nearest neighbor algorithm (KNN) or Support vector machines (SVM).

Contribution

This thesis introduces a new method for classifying Pap smears cells, which are not previously segmented. This classification is carried out using global color and texture features. Yet the strategy was devised for noisy images, the method was evaluated using the available databases that are captured from monolayer preparations. The performance of the proposed descriptors was tested using two databases of images of Pap smears cells, one database with 500 cells and the other with 917 cells. Both databases were built by the Herlev University Hospital, Denmark. These databases are distributed unevenly into two categories, namely 2-class (normal and abnormal cells) and 7-class. The advantages of the proposed strategy are as follows: segmentation is not required and global features are chosen to capture most important characteristics of a cell, considered as a whole.

This work is organized as follows. Chapter 2 presents the description of cervical cell and reviews the most representative methods used to segment and classify cervical cytology images. Chapter 3 presents the visual features used to classify cells together with their performance for a two class problem: normal and abnormal cells, or the seven class problem.

Publications

Two articles were published as a result of the present work. In the first article, the MPEG-7 a descriptor was used, and the KNN classifier was assessed (annex A); subsequently the SVM classifier was evaluated. These latter results were presented in the second article (annex B).

Camargo LH, Díaz G, Romero E. Clasificación de células escamosas en citología cervical usando descriptores MPEG-7. VII Seminario Internacional de Procesamiento y Análisis de Información Médica - SIPAIM 2011. Universidad Industrial de Santander, Bucaramanga, Colombia.

Camargo LH, Díaz G, Romero E. Pap smear cell image classification using global MPEG-7 descriptors. 11th European Congress on Telepathology and 5th International Congress on Virtual Microscopy. 2012, Venice, Italy.

Grants

The present research was partially supported by the Research Direction, Bogotá - DIB, Universidad Nacional de Colombia, QUIPU code: 202010011343, HERMES code: 8521.

2. Description and classification system of cervical cytology cells

This chapter introduces the cervical cell. In section 2.1, the histology of the squamous epithelium is described. Section 2.2 presents the Papanicolaou technique, while section 2.3 explains the main cellular elements and the morphological characteristics of the superficial, intermediate, parabasal and basal cells. Section 2.5 tackles with the diagnosis classification systems and finally section 2.6 presents some previous cervical cell classification works.

2.1. Squamous epithelium

The epithelium that covers the vagina, and the cervical vaginal portion of the sexually mature woman, is a stratified squamous epithelium, also called the squamous epithelium. It measures approximately 0.2 mm thick and is composed of four layers [56] or strata histological, that can be well differentiated, as shown in figure 2-1:

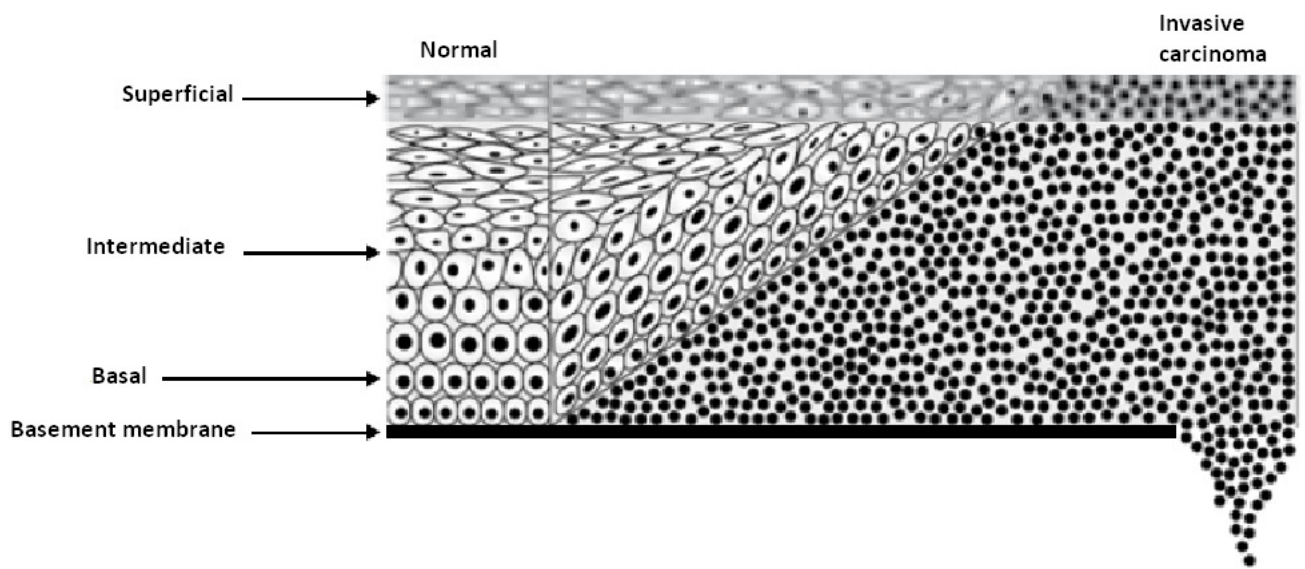


Figure 2-1.: Histological location of the type of cells observed in cytology [56]

- Basal Layer - basal cell: Is the innermost layer. This tissue is formed by a single layer of cylindrical palisade-wise arranged cells at the base membrane. Germ cells (those cells that spring the others) are either round or oval, with large nucleus (oval center), which occupy much of the cell. The process of regeneration of the epithelium starts at this cell layer. The remaining layers represent the stages of cell maturation and differentiation that migrate to the surface; it is constantly germinating new cells.
- Parabasal Layer - parabasal cells: This tissue consists of several layers of rounded or polyhedral cells, which have a central nucleus and cytoplasm thick with intercellular bridges.
- Intermediate Layer - intermediate cells: It is composed of many layers of flattened cells with glycogen-rich cytoplasm, relatively small, central and vesicular nucleus.
- Surface Layer - surface cell: This part consists of several layers of flattened, polygonal-shaped and large cells with thin and bright cytoplasm, without intercellular bridges and pyknotic nucleus.

The squamous epithelium is able to exfoliate various cell types, depending on the maturation degree. Typically, cells are exfoliated from the upper strata, unless either ulcers are present or the sample is taken with excessive energy. When smearing this sample, cells are flattened, allowing visual evaluation of four types of cells: superficial, intermediate cells, parabasal cells and basal cells.

2.2. Papanicolaou staining

The smeared specimen of female genital tract is almost exclusively stained using the Papanicolaou technique. It is a polychrome staining method, consisting in a nuclear staining with good cytoplasmic contrast [15, 59, 61].

The method presents many advantages since it allows a proper nucleus definition, detection of chromatin patterns and a well delimited cytoplasm, thereby permitting observation of maturation and metabolic activity. Papanicolaou staining utilizes three types of dyes, namely hematoxylin, Orange G and eosin alcohol (EA) [15].

Hematoxylin selectively stains the nuclei, clearly highlighting chromatin. Orange G is a monochromatic stain that colors cytoplasm, staining keratin with bright orange. Keratin is not normally present in the cervical and vaginal epithelium, however, it is usually present in keratinized carcinomas, so the presence of a brighter orange color in the cytoplasm is an important diagnostic sign. Eosin alcohol (EA) is a polychrome stain composed of eosin, light green and Bismarck brown. Eosin stains the cytoplasm of mature squamous cells, nucleoli and cilia. Superficial cells are pink stained with eosin and are therefore eosinophilic. Parabasal

and intermediate cells are stained in green or blue, depending on the EA staining time, and therefore are known as cyanophilic. These cells are metabolically active [15].

2.3. Cellular elements

The squamous epithelium cell types that are recognized in the cytological smear are as follows [56]:

2.3.1. Superficial cells

These cells present the greatest maturation degree; they have a diameter of 50 to 60 μm , they are also eosinophilic (pink stained with eosin) or basophils (blue stained with hematoxylin), as illustrated in figure 2-2. In this case, the cytoplasm has polygonal shape; it is also clear and homogeneous, and clearly delimited. The cells appear on a one-by-one basis or in groups. In the cytoplasm, occasionally, keratohyaline granules are identified. These cells may have a large nucleus, of approximately 7 μm (in which the chromatin structure is easily recognizable, being also regularly and finely distributed), otherwise, cells are small and may present a pyknotic nucleus (in which chromatin is condensed by degeneration and is dark-colored uniformly).

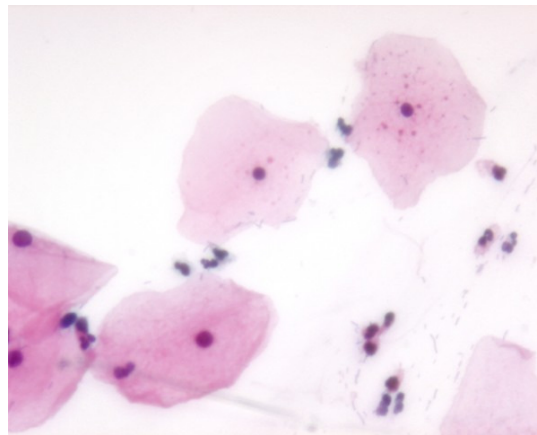


Figure 2-2.: Superficial squamous cervical cells. Pap staining.

2.3.2. Intermediate Cells

They represent the most consistent and largest cell in the vaginal smear. They have a size of about 30 a 50 μm and nuclei around 8 μm (Figure 2-3). They are basophilic, polygonal, rounded shaped and their nuclei are usually vital, but occasionally pyknotic. Intermediate cells can differentiate in large (with glycogen) and small (without glycogen) cells. During

pregnancy and under the strong influence of progestins or androgens, these cells show a canoe or boat and wide margin, reason why they are called navicular cells.

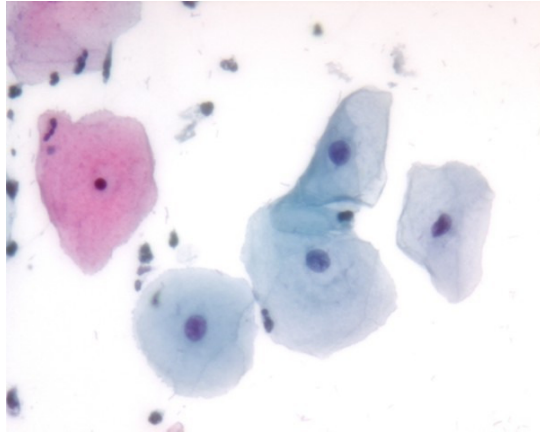


Figure 2-3.: Intermediate squamous cervical cells. Pap staining.

2.3.3. Parabasal cells

These cells have a diameter of approximately 20 μm and nuclei of approximately 9 μm . Intensely stained, they are basophilic and exhibit an elongated, oval-like and also rounded shape (Figure 2-4). The parabasal cells can be classified as being either large or small.

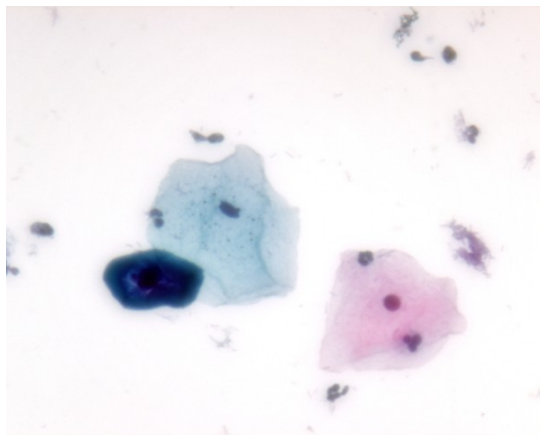


Figure 2-4.: Parabasal cells. Pap staining.

2.3.4. Basal cells

These cells have a cytoplasmic diameter of 12 - 14 μm and a nucleus diameter varying between 8-10 μm . The cells have intense basophilic staining and frequently a deteriorated nucleus. If chromatin is preserved, nucleoli may be visible.

2.4. Cellular alterations

2.4.1. Benign cellular alterations

In the cell nucleus, some of the following changes may occur:

- Nucleus swelling: A loss of cellular osmotic regulation causes cell death by nuclear edema formation, resulting in swelling of the cell nucleus.
- Karyolysis: Presence of large, hyperchromatic and without structure nuclei.
- Karyorrhexis: nuclear destruction occurs, with thickening of the chromatin structure, which leads to the rupture of chromatin bridges.
- Karyopyknosis: nuclear condensation caused by loss of water. The dead nuclei then retract, reducing their size and becoming hyperchromatic.
- Nucleus hyperchromasia wall: It occurs if the particles are deposited on the chromatin nuclear membrane.

Some cytoplasm changes:

- Perinuclear halo formation: This causes nuclear shrinkage.
- Pseudoeosinophilic: As a consequence of cytoplasm degradation, a pseudoeosinophilic staining may appear.
- Vacuolization: The cellular uptake of water causes an increase in cell size and vacuole formation.
- Hyalinization: Condensation of the cytoplasm which is stained in a deep red.
- Shape changes: Parabasal cells are influenced by degradation processes. A possible reaction of the surface can be caused by cytoplasm elongation, i.e., it takes a spindle-like shape.

Causes

- Reactive cellular inflammatory changes
- Reactive cellular changes associated with radiation
- Cellular changes associated with Intrauterine device (IUD). When patients are treated with an intrauterine spiral device, reactive changes are observed: enlarged nuclei and vacuolization of the endocervical epithelium. Metaplastic cells and endometrial cells may appear together with an inflammatory smear.

- Post - hysterectomy glandular cells: During the first four weeks after intravaginal surgery, like conization or hysterectomy, a granulation tissue is observed in the region of the surgical wound. This tissue consists of squamous and endocervical epithelium with inflammatory disorders, regenerative epithelium, leukocytes, histiocytes, fibroblasts and protein-rich debris.
- Atrophy.

2.4.2. Malignant cellular alterations

The squamous epithelium cancerization involves not cornified epithelium. The histological criteria of malignancy depend on the degree of vertical extent of atypical changes of the basal cell layer. In addition to the tissue widening, there appear nuclear polymorphism, enlargement and hyperchromasia, as well as atypical chromatin and a large number of mitosis.

2.5. Diagnosis classification systems

There are several systems for classifying premalignant cervical tumors, including cytology and histology. The first system corresponds to the Papanicolaou classification, which has 5 numerical classes or grades (I, II, III, IV, V); this system is purely cytological [60].

In 1961, in the first International Congress of Cytology, held in Vienna, experts agreed about the terms for the three most crucial cytological cervical lesions, namely invasive carcinoma, carcinoma in situ (CIS) and dysplasia, the latter classified as mild, moderate, and severe or severe [32]. In 1967, in order to propose the classification of cervical intraepithelial neoplasia (NIC), this system considered various changes observed in three different degrees of dysplasia, including Grade III severe dysplasia and the previous CIS [66,67], see table **2-1**.

In 90s, a system was proposed at the Bethesda National Cancer Institute U.S. (United States National Cancer Institute). In this system, NIC II and III were assembled in a single group, "High-grade squamous intraepithelial lesion" (HSIL). In 2001, atypical cells were divided into ASC-US (atypical squamous cells of undetermined significance) and ASC-H (atypical squamous cells cannot exclude a squamous intraepithelial high grade lesion). The Bethesda system is recommended by OMS for cytological reports [14]. Table **2-1** is a summary of the different classification systems.

Papanicolaou [60]	NIC [66, 67]	OMS [32]	Bethesda [6, 55, 72]
Class I	Normal	Normal	Normal
Class II	Atypia	Atypia	ASC
Class III	NIC I	Mild dysplasia	LSIL
Class III	NIC II	Moderate dysplasia	HSIL
Class III	NIC III	Severe dysplasia	HSIL
Class IV	NIC III	in situ cancer	HSIL
Class V	Invasive carcinoma	Invasive carcinoma	Invasive carcinoma

Table 2-1.: Classification systems. Modified from OMS 2007

2.5.1. Histological classification: WHO-NIC

In many countries, the NIC classification and the Classification of the World Health Organization (WHO) are widely used in cytological reports, although these systems should only be used in histological reports, including various tissues as follows:

- Mild dysplasia (NIC I, cervical intraepithelial neoplasia-grade I): The extension of the atypical epithelium is limited to the lower third of the epithelium thickness.
- Moderate dysplasia (NIC II, grade II): It affects the lower two thirds of the epithelium thickness.
- Severe dysplasia (NIC III, grade III): The upper third is also involved in the atypia.
- Carcinoma in situ: There are no mature cells migrating towards the surface. The biological significance of these alterations does not differ from severe dysplasia and may therefore be grouped according to the concept on the lesion (NIC III).

2.5.2. Cytological classification: Bethesda system

The Papanicolaou classification system and the Bethesda System diagnostic reports are used for vaginal or cervical cytology; the latter was introduced in 1988 [55], revised in 1991 [6], and updated in 2001 [72] in order to establish a uniform terminology and standardize diagnostic reports. This system includes epithelial cell abnormalities in squamous cells, glandular and other malignancies.

Atypical squamous cell (ASC)

This category includes atypical squamous cell of undetermined significance (ASC-US), also atypical squamous cells of undetermined significance, suggestive of high-grade squamous

intraepithelial lesion (ASC - H) [72]. Among the morphological changes that characterize the immature metaplastic cells with atypia (in [28]):

- Cells isolated or grouped with slight loss of cohesiveness.
- Round or oval.
- With scant cytoplasm, also dense and anaphilic or cyanophilic.
- Loss of nuclear/cytoplasmic ratio.
- With enlarged nuclei, and showing a hyperchromasia trend, with slight irregularity of the nuclear membrane which can also be smooth.
- Thin granular chromatin either evenly distributed or slightly irregular.
- Variations in nuclear size and shape.

Atypia squamous cells of undetermined significance (ASC-US)

Most changes occur in the mature squamous cells of intermediate superficial type. Changes can also be observed in cells with metaplastic morphology (atypical metaplasia), in cells being repaired (repair atypical), as well as in cells with morphology atrophy and parakeratotic changes.

Changes occur in the form of slight increases in the nuclear/cytoplasmic ratio. The cytoplasm presents a cyanophilic-acidophilic appearance marked with halos; the nuclear membrane looks smooth or slightly irregular, the nucleus is moderately enlarged up to 2-3 times the nucleous size of an intermediate cell (75 to 120 μm^2). A normal or slightly hyperchromatic, usually negative chromocenters, and absent nucleoli, can be associated to binucleation [33,65]. Similarly, nuclear abnormalities are associated with orange and dense cytoplasm (atypical parakeratosis) [71].

Atypical squamous cells of undetermined significance cannot exclude HSIL (ASC - H)

These cells include cases in which the cellular changes are quite important, but not enough to be considered as conclusive, due to a shortage of these cells or due to other reasons, such as inflammation or bleeding [71].

Low-grade squamous intraepithelial lesion (LSIL) (including changes associated with infection by Human Papilloma Virus (HPV) or mild dysplasia (NIC I))

Because HPV infects only the cell nucleus and produces a multiplication of DNA and chromosomes, the criteria of malignancy is related to nuclear changes, including the following:

- Structure of atypical chromatin: irregular structure of chromatin in the form of lumps. The hematoxylin used in the method of Papanicolaou stains the nucleus; the biologically active parts, which constitute euchromatin and interface regions at rest, are active genetically- heterochromatin-like dark regions. In a normal nucleus these particles are distributed uniformly and take a finely granulated aspect. An atypical nucleus shows an increase of heterochromatin with a degree of atypia.
- Hyperchromasia nuclear and formation of the chromocenter.
- Polymorphism.
- An absolute increase of nucleus size.
- Increase the nucleus/cytoplasm ratio.
- Formation of the nucleolus.

High-grade squamous intraepithelial lesions (H-SIL), (comprising moderate dysplasia, severe NIC II, NIC III, Ca in situ)

This lesion is characterized by cells, smaller than those of low-grade lesion, usually isolated cells non-cohesive. The nuclear size is comparable to the low grade, but with a marked increase in the nuclear/cytoplasmic ratio. This finding is related to immature cytoplasm of occasionally dense metaplastic type; there is nuclear hyperchromatism with fine or granular chromatin. The nuclear membrane is irregular, without nucleoli [71].

Carcinoma

The affected cells show irregular distribution of chromatin and paracromatina (clear spaces between thick clumps of chromatin), there is also a prominent nucleoli. This case presents keratinizing differentiation criteria (orangofilia, cannibalism, elongated cells with pyknotic nuclei) and also keratinizing criteria (syncytial, dense cytoplasm with irregular, angulated and central nuclei with nucleoli and irregular chromatin with highly parachromatic clarification) [71].

Table 2-2 shows the morphologic characteristics according to the Bethesda system.

Cell structure	ASC-US	ASC-H	L SIL	HSIL
Cytoplasm	cyanophilic - acidophilus			
Cytoplasm	binucleate			
Nuclear membrane	irregular		slightly irregular	
Nuclear hyperchromasia	low		hyperchromatic	
Nucleus	75-120 μm			
Nucleolus	asent		small o absent	absent
Chromatin	uniform	irregular	irregular	Fine or granular

Table 2-2.: Morphological characteristics of abnormal cells

2.6. Computer-assisted diagnosis for the analysis of cervical cytology

Manual screening of cutology smears is an intensive and demanding labor, for which the cytotechnologist must be concentrated during extended periods of time [5]. Different sources of error are inherent of this kind of examination, including sampling error, blood, inflammatory cells or overlapping cells, as illustrated in figure 2-5.

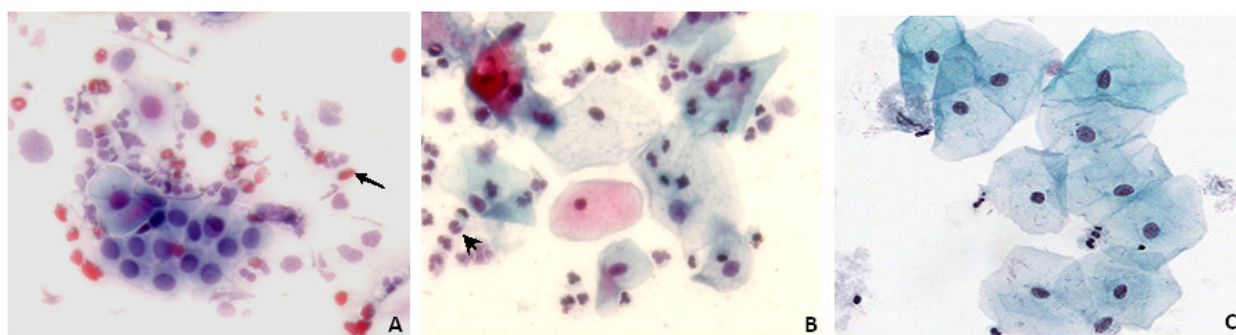


Figure 2-5.: Cervical cells: Note the different sources of noise. A. Presence of blood cells. B. Inflammatory cells. C.Overlapping cells.

Current available commercial softwares use more or less the same scheme for automatic detection and classification of cervical cells, starting by a cell segmentation which is then binarized for somehow extracting different morphological characteristics [7], (see figure 2-6). Traditionally, the captured cytology image is firstly pre-processed to improve the segmentation quality and then cells are segmented and classified.

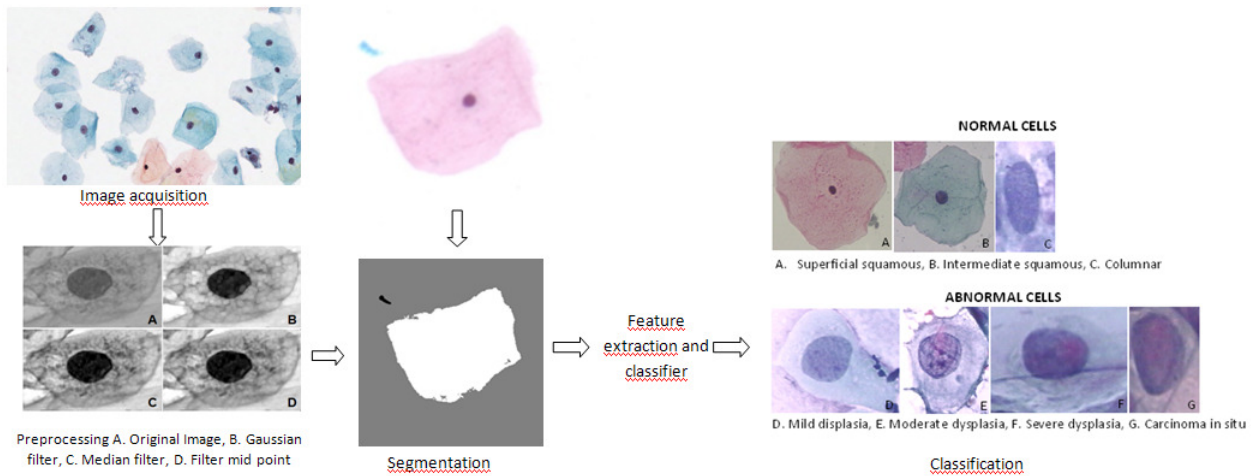


Figure 2-6.: Computer-assisted diagnosis for detection and classification of cervical cell, including segmentation, feature extraction and classification

2.6.1. Image Preprocessing

The image preprocessing stage is required when segmenting the cell for background extraction [64]. Every captured image exhibits a certain percentage of noise [20] and may have low contrast [26]. Overall, most methods start by reducing noise and increasing contrast. Median filters are traditionally used to reduce noise, where the gray level of every pixel is replaced by the median of the intensity levels of the pixel neighborhood [26].

On the other hand, different methods have been proposed to improve the image contrast. These methods are usually divided as point-, spatial- and frequency- domain strategies. Plisiti et al. achieved contrast enhancement and edge sharpening by applying contrast limited adaptive histogram equalization to the red, green and blue image components. A global threshold is then extracted for each image, using the Otsu method [58]. Finally, a single image is obtained by using a logical or-operator on the three components [64]. Li et al., proposed to map the Pap image into a gray scale, removing the small white pores. Then, an edge detector is used to find the cell and nucleus boundaries. Finally, these borders were smoothed [36].

Among the studies that use methods of spatial filters, Li and Najarian used a weighted median filter, introduced by Brownrigg [17], who preserved the edge information. In this work, the image was preprocessed. First, the image color was transformed into gray-scale, then a modified weighted filter was calibrated below threshold gray-scale; finally, the modified weighted filter was used [37]. Mat Isa implemented three algorithms in the preprocessing stage, namely a median filter, histogram normalization and histogram equalization [46, 47]. Yang et al. proposed a trim-meaning filter (designed by the authors) in order to remove

noise, remove impulses, Gaussian noise and also to preserve the edge sharpness [76].

Figure 2-7 illustrates an example of preprocessing with transformation from RGB space into HSI. The RGB cube becomes a cylindrical shape, thus saturation corresponds to a value of radial distance, hue becomes a function of the angle with respect to the polar-coordinate system, and intensity is the distance along the axis perpendicular to the plane of polar coordinates. Figure 2-7 (a) is the original image, and 2-7 (b, c, d) is the image created by the transformation form RGB to HSI.

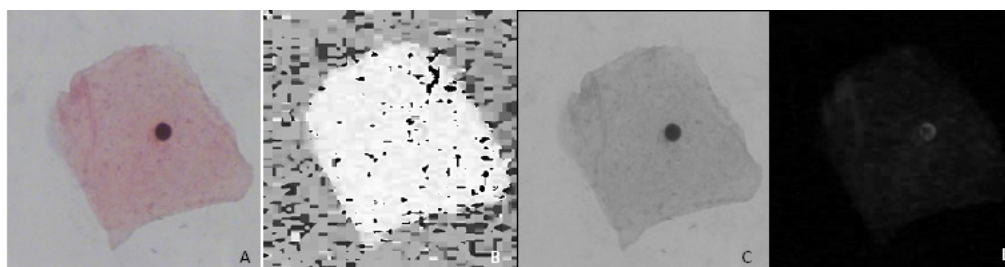


Figure 2-7.: RGB space transformation to HSI. Component of A. Original Image, B. Saturation, C. Intensity, D. Hue.

For the image obtained in the component of intensity, there are several noise-filtering techniques, such as Gaussian filter, median filter, and midpoint filter. Figure 2-8 (a) shows the Gaussian filter. Figure 2-8 (b) shows the Median filter and figure 2-8 (c) shows the Midpoint filter.

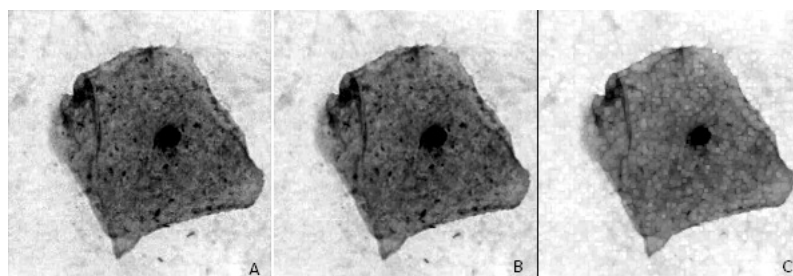


Figure 2-8.: Preprocessing A. Gaussian filter, B. Median filter, C. Midpoint filter.

2.6.2. Classification of cervical cells

Neural Networks

Among the most widely-used classifiers for the classification of Pap smear we find neural networks [68] [49]. The multilayer feed-forward back-propagation network has been implemented in the trading system Papnet. The classification accuracy obtained was between 80 % and 90 % [39].

Li et al. employed the multilayer feed-forward back-propagation network, with features such as percentage of cell coverage, cell coloration, distribution of cell size, average of cell size, existence of multinuclear cells, existence of halos in cells, distribution of nuclear size, average of nuclear size, nuclear/cytoplasm ratio, and percentage of cells without nucleus, finding proper classification accuracy of cervical precancerous cells in normal and abnormal classes almost in 99 % of the times [36].

The multilayer perceptron neural network has also been used to classify normal cells, LGSIL and LIEag, with accuracy levels of 76 % [45]. The performance of the mentioned systems might be dependent on the ability of human experts to extract the characteristics of cervical cells [46].

[In 1991, Lee et al. used decision tree classifiers, neural networks classifiers, hybrid classifiers and multiple classifier integration as learning and classification algorithms. This led to the best classification results. The error rate was 18.7 % when implementing a three-classifier integration involving the neural network classifier, the hybrid classifier and the quadratic tree classifier, obtaining the best performance [35]].

Tabu search (TS)

The Tabu Search (TS) is a meta-heuristic memory-adaptive procedure proposed by Glover in 1986 [25]. In 2006 Marinakis et al. used this local algorithm and combined three binders of the nearest neighbor, finding considerable accuracy levels (91.93 % to 98.47 %) regarding the new database in the problem of two classes (normal and abnormal). Also for the new database, accuracy levels in the problem of seven classes ranged from 91.38 % to 96.95 %, while in the old data set, for the problem of two classes, accuracy ranged from 97.6 % to 100 % , and for the problem of seven classes it fluctuated between 95.2 % and 100 % [40].

Ant colony

Another algorithm was based on the optimization of ant colony for the problem of feature selection. Such an algorithm uses two different types of the binder of nearest neighbor classification [41]. In this approach, accuracy in the new set of data, for the problem of two classes, varied from 93.78 % to 99.78 % , and for the problem of seven classes it fluctuated between

93.34 % to 97.92 %. On the other hand, in the old set of data, and for the problem of two classes, accuracy fluctuated between 98.8 % 100 % , and for the problem seven classes these figures oscillated between 98.2 % and 100 %. The main characteristics to correctly classify new cells were the following: cytoplasm and nucleus brightness, and N/C ratio [43].

Genetic algorithms

Genetic algorithms were subsequently used to solve the problem of feature selection. These algorithms were complemented with three types of nearest-neighbor classifiers. The accuracy regarding the new set of data, for the problem of two classes, fluctuated from 94.11 % to 99.67 % , and for the seven classes is ranged from 93.45 % to 98.36 % . Regarding the old data set, for two classes, accuracy ranged from 98.8 % to 100 % , and for the problem of the seven classes, accuracy oscillated between 98 % and 100 % [42]. Nanni also used the previous database and applied it for a cross-validation of 5 partitions, namely texture descriptor local binary patterns (DL), which is defined as a measure of invariant texture in grayscale; then, the procedure was performed using SVM classification [54].

In 2010, Samsudin used a public database, and addressed the two-class problem with the following 10 features: area of the nucleus and the cytoplasm, nuclear/cytoplasm ratio, brightness of the nucleus and the cytoplasm, smaller diameter of the nucleus, perimeter of the nucleus and position of the nucleus. This study also deployed classification in a group of classifiers, namely GB k-NN, GDW k-NN and GLMV; the results outperformed those of individual techniques [69].

A hybrid algorithm that combines genetic algorithms with ID3 algorithms (Interactive Dichotomizer Version 3) was also used as binder by [50]. Liu, experimented with a database of 149 cells, comparing four types of Wavelet transforms to classify pixels. The best performance occurred when using Daubechies 16 and Gabor. The authors also studied three types of features, with and without multi-spectral or texture information, finding better results when using both types of information [38].

Mat et al. proposed a feature extraction algorithm (through growth of regions) based on seeds. This method chooses an $N \times N$ neighborhood. Once an initial seed has been located within the neighborhood, the study of the vicinity of pixels begins. Each pixel is compared with the average gray levels and the standard deviation. The pixels that belong to the new neighborhood are marked, and then the growth of the region on the edge of the object is stopped. The method uses the size and intensity of the nucleus and cytoplasm as input data for the classification of the precancerous cervical cells. The size of the regions of nucleus and cytoplasm is calculated from the number of pixels. The authors classified 211 normal cases, 143 LSIL and 196 HSIL, reaching accuracy levels over 96.50 %, specificity was 100 % and

sensitivity was 95.33% [46].

Table **2-3** compares the different classification algorithms. Note that results obtained with tabu search, ant colony and genetic algorithms are over monolayer preparations and they use characteristics that completely dependent on the cell segmentation. Importantly, because of the intrinsic nature of these algorithms, they are computationally very expensive and require an heavy computational infrastructure, an unthinkable situation at any real cytology laboratory.

Algorithm	Old database	New database
	2 Classes — 7 Classes(%)	2 Classes — 7 Classes (%)
Tabu Search [40]	97.6-100 — 95.2-100	91.93-98.47 — 91.38-96.95
Colony of ants [43]	98.8-100 — 98.2-100	93.78-99.78 — 93.74-97.92
Genetic Algorithms [43]	98.8-100 — 98-100	94.11-99.67 — 93.45-98.36

Table 2-3.: The numbers of cases are maintained in each of the methods. (%) indicates the percentage of accuracy of detection of each method

3. Classification of cervical smear images using MPEG-7 descriptors

This chapter presents the classification of cervical cell. Section 3.1, describes system development. This section also introduces the main results 3.2.

Several strategies have been previously applied for classifying cervical cytology cells, all pursuing nucleus segmentation. Sanchez sets regions [70] using a simple threshold [62], a procedure broadly adapted to different techniques: namely local adaptive segmentation nuclei procedures [37], seed growing (Romberg y col. en [18]), mathematical morphology [34], a Hough transform [23, 52], and active contours [2, 3, 9, 30, 64]. Jantzen and Dounias proposed several cell features as morphometric descriptors; these features include the nucleus and cytoplasm areas, the nucleus/cytoplasm proportion, the nucleus and cytoplasm brightness, smaller and larger nucleus/cytoplasm diameters, nucleus and cytoplasm roundness, nucleus and cytoplasm perimeters, nucleus position, nucleus/cytoplasm maxima and minima. Nevertheless, these morphometric characteristics require prior accurate segmentation, which is hard to achieve by human intervention using commercial software packages such as CHAMP (Cytology and Histology Modular Analysis Package, Aarhus, Denmark) of DIMAC (Digital Image Company) [8, 29, 44, 57].

3.1. Methods

Rather than attempting to detect some of the previously reported morphometric features, the present investigation used two global MPEG-7 color descriptors, namely Color Layout and Scalable Color. It also used one texture descriptor, the Edge Histogram descriptor, which was used as the representation space together with two supervised classification algorithms (SVM and KNN) that divide the different classes, this is shown in figure 3-1.

3.1.1. Classification based on global MPEG-7 descriptors

The cell classification approach is carried out using color and texture MPEG-7 descriptors, thereby attempting to capture information related with the particular color spatial location and global color distribution of both the nucleus and cytoplasm. The texture descriptor

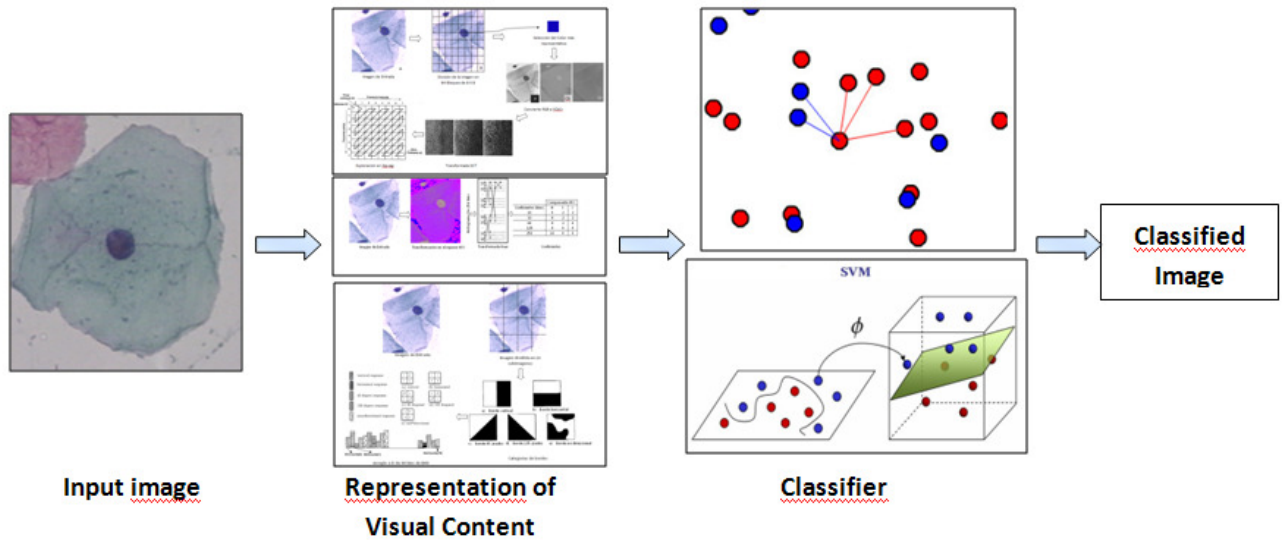


Figure 3-1.: Overview of proposed method.

stands for the particular borders of both nucleus and cytoplasm and their intrinsic relationships. These global characteristics are not evaluating the classical morphometric features, but they are using nucleus and cytoplasm visual primitives as discriminate factors.

Color layout

This descriptor, typically used in the YCrCb color space, captures the spatial color distribution in an image or an arbitrary region. Basically, the color layout descriptor uses representative colors on a grid, followed by a Discrete Cosine Transform (DCT) and an encoding of the resulting coefficients. The feature extraction process consists of two parts; grid based representative color selection and DCT transform with quantization. Specifically, an input image is divided into 64 blocks, their average colors are derived and transformed into a series of coefficients by performing a conventional DCT. A few low-frequency coefficients are selected using zigzag scanning and quantized to form a Color Layout Descriptor, is show in figure 3-2 [10].

Scalable color

This descriptor is extracted from a color image histogram in the hue-saturation-value (HSV) color space. This histogram, constructed with fixed color space quantization, is projected into a set of Haar bases so that the obtained coefficients constitute a scalable color representation. The histogram values must be normalized and non linearly mapped into a 4-bit integer representation, giving higher weight to small values. The Haar transform is applied then to this histogram version with two basic operators: sum and difference bin neighbor,

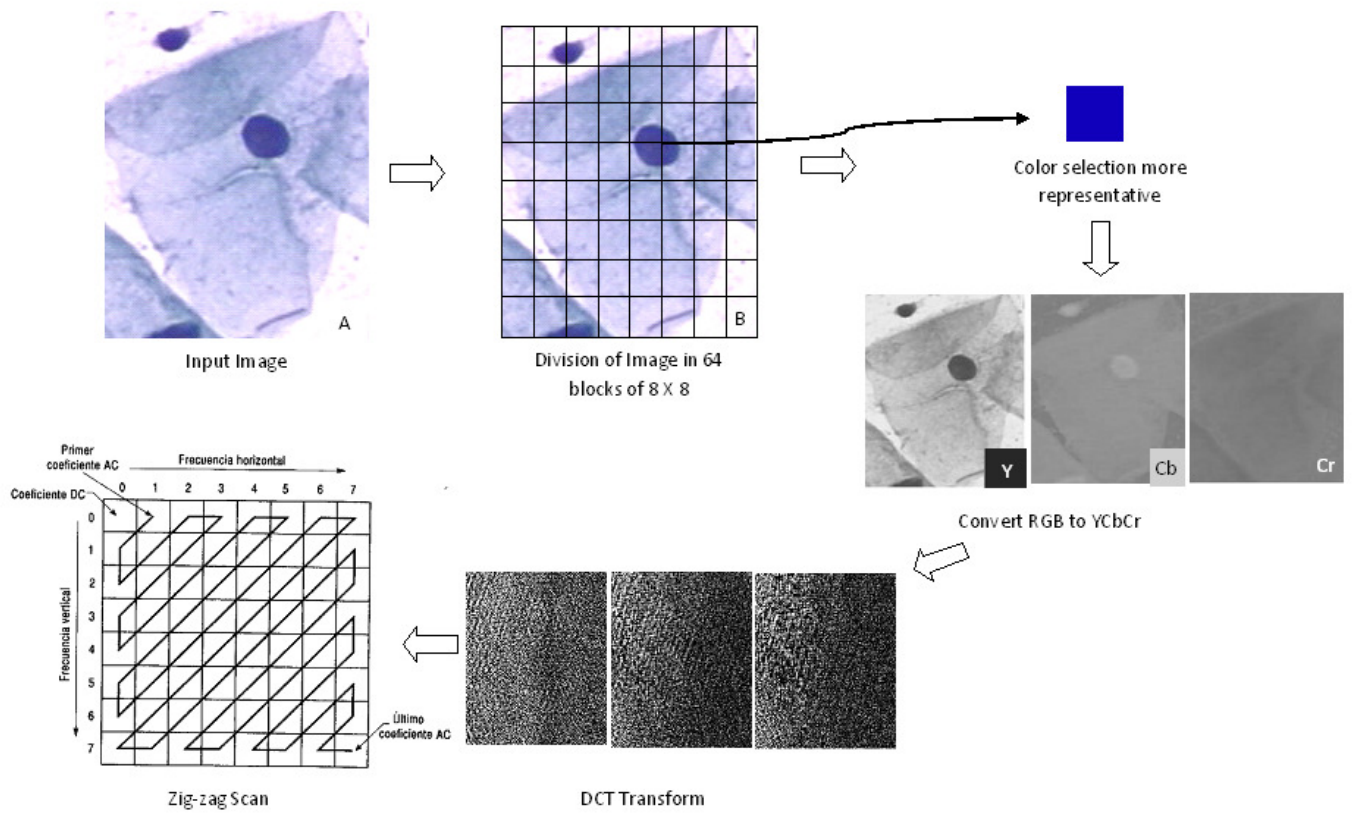


Figure 3-2.: Extraction process descriptor of spatial distribution of color

decomposing the histogram into low and high frequency subbands, as is showed in figure 3-3 [10].

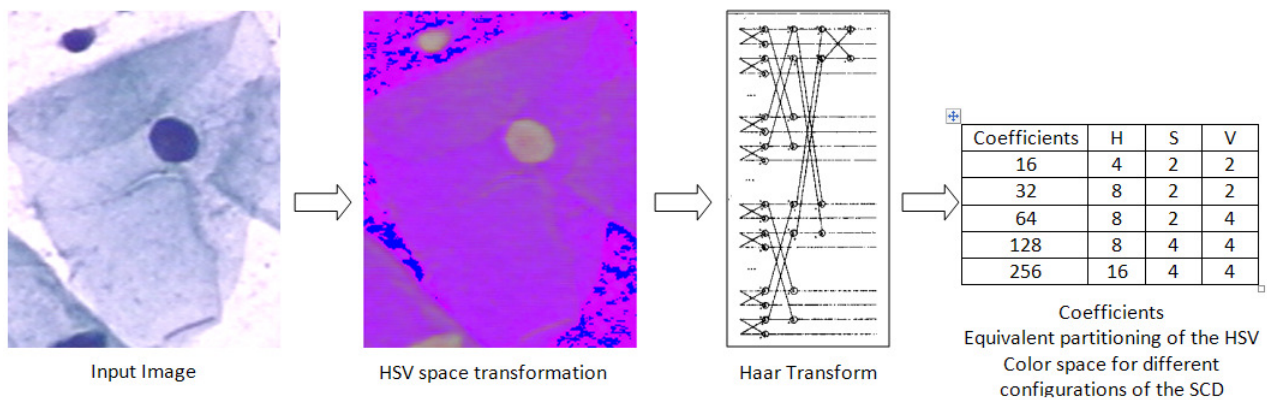


Figure 3-3.: Extraction process descriptor of color distribution

Edge histogram

This descriptor captures the spatial edge distribution, a very useful feature for image matching, even though the underlying texture may not be homogeneous. A given image is first sub-divided into sub-images and local edge histograms, for each of these sub-images, are computed. Edges are then coarsely grouped into five categories: vertical, horizontal, 45 diagonal, 135 diagonal, and isotropic (nonorientation specific). Thus, each local histogram has five bins corresponding to the above five categories. The image partitioned into 16 sub-images results in 80 bins. These bins are nonuniformly quantized using 3 bits/bin, resulting in a descriptor with size of 240 bits (see figure 3-4), [75].

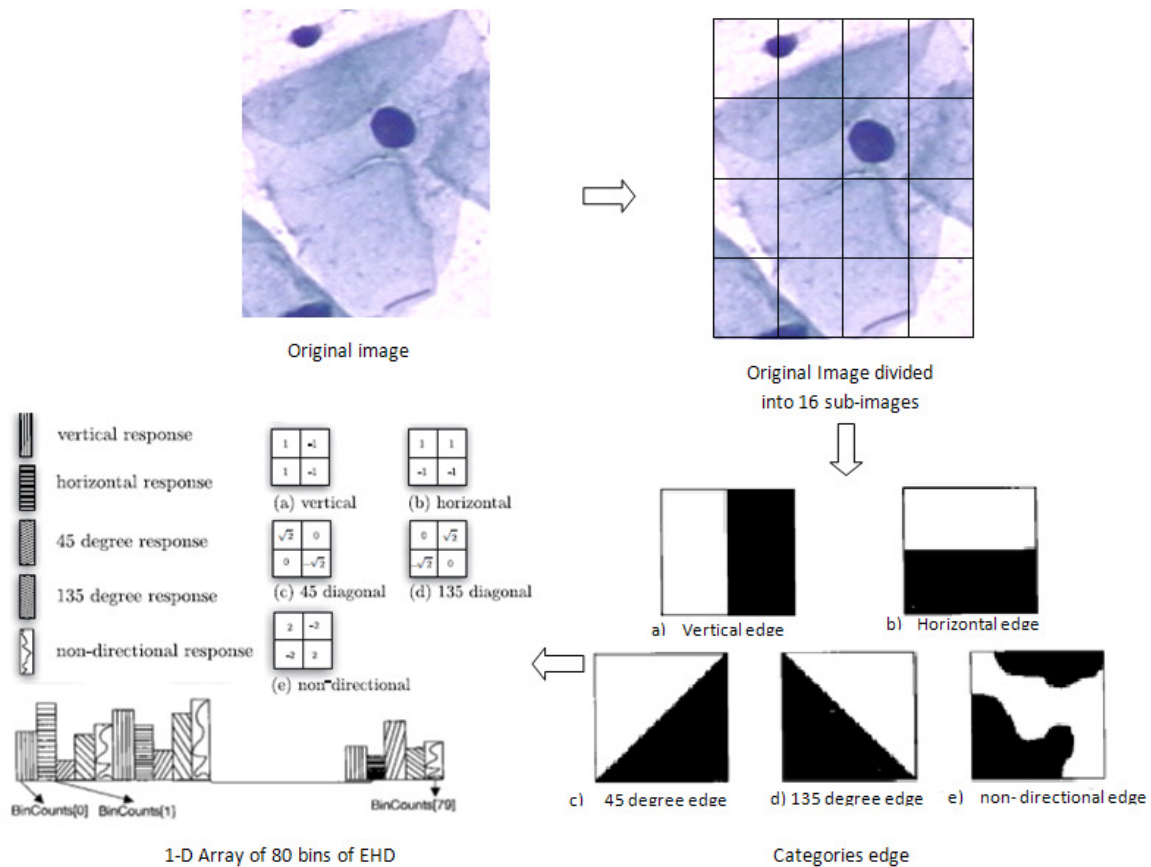


Figure 3-4.: The edge histogram descriptor identifies five types of edges and stores the values in an array 1D

3.1.2. Classification models

The classification method used a classical K-Nearest Neighbor algorithm and a Support Vector Machine. The proposed approach was evaluated under a 10-fold experimental setup.

The k-NN decision rule

The k-nearest neighbor method is an intuitive method that classifies unlabeled samples based on their similarity with samples in the training set. Given the knowledge of N prototype features (vectors of dimension k) and their correct classification into M classes, the k-NN rule assigns an unclassified pattern to the class that is most heavily represented among its k neighbors in the pattern space, under some appropriate metric [13]. In this work euclidean distance was used.

The SVM algorithm

A support vector machine (SVM) is a classification model that finds an optimal separating hyperplane that discriminates two classes. A SVM is a linear discriminator; however it can perform non-linear discriminations thanks to the fact that this is a kernel method. In this work, we used a SVM version that uses sequential minimal optimization algorithm. The multi-class classification problem is solved using a one vs. all strategy: a binary classifier for each class labeling the class samples as positive examples and other class samples as negative ones. The final decision is set to the class having the largest decision function among all classes [73].

3.1.3. Database

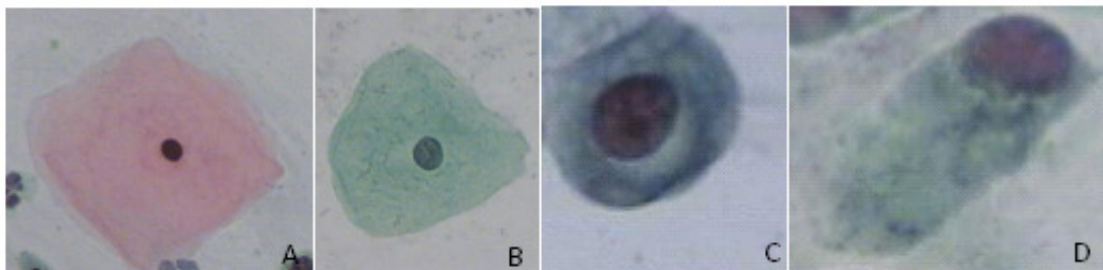
Two databases, which consist of images with single cells from the Herlev University Hospital in Denmark, were used (<http://labs.fme.aegean.gr/decision/downloads>). Skilled cyto-technicians and doctors manually classified each cell into 2 classes: abnormal and normal and then sub-classified them into seven classes. Each cell was examined by two cyto-technicians, and difficult samples also by a doctor. In case of disagreement, the sample was simply discarded [8, 44, 57]. Finally there were two databases: B_1 data contains 500 cells with the following distribution: 1. Normal: columnar epithelial, parabasal squamous epithelial, intermediate squamous epithelial, superficial squamous epithelial. 2. Abnormal: mild squamous non-keratinizing dysplasia, moderate squamous non-keratinizing dysplasia, severe squamous non-keratinizing dysplasia (Table.3-1). Examples of images of this database are shown in Figure 3-5

The B_2 data contain 917 cells with the following distribution: 1. Normal: superficial squamous epithelial, intermediate squamous epithelial, columnar epithelial. 2. Abnormal: mild squamous non-keratinizing dysplasia, moderate squamous non-keratinizing dysplasia, severe squamous non-keratinizing dysplasia, squamous cell carcinoma in situ intermediate (Table.3-1). Examples of images of this database are shown in Figure 3-6.

Clase	B-1	B-2
Normal	1. Superficial (50)	1. Superficial (74)
N.	2. Intermediate (50)	2. Intermediate (70)
N.	3. Parabasal (50)	
N.	4. Columnar (50)	3. Columnar (98)
Total N.	200	242
Anormal	5. Mild Dysplastic (100)	4. Mild Dysplastic (182)
A.	6. Moderate Dysplastic (100)	5. Moderate Dysplastic (146)
A.	7. Severe Dysplastic (100)	6. Severe Dysplastic (197)
A.		7. Carcinoma in situ (150)
Total A.	300	675
Total	500	917

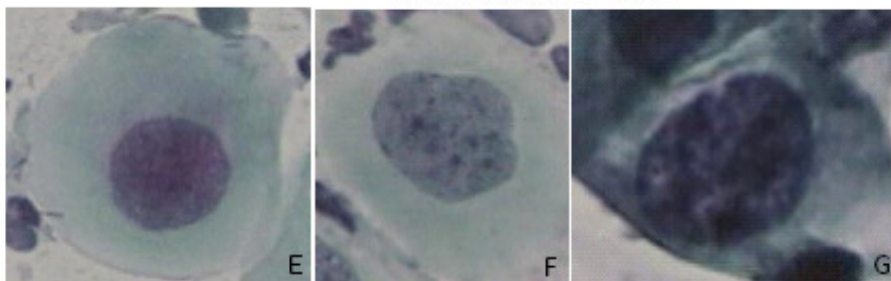
Table 3-1.: Database

NORMAL CELLS



A. Superficial squamous, B. Intermediate squamous, C. Parabasal, D. Columnar

ABNORMAL CELLS



E. Mild dysplasia, F. Moderate dysplasia, G. Severe dysplasia

Figure 3-5.: Cells from the database B.1

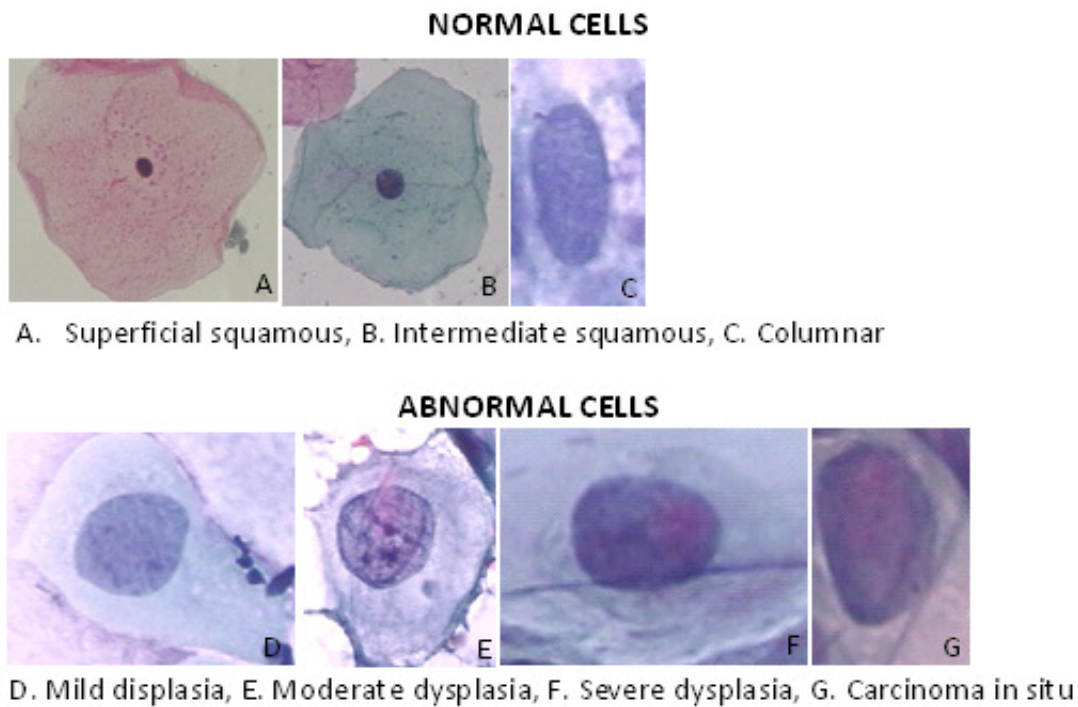


Figure 3-6.: Cells from the database B.2

3.1.4. Experimental setup

Two classification algorithms were assessed, namely KNN and SVM. Each classification model was tuned independently for its own particular set of parameters as follows: k-NN was assessed by varying the k nearest neighbors between 1 and 15 with increment steps of two. For SVM, two kernel types were evaluated, namely radial basis function (RBF) and polynomial functions. For the RBF kernel, the γ parameter was varied from 0,00 to 0,90 with increment steps of 0.10, while the polynomial kernel degree was set at 1, 2 and 3. The regularization parameter Complexity C, was varied between 1,0 to 20 with steps of 0.2. Evaluation was carried out when solving both two-class problems (namely normal and abnormal) and seven-class problems. A conventional 10-fold cross validation was performed for every parameter combination.

3.2. Results

The effects of different levels of complexity were evaluated.

Two class problem

Figure 3-7 shows the error percentage of the proposed method in the two-class problem for datasets B_1 and B_2. In the case of B_1, the best performance was achieved with the edge histogram descriptor and the KNN classifier with $k=9$. The proposed algorithm reached averaged accuracy levels of 72 % in this dataset. Database B_2 shows an accuracy of 83 %, using the three-descriptor integration (Color Layout, Scalable Color Descriptor and Edge Histogram Descriptor) and the KNN classifier. The average accuracy level in this database, in case of the Scalable Descriptor, the Edge histogram, and the three-descriptor integration, was 82 %, except for the Color Layout Descriptor which reached an accuracy of 72 %.

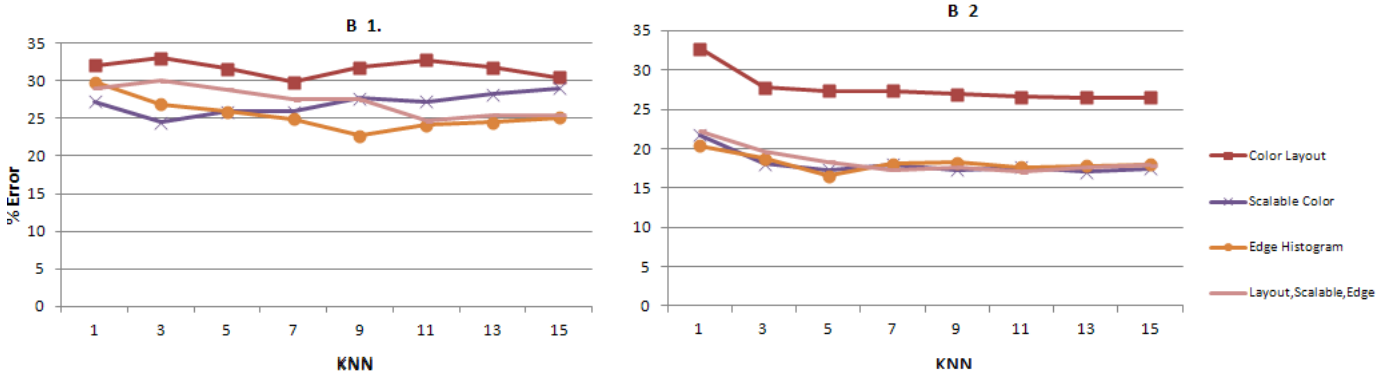


Figure 3-7.: Error classification for the B_1 and B_2 databases, in the two class problem.

Figure 3-8 shows the contour plot of the error measurement for SVM. Lighter regions correspond to the values of the parameters with best performance. Figure 3-8 a) shows the performance of the SVM with color-layout descriptor and polynomial kernel. The best performance was found when using a degree of 1.0 together with 4 complexities values (an accuracy of 75 %). For the SVM with radial base kernel (figure 3-8 b)), the best performance was achieved for $\gamma = 0,01$ and complexities values between 12.4 and 13.8 (a 76 % of accuracy). In figure 3-8 c), using scalable color and polynomial kernel, the best performance was achieved when applying degrees 3.0 and 7 and 16.8 complexities values (74 % accuracy). Additionally, when applying polynomial kernels, the best performance was reached for γ values between 0.60 to 0.90 (76 % accuracy), this is shown in figure 3-8 d). In case of the edge-histogram descriptor with polynomial kernel, the best performance was achieved by using degrees 3 and 1 with 2 complexities values (72 % accuracy); see figure 3-8 e). The three descriptor integrators, namely Color Layout, Scalable Color Descriptor and Edge Histogram Descriptor, obtained accuracy levels of 80 %.

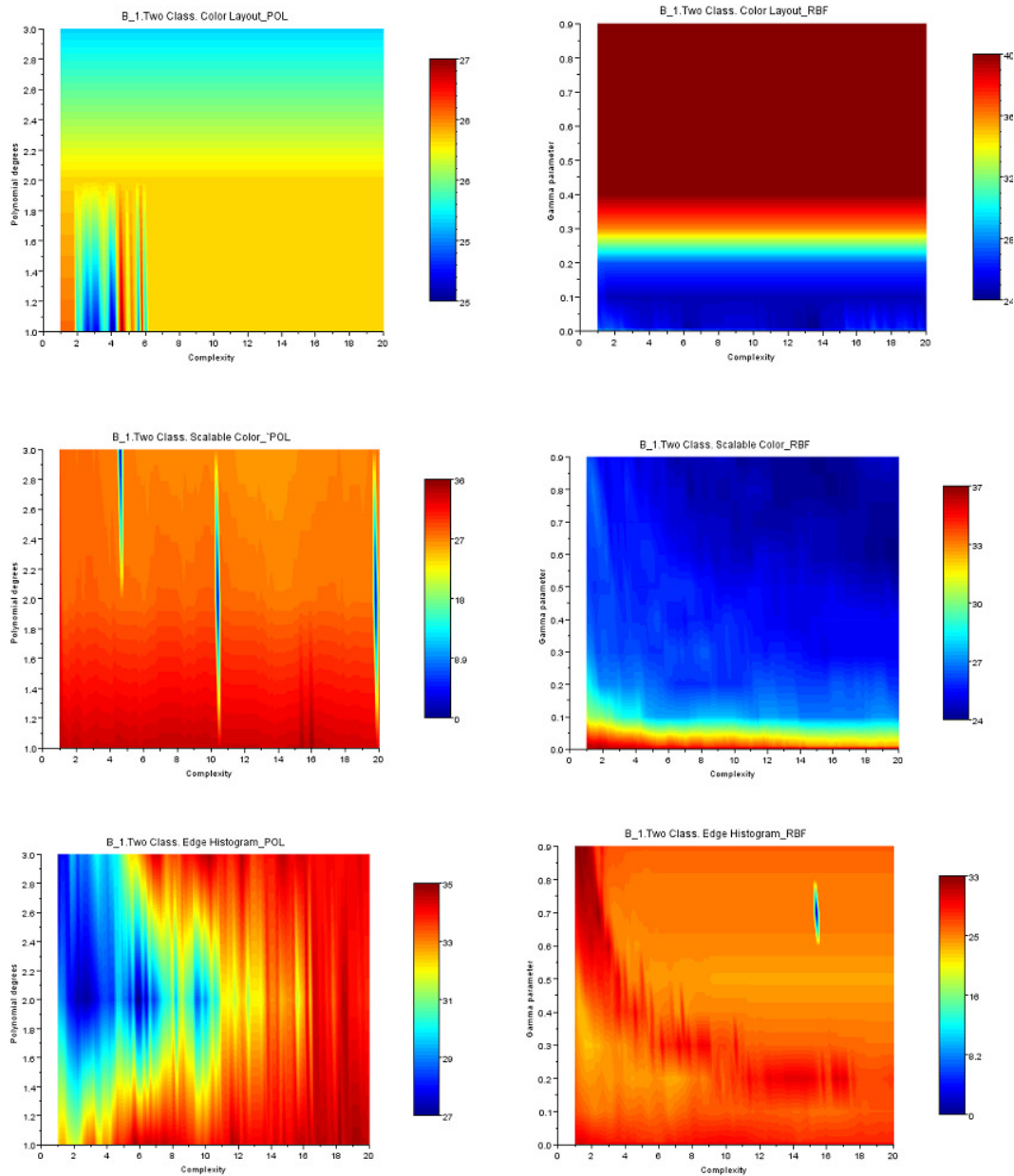


Figure 3-8.: Error classification for the B.1 database, together with the two-class problem. The panels are as follows: Color Layout (upper panel), Scalable Color (middle panel) and Edge Histogram descriptor (lower panel). Left plots correspond to the SVM polynomial function, whilst SVM radial-basis function kernel is depicted in the right panel. The blue color represents the parameters with the best performance for each case.

Seven class problem

Figure 3-9 shows the error percentage for the method proposed in the 7-class problem for datasets B.1 and B.2. In B.1, the best performance was achieved when using the three-descriptor integration by means of the KNN classifier with $k=1$ (accuracy 43%). In average, this database showed an accuracy of 42% when using three-descriptor integration.

The B.2 database exhibited a performance in terms of accuracy of 41% with the edge histogram descriptor, $k=13$, and with the three-descriptor integration (Color Layout, Scalable Color Descriptor and Edge Histogram), involving the KNN classifier ($k=11$). The average accuracy for the Scalable Descriptor, Edge histogram, and the three-descriptor integration was 34%.

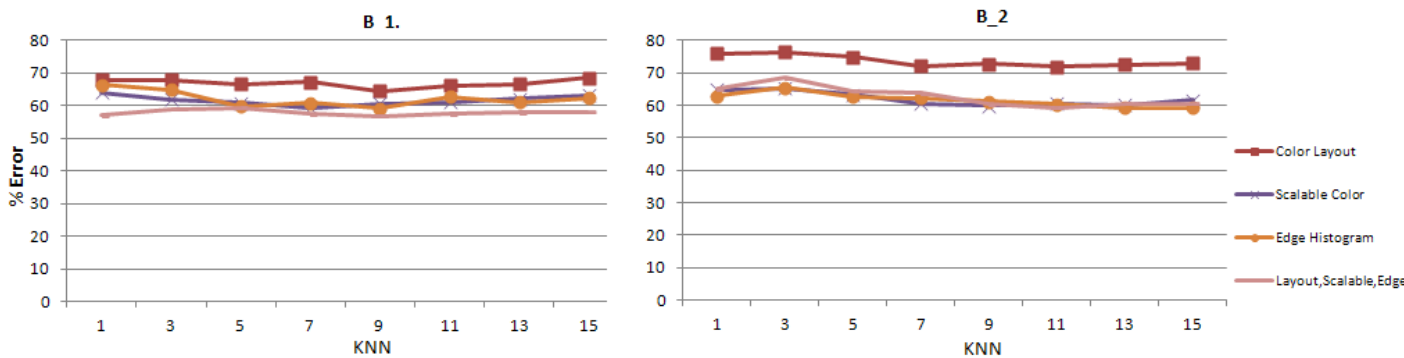


Figure 3-9.: Classification error reported in the B.1 and B.2 database, for the seven Class.

Figure 3-10 shows the contour plot of error measurement for SVM. Brighter regions correspond to the values of parameters with the best performance. Figure 3-10 a) show the performance of color layout and SVM when using the polynomial kernel, the best performance was reported when using degree 3 for every complexity value. In SVM with radial-base kernel, (see figure 3-10 b)), the best performance was reached for complexity values equal to 0.10.

Figure 3-10 c) shows the results with scalable color and polynomial kernel, the best performance was reported when using degrees 1.2 to 1.8, and 2 complexities values. As shown in figure 3-10 (d)), the best performance with radial-base kernel was reached for γ values between 9.0 to 11.9 and complexities values equal to 0.20.

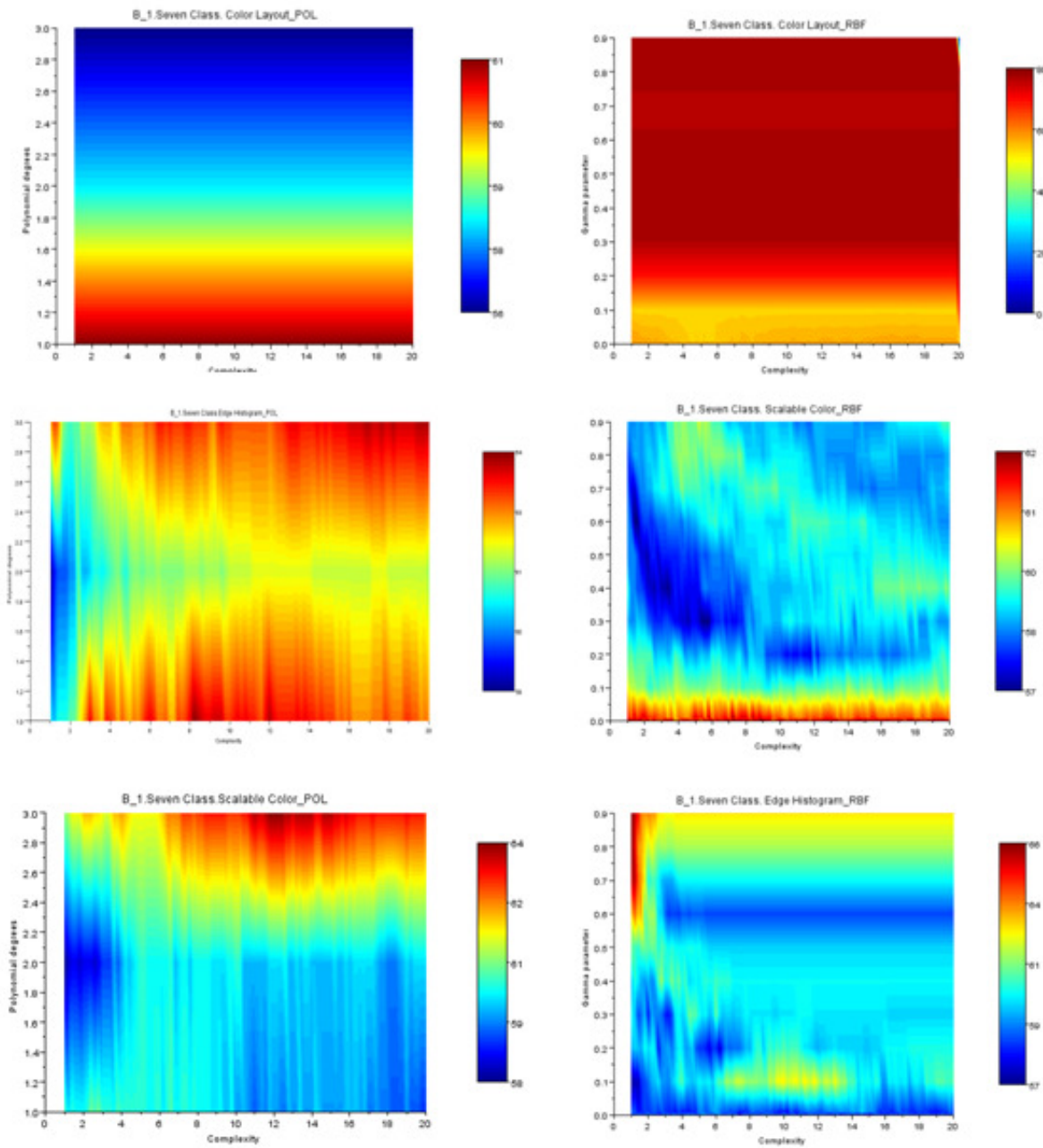


Figure 3-10.: Error classification for the B_1 database, together with the seven class problem. Left plots correspond to the SVM polynomial function, whilst SVM Radial-Basis function kernel is depicted on the right panels as follows: Color layout (upper panel), Scalable Color (middle panel) and Edge Histogram descriptor (lower panel). Clear blue color represents parameters with the best performance for each case.

3.3. Discussion

A new strategy for assisting the diagnosis of pap smears was proposed, implemented and evaluated. Unlike classical approaches, which aim to segment the cell since they use mostly anatomical features to describe the shape and therefore are quite prone to failure, the described strategy globally characterizes the object of interest using its very intrinsic features. The present investigation aimed to devise a strategy able to capture valuable information to perform a proper classification of conventional cytological images. A first evaluation was carried out in monolayer images and the performance obtained with the presented method reached a sensitivity of 85 %, a figure that demonstrates that the technique can be used in real scenarios .

The intentional fact of attempting to avoid any segmentation step led us to propose the use of global descriptors. This point is crucial towards implementing the type of techniques that can be robust to high levels of noise and to mount a functional screening process in Colombia. In particular we used a color histogram extracted from the HSV color space. Global descriptors included the spatial distribution of in the YCbCr color space, and the spatial distribution of five types of edges to classify normal and abnormal cells in two types of classification problems, two-classes or seven-classes. This strategy achieved a precision of 90 % and a sensitivity of 83 % in the two clas classification problems, similar to what was reported by Lee et al, 1991 who included a segmentation step [35]. Overall, precision reported in the literature is completely dependent on the segmentation quality and most works use commercial software such as CHAMP [8, 44], for instance Byriel et al., or Martin et al. set a large number of morphometric features that are then used to classify cervical cytology images either in two or seven classes [40–44, 57]. Yet their results outperform what was herein presented, they usually need to correct the obtained segmentations by hand whereby the the procedure results very expensive and a burden work. Discussion remains open since some evidence points out to demonstrate the relativity of those results, for instance Yang-Mao et a. have shown that CHAMP software was unable to provide satisfactory segmentation performance in different challenges, especially for the case of abnormal cervical cells, whose cytoplasm and nucleus contours are often very blurred [76].

Cells obtained by Pap staining the cytology smear are difficult to segment because of the large associated variability, namely the intrinsic biological tissue variability, dye conservation differences, fixation procedure differences, cell staining [74], and clustering [64], a set of factors that blur identification of the cell boundaries [18]. Most works in the literature have been performed in monolayer preparations, a technique that is not applicable in our country, as has been largely discussed before in this thesis, basically because this is very ex-

pensive. From a technical standpoint, another difficulty arises when one attempts to extend this technique to the conventional cytological image: the segmentation quality. This issue is illustrated in figure 3-11 that shows two images, the upper is a single isolated cell along with its luminance histogram and the lower panel shows a set of superimposed cells from a cytological image as well as its luminance histogram.

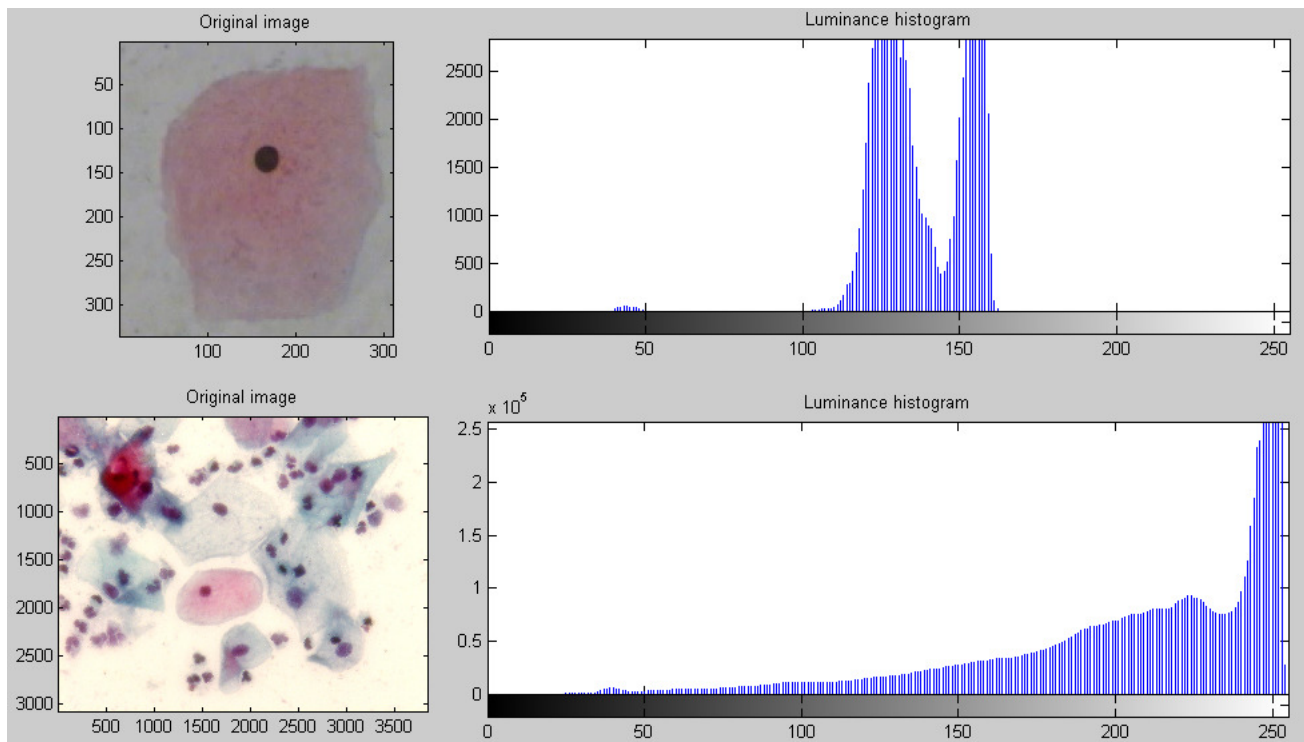


Figure 3-11.: Histogram for isolated and complete images

The two histograms illustrate the segmentation problem, while the upper one is very easy to segment since the bimodal distribution presents a clearly defined valley between the two peaks, The lower histogram is practically dominated by a large peak in the high frequency band. The segmentation problem is much more difficult in this case and practically the distinction between background and cels, so well defined in the upper histogram, is practically inexistent in the lower one.

4. Conclusions

In this work, an original method was proposed, combining MPEG-7 and classifier for discrimination of two-classes: normal and abnormal or seven classes. Instead of attempting to segment cells and its structural components, we proposed to characterize the whole internal cell structure using the well known MPEG-7 global descriptors, without segmentation techniques, avoiding the high dependency on image quality of the existent methods, the particular parameterization for each image and high associated costs.

Most image segmentation methods work well when the images are obtained with a good quality, an issue that is not a problem in most monolayer staining techniques but that constitutes a fundamental drawback for the conventional cytology. In particular, images obtained from specimen processed as conventional cytology, samples are frequently contaminated with different species, i.e. bacillus or fungi,, and the cytoplasm and nucleus contours are often blurred, especially for abnormal cervical cells. This work has presented on the contrary, an original method for classification of normal and abnormal cell, using structural information which is much less susceptible of being contaminated by any type of noise. This kind of Computer aided decision systems for cervical cytology would not require segmentation or large software resources. Results are very promising using two global descriptors of the MPEG-7 standard that were chosen because they were devised to globally capture semantic information, however implementation with more MPEG-7 descriptors and classification algorithms can definitely improve the classification performance and will be considered as the future work of the present investigation.

A. Clasificación de células escamosas en citología cervical usando descriptores MPEG-7

Camargo LH, Díaz G, Romero E. Clasificación de células escamosas en citología cervical usando descriptores MPEG-7. VII Seminario Internacional de Procesamiento y Análisis de Información Médica - SIPAIM 2011. Universidad Industrial de Santander, Bucaramanga, Colombia.

Clasificación de células escamosas en citología cervical usando descriptores MPEG-7

Luz Helena Camargo^{a,b}, Gloria Díaz^a, Eduardo Romero^{a,*}

^a*Grupo de investigación BioIngenium, Universidad Nacional de Colombia, Bogotá, Colombia*

^b*Universidad Distrital Francisco José de Caldas, Bogotá, Colombia*

Resumen

El cáncer de cérvix afecta a gran parte de la población femenina en países en vía de desarrollo. La citología convencional es la técnica más empleada para determinar si las células cervicales son precancerosas, sin embargo su lectura es costosa en tiempo de especialistas. Este artículo presenta una estrategia para la clasificación automática de células cervicouterinas usando descriptores globales de color y textura, la cual alcanzó una precisión del 90% y una sensibilidad del 83%. A diferencia de otras propuestas, el método desarrollado no requiere segmentación del núcleo y el citoplasma de la célula, lo que facilita su aplicación en herramientas de cribado.

Palabras clave: Célula Cervical, Análisis de imágenes citológicas, Cáncer cervical, Papanicolaou.

1. Introducción

El cáncer cervicouterino es una de las neoplasias de más fácil abordaje terapéutico y preventivo. Aún así, sigue siendo el segundo tipo de cáncer como causa de muerte en la población femenina[1]. Una de las principales razones para que esto suceda es la insuficiencia de sistemas de detección precoz que permitan el diagnóstico cuando la enfermedad aún tiene buen pronóstico. La técnica de detección precoz más empleada es la citología cervical, que consiste en la identificación al microscopio de células precancerosas y cancerosas en muestras citológicas obtenidas del cuello uterino. Aunque este método es considerado como un método eficaz, sencillo, rápido

*Correspondencia: Eduardo Romero, Carrera 30 45-03 Edificio 471 Medicina, Piso 1, Tel/Fax (1) 3165000 Ext. 15025, <http://www.bioingenium.unal.edu.co>

y económico, no sólo para detectar tempranamente el cáncer del cuello del útero sino también infecciones por microorganismos, su aplicación requiere tiempo y experticia del cito tecnólogo en el diagnóstico preliminar y del cito patólogo en el diagnóstico definitivo. Por lo anterior, se han realizado esfuerzos importantes en el desarrollo de métodos de apoyo al diagnóstico citológico que faciliten y/o mejoren la detección de anormalidades. Actualmente existen en el mercado diferentes métodos como AutoPap Primary Screening System (aprobado por la FDA), ThinPrep Imaging System (Cytoc) y Path (Molecular Diagnostic). Estas herramientas se caracterizan por identificar regiones con probabilidad de contener células patológicas, pero no realizan una identificación real de las células en cada una de las tipologías patológicas posibles, adicionalmente su implementación tiene costos elevados, que reducen la posibilidad de ser usados en sistemas de cribado en países en desarrollo[2]. Otros métodos para la clasificación de imágenes de células de citopatología han sido propuestos [3, 4], los cuales requieren una segmentación precisa de cada célula, para a partir de ella calcular un conjunto grande de características morfométricas (áreas, perímetros, redondez, elongación de núcleo y citoplasma, entre otras), aunque su desempeño es adecuado, el depender del resultado de una segmentación previa, hace que su aplicación tenga limitaciones debido a la diversidad de estructuras contenidas en las imágenes, la intensa variación del fondo, agrupación y alto grado de superposición que dificultan la identificación de los límites de los núcleos de las células[5]. Nanni y otros [6], proponen el uso de un descriptor de texturas conocido como patrones locales binarios para clasificar células sin requerir segmentación precisa, los resultados reportados para el área bajo la curva ROC, son de entre 0.72 y 0.88. Este trabajo, presenta un enfoque sencillo para la clasificación de células de citología cervical que hace uso de los descriptores MPEG-7, para describir una región conteniendo la célula pero que no requiere una segmentación previa de núcleo y citoplasma. Esto permite, no sólo reducir el costo computacional, sino que además reduce el tiempo empleado por los cito tecnólogos para revisar manualmente la exactitud de la segmentación resultante.

2. Método Propuesto

2.1. Descriptores de características visuales

En este trabajo se consideraron como principales características para clasificar correctamente células normales y precancerosas la distribución del color y la textura. Estas características permiten implícitamente capturar información relacionada con la relación núcleo/citoplasma, la presencia de

variaciones en forma y tamaño en el núcleo asociadas a células patológicas y otras características morfológicas citoplasmáticas de cada componente celular. Tres descriptores de características propuestos en el estándar MPEG-7 fueron extraídos para cada región: diseño o distribución espacial del color (color layout descriptor), color escalable e histograma de bordes. Adicionalmente el efecto de combinar estos descriptores fue también evaluada (distribución espacial del color con histograma de bordes y color escalable con histograma de bordes), resultando en cinco descriptores evaluados. Estos descriptores consideran características de color [7] y de textura [8], las cuales son consideradas como relevantes en la clasificación de malignidad de las células cervicouterinas[9].

Diseño o distribución espacial del color

Este descriptor ofrece una manera compacta de representar la distribución espacial de las características de color de una imagen en el dominio frecuencial. La representación se basa en los coeficientes de la Transformada Discreta del Coseno (DCT) sobre los valores de las componentes Y, Cb y Cr de la imagen. El cálculo de este descriptor es ilustrado en la figura 1. El proceso inicia con la partición de la imagen en 64 bloques (cada bloque de tamaño $W/8 \times H/8$, donde W y H representan la anchura y la altura respectivamente de la imagen). Cada bloque es representado por la media de los colores contenidos en él, obteniendo así una imagen de 8×8 de apariencia borrosa que describe la apariencia general de la distribución del color en la imagen. Luego, la imagen es transformada al espacio de color YCbCr y, los primeros 64 coeficientes, resultantes de transformar las matrices correspondientes a cada uno de los componentes de color usando la DCT, son usados para describir la imagen, obteniendo así un descriptor de 192 características[7].

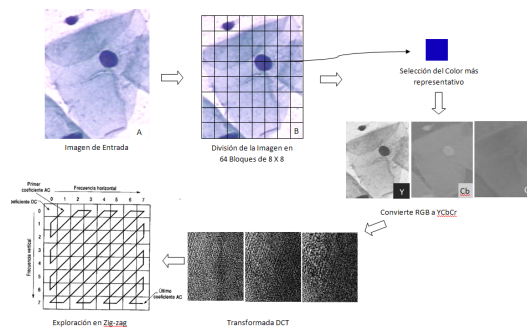


Figura 1: Proceso de extracción del descriptor de distribución de color

Color escalable

Este descriptor representa la distribución de probabilidad de los colores en la imagen mediante un histograma de color HSV codificado usando la transformada de Haar. El concepto de escalabilidad se refiere al número de bins usados para calcular el histograma. Como se ilustra en la figura 2, inicialmente un histograma de 256 bins es calculado (incluyendo 16 niveles en H, 4 en S y 4 en V), este histograma es luego codificado usando la transformada de Haar, esta transformada permite codificar el histograma con la mitad de bins en cada iteración, por la aplicación sucesiva de operaciones de suma y diferencia que se relacionan con filtros pasa bajos y pasa altos del histograma[7]. En este trabajo el descriptor color escalable fue reducido a 64 bins con 8 niveles para el componente H, 2 para S y 4 para V.

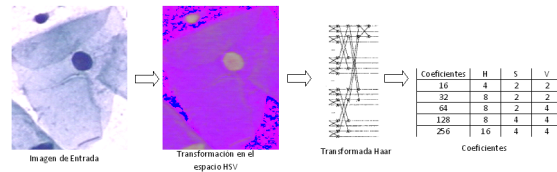


Figura 2: Proceso de extracción del descriptor color escalable

Histograma de bordes

Es un descriptor de textura que representa la distribución espacial de cinco tipos de bordes en regiones locales de la imagen, cuatro bordes direccionales generados a 0° (vertical), 45° , 90° (horizontal) y 135° y un borde no direccional. La imagen es dividida en 16 sub-imágenes en las que se genera un histograma de borde local con 5 bins, como se observa en la figura 3, obteniendo finalmente 80 histogramas como resultado de las 16 subimágenes por 5 bins. La información de borde es extraída con operadores de detección de bordes [8].

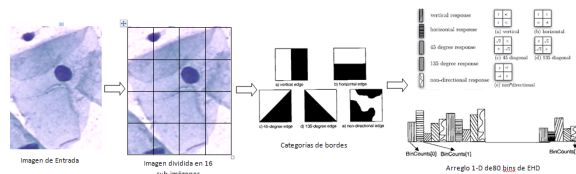


Figura 3: Descriptor histograma de borde, distingue cinco tipos de bordes y almacena los valores en un arreglo 1D

2.2. Clasificación

Para la clasificación de cada célula, fue empleado un modelo de clasificación por vecindad conocido como k vecinos más cercanos (Knn), el cual se basa en la suposición de que los prototipos más cercanos tienen una probabilidad a posteriori similar de pertenecer a la misma clase, para ello se emplea la métrica de distancia euclidiana. Este método es sencillo de implementar y no depende de muchos parámetros para su entrenamiento y se reporta buen desempeño en tareas similares [10].

3. Resultados Experimentales

3.1. Base de Datos

Dos bases de datos creadas en el Departamento de patología del Hospital universitario Herlev¹ fueron usadas para evaluar el método propuesto. Estas imágenes tienen una resolución de 0.201um/pixel. Dos citotecnólogos especializados clasificaron manualmente cada una de las células y las muestras difíciles fueron examinadas también por un médico especialista[3, 11, 12]. La primera base de datos (BD_1) contiene 200 imágenes normales y 306 anormales[3], ejemplos de estas imágenes son mostrados en la figura 4. La segunda base de datos (BD_2) contiene 917 células, 242 normales y 675 anormales [11, 12], Ejemplos de imágenes de esta base de datos son mostrados en la figura 5.

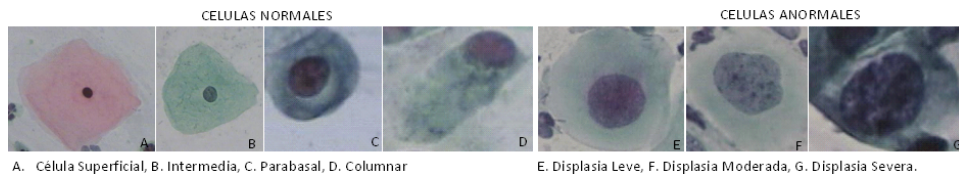


Figura 4: Ejemplos de células de la primera base de datos

El método propuesto fue evaluado usando un modelo de validación cruzada de 10 particiones (10-fold). El desempeño de cada descriptor propuesto y la diferentes combinaciones de éstos fue evaluado independientemente. Por otro lado, el valor de K , para el algoritmo de clasificación KNN fue ajustado usando valores impares entre 1 y 15. La figuras 6 y 7 presentan la precisión y sensibilidad, respectivamente, de cada uno de los descriptores para los diferentes valores de K en las dos bases de datos evaluadas.

¹<http://fuzzy.iau.dtu.dk/downloads/smear2005/>

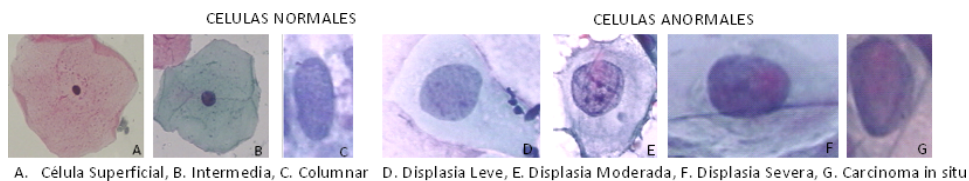


Figura 5: Ejemplos de células de la segunda base de datos

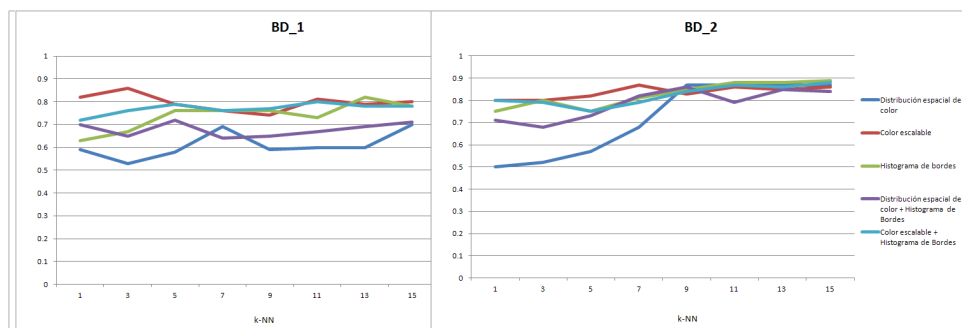


Figura 6: Precisión para diferentes valores de K en cada una de las bases de datos evaluadas

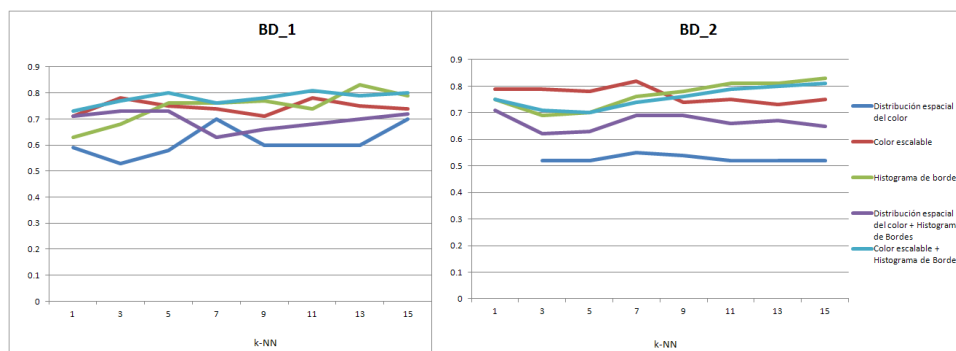


Figura 7: Sensibilidad para diferentes valores de K en cada una de las bases de datos evaluadas

En promedio el mejor desempeño para la base de datos 1 fue obtenido por el descriptor histograma de bordes con un modelo de clasificación 13-NN, con una precisión de 0,82 y una sensibilidad de 0,83. Similarmente, el mejor desempeño para la base de datos 2 fue obtenido por el mismo descriptor histograma de bordes pero esta vez con un modelo de clasificación 15-NN, en este caso la precisión fue de 0,89 y la sensibilidad de 0,83. Resultados similares fueron obtenidos cuando este descriptor fue combinado con

el descriptor color escalable con una sensibilidad de 0,81 para las dos bases de datos y una precisión de 0,8 y 0,88 para la primera y segunda base de datos, respectivamente. Cabe resaltar que el valor del parámetro K no resulta ser definitivo en el desempeño, salvo para el caso histograma de bordes en la base de datos 2, caso en el cual se observa una mejora significativa del desempeño cuando K es mayor a 11.

4. Conclusiones y Trabajo Futuro

Este trabajo presenta un modelo de clasificación de células cervicouterinas basado en el uso de descriptores de color y de textura del estándar MPEG-7. La principal ventaja del modelo propuesto es que éste no requiere un proceso de segmentación previa, lo que hace factible su utilización en sistemas masivos de cribado. El modelo propuesto fue evaluado en bases de datos de dominio público. De los resultados se puede observar que la información de texturas descrita por un histograma de bordes presenta mayor poder de discriminación que los demás descriptores. Sin embargo, el uso de modelos de clasificación más robustos puede mejorar estos resultados. Así mismo, el uso de otros descriptores o una combinación lineal de estos deberá ser evaluada para determinar el aporte real de cada uno de ellos en la caracterización de este tipo de células.

El trabajo futuro se enfocará en el uso de estos descriptores para realizar una clasificación de las células en sus diferentes tipologías como son: superficial, intermedia y columnar para las células normales y displasia leve, moderada, severa y Carcinoma in situ para las células anormales. Adicionalmente, se espera que este modelo de clasificación pueda ser usado en un sistema completo que permita la detección y clasificación de células anormales en imágenes de citología cervicouterina.

Referencias

- [1] O. M. de la Salud, Control integral del cáncer cervicouterino. Guía de prácticas esenciales, 2007.
- [2] J. Giménez, P. Sanz, J. Torres, C. Hörndler, E. Urbiola, Evaluación de dispositivos automatizados para diagnóstico citológico en la prevención del cáncer de cérvix, *Rev. Esp. Patol.* 35 (2002) 301–314.
- [3] J. Byriel, Neuro-fuzzy classification of cell in cervical smears, Master's thesis, Technical University of Denmark (1999).

- [4] Y. Marinakis, G. Dounias, J. Jantzen, Pap smear diagnosis using a hybrid intelligent scheme focusing on genetic algorithm based feature selection and nearest neighbor classification, *Computers in Biology and Medicine* 39 (2009) 69–78.
- [5] C. Duanggate, B. Uyyanonvara, T. Koanantakul, A review of image analysis and pattern classification techniques for automatic pap smear screening process, in: *The 2008 International Conference on Embedded Systems and Intelligent Technology February 27-29, 2008*, Grand Mercure Fortune Hotel, Bangkok, Thailand, 2008.
- [6] L. Nanni, A. Lumini, S. Brahnam, Local binary patterns variants as texture descriptors for medical image analysis, *Artificial Intelligence in Medicine* 49 (2010) 117–125.
- [7] L. Cieplinski, Mpeg-7 color descriptors and their applications, in: W. Skarbek (Ed.): *CAIP 2001, LNCS 2124*, pp. 11-20. Springer-Verlag Berlin Heidelberg, 2001.
- [8] X. F. Z. Y-J, Evaluation and comparison of texture descriptors proposed in mpeg-7, *J. Vis. Commun. Image R* 17 (2006) 701–716.
- [9] L. Na, L. J. Xin., Automatic color image segmentation based on anisotropic diffusion and the application in cancer cell segmentation, in: *The 1st International Conference on Bioinformatics and Biomedical Engineering, Vol. 992-994*, 2007.
- [10] F. Narvaez, Sistema de anotación para apoyo en el seguimiento y diagnóstico de cáncer de seno, Master’s thesis, Universidad Nacional de Colombia (2010).
- [11] E. Martin, Pap-smear classification, Master’s thesis, Technical University of Denmark (2003).
- [12] J. Norup, Classification of pap-smear data by transductive neuro-fuzzy methods, Master’s thesis, Technical University of Denmark (2005).

B. Pap smear cell image classification using global MPEG-7 descriptors

Camargo LH, Díaz G, Romero E. Pap Smear Cell Image Classification Using Global MPEG-7 Descriptors. 11th European Congress on Telepathology and 5th International Congress on Virtual Microscopy. 2012, Venice, Italy.

Pap smear cell image classification using global MPEG-7 descriptors

Abstract

Background

Although cervical cancer is fully curable if diagnosis is early achieved, it is still the most frequent female fatal disease worldwide, the second death cause in female population. On the other hand, Papanicolaou smear test is the screening method for detecting abnormalities in the uterine cervix cells, including the precancerous cells. Provided that this test is extremely labor intensive and reader dependant, automated screening has been pursued during last decades, resulting in some commercial developments that identify suspicious regions i.e. regions with the possibility of containing malignant cells, but without an accurate identification of pathological cells. As a consequence, cell classification is still an open research problem, usually approached by characterizing morphological features that are highly dependent on the accurate segmentation of the boundaries. In this paper we propose a cell classification method that instead of attempting to segment the cell cytoplasm and nucleus, it characterizes the very inner cell relationships using global features and a standard classifier, the k nearest neighbour, that learns a particular data partition.

Findings

The cell classification approach is carried out using color and texture MPEG-7 descriptors, specifically the color layout, color scalable and Edge Histogram descriptors. The proposed approach raises a mean sensitivity of 84 % and specificity between 82 % and 89 %, which is really promising since no image preprocessing was carried out.

Conclusions

The most relevant feature was the edge histogram, with which the best results were reported, whilst combination of this feature with the color scalable feature reported the poorer performance. On the other hand, modification of the learning parameters, did not significantly change results w.r.t. the classification task.

Keywords: Papanicolaou; MPEG-7; cervical cancer; color layout descriptor; edge histogram descriptor; scalable color descriptor

Background

Several strategies have been previously applied for classifying cervical cytology cells, all pursuing a nucleus segmentation. Sanchez sets regions [1] using a simple threshold [2], a procedure broadly adapted to different techniques: a local adaptive segmentation nuclei procedure [3], seed growing (Romberg y col. en [4]), mathematical morphology [5], a Hough transform [6], and active contours [7]. Jantzen and Dounias propose several cell features as morphometric descriptors, including the nucleus and cytoplasm areas, nucleus/cytoplasm proportion, nucleus and cytoplasm brightnesses, smaller and larger nucleus/cytoplasm diameters, nucleus and cytoplasm roundness, nucleus and cytoplasm perimeters, nucleus position, nucleus/cytoplasm maxima and minima. Nevertheless, these morphometric characteristics require a previous accurate segmentation, hardly achieved by human intervention using commercial software such as CHAMP (Cytology and Histology Modular Analysis Package, Aarhus, Denmark) of DIMAC (Digital Image Company) [8,9].

Methods

Rather than attempting to detect some of the previously reported morphometric features, the present investigation used two global MPEG-7 color descriptors, Color Layout and Scalable Color, and one texture descriptor, the Edge Histogram descriptor, as the representation space and two supervised classification algorithms (SVM and KNN) that divide the different classes.

Classification based on global MPEG-7 descriptors

The cell classification approach is carried out using color and texture MPEG-7 descriptors, thereby attempting to capture information related with the particular color spatial location and global color distribution of both the nucleus and cytoplasm. The texture descriptor stands for the particular borders of both nucleus and cytoplasm and their intrinsic relationships. These global characteristics are not evaluating the classical morphometric features, but they are using nucleus and cytoplasm visual primitives as discriminant factors.

Color layout

This descriptor, typically used in the YCrCb color space, captures the spatial color distribution in an image or an arbitrary region. Basically, the color layout descriptor uses representative colors on a grid, followed by a Discrete Cosine Transform (DCT) and an encoding of the resulting coefficients. The feature extraction process consists of two parts; grid based representative color selection and DCT transform with quantization. Specifically, an input image is divided into 64 blocks, their average colors are derived and transformed into

a series of coefficients by performing a conventional DCT. A few low-frequency coefficients are selected using zigzag scanning and quantized to form a Color Layout Descriptor [10].

Scalable color

This descriptor is extracted from a color image histogram in the hue-saturation-value (HSV) color space. This histogram, constructed with fixed color space quantization, is projected into a set of Haar bases so that the obtained coefficients constitute a scalable color representation. The histogram values must be normalized and non linearly mapped into a 4-bit integer representation, giving higher weight to small values. The Haar transform is applied then to this histogram version with two basic operators: sum and difference bin neighbor, decomposing the histogram into low and high frequency subbands [10].

Edge histogram

This descriptor captures the spatial edge distribution, a very useful feature for image matching, even though the underlying texture may not be homogeneous. A given image is first sub-divided into sub-images and local edge histograms, for each of these sub-images, are computed. Edges are then coarsely grouped into five categories: vertical, horizontal, 45 diagonal, 135 diagonal, and isotropic (nonorientation specific). Thus, each local histogram has five bins corresponding to the above five categories. The image partitioned into 16 sub-images results in 80 bins. These bins are nonuniformly quantized using 3 bits/bin, resulting in a descriptor with size of 240 bits [10].

Classification models

The classification method used a classical K-Nearest Neighbor algorithm and a Support Vector Machine. The proposed approach was evaluated under a 10-fold experimental setup.

The k-NN decision rule

The k-nearest neighbor method is an intuitive method that classifies unlabeled samples based on their similarity with samples in the training set. Given the knowledge of N prototype features (vectors of dimension k) and their correct classification into M classes, the k-NN rule assigns an unclassified pattern to the class that is most heavily represented among its k neighbors in the pattern space, under some appropriate metric. In this work euclidean distance was used.

The SVM algorithm

A support vector machine (SVM) is a classification model that finds an optimal separating hyperplane that discriminates two classes. A SVM is a linear discriminator, however it can

perform non-linear discriminations thanks to the fact that this is a kernel method. In this work, it is used a SVM version that uses sequential minimal optimization algorithm. The multi-class classification problem is solved using a one vs. all strategy: a binary classifier for each class by labeling the class samples as positive examples and other class samples as negative ones. The final decision is set to the class having the largest decision function among all classes.

Results and discussion

Database

Two databases composed of images with single cells, from the Herlev University Hospital, Denmark, were used (<http://fuzzy.iaa.dtu.dk/downloads>). Skilled cyto-technicians and doctors manually classified each cell in 2 classes: abnormal and normal and then subclassified into seven classes. Each cell was examined by two cyto-technicians, and difficult samples also by a doctor. In case of disagreement, the sample was simply discarded [8,9]. Finally there are two database: B_1 data contains 500 cells with the following distribution:

1. Normal: columnar epithelial, parabasal squamous epithelial, intermediate squamous epithelial, superficial squamous epithelial.
2. Abnormal: mild squamous non-keratinizing dysplasia, moderate squamous non-keratinizing dysplasia, severe squamous non-keratinizing dysplasia.

The B_2 data contains 917 cells with the following distribution:

1. Normal: superficial squamous epithelial, intermediate squamous epithelial, columnar epithelial.
2. Abnormal: mild squamous non-keratinizing dysplasia, moderate squamous non-keratinizing dysplasia, severe squamous non-keratinizing dysplasia, squamous cell carcinoma in situ intermediate.

Experimental setup

Two classification algorithms were assessed (KNN and SVM). Each classification model was tuned independently for its own particular set of parameters as follows: k-NN was assessed by varying the k nearest neighbors between 1 and 15 with increment steps of two. SVM used two kernel types were evaluated; radial basis function (RBF) and polynomial functions. For the RBF kernel, the γ parameter was varied from 0,00 to 0,90 with increment steps of 0.10, while the polynomial kernel degree was set at 1, 2 and 3. The regularization parameter Complexity 1,0. Evaluation was carried out with both 2-class (normal and abnormal) problems. A conventional 10-fold cross validation was performed for every parameter combination.

Results

The effects of different levels of complexity were evaluated, but do not show important variations. Figure **B-1** shows the percentage error for the method proposed in the 2-class problem for B_1 datasets. Better performance was achieved by the KNN classifier with k 15. The B_2 database shows the performance of SVM with radial base kernel, for values of 0.01 (figure **B-2**). In both database these results were obtained with the edge histogram descriptor.

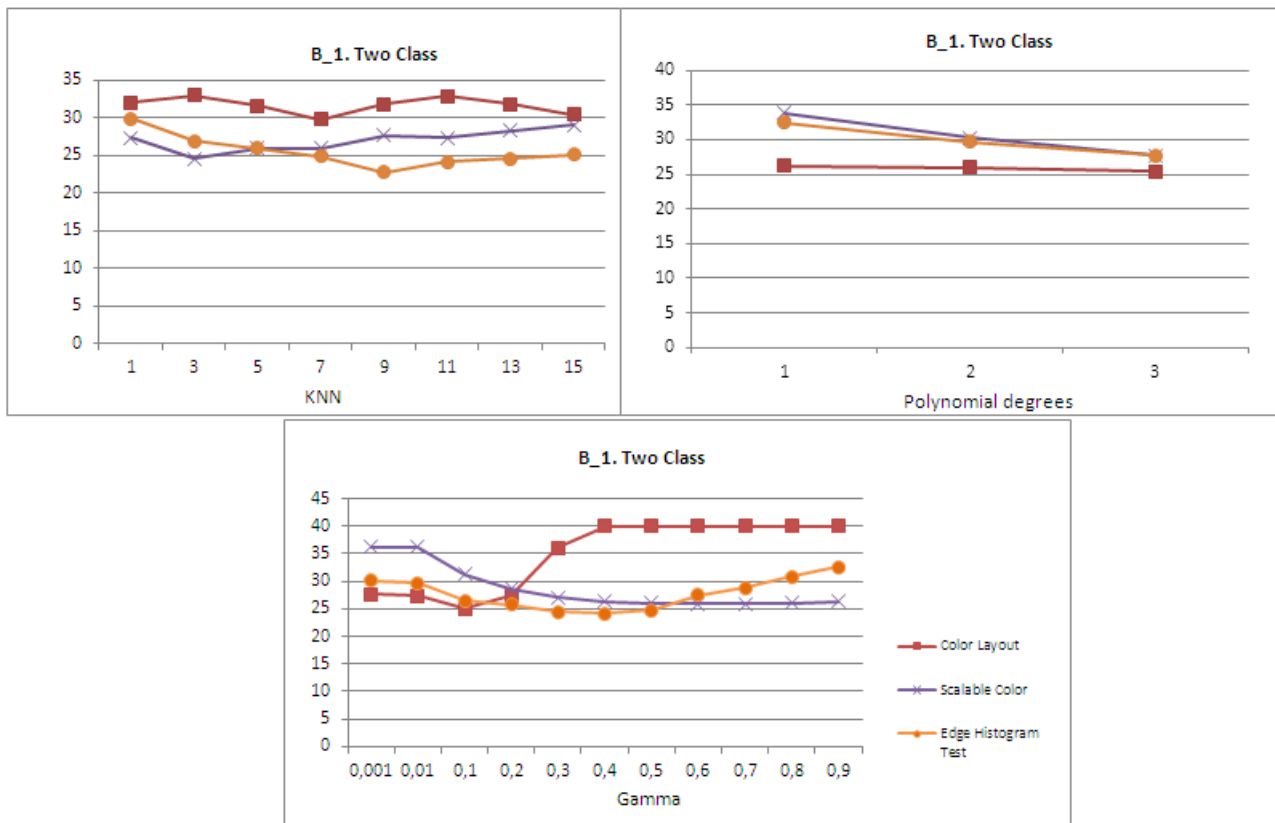


Figure B-1.: Error for different methods in B_1

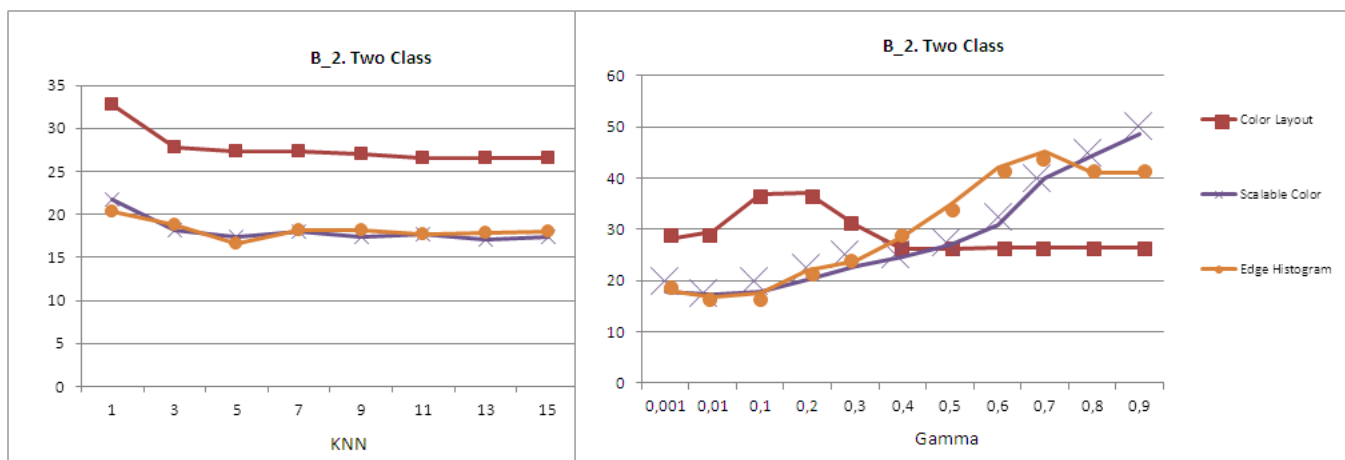


Figure B-2.: Error for different methods in B_2

Discussion

A new strategy for assisting the diagnosis of pap smears which requires no segmentation was proposed, implemented and evaluated. Instead of attempting to segment cells and its structural components, we propose to characterize the internal cell structure using well known MPEG-7 global descriptors. Results are very promising, however evaluation of other specific features and classification algorithms can improve the classification performance.

Conclusions

This work presented an original method for discrimination of each class: normal and abnormal. Future work is to evaluate the classification of the seven class.

Competing interests

The authors declare that they have no competing interests.

Authors' contributions

All authors read and approved the final manuscript.

Acknowledgements

This study was supported partially by the Research Headquarters Address Bogota - DIB, Universidad Nacional de Colombia, QUIPU code 202010011343, HERMES code: 8521.

References

- [1] Sánchez, A. Aplicaciones de la visión artificial y la biometría informática, Dykinson S.L.,(Madrid), U. R. J. C. (ed.) 2005.
- [2] Passariello. Imágenes médicas. Adquisición, Analisis, procesamiento e Interpretación Ediciones de la Universidad Simón Bolívar, 1995.
- [3] Li, Z. and Najarian, K. Biomedical image segmentation based on shape stability. IEEE International Conference on Image Processing, 2007, 281-284.
- [4] Romberg, J.; Akram, W. and Gamiz, J. Image segmentation using region growing *http://www.owlnet.rice.edu/elec539/Projects97/WDEKknow/index.html*, 1997.
- [5] Lassouaoui, N. and Hammami, L. Genetic algoritms and multifractal segmentation of cervical images. Proc. Of. IEEE-EURASIP 7th International Symposium on Signal Processing and its Applications, 2003.
- [6] Garrido, N. P. Applying deformable templates for cell image segmentation Pattern Recognition Pattern Recognition, 2000, 33, 821-832.
- [7] Bamford, P. and Lovell, B. Method for Accurate Unsupervised Cell Nucleus Segmentation. Engineering in Medicine and Biology Society, 2001. Proceedings of the 23rd Annual International Conference of the IEEE, 2001.
- [8] Byriel J: Neuro-fuzzy classification of cell in cervical smears, Master's thesis, Technical University of Denmark 1999.
- [9] Martin E: Pap-smear classification, Master's thesis, Technical University of Denmark 2003.
- [10] Manjunath, B.; Ohm, J.; Vasudevan, V. and Yamada, A. Color and Texture Descriptors. IEEE Transactions on circuits and systems for video technology, 2001, 11, 703-712.

Bibliografía

- [1] *Fusing MPEG-7 visual descriptors for image classification*, ICANN 05, Warsaw, Poland, September 2005.
- [2] P. Bamford and B. Lovell, *A water immersion algorithm for cytological image segmentation*, Proceedings of the APRS Image Segmentation Workshop, Sydney, 1996., pp. 75–79.
- [3] ———, *Unsupervised cell nucleus segmentation with active contours*, Signal Processing **71 (2)** (1998), 203–213.
- [4] ———, *Method for Accurate Unsupervised Cell Nucleus Segmentation*, Engineering in Medicine and Biology Society, 2001. Proceedings of the 23rd Annual International Conference of the IEEE, 2001.
- [5] G. Birdsong, *Automated Screening of Cervical Cytology Specimens*, Human Pathology **27 (5)** (1996), 468–481.
- [6] S. Broder, *From the National Institutes of Health. Rapid Communication - The Bethesda System for Reporting Cervical / Vaginal Cytologic Diagnoses - Report of the 1991 Bethesda Workshop*, JAMA **267** (1992), 1892.
- [7] C. Bustaraca and S. Huertas, *Reconocimiento de Patronos en dos dimensiones*, Tesis, 1994.
- [8] J. Byriel, *Neuro-Fuzzy Classification of cell in Cervical Smears*, Master's thesis, Technical University of Denmark, 1999.
- [9] C. Chang, M. Lin, H. Harn, Y. Harn, C. Chen, K. Tsai, and C. Hwang, *Automatic Segmentation of Abnormal Cell Nuclei from Microscopic Image Analysis for Cervical Cancer Screening*, Proceedings of the 2009 IEEE 3rd International Conference on Nano/Molecular Medicine and Engineering October 18-21, Tainan, Taiwan, 2009.
- [10] L. Cieplinski, *MPEG-7 Color Descriptors and Their Applications*, W. Skarbek (Ed.): CAIP 2001, LNCS 2124, pp. 11-20. Springer-Verlag Berlin Heidelberg, 2001.
- [11] Liga Colombiana contra el Cáncer, *Normas para la garantía de la calidad en citología cervico-uterina*, 2005.

-
- [12] Cytoc Corporation, *Operation Summary and Clinical Information, ThinPrep Imaging System*, Patent no. 86093-001. rev. e00. marlborough, mass: Cytoc corporation;, 2006.
- [13] T.M. Hart P.E. Cover, *Nearest neighbor pattern classification*, IEEE Transactions on Inf. Theo. **IT -13** (1967), 21–27.
- [14] Organización Mundial de la Salud, *Control integral del cáncer cervicouterino. Guía de prácticas esenciales*, 2007.
- [15] Organización Panamericana de la Salud., *Manual de Procedimientos del Laboratorio de Citología*, Organización Panamericana de la Salud., 2002.
- [16] P. Disaia and W. Creasman, *Oncología Ginecológica Clínica*, Elsevier Science, 2002.
- [17] Brownrigg DRK, *The weighted median filter*, Communications of the Association of Computer Manufacturers (CACM). **27** (1984), 807–818.
- [18] C. Duangate, B. Uyyanonvara, and T. Koanantakul, *A Review of Image Analysis and Pattern Classification Techniques for Automatic Pap Smear Screening Process*, The 2008 International Conference on Embedded Systems and Intelligent Technology February 27-29, 2008, Grand Mercure Fortune Hotel, Bangkok, Thailand, 2008.
- [19] Mathias Eitz, Kristian Hildebrand, Tamy Boubekeur, and Marc Alexa, *An evaluation of descriptors for large-scale image retrieval from sketched features lines*, Computers & Graphics **348** (2010), 482–498.
- [20] Jose Esqueda and Luis Palafox, *Fundamentos de Procesamiento de imágenes*, Universidad Autónoma de Baja California, 2005.
- [21] J Ferlay, F Bray, P Pisani, and D Parkin, *GLOBOCAN 2002: Cancer Incidence, Mortality and Prevalence Worldwide*, Tech. report, Lyon: IARC; 2004. Report No. :, 2002.
- [22] J Ferlay, H Shin, F Bray, D Forman, C Mathers, and D Parkin, *Estimates of worldwide burden of cancer in 2008: GLOBOCAN 2008*, Tech. report, Int J Cancer 2010, Jun 17., 2010.
- [23] A. & de la Blanca N. P. Garrido, *Applying deformable templates for cell image segmentation Pattern Recognition*, Pattern Recognition **33** (2000), 821–832.
- [24] J. Giménez, P. Sanz, J. Torres, C. Hörndler, and E. Urbiola, *Evaluación de dispositivos automatizados para diagnóstico citológico en la prevención del cáncer de cérvix*, Rev. Esp. Patol. **35** (2002), 301–314.
- [25] F Glover, *Future Paths for Integer Programming and links to artificial intelligence*, Computers & Operational Research, 1986.

- [26] R. Gonzalez and R. Woods, *Tratamiento Digital de Imágenes*, Reading, MA., 1996.
- [27] INC Grupo Vigilancia Epidemiológica del Cáncer, *Grupo Vigilancia Epidemiológica del Cáncer*, INC, Tech. report, Instituto Nacional de Cancerología, 2008.
- [28] J. Guaithero, de M. Toro, and de S. López, *Células Metaplásicas Inmaduras Atípicas como predictor de Lesión Intraepitelial Escamosa de Alto Grado. A propósito de un caso*, Revista de la Facultad de Farmacia **42** (2001), 63–66.
- [29] J. Jantzen and G. Dounias., *Analysis of Pap- Smear image Date*, Proc. of the Nature-Inspired Smart information System. 2nd Annual Symposium, NISIS, 2006.
- [30] A. Kale, S. Aksoy, and S. Onder, *Cell Nuclei Segmentation in Pap Smear Test Images*, Signal Processing and Communications Applications Conference, 2009., 2009.
- [31] E. Kasutani and A. Yamada, *The MPEG-7 Color Layout Descriptor: A Compact Image Feature Description for Highspeed Image/Video Segment Retrieval*, Proceedings of the 2001 International Conference on Image Processing (ICIP 2001), 2001. pp.7-10.
- [32] C. Lacruz, *Nomenclatura de las lesiones cervicales (de Papanicolau a Bethesda 2001)*, Rev. Esp Patol. **36** (2003), 5–10.
- [33] C. Lacruz and J. González, *Citología ginecológica: de Papanicolaou a Bethesda*, Editorial Complutense S.A., 2003.
- [34] N. Lassouaoui and L. Hammami, *Genetic algorithms and multifractal segmentation of cervical images*, Proc. Of. IEEE-EURASIP 7th International Symposium on Signal Processing and its Applications, 2003.
- [35] J. Lee, J. Hwang, D. Davis, and A. Nelson, *Integration of Neural Networks and Decision Tree Classifiers For Automated Cytology Screening*, IEEE (1991), 257–262.
- [36] Z. Li and K. Najarian, *Automated classification of Pap smear tests using neural networks*, In: Proceedings of international joint conference on neural networks **4** (2001), 2899–2901.
- [37] ———, *Biomedical image segmentation based on shape stability*, IEEE International Conference on Image Processing (2007), 281–284.
- [38] Y. Liu, T. Zhao, and J. Zhang, *Learning multispectral texture features for cervical cancer Detection*, IEEE (2002), 169–172.
- [39] L. Mango, *Computer-assisted cervical cancer screening using neural networks*, Cancer Letters **7** (1994), 155–162.

-
- [40] Y. Marinakis and G. Dounias, *Nearest neighbor based pap-smear cell classification using tabu search for feature selection*, *The pap smear benchmark. Intelligent and nature inspired approaches in pap smear diagnosis*, Special session proceedings of the NISIS -2006. Symposium (25-34) November 29 - December 1, Puerto de la Cruz, Tenerife, Spain., 2006.
- [41] ———, *Nature inspired intelligence in medicine: ant colony optimization for pap-smear diagnosis*, *International Journal on Artificial Intelligence Tools (IJAIT)* **17** (2008), 279–301.
- [42] Y. Marinakis, G. Dounias, and J. Jantzen, *Pap smear diagnosis using a hybrid intelligent scheme focusing on genetic algorithm based feature selection and nearest neighbor classification*, *Computers in Biology and Medicine* **39** (2009), 69–78.
- [43] Y. Marinakis, M. Marinak, and G. Dounias, *Particle swarm optimization for pap-smear diagnosis*, *Expert Systems with Applications* **35** (2008), 1645–1656.
- [44] E. Martin, *Pap-Smear Classification*, Master’s thesis, Technical University of Denmark, 2003.
- [45] N. Mat-Isa, M. Mashor, and N. Othman, *Early diagnostis system for cervical cancer based on neural networks*, Ph.D. thesis, Universiti Sains Malaysia, 2003.
- [46] ———, *An automated cervical pre-cancerous diagnostic system*, *Artificial Inteligence in Medicine* **42** (2008), 1–11.
- [47] NA Mat-Isa, *Automated detection technique for pap smear images using moving K-means clustering and modified seed based region growing algorithm.*, *Int J Comput Internet Manage* **13** (3) (2005), 45–49.
- [48] J Mayer, I Bruchim, SV Blank, and P Petignat, *Advancing women’s cancer care. Report from the 37th Annual Meeting of the Society of Gynecologic Oncologists*, *Gynakol Geburtshilffiche Rundsch.*, 2006.
- [49] S. McKenna, I.Ricketts., A. Caim, and K. Hussein, *A comparison of neural network architectures for cervical cell classification*, *Proceedings of the 3rd International Conference on Artificial Neural Networks Brighton*, 1993.
- [50] P. Mitra, S. Mitra, and S. Pal, *Staging of cervical cancer with soft computing*, *IEEE Trans Biomed Eng* **47** (2000), 934–40.
- [51] A Moscicki, M Schiffman, S Kjaer, and L Villa, *Updating the natural history of HPV and anogenital cancer.*, *Vaccine*. **24 Suppl 3**: (2006), S42–S51.

- [52] M. Mouroutis and S. Roberts, *Robust cell nuclei segmentation using statistical modeling*, IOP Bioimaging (1998), 79–91.
- [53] N Muñoz, X Bosch, de Sanjosé S., R Herrero, and Castellsagué X y col., *Epidemiologic Classification of Human Papillomavirus Types Associated with Cervical Cancer*. International Agency for Research on Cancer Multicenter Cervical Cancer Study Group, N Engl J Med **348** (2003), 518–527.
- [54] L. Nanni, A. Lumini, and S. Brahnam, *Local binary patterns variants as texture descriptors for medical image analysis*, Artificial Intelligence in Medicine **49** (2010), 117–125.
- [55] Bethesda. National Cancer Institute, *The 1988 Bethesda System for reporting cervical/vaginal cytologic diagnoses*, JAMA **262** (1989), 931–4.
- [56] H. Friedrich Nauth, *Citodiagnóstico Ginecológico*, Hans Friedrich Nauth, 2005.
- [57] J. Norup, *Classification of pap-smear data by transductive neuro-fuzzy methods*, Master's thesis, Technical University of Denmark, 2005.
- [58] N. Otsu, *A Threshold Selection Method from Gray-Level Histograms*, vol. 9, IEEE Transactions on Systems, Man, and Cybernetics, 1979.
- [59] G. Papanicolaou, *A new procedure for staining vaginal smears*, Science **95** (1942), 438–439.
- [60] ———, *Atlas of Exfoliative Cytology*, Cambridge, Mass., 1954.
- [61] G. Papanicolaou and E. Bridges, *Simple method for protecting fresh smears from drying and deteriorations during mailing*, Journal of the American Medical Association **164** (1957.), 1330–1331.
- [62] Passariello, *Imágenes médicas. Adquisición, Análisis, procesamiento e Interpretación*, Ediciones de la Universidad Simón Bolívar, 1995.
- [63] Piñeros, J. Ferlay, and R. Murillo, *Cancer incidence estimates at the national and district levels in Colombia*, Salud Pública Mex. **48** (2006), 455–465.
- [64] M. E. Plissiti, A. Charchanti, O. Krikoni, and D.I. Fotiadis., *Automated segmentation of cell nuclei in PAP smear images*, IEEE International Special Topic Conference on Information Technology in Biomedicine. Greece, Oct. 26-28., 2006.
- [65] M. Richard, *Art. And Science of Citopathology*, Chicago, American Society of Clinical Pathologists, 1996.
- [66] R. Richart, *Theory of Cervical Carcinogenesis*, Obstet and Gynec Surv **24** (1969), 874.

-
- [67] ———, *Cervical intraepithelial neoplasia*, *Pathol Ann.* **8** (1973), 301.
- [68] I. Ricketts, *Cervical cell image inspection- a task for artificial neural networks*, *Network* **3** (1992), 15–18.
- [69] N. Samsudin and A. Bradley, *Nearest neighbour group-based classification*, *Pattern Recognition* **43** (2010), 3458–3467.
- [70] A. Sánchez, *Aplicaciones de la visión artificial y la biometría informática*, DYKINSON S.L., 2005.
- [71] J. Santamaría, J. Rodríguez, and D. Agustín, *Cuadernos de Citopatología*, Díaz de Santos, 2006.
- [72] D. Solomon, D. Davey, R. Kurman, A. Moriarty, D. O'Connor, M. Prey, and et al., *The 2001 Bethesda System. Terminology for reporting results of cervical cytology*, *JAMA* **287** (2002), 2114–9.
- [73] Vapnik V., *The Nature of Statistical Learning Theory*, 1989.
- [74] R. Walker, P. Jackway, B. Lovell, and I. Longstaff, *Classification of cervical cell nuclei using morphological segmentation and textural feature extraction*, *Intelligent Information Systems, Proceedings of ANZIIS* (1994), 297–301.
- [75] Xu F. and Zhang Y-J, *Evaluation and comparison of texture descriptors proposed in MPEG-7*, *J. Vis. Commun. Image R* **17** (2006), 701–716.
- [76] S. Yang, Y. Chan, and Y. Chu, *Edge Enhancement Nucleus and Cytoplasm Contour Detector of Cervical Smear Images*, *IEEE Transactions On Systems, Man, And Cybernetics- Part B: Cybernetics* **38** (2008), 353–366.