

# **FEATURES EXTRACTION AND NUCLEI CLASSIFICATION IN TISSUE SAMPLES OF COLORECTAL CANCER**

By

**Sameer Akhtar Syed**

A Thesis

Submitted to the Faculty of Science  
in Partial Fulfillment of the Requirements for  
the Degree of Master of Computer Science  
at the University of Windsor

Windsor, Ontario, Canada

2020

**FEATURES EXTRACTION AND NUCLEI  
CLASSIFICATION IN TISSUE SAMPLES OF  
COLORECTAL CANCER**

by

Sameer Akhtar Syed

APPROVED BY:

---

Dr. Mohamed Belalia  
Faculty of Science

---

Dr. Imran Ahmad  
School of Computer Science

---

Dr. Boubakeur Boufama, Advisor  
School of Computer Science

# Declaration of Originality

I hereby certify that I am the sole author of this thesis and that no part of this thesis has been published or submitted for publication.

I certify that, to the best of my knowledge, my thesis does not infringe upon anyone's copyright nor violate any proprietary rights and that any ideas, techniques, quotations, or any other material from the work of other people included in my thesis, published or otherwise, are fully acknowledged in accordance with the standard referencing practices. Furthermore, to the extent that I have included copyrighted material that surpasses the bounds of fair dealing within the meaning of the Canada Copyright Act, I certify that I have obtained a written permission from the copyright owner(s) to include such material(s) in my thesis and have included copies of such copyright clearances to my appendix.

I declare that this is a true copy of my thesis, and that this thesis has not been submitted for a higher degree to any other University or Institution.

# Abstract

Cancer is an important public health problem and the third most important cause of death in North America. The American Cancer Society (ACS) estimates that by the end of 2020, 1,806,590 new cases and 606,520 deaths will occur in the US. Among the highest impact types of cancer are colorectal, breast, lung, and prostate.

This thesis addresses the features extraction by using different artificial intelligence algorithms that provide distinct solutions for the purpose of Computer-Aided Diagnosis (CAD). For example, classification algorithms are employed in identifying histological structures, such as lymphocytes, cancer-cells nuclei and glands, from features like existence, extension or shape. The morphological aspect of these structures indicates the degree of severity of the related disease. Unsupervised-learning algorithms, such as thresholding and c-fuzzy means, are used to segment cells or nuclei cells in histological images. Also, dimensionality reduction algorithms are used to simplify and extract hidden but important information on features. In this paper, we use a large dataset of 5000 images to classify eight different tissue types in the case of colorectal cancer. We employ a dimensionality reduction algorithm to extract important information about morphological and color features. At the time, deploying algorithm of Watershed, we perform image segmentation and extract statistical information about the area, perimeter, circularity, eccentricity and solidity of the interest points in the image.

Then, we use and compare four popular machine learning techniques, i.e., Naive Bayes and Random Forest, Support Vector Machine and Multilayer Perceptron to classify and to improve the precision of category assignation. The performance of each algorithm was measured using 3 types of metrics: Precision, recall and F1-Score representing a huge contribution to the existing literature complementing it in a quantitative way. The large

number of images has helped us to circumvent the overfitting and reproducibility problems with different datasets. The main contribution is the use of new characteristics different from those already studied, this work researches about the color and morphological characteristics in the images that may be useful for classifying the classes in the dataset of colorectal cancer.

# Acknowledgements

I'd like to acknowledge everyone here.

# Table of Contents

<b>Declaration of Originality</b>	<b>iii</b>
<b>Abstract</b>	<b>iv</b>
<b>Acknowledgements</b>	<b>vi</b>
<b>List of Figures</b>	<b>ix</b>
<b>1 INTRODUCTION</b>	<b>1</b>
1.1 Computer Vision . . . . .	1
1.2 Machine Learning . . . . .	4
1.3 Deep Learning . . . . .	8
1.4 Digital Pathology . . . . .	11
1.4.1 Histopathology . . . . .	12
1.4.2 Features Extraction . . . . .	13
1.5 Problem Statement . . . . .	14
1.6 Applications . . . . .	16
1.7 Organization . . . . .	19
<b>2 LITERATURE REVIEW</b>	<b>20</b>
2.1 Features Extraction . . . . .	20
2.1.1 Texture Features . . . . .	21
2.1.2 Color Features . . . . .	21
2.1.3 Topological Features . . . . .	22
2.1.4 Morphological Features . . . . .	22
2.1.5 Similarity Features . . . . .	24
2.1.6 Extracting Methods Review . . . . .	26
2.2 Processing Challenges . . . . .	28
2.2.1 Relation to our Work in Nuclei Classification . . . . .	34
<b>3 METHODOLOGY</b>	<b>36</b>
3.1 Motivation . . . . .	36
3.2 Proposed Methodology . . . . .	38
3.2.1 Dataset Collection . . . . .	38
3.2.2 Image Features . . . . .	40
3.2.3 Cancer Classification Algorithms . . . . .	44
3.2.4 Assessment Metrics . . . . .	46

<b>4 RESULTS AND DISCUSSION</b>	<b>48</b>
4.1 Cancer Classification Results . . . . .	49
4.1.1 Experimental dimensionality reduction algorithm . . . . .	49
4.1.2 8-Classes Dataset - Clasification Results . . . . .	49
4.1.3 Binary Dataset - Clasification Results . . . . .	55
4.2 Comparison . . . . .	57
<b>5 CONCLUSIONS AND FUTURE WORK</b>	<b>59</b>
 <b>Bibliography</b>	 <b>61</b>

# List of Figures

1.1	Optical character recognition (OCR). Image acquired from <a href="http://yann.lecun.com/exdb/lenet/">http://yann.lecun.com/exdb/lenet/</a> . . . . .	2
1.2	Surveillance and traffic monitoring (Image acquired from [7]). . . . .	3
1.3	Face detection algorithms can locate and recognize the individuals in this image. (Image acquired from [9]). . . . .	3
1.4	Medical imaging - radiological evaluation of dynamic contrast-enhanced from a magnetic- resonance imaging (MRI) of the breast (Image acquired from [8])	4
1.5	Types of Machine Learning algorithms (Image acquired from <a href="https://www.geeksforgeeks.org/">https://www.geeksforgeeks.org/</a> ) . . . . .	5
1.6	This is an example where we know there are three classes. There are three hypotheses induced, each one covering the instances of one class and leaving outside the instances of the other two classes. (Image acquired from <a href="https://www.javatpoint.com/supervised-machine-learning">https://www.javatpoint.com/supervised-machine-learning</a> ) . . . . .	6
1.7	Neural Networks consist of the following components An input layer $x$ . An output layer $y$ . An arbitrary amount of hidden layers. A set of <i>weights</i> ( $W$ ) and <i>biases</i> ( $b$ ) between each layer. A choice of activation function for each hidden layer, $\sigma$ . This diagram shows the architecture of a 2-layer Neural Network (input layer excluded). (Image acquired from <a href="https://towardsdatascience.com/how-to-build-your-own-neural-network-from-scratch-in-python-1-1-1">https://towardsdatascience.com/how-to-build-your-own-neural-network-from-scratch-in-python-1-1-1</a> )	
1.8	Common activation Functions. The Rectified Linear Unit (RELU) function is the most widely-used activation function in neural networks today. One of the greatest advantages RELU has over other activation functions is that it does not activate all neurons at the same time it converts all negative inputs to zero. This makes it very computationally efficient as few neurons are activated at any given time. In practice, RELU converges six times faster than other functions like <i>tanh</i> and <i>sigmoid</i> [18]. (Image acquired from <a href="https://7-hiddenlayers.com/deep-learning-2/">https://7-hiddenlayers.com/deep-learning-2/</a> ) . . . . .	10

1.9	Disease classification task. Input image is an abnormal axial slice MRI brain going through a schematic of a Convolution, RELU and pooling layers, before classification by the fully connected layers (Image aquired from [16]). . . . .	11
1.10	Sample image of epithelium tumour and stroma tumour along with the feature values from a colorectal cancer database consists of 1332 image of epithelium and stroma tissue sample extracted from the patients (available on <a href="http://fimm.webmicroscope.net/supplements/epistroma">http://fimm.webmicroscope.net/supplements/epistroma</a> [17]) (Image acquired from [24]) . . . . .	15
1.11	Examples of H&E and IHC images (Image acquired from [23]) . . . . .	17
1.12	Acquisition workflow diagram in histopathology images [46]. . . . .	17
2.1	Example of color and texture differences between histopathological images of colon tissues: (a), (b) Normal and (c), (d) cancerous. Nonglandular regions in images are shaded with gray. Image aquired from [60] . . . . .	22
2.2	a) A sample image, (b) the Voronoi diagram of the image (black dotted lines) and its Delaunay triangulation (red solid lines), and (c) a cell-graph of the image. Image acquired from [32] . . . . .	23
2.3	(a) A sample of a nucleus with its boundary points and centroid, (b) line segments used in symmetry computation, (c) chords used in concavity computation. Image acquired from [32] . . . . .	23
2.4	Similarity relationship from a histopathological images of biopsy samples: (a) a healthy breast tissue, (b) a cancerous breast tissue, (c) a healthy brain tissue, and (d) a cancerous brain tissue. Image acquired from [32] . . . . .	23
2.5	Flowchart of a proposed BCNN-based classification method for histopathological images [57]. . . . .	30
2.6	Schematic overview of the resampled-based Markovian model (RMM) to classify given images [60]. . . . .	30
2.7	Stages defined for multidimensional and fuzzy approaches [47]. . . . .	31
3.1	Object boundary detection using morphological filter. a) Image without filter b) Image with Filter [42]. . . . .	37
3.2	a) Results of corner detector b) corners projected on input image [43] . . . . .	38
3.3	Proposed Methodology divided by four steps . . . . .	38
3.4	Selected types of cancer and non cancer cells for the study . . . . .	39
3.5	Decomposition of an image in the RGB components. . . . .	42
3.6	Decomposition of a image in the HSV components. . . . .	42

3.7	Decomposition of a image in the YUV components. . . . .	43
3.8	Nuclei segmentation. . . . .	43
4.1	General Performance of the 4 algorithms - 8-Classes Dataset. . . . .	50
4.2	Performance of the 4 algorithms using Precision - 8-Classes Dataset. . . . .	53
4.3	Performance of the 4 algorithms using Recall - 8-Classes Dataset. . . . .	53
4.4	Performance of the 4 algorithms using F1-Score - 8-Classes Dataset. . . . .	54
4.5	Confusion matrix of four algorithms - 8-Classes Dataset. . . . .	55
4.6	Performance of the 4 algorithms - Binary Dataset . . . . .	56
4.7	Confusion matrix of four algorithms - Binary Dataset. . . . .	57

# Chapter 1

## INTRODUCTION

### 1.1 Computer Vision

As humans, we perceive the three-dimensional structure of an image taken by the world around us, otherwise as we can perceive evident elements, we can disconcert for some common optical illusions that makes it difficult to detect useful information. As mentioned in [11], visual information is being attracted by our eyes at great speed. Much of this information is redundant and compressed by several layers in the visual cortex, so that the higher centers of the brain only interpret a small fraction of the information.

Computer Vision is the study of methods and techniques to extract more information than we can observe, this system can be efficiently used in practical applications. Therefore, it encompasses both the science and engineering of vision. Computer vision focuses on extracting relevant information from images at a high level to perform vision-based tasks [5].

Research topics under computer vision include motion detection, autonomous navigation, scene reconstruction and recognition, augmented reality (AR), object recognition, object tracking and many more vision-based tasks. Consequently, there are factors involved that make vision such a difficult task for machines to accomplish. Computer vision systems often face challenges such as scale change, variations in lighting conditions, point of view changes, partial occlusion, deformation of non-rigid objects, or intra-class variation of visual

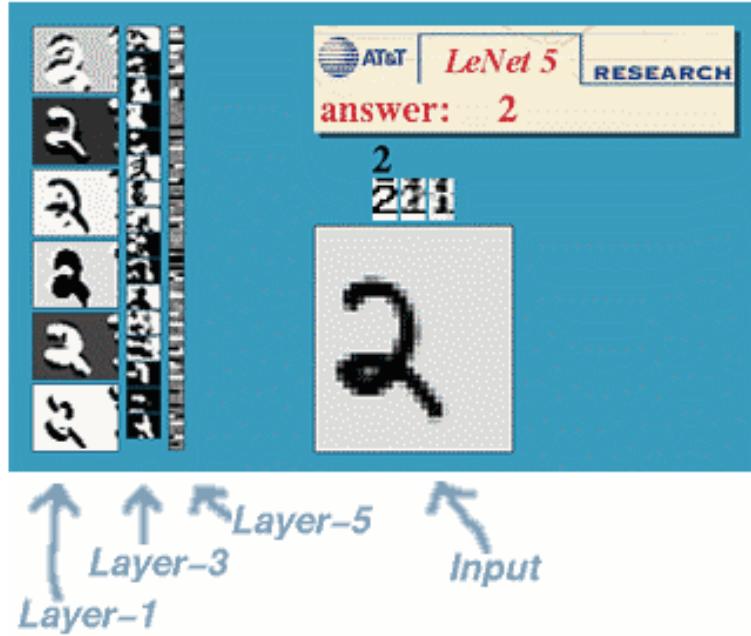


FIGURE 1.1: Optical character recognition (OCR). Image acquired from <http://yann.lecun.com/exdb/lenet/>

perspectives [10].

Some typical applications of computer vision can be found in optical character recognition (OCR) which extract printed or handwritten text from images (Figure. 1.1). In other field, with the help of computer vision it is possible to detect the movement of cars and to track and count the different vehicles by analyzing a camera picture (Figure 1.2). Another example of a successful application of computer vision is face detection, known as the first step towards many face-related technologies, with recognition or verification purpose (Figure 1.3).

Figure 1.4 is an example of a medical assistance, oriented to Computer-aided diagnosis (CADx). CADx is known as one of the major research subjects in medical imaging and histological diagnostic, that the performance by computers complement physicians view [4].

Computer vision methods for assisting medical diagnosis generally follow a similar structure: it relies on the extraction and combination of several features obtained from pre-processed images, and uses them to build models that can be generalized to invisible data. This image processing methodology can be used to extract out of sight information, thereby providing valuable outcomes to perform related tasks in CADx.

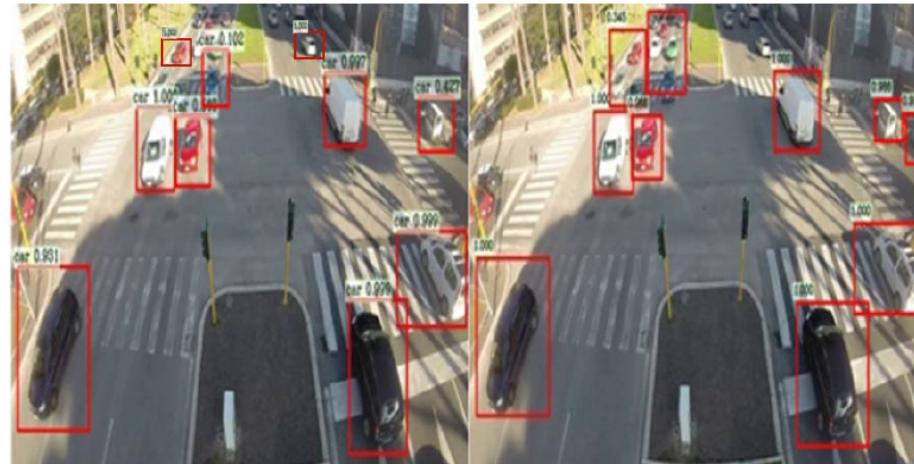


FIGURE 1.2: Surveillance and traffic monitoring (Image acquired from [7]).

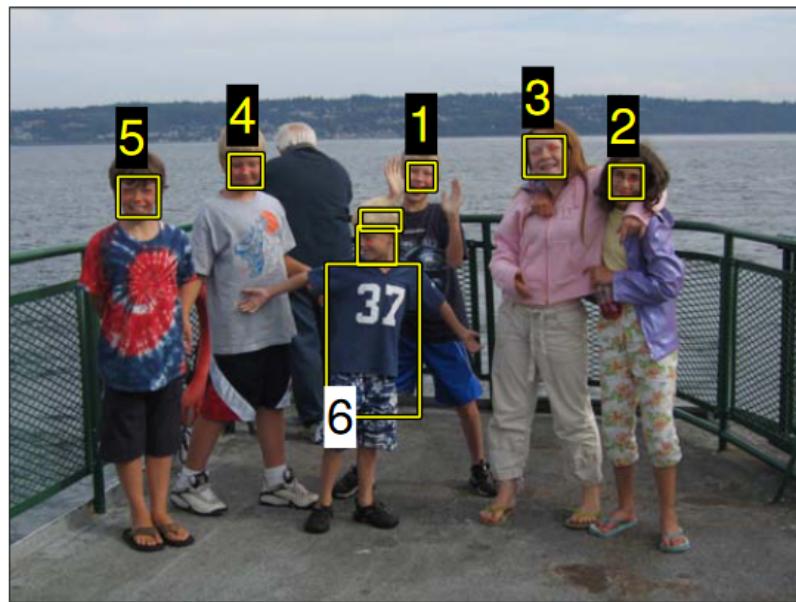


FIGURE 1.3: Face detection algorithms can locate and recognize the individuals in this image. (Image acquired from [9]).

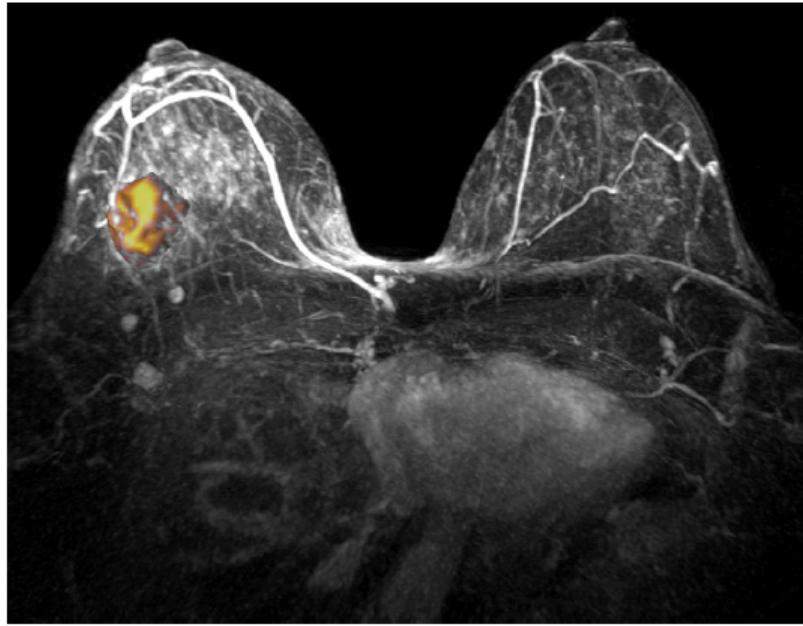


FIGURE 1.4: Medical imaging - radiological evaluation of dynamic contrast-enhanced from a magnetic- resonance imaging (MRI) of the breast (Image acquired from [8])

We further explain machine learning and supervised techniques to features extraction in section 1.2 and deep learning in computer vision in section 1.3. We present an overview of Digital Pathology in section 1.4. Then, the problem statement in section 1.5, application in section 1.6, and finally the thesis organization in section 1.7.

## 1.2 Machine Learning

In order to solve a problem on the computer, we need an algorithm. An algorithm is a set of instructions that should be executed to convert input to output. There may be various algorithms, we maybe interested in finding the most efficient. *Machine learning* is the programming of computers to optimize performance, using sample data or past experience as a standard to solve a given problem.

Machine learning uses statistical theory to build mathematics Model, due to make inferences based on samples. The core task is twofold: First, in training, we need to write efficient algorithms to solve optimization problems and storage and processing data. Two, a model is learning, its representation and reasoning needs algorithm solutions. We build a model depending on some parameters, learning is the execution of a computer program

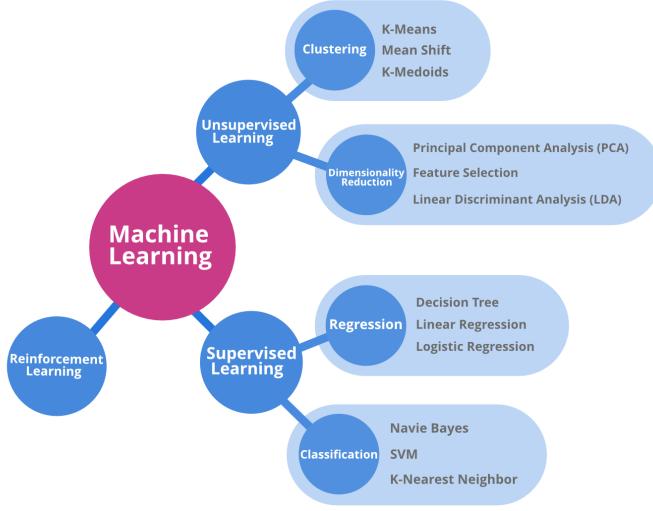


FIGURE 1.5: Types of Machine Learning algorithms (Image acquired from <https://www.geeksforgeeks.org/>)

using training data to optimize model features to get a model that may be predictive and acquire knowledge from data.

There are 3 different types Machine Learning (ML) techniques named:

- **Unsupervised ML Algorithms:** refer to drawing inferences from unlabelled data.
- **Supervised ML Algorithms:** refer to training and testing the algorithms over labelled data.
- **Reinforcement ML Algorithms:** refer to learning a lesson by optimal actions through trial and error. This means that the algorithm decides the next action by learning behaviors that are based on its current state and that will maximize the reward in the future [25]. We summarize ML types on Figure 1.5

Furthermore, we discuss along this thesis, a classification algorithm trained by *supervised learning*. In supervised learning, classification algorithms are used when the output variable is categorical, which means there are two or more classes such as Yes-No, True-false, etc. Figure 1.6 shows a model trained using labeled data sets, where the model learns for each data type. After the training process is completed, the model will be tested on the test data (a subset of the training set), and then the output will be predicted.

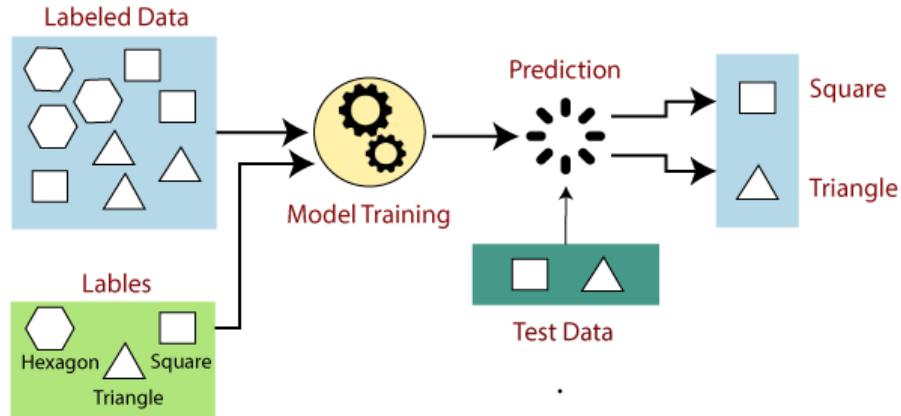


FIGURE 1.6: This is an example where we know there are three classes. There are three hypotheses induced, each one covering the instances of one class and leaving outside the instances of the other two classes. (Image acquired from <https://www.javatpoint.com/supervised-machine-learning>)

On the classification Machine learning field, we find intelligent application in eHealth, such as early diagnosis and prognosis of cancer types. Cancer research has been popularized due to the importance of classifying cancer patients into high-risk or low-risk categories has led many research teams.

Therefore, these technologies have been used for the purpose of simulating the progress and treatment of cancer conditions. In addition, the ability of machine learning tools to detect key features from complex data sets reveals their importance. A variety of these techniques, including Bayesian Networks (BNs), Support Vector Machines (SVMs), Decision Trees (DTs), and Multilayer Perceptrons (MLP) have been widely applied in cancer research for the development of predictive models, resulting in effective and accurate decision making [26].

Some machine learning techniques are majorly related with colorectal cancer diagnosis and prognosis. This thesis focuses mainly on detecting the primary site of cancer. Then, there is need of popular ML techniques required for the early detection of this type of cancer as we present:

- **Naive Bayes**

**Naive Bayes** classifiers are statistical classifiers that can predict class membership probabilities, such as the probability that a given sample will belong to a particular

case. The naive Bayes classifier depends on Bayes' theorem. This classifier assumes that the influence of attribute values on a given class is independent of the values of other attributes. [29].

- **Random forest**

**Random forest** is a collection (ensemble) of decision trees. It is a popular ensemble technique in pattern recognition. Random forest for cancer classification is based on gene expression and generated with the different number of features per node split [29].

- **Support Vector Machines (SVMs)**

**Support vector machines** are supervised learning models. In the SVM algorithm, each data item is drawn as coordinates in an n-dimensional space, where n is the total number of elements used for classification, and the value of each element is represented by the coordinates of the data point. SVM contains a decision hyperplane, which is used to divide different types of data points using maximum margin. The data points located near the hyperplane are called support vectors. This classification process generates a nonlinear decision boundary and classifies data points that are not represented in a vector space [27].

- **Multilayer Perceptron (MLP)**

The **multilayer perceptron** is a feedforward artificial neural network. MLP is a supervised learning algorithm used to learn functions by training on the basis of a given data type. MLP can learn a nonlinear function approximator for any classification. [29]

These techniques allow to classify important characteristics of pathology diagnosis, but we need to compare each other. With this in mind, we use 4 performance measures for classification algorithm [30]:

- **Recall:** It is the relation between correct matches and correspondences. It is the division of relevant instances that are retrieved as result.

- **Precision:** It is defined as the number of true positive cases over the total number of everything classified as positives. It evaluates how much of results were essentially true.
- **Accuracy:** It is used to express the closeness of a measurement to the true value.
- **F1 score:** Also known as balanced F-score or F-measure. It can be interpreted as a weighted average of the precision and recall, where an F1 score reaches its best value at 1 and worst score at 0.

### 1.3 Deep Learning

Deep learning is a subset of machine learning, consisting of algorithms inspired by the structure of a human brain. It is particularly useful for working with unstructured data which nevertheless follows some internal rules: images, videos, text, etc. Through the use of deep learning, computer vision has improved rapidly. Using deep learning models, such as Convolutional neural networks (CNN), image processing has become a very broad field, covering a variety of methods to reproduce human visual capabilities, such as recognizing human faces, objects or patterns [15].

CNNs are developed based on Artificial neural networks (ANNs). ANNs is one of the most famous machine learning models. It was introduced as early as the 1950s and has been actively researched since then [15]. Roughly speaking, a neural network consists of many connected computing units (called neurons), which are arranged in layers. There is an input layer, the data enters the network, and then there are one or more hidden layers, which transform the data as it flows through the data, and then terminate in the output layer that produces the neural network predictions.

In computer vision, this methods begins from a random value, training repeatedly until the network learns the most appropriate weights for each neuron and the output layer provides an accurate prediction. The network is adjusted according to update each weight of each neuron until the pattern recognized by the network produces a good prediction on the training data. However, we need to evaluate the “goodness” of these predictions, i.e., how far off are our predictions?. Hence, we use a *loss function* that allows to calculate

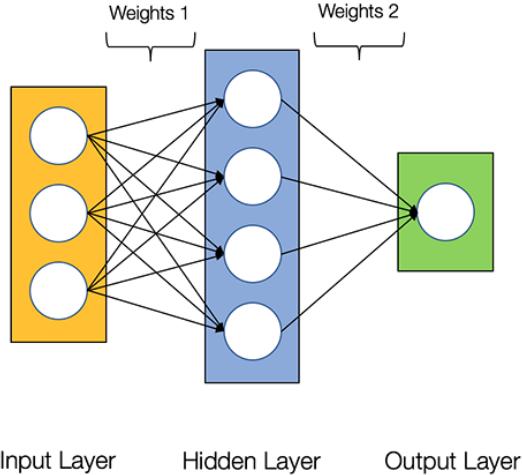


FIGURE 1.7: Neural Networks consist of the following components An input layer  $x$ . An output layer  $y$ . An arbitrary amount of hidden layers. A set of *weights* ( $W$ ) and *biases* ( $b$ ) between each layer. A choice of activation function for each hidden layer,  $\sigma$ . This diagram shows the architecture of a 2-layer Neural Network (input layer excluded). (Image acquired from <https://towardsdatascience.com/how-to-build-your-own-neural-network-from-scratch-in-python-68998a08e4f6>)

an error rate which show us the difference between each predicted value and the actual value. There are many available loss functions, we choose the most convenient according to the nature of our problem, once we determine one, the goal in training is to find the best set of weights and biases that minimizes the loss function. Furthermore, we need to update our weights and biases based on *Backpropagation*. Backward propagation of errors (*Backpropagation*), is an algorithm that calculates the gradient of the error function with respect to the neural network's weights [18].

At this part, a standard neural network can be constructed and trained. However, deep learning is not possible because of training process involves a different activation function in contrast to classical bounded activations like  $\text{sign}(x)$ ,  $\sigma(x)$ , and  $\tanh(x)$ . These functions are associated to each node, and are applied to node inputs to produce node outputs. In this case we use the *rectified linear activation function (RELU)*. This is a non-linear function that will output the input directly if it is positive, otherwise, it will output zero. It has become the default activation function for many types of CNNs because a model that its derivative with respect to its input, is always 1 for positive input value, hence it solves the problem of vanishing gradient and makes easier to train and often achieves better performance [17].

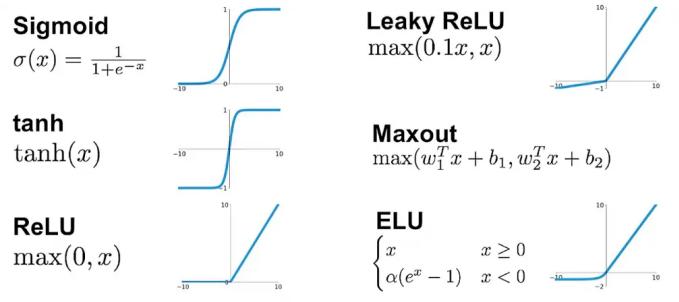


FIGURE 1.8: Common activation Functions. The Rectified Linear Unit (RELU) function is the most widely-used activation function in neural networks today. One of the greatest advantages RELU has over other activation functions is that it does not activate all neurons at the same time it converts all negative inputs to zero. This makes it very computationally efficient as few neurons are activated at any given time. In practice, RELU converges six times faster than other functions like *tanh* and *sigmoid* [18]. (Image acquired from <https://7-hiddenlayers.com/deep-learning-2/>)

Returning to computer vision, CNNs are the most researched machine learning algorithms in medical image analysis for preserving spatial relationships when filtering input images. A CNN takes an input image of raw pixels, and transforms it via *Convolutional Layers*, *Rectified Linear Unit (RELU) Layers* and *Pooling Layers*. Then, this output feeds into a final Fully Connected Layer which assigns class scores or probabilities, thus classifying the input into the class with the highest probability [19].

The *convolutional layer* is the core component of CNN. It bears the main part of the network computing load. The main purpose of convolution is to extract features such as edges, colors, and corners from the input. As we go deeper into the network, the network also begins to recognize more complex features, such as shapes, numbers, and faces.

*Convolution* is defined as the operation of two functions. In image analysis, one function composed of the input value (such as pixel value) at a certain position in the image, and the second function is a filter (or kernel); each can be represented as an array of numbers. Calculating the dot product between two functions to get the output. This process is repeated until cover the entire image, generating a *feature map*, a map where the filter is strongly activated and highlight a feature such as a straight line, a dot, or a curved edge [16].

Moreover, the *RELU layer* uses an activation function that turns into negative input values to zero. Thus, simplifies and accelerates computing and training, and helps to evade vanishing gradient problem. The *pooling layer* is inserted between the convolutional layer

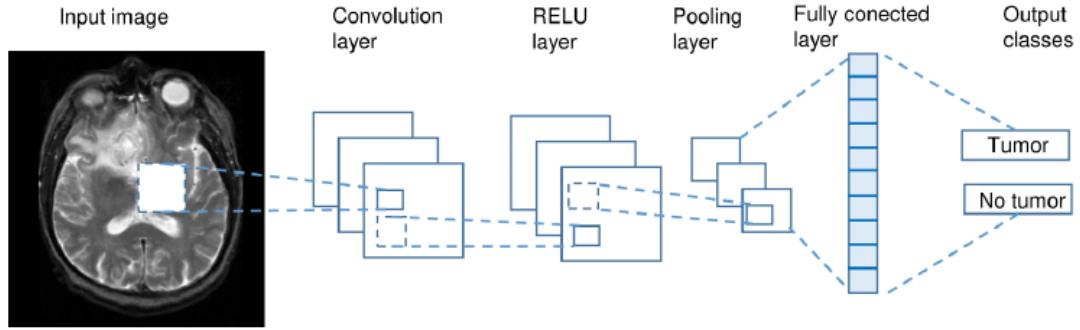


FIGURE 1.9: Disease classification task. Input image is an abnormal axial slice MRI brain going through a schematic of a Convolution, RELU and pooling layers, before classification by the fully connected layers (Image aquired from [16]).

and the RELU layer to reduce the number of parameters to be calculated and the size of the image (width and height, but not depth). Finally we find in a CNN, the *Fully Connected Layer*. This layer takes the output of the previous layer (convolution, RELU or pooling) as input and calculates the probability score of the classification from different classes available. As shown in Figure 1.9, we notice that final connected layer looks at the combination of the most strongly activated functions that will determine that the image belongs to a particular category. Other example where we can apply this CNN, is on histological slides. Cancer cells have a higher ratio of DNA to cytoplasm than normal cells, if the DNA signature is strongly detected from convolution, RELU or pooling layer, CNN is more likely to predict its existence cancer cell.

## 1.4 Digital Pathology

In the last century, the basic process of pathologists for diagnosis has remained relatively unchanged, but advances in information technology have provided huge opportunities for image-based diagnosis and research applications. Pathology has lagged behind other health care practices, such as histology, where this technology is widely adopted in digital medicine. As the equipment that generates full slide images becomes more practical and affordable, practice will increasingly adopt the technology and eventually generate new data that will quickly make the present large amount of histological imaging data overshadowed [20].

This work is framed into what is called *Digital Pathology*, which is a relatively recent

area of research which is dedicated to providing accurate and efficient computational methods to support quantitative detection, diagnosis, and prognosis in pathology. It presents various computational learning methods and frameworks for the automatic representation of histopathology images by learning from databases of different digital pathology tasks including the detection, localization and quantification of tumours and tissues in various types of cancer.

The role of machine learning in computer vision learning is very significant. There are many supervised learning techniques which are capable of segment, measure and classify images to applied to routine pathological tasks, including the quantification of antibody staining, the identification and classification of cells, and the characterization of essentially multi cellular or regional microstructures [20]. The development and evolution of computational techniques to classify categories such as K-means algorithms, Support Vector Machines (SVM), Multilayer Perceptron (MLP), Random Forest Algorithm among others have allowed the appearance of some solutions to apply on diagnosis assistance of malignant tumours and tumour vs non-tumour identification.

#### 1.4.1 Histopathology

Delving deeper into the field of pathology and histopathology, we find Histopathology slides provide a more comprehensive view of the disease and its impact on the tissue, because the preparation process preserves the underlying tissue structure. In this way, certain disease characteristics, such as lymphocytic infiltration of cancer, can be inferred only from histopathological images [21].

In the last decade, technological advances in software and hardware have made possible the digitization of histological slides and the creation of histopathological images, through robotic microscopes, known as operating tissue scanners. Digitized slides are generated from image data of the entire sheet in high resolution. These are visualized through computer platforms that emulate the functionalities of a microscope, increasing with information technologies through the internet and accessed through the use and exchange of data and images online through computers or mobile devices. These technological developments offer great advantages over the traditional treatment of using microscopes, by allowing in a

very expeditious way ease of access, interconsultation, second opinion, as well as multiple applications and useful tools in clinical diagnosis, research and education [22].

With the arrival and mass use of digitized slides, great interest has arisen in the development and use of image analysis algorithms that allow the quantification of various biological markers. The benefits provided by such algorithms include improvements in accuracy and precision in the detection, classification and measurement of morphometric patterns. In this last aspect, these analyzes are currently being applied mainly in the quantification of the expression of tumour markers in cancer.

In recent years, to analyze these characteristics and find out abnormalities of some type of cancer, computerized systems have applied segmentation algorithms by thresholding have been used (unsupervised machine learning). Other studies have focused on computer classification of the type of tissue present, in order to make a diagnosis based on it. The fundamental problem in most of the algorithms developed is the large number of parameters that need to be tuned, which is why work has been carried out to adjust them automatically. The process of segmentation allows to distinguish between different regions in the image and thus finding regions of interest.

#### 1.4.2 Features Extraction

The systems used for digital pathology are designed to process images obtained from different microscopes. Multiple visual characteristics (colors, shapes, etc.) are obtained from the regions of interest (ROI) to identify and make a better diagnosis from tissue samples, through image segmentation and feature extraction [23].

Additionally, Morphological characteristics provide information about the size and shape of the described region, object or image. For example, the glandular area can be used as a discriminatory criterion to classify between cancer and non-cancer, the perimeter can also be used as a descriptor of the traits to characterize the size of the cells in the segmentation process. Taking advantage of the pathologist's experience in the diagnosis of Oral Submucosal Fibrosis (FSO), Muthu and the other researchers represented the image using morphological characteristics such as the perimeter of eccentricity, thus being the diameter of the equivalent area to describe the nucleus of the cell. Not only can classification tasks

be performed using these types of characteristics, there is also an automated system for differential white blood cell (GB) count based on 19 characteristics such as area, perimeter, convex area, solidity, orientation and eccentricity. In addition, these features have been used to construct Content Based Image Retrieval (CBIR) to find histopathology images of the prostate based on morphological similarity. Finally, most of the estimation of morphological measurements is made based on a previous segmentation, so its performance depends on the precision of said segmentation.

The intensity characteristics provide information on the gray level or color of the pixels located in the regions of interest. This feature extraction approach uses different color spaces where the hue channel of the HSV (Hue-Saturation-Value) color space conversion of the original image is used. It can also work in the white / pink / purple dimension with the objective of evaluating color models. It allows to compare the incidence on the performance in a classification task and make sure that there is no single model that works better than others in all cases. Figure 1.10, presents an example of the perception-based feature values of an epithelium and stroma tumour from a histopathological slide.

## 1.5 Problem Statement

All these studies are motivated due to the fact that Cancer is a global public health problem in the 21st century. According to WHO (World Health Organization) is ranked as the third cause of death in the world and the second in developed countries, only surpassed by cardiovascular diseases. The WHO GLOBOCAN developments aim to provide updated estimates on the incidence, mortality and prevalence of the main types of cancer, at the national level, for 184 different countries. These estimates are based on the most recent data available from the International Agency for Research on Cancer (IARC) and on information available to the public on the Internet. For the year 2012, the GLOBOCAN project estimated that there were 14.1 million new cases of cancer worldwide, 8.2 million deaths from cancer and 32.6 million people living with this disease (in the 5 years after the diagnosis).

According to estimates of cancer treatment and survival rates, on january 2019, more than 16.9 million Americans had a history of cancer, only to population growth and aging,

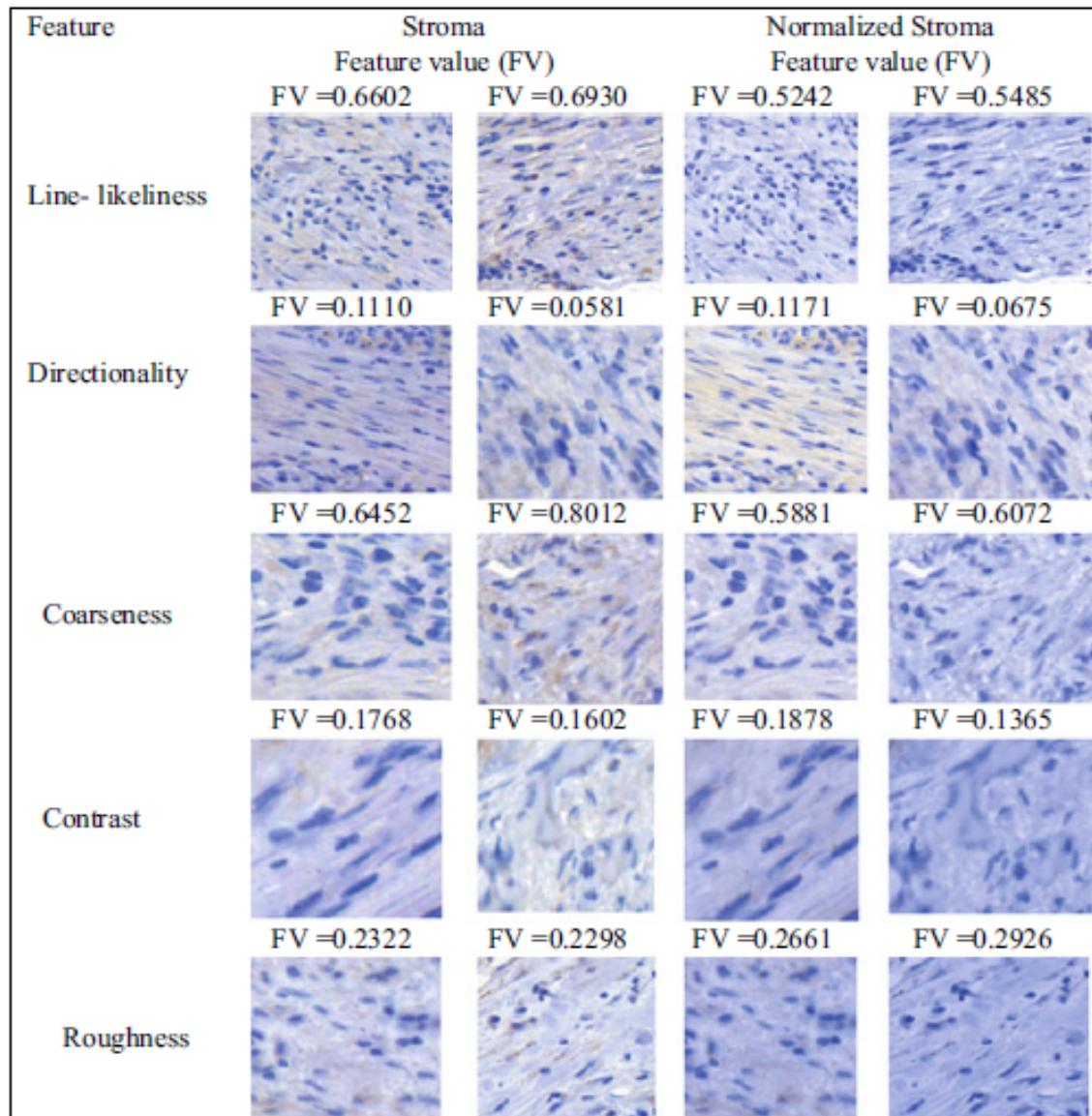


FIGURE 1.10: Sample image of epithelium tumour and stroma tumour along with the feature values from a colorectal cancer database consists of 1332 image of epithelium and stroma tissue sample extracted from the patients (available on [http://fimm.webmicroscope.net/supplements/epistroma\[17\]](http://fimm.webmicroscope.net/supplements/epistroma[17])) (Image acquired from [24])

this number is expected to exceed 22.1 million statistics by 2030. The report is produced every three years by the American Cancer Society and the National Cancer Institute to help the public health community better serve this growing population [41]. The current numbers mean that one in five men and one in six women will eventually get cancer by the time they reach 75. One in eight and one in 12 women will pass away from the disease [64]. One of the most concerning factors was that 57 percent of new cancer cases and 65 percent of deaths came from the developed world.

The three most common cancers among men in 2019 are prostate cancer (3,650,030), colon and rectal cancer (776,120) and skin melanoma (684,470). Among women, the three most common cancers are breast cancer (3,861,520), endometrial (uterine body) (807,860), and colon and rectal cancer (768,650). The author's estimate of the number of cancer survivors in 2030 (22.1 million) is based on population projections [41]. Therefore, the less developed regions and that have few resources will occur 57% (8 million) of new cancer cases, 65% (5.3 million) of cancer deaths and 48% (15.6 million) of people with the disease (5 years after diagnosis) [64].

## 1.6 Applications

The histopathology acquisition process follows a well-defined methodology in its first steps, according to Figure 1.12 [46]: First, the biological sample is taken from an organ. Then, a fixation process is done over the biopsy to ensure chemical stability on the tissue and to avoid post-mortem changes. After this, it must be cut into sections that can be placed onto glass slides. The sections are stained to reveal cellular components by chemical reactions. The most common dyes used are Hematoxylin - Eosin (H&E), which stain cell nuclei in a dark blue or purple and cytoplasm and connective tissue in bright pink (i.e. as shown in Figure 1.11). Finally, the section covers slipped into being viewed and digitized with a microscope.

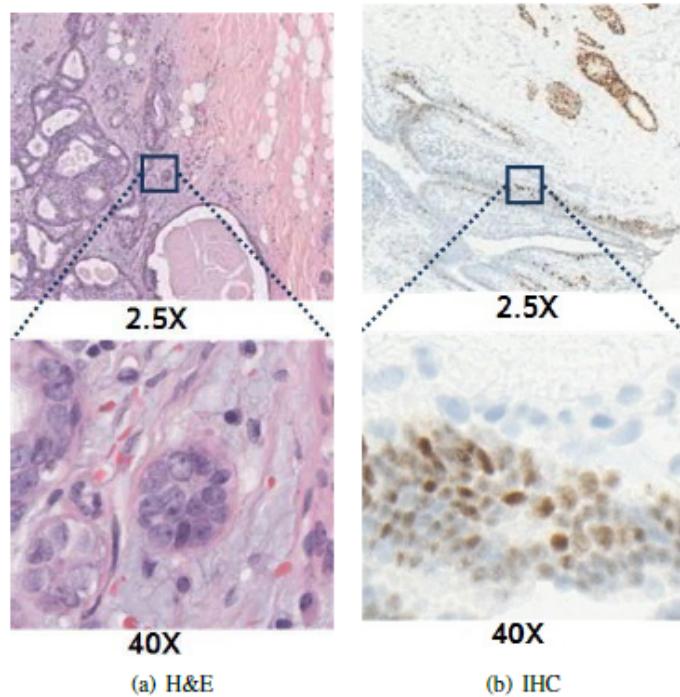


FIGURE 1.11: Examples of H&E and IHC images (Image acquired from [23])

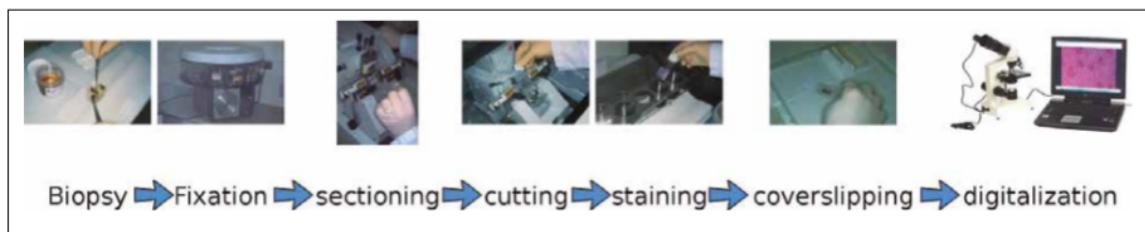


FIGURE 1.12: Acquisition workflow diagram in histopathology images [46].

Human factors in this process constitute sources of variability. This visual heterogeneity is due to three factors at least [46]: (i) **Magnification**, referring to increasing the proportion of biological structures which are visible under the microscope according to the set of lenses; (ii) **Staining**, that is, improve the contrast in a biological sample seen under a microscope; (iii) **Slice Orientation**: how the tissue appearance changes when the cut is done in a longitudinal or cross-sectional orientation in smooth muscle tissue using the same staining (H&E).

The usual histopathology image workflow usually goes as follows [46]:

- **Image preprocessing:** Raw image data is transformed to reduce visual variability and noise, as well as making it more suitable for the subsequent steps.
- **Feature extraction:** whose purpose is to produce a more descriptive representation of the image-making important explicit information which is not directly manifest from the raw pixels.
- **Pattern recognition:** In which interesting visual patterns are detected and identified through supervised and unsupervised algorithms.

As described above, the development of techniques to solve the public health problem caused by cancer has evolved significantly in recent years. Even if this has been caused by the emergence of new algorithms due to the boom in computer science, the growth and development of the applications of these techniques in cancer detection has been significantly high than others areas due to the high investment that countries have invested in the study for the treatment and cure of cancer. For example, the European Union in recent times has been increasing its budget for cancer research and treatment, measures that have been taken by each and every one of the member countries in proportion to their economic capacities, which shows the importance given to the subject in this region and the understanding that prevention is cheaper than treatment.

To summarize, research in the field of computer vision for assistive medical diagnosis is increasing because it provides objective goals about the patient's condition. Therefore, a non-invasive way to detect colorectal cancer, is essential for medical diagnosis, because developed detection algorithms are still rarely used in clinical practice, their reliability is still worthy of attention due to lack of research on this type of cancer.

As a result, this thesis analyzes histopathological images to decide whether a tissue contains cancer or not. Using computational learning tools which automatically find patterns for healthy and abnormal tissues, we contribute a fundamental support for the diagnosis of colorectal cancer with this research. It is also embedded into histology which studies the microscopic anatomy of biological tissues. In this study histopathological cancer images are the cornerstone to understand the state of biological structures, provide the diagnosis and analyze the state of diseases [46].

## 1.7 Organization

This thesis es organized as follows: Chapter 2 reviews the literature in image processing and computational algorithms used in cancer detection. Chapter 3 will go over the methodology used to address the problem of deciding whether a tissue is cancer positive or not as stated in subsection 1.6. Chapter 4 will go over the results showing the performance of the algorithms used and tested to classify the histopathological images. Finally Chapter 5 Resumes the mains aspects and conclusion of this work.

# Chapter 2

## LITERATURE REVIEW

This chapter reviews the literature related to classification of histopathological images for diagnosis/procedures. We reviewed work focused on an increasing amount of literature with different datasets and objectives, having in common a high degree of accuracy. The common feature in CADx works is that specific image analysis problems require specific image representation schemes. Therefore, we make an overview of the most common features extraction scheme for histopathological image representation listed on the visual features which describe and provide the most relevant information for specific machine learning algorithms, depending on the type of image, some features may be more significative than others.

Thus, we review features extraction in section 2.1. This section reviews of machine learning and deep learning techniques in visual features extraction. Subsequently, we notice some challenges from high-level semantic concepts as connected areas and objects which could determine the visual pattern of disease if it exists, or, a false positive case. Thus, we focus on processing histopathological challenges for CADx in section 2.2.

### 2.1 Features Extraction

Numerous methods have been used to obtain a host of features. The latter represent clinically meaningful information after the preprocessing of images. Indeed, according to

[46] *Feature extraction methods* are usually performed on the following feature types:

### 2.1.1 Texture Features

Texture features provide data about variance of intensities inside of a specified region. A texture is a set of connected pixels that repeatedly appear in an image. It provides information about changes that determine the surface strength by quantifying properties such as smoothness, roughness, and regularity. To extract information from these features one can use statistical properties like correlations, means, etc, to relate between each others to tissue identification. Similarly, after such statistical processing, these features can be used to classify sub-images. For example, the image is divided into squares of sub-images and texture features are extracted from these squares, and then the squares containing malignant cells are distinguished from the squares composed [60].

### 2.1.2 Color Features

Color features or visual features based on intensity provide information on the gray level or color of pixels located in the ROI. This type of function does not provide any information about the spatial distribution of pixels. Intensity histogram units are used to define characteristics. For example, by using the gray value of the pixel to define the optical density of the pixel and by using the pixel values in a single color channel may establish a relationship between color values in different channels.

Figure 2.1 shows an example of color differences between a normal and a cancerous tissue from a histopathological sample. Also, at the same figure we can distinguish patterns like texture that can be extracted.

In Welsh approach, the transfer result depends on the luminance information of reference image. Since the process of Gupta approach is on the purpose of propagating color information using the least-squares optimization method, the result shows limited relevance to the reference image luminance

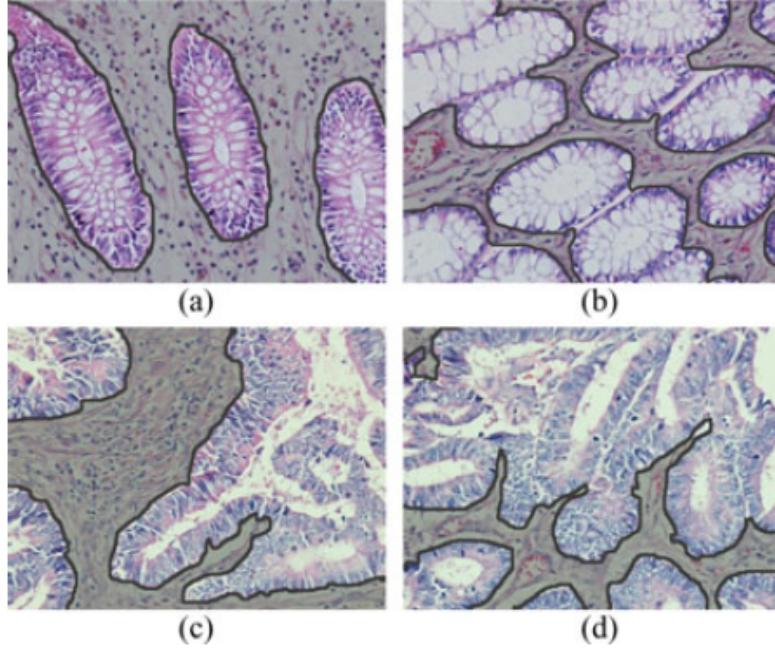


FIGURE 2.1: Example of color and texture differences between histopathological images of colon tissues: (a), (b) Normal and (c), (d) cancerous. Nonglandular regions in images are shaded with gray. Image aquired from [60]

### 2.1.3 Topological Features

Topological features provides information about some structure within the image. The structure of a tissue by quantifying the spatial distribution of its cells. Thus, graph construction, local and global graph features is a usual method to extract patterns that consist of a bunch of cells. Therefore, the local metrics correspond to the properties of segmented images or sub-images. Moreover, the global ones are the giant connected component ratio, spectral radius. As example, Figure 2.2, shows a sample image and its respective graph model.

### 2.1.4 Morphological Features

Morphological features provide information about the size and shape of the described region, object or image. This information is particularly useful on the segmentation task, where graph theory is commonly the most used in this group of features. Area, perimeter, angle and other geometrical features are calculated. On the other hand, the shape consists of

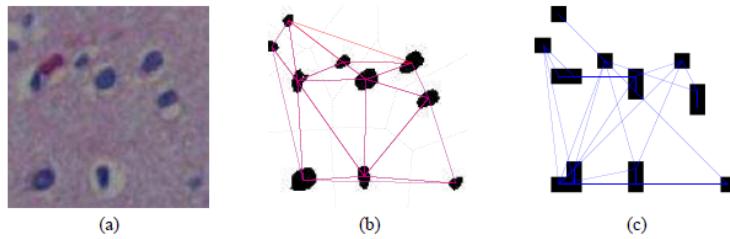


FIGURE 2.2: a) A sample image, (b) the Voronoi diagram of the image (black dotted lines) and its Delaunay triangulation (red solid lines), and (c) a cell-graph of the image. Image acquired from [32]

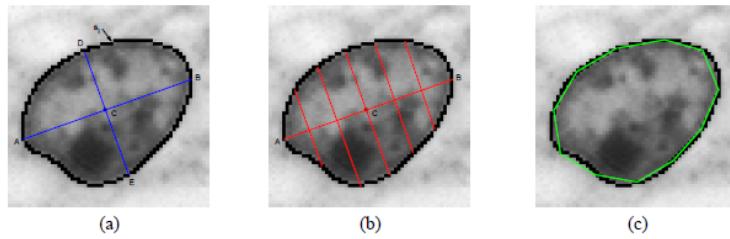


FIGURE 2.3: (a) A sample of a nucleus with its boundary points and centroid, (b) line segments used in symmetry computation, (c) chords used in concavity computation. Image acquired from [32]

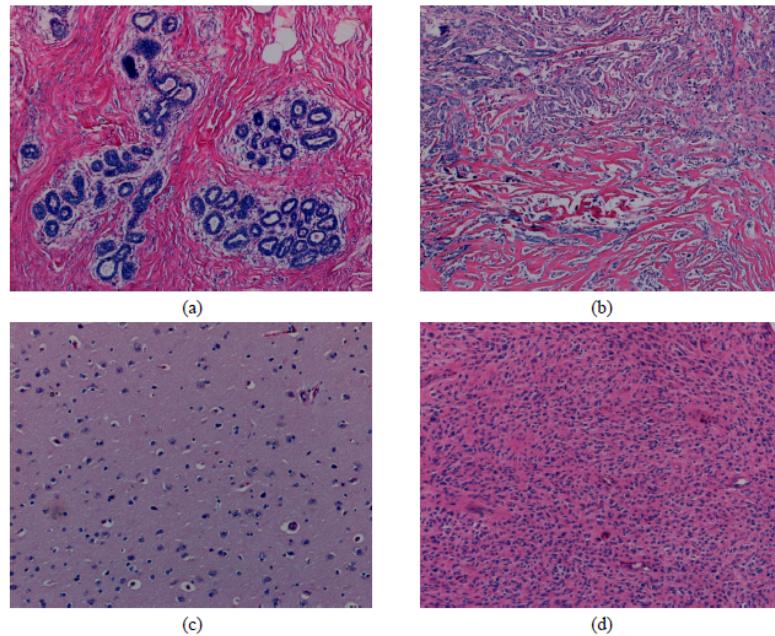


FIGURE 2.4: Similarity relationship from a histopathological images of biopsy samples: (a) a healthy breast tissue, (b) a cancerous breast tissue, (c) a healthy brain tissue, and (d) a cancerous brain tissue. Image acquired from [32]

compactness, roundness, smoothness, length of major and minor axes, symmetry, concavity and circumference [32]. These characteristics can be found in Figure 2.3 and a sample image to describe which segments and patterns may be calculated.

### 2.1.5 Similarity Features

Similarity measures provide information about similar image sequences. Commonly used to compare characteristics from different samples. Proof of this, Figure 2.4, shows a couple of different sample images that can be related. Furthermore, this review can summarize the above features in Table 2.1.

Feature Type	Feature Extraction Method or Metric	Authors
Texture	First-Order Statistics	[51]
	Lower Order - High Order Histograms	[60]; [52]
	Local Binary Patterns (LBP)	[54], [52], [56]
	Run Length Matrix (RLM)	[51]
	Gray-level co-matrix (GLCM)	[44], [52], [51], [56]
	Haralick Descriptor	[54], [44], [50], [56]
	Gabor filters	[54], [52], [56]
Color	Perception-like features: coarseness, contrast, directionality, line-likeness, and roughness.	[52]
	Sample Entropy	[47]
	First Order Statistics	[50], [51]
	Contrast Measure	[54]
Morphological	Quantile Normalization	[53]
	Perimeter (Size)	[48]
	Ratio of major to minor axis (Shape)	[48]
	Mean of the gray-level intensity (Nuclear Appearance)	[48]
Topological	Axis of Least Inertia	[59]
	Graph Analysis: Diameter, Degree, Clustering Coefficient	[45]
	Delaunay Triangulation	[60], [45]
Similarity	Grid-Based Approach	[60]
	Voting Approach	[60]
	Bag of Words Approach	[60]

TABLE 2.1: Feature extraction method for different feature types

### 2.1.6 Extracting Methods Review

According to Table 2.1, most used computational learning models are less heterogeneous due to low sample sizes. While supervised learning techniques are intended to classify new observations in the defined categories, Support Vector Machine (SVM) is the most common method and provides the most accurate estimations in general. Unsupervised learning is used mostly in preprocessing. To detect usual ductal hyperplasia (UDH), [48] segmented cell regions by clustering the pixel data with Gaussian Mixture Models (GMMs) in four cytological regions (i.e., cellular: nuclear and cytoplasmic, extracellular, regions with hues and Illumina). The pixels classified as cellular components are further clustered by using dynamic thresholding to eliminate pixels with less luminance. In the end, individual cells are segmented by converting them to gray-level and using watershed. The idea is an amount of waterflows along with a topographic relief following a certain descending path to eventually reach a catchment basin. Blobs in the image can be separated by identifying the limits of adjacent catchment basins and then separating them. Even so, [51] classify regions of interest containing benign/malignant cancerous colon tissues. RGB intensities of the samples are converted to optical density values and, therefore, the saturation of each stain. 2-Means clustering was used to extract nuclei structures by identifying pixels with high hematoxylin concentration. The only exception is the one found in [59], which used Self Organizing Maps to split normal and tumorous epithelium tissues.

Computational Learning Class	Method	Authors
Supervised Learning	Support Vector Machine	[45], [48], [54], [60], [52], [44], [50] & [47]
	Logistic Regression	[47]
	Neural Networks	[44], [52], [47]
	Linear Discriminant Analysis	[53]
	Ensemble Trees	[52]
	Decision Trees	[47]
	Random Forests	[51], [47]
	RBF	[47]
	Multiple Instance Learning	[48]
	Resampling Based-Markovian Models	[60]
Unsupervised Learning	Bilinear Convolutional Neural Networks	[57]
	Convolutional Neural Networks	[58]
	k-Means	[51], [44], [60]
	C-Fuzzy Means	[59]
Self-Organizing Maps (SOMs)	Self-Organizing Maps (SOMs)	[59]
	Gaussian Mixture Models	[48]

TABLE 2.2: Different types of methods for computational learning classes

The most popular state-of-the-art study of the reviewed literature is [58]. Using 86 H&E slides to obtain 100.000 image patches and data augmentation, they identify nine tissue types with several CNNs (Convolutional Neural Networks): VGG19, AlexNet, SqueezeNet, GoogLeNet, and Resnet50, with a 70/15/15 division of the dataset into train/validation/test sets, where the first one obtained a 98% accuracy. Visual representations of tissue classes are obtained through t-SNE on deep layer activations, with an almost perfect separation

of the classes in the testing set; visualization of morphological features are derived from a Deep-Dream approach. Their excellent performance let them establish four categories of consensus molecular subtypes (CMSs) and calculate a deep stroma score.

The performance of the algorithms depends on the ways the data are incorporated into them. As shown in [51] the performance of eight algorithms is tested for different ways of incorporating the raw data. see Table 2.3:

Classifiers	Vector with 24 Attributes	Vector with 13 Attributes	Vector with 37 Attributes (24+13)
Random Forest	0.742	0.769	0.79
J48	0.72	0.651	0.742
SMO	0.603	0.687	0.7
Rotation Forest	0.923	0.896	0.913
MultiClassClassifier	0.889	0.94	0.886
Multilayer Perceptron	0.891	0.917	0.907
Logistic	0.889	0.94	0.866
RBFNetwork	0.641	0.682	0.675
Average $\pm SD$	0.787 $\pm$ 0.111	0.810 $\pm$ 0.113	0.807 $\pm$ 0.081

TABLE 2.3: Performance of eight algorithms according to [51].

According to Table 2.3, the best algorithm is Rotation Forest allowing to have the best result with 24 and 37 attributes, while better results are obtained using 13 vectors with MultiClassClassifier and logistic regression algorithms.

## 2.2 Processing Challenges

A significant methodology related to computer vision systems in medical imaging involves a stage of preprocessing, for example a segmentation algorithm to easily capture the variability induced by texture or color inside histopathological samples, among others.

Most preprocessing steps aims to reduce noise when they are properly applied. Nonetheless there are many challenges in high precision algorithms of tissue classification and high-efficiency algorithms in computation time. Specifically, the number of samples in different studies is too small to draw a general conclusion [52]. The natural high-dimensionality of images increases the complexity of analysis and the over-division into patches could increase the number of false-positive observations, even though the overall accuracy degree is high. In most practical applications the number of labeled images is too low and the staining process depends on laboratory procedures.

Hence, in [57] aiming at internalizing the transformation processes, the authors tried to solve these problems by proposing a Bilinear Convolutional Neural Network (BCNN) for the classification of histopathological images of tissues with colorectal cancer. Also, enhancing the feature representation, we find that the categories specified by [52], as shown on Figure 2.5, were learned from hematoxylin-eosin (*H&E*) components jointly and the classification of the images. They applied a decomposition algorithm to obtain *H* and *E* representations through a stain decomposition algorithm. Moreover, with the feature functions associated, they use two parallel CNNs, which are a hierarchy of neural units, including a convolutional, pooling, and non-linear layer. The Feature outputs are combined at each location using the outer matrix product before the classification is done. Similarly in [52], the authors used CNNs for category classification, combining a DeepNet architecture with a focal loss function. Still, they carried out magnification and color normalization to images.

Along with the authors in [60], they proposed to generate perturbed images from training data and modeling them by a Markov process. The classification of the images is done using their perturbed samples and thereby, reducing the negative outcomes due to the large variance in tissue images (Figure 2.6). They present their solution as follows. They randomly select points of the stained images and locate a window at the center of each point to extract four-color and texture features. Following, a k-Means algorithm is run to quantify the pixels into dominant cluster-colours of H&E staining (white, pink and purple). The features are then discretized into observation symbols. Each new data point is labeled into the cluster of the nearest neighbour, so each observation has a set of observation symbols that can be ordered by the shortest distance. A sequence can be therefore defined, and the noise is reduced due to the existence of large sequences, modeling images in a

Markov framework.

In another related work [47]. In this work, the authors start dividing stained images into equal-size cells and a rolling squared-window of shape  $m$  is defined by moving from upper left to lower right (see Figure 2.7). Then, they define a similarity measure based on comparing maximal individual distances to average intensity in each window and channel. A probability of similarity of each window and the other ones is calculated, which lets them calculate a sample entropy measure from the average of these probabilities as the window moves. The sample entropy curves are described as a function of a constant tolerance and the window size, which let them calculate the area under the curve, obliquity, area ratio, maximum point value and maximum point scale metrics under different configurations. These calculated entities are used as 13 features for classification between benign and malignant images of colorectal adenocarcinoma.

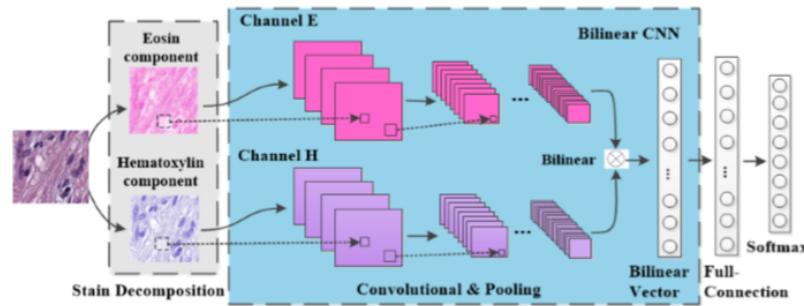


FIGURE 2.5: Flowchart of a proposed BCNN-based classification method for histopathological images [57].

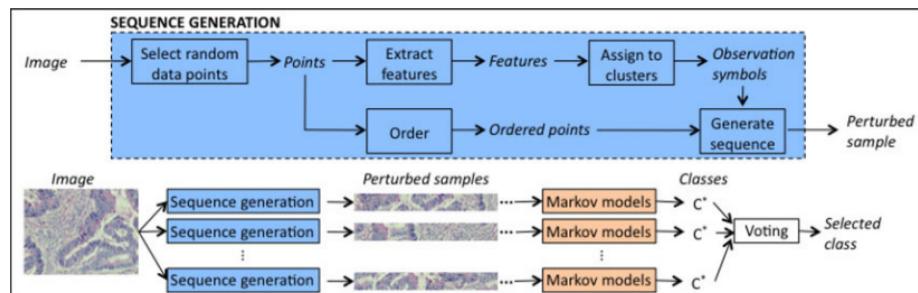


FIGURE 2.6: Schematic overview of the resampled-based Markovian model (RMM) to classify given images [60].

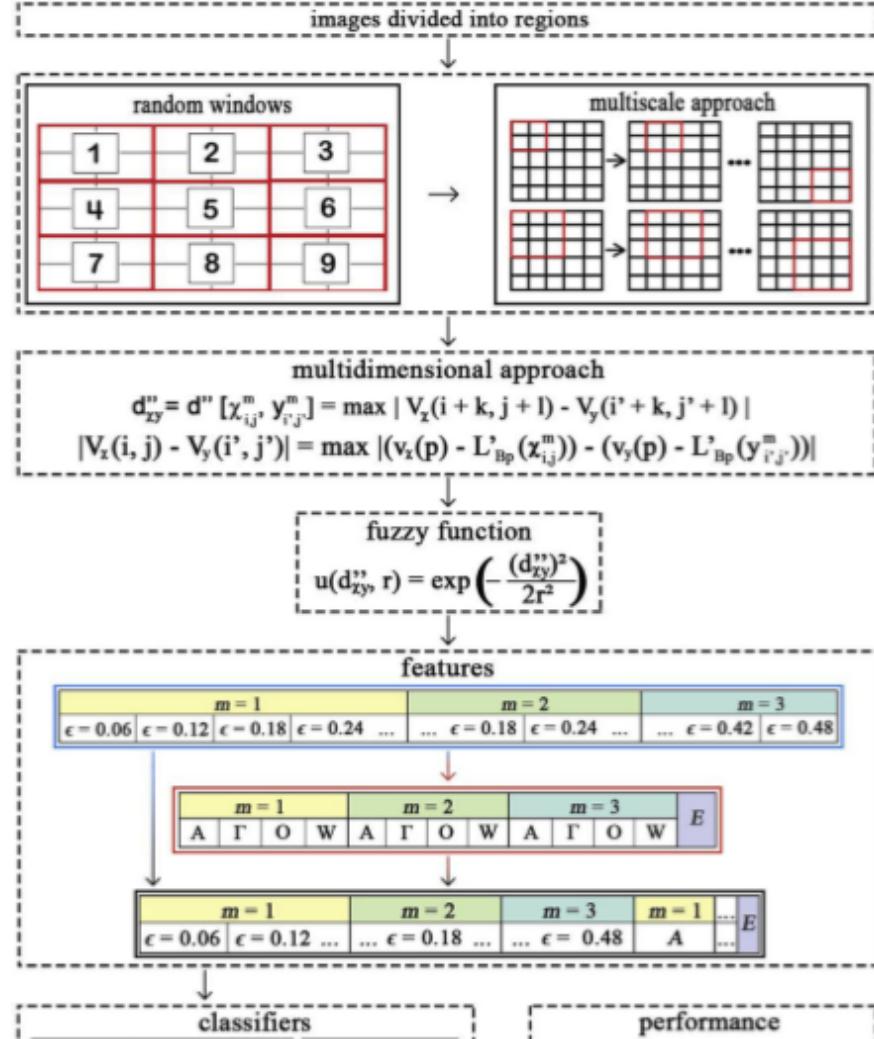


FIGURE 2.7: Stages defined for multidimensional and fuzzy approaches [47].

Finally, we present in Table 2.4 a record of literature which exhibits components of relatives works with a performance measure and a sample from a database of histopathological images to establish a complete medical diagnosis/prognosis context.

Reference	Classification Objective	Dataset Type	Sample Size	Accuracy
[52]	Tumor epithelium; Simple and complex stromas; Immune cell groups; debris and mucus; mucosal glands; adipose tissue; background.	Colorectal Cancer	5000 images (625 samples for each class)	0.85
[53]	Basophilic structures; Eosinophilic intra & extra cellular proteins; Lumen of glands; Red blood cells.	Renal Tumor; Glioblastoma; Ovarian.	Two renal tumors (RCC1 and RCC2 with 55 and 47 images, respectively), one glioblastoma (Gbm, 52 images), and one ovarian (Ov, 50 images).	0.87
[48]	Binary classification of Usual Ductal Hyperplasia (UDH)	UDH	327 regions of interest are used for training; 149 for testing.	0.879
[54]	Epithelium Stroma Tissues	Epithelium Stroma	576 images of regions of interest were used for training and 720 images for testing.	0.995
[47]	Benign and malignant colorectal adenocarcinoma.	Colorectal Cancer.	50/50 images of benign/malignant cancer for training and 17/34 for testing.	AUC: 0.983
[45]	Normal; Low-Grade Adenocarcinomatous; High-Grade Adenocarcinomatous Colon Tissues.	Colon Cancer Diagnosis	213 microphotographs 115 images of tissues are used for training and 98 images for testing.	0.8265
[51]	Benign and malignant colon tissues	Colon Cancer Diagnosis	44 benign and 43 malignant H&E stained tissue samples.	0.91

Reference	Classification Objective	Dataset Type	Sample Size	Accuracy
[50]	Well, intermediate, poor cancer-degree	Colon Cancer Diagnosis	92 malignant colon biopsy samples (23/44/25 samples associated with poor/mid/well-degrees of cancer)	0.9613
[57]	Tumor epithelium, stromas, immune cell groups, debris and mucus, mucosal glands, adipose tissue, background	Colon cancer	125 images for each class	0.985
[56]	Normal, Hyperplastic polyp (HP), Tubular Adenoma with low-grade dysplasia (TA_LG) and Carcinoma (CA)	Colorectal Cancer Diagnosis	200 images, containing 50 images from each of the four classes	0.9317
[60]	Normal; Low-grade cancerous; High-grade cancerous.	Colorectal Cancer Diagnosis	3226 images of colon tissues. 1644 images for training (510 normal, 859 low-grade cancerous, 275 high cancerous) and 1592 for testing (491 normal, 849 low-grade cancerous, 257 high-grade cancerous).	0.9522 0.8945 0.8646 (High).
[44]	Cancerous; Normal	Colorectal Cancer Diagnosis	113 colon images: 64 cancerous & 49 normal.	0.833
[59]	Tumour and Normal Epithelium	Colorectal Cancer	Diagnosis 134 images of normal and tumorous regions; 14 whole slide images stained for different biomarkers.	0.81-0.96

Reference	Classification Objective	Dataset Type	Sample Size	Accuracy
[61]	Background; Adipose; Mucus; Tumor Epithelium; Mucosal glands; Muscle; Stroma; Blood vessel; Immune cell; Necrosis	Colorectal Cancer Diagnosis	660 digitalized colorectal cancer specimens	0.72-0.96
[58]	Adipose tissue; Background; Debris; Lymphocytes; Mucus; Smooth muscle; Normal colon mucosa; Cancer-associated stroma; Colorectal adenocarcinoma epithelium.	Colorectal Cancer.	86 H&E slides of human cancer tissue to create 100.000 image patches for training; 25 H&E slides to create 7.180 patches for testing.	0.943

TABLE 2.4: Literature Review

### 2.2.1 Relation to our Work in Nuclei Classification

Our work in Nuclei Classification uses components from computer vision and classification algorithm. In particular we start at cellular-level feature extraction. We build a nuclei classification algorithm to determine the locations of the nuclei/cells in a tissue. The type of the feature extraction method to be deployed depends on the sample. In this case, computer vision allows us to segment a sample into an specified area of interest. On the other hand, the complex nature of histopathological tissue images means a challenge to characterize an individual nucleus as well as an entire tissue by aggregating the features of its nuclei. For this reason, we use color and morphological features to determine the exact locations of nuclei beforehand. Briefly, our objective is to establish a relationship between color characteristics and morphology for the classification of tissues with colorectal cancer

Furthermore, color and morphological methods allow us to analyze changes and quantify the size and shape characteristics of cell nuclei in order to construct a nuclei classification algorithm. In this case, the dataset consists of 150 morphological features, and 15 color features.

Finally, to solve the exact details of cells and the success, we detailed in section 3 the next steps to show how we treat the difficulty of the complex nature of image scenes, i.e., stain related problems including lack of dark separation lines between a nucleus and its surroundings, inhomogeneity of the interior of a nucleus, and occurrence of non-nuclei stain artifacts in a tissue [34].

# Chapter 3

## METHODOLOGY

This chapter displays our methodology for the extraction of color and morphological features and its implementation of nuclei classification in Colorectal Cancer (CRC) digital histopathological images using machine learning. In section 3.1, we explain the motivation behind the feature extraction techniques with a brief overview. Then, we explain the details of the methodology in section 3.2. This section is divided into four main phases. Subsection 3.2.1 reviews the *dataset collection*. Then subsection 3.2.2 details the extraction of *image classification features*. In subsection 3.2.3, we emphasize on *modelling of cancer classification algorithms* and finally subsection 3.2.4, we explain the *assessment metrics* used to measure the performance of classification algorithms, to compare it later.

### 3.1 Motivation

In this approach, we apply artificial intelligence in the biomedical area to support colorectal cancer diagnosis using digitized histopathological images. The analysis of digital histopathological images often uses general image recognition technology but this has some problems, such as texture analysis since this approach is used in pixels, it is sensitive to noise in the values of pixels. To address this issue, this research focuses on the relationship between *morphological* and *color* characteristics. We take advantage of this type of information to classify images of histopathology.

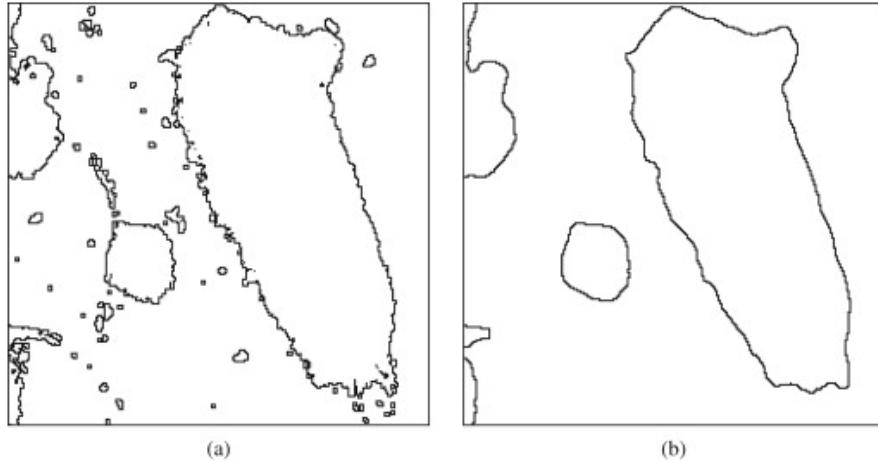


FIGURE 3.1: Object boundary detection using morphological filter. a) Image without filter  
b) Image with Filter [42].

Overall, *morphological* features are commonly used to boundary detection and highlight its importance to preserve, uncover, and detect the geometric structure of image objects. Also, morphological features are quite relevant at shape analysis offering efficient solutions to apply denoising-filters like median-type and stack filters. In short, boundary detection becomes quite sensitive to small noise artifacts, for example, we present Figure as boundary detection and its properties for noise reduction, detail preservation, and artifact-free images [42].

On the other hand, *color* features are important for expressing saliency, and color uniqueness can be clearly incorporated into the design of image feature detectors. For the most part detection methods are based on the statistical analysis of color derivatives [43]. Figure 3.2 presents results of a corner detector algorithm based on color feature trained according where more pixels are located. Briefly, color is determined by concentration of pixels, it means intensity of each color channel (Red,Green,Blue for RGB images).

To achieve this, we work with a method split into four main steps. First, we obtain the digitized histopathological images dataset in order to select the types of tissue for analysis based on previous works [46] [36]. The second step is construct a classification model able to determine if a nuclei is cancerous or not, from color and morphological characteristics. Then, on third step, we apply the constructed model using machine learning algorithms to our dataset. Finally, at fourth step, we evaluate performance of algorithms

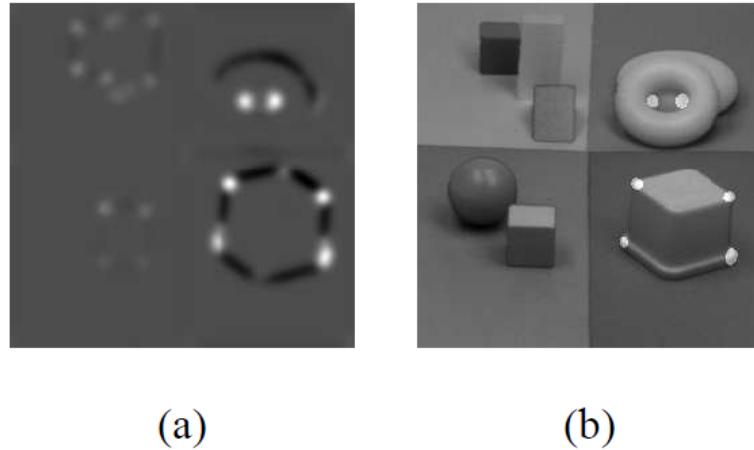


FIGURE 3.2: a) Results of corner detector b) corners projected on input image [43] .

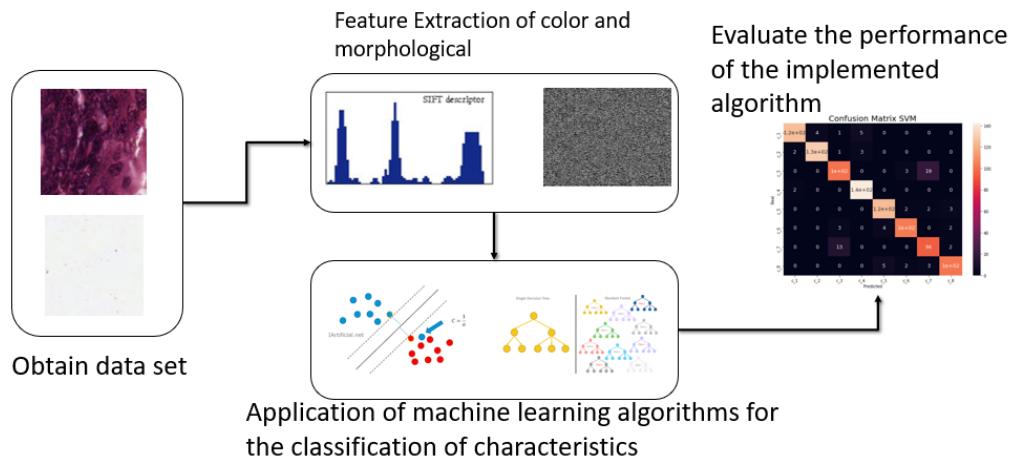


FIGURE 3.3: Proposed Methodology divided by four steps

through a comparison of four popular assessment metrics to detect changes and test the best performing algorithms. See Figure 3.3 where we provide the scheme of proposed methodology.

## 3.2 Proposed Methodology

### 3.2.1 Dataset Collection

First phase consists of obtaining the dataset. *The Medical Center of Mannheim* (University of Heidelberg, Germany) provides us a database which contains 8 types of labels separated

by folders as Figure 3.4 on images *.tif*. These classes were selected from the most representative images from our dataset. These ones shows the wide variation of illumination, stain intensity and tissue textures present in routine histopathological images [36]. Images were extracted from 10 independent samples of colorectal cancer (CRC) primary tumours.

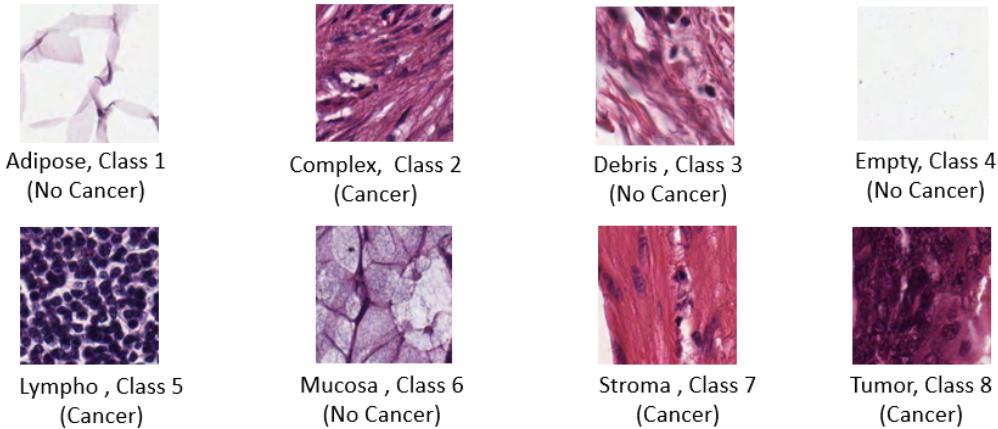


FIGURE 3.4: Selected types of cancer and non cancer cells for the study

Labeling histopathological images as cancerous or non-cancer regions is a key task in cancer diagnosis. It is also important to divide the cancer tissue into different categories to give a medical context. In this thesis, we analyze a dataset compound by anonymized H&E stained Colorectal Cancer (CRC) tissue slides. Each tissue image in Figure 3.4 is a label that explains if an image is cancer diagnosed or not, if it is the only tissue, or a part of the sheet, among others. However, supervised approaches require intensive work and time to obtain this labels. In this case, dataset used is already labelled and exhibited on Table 3.1. The slides were first digitized [37]. Moreover, contiguous tissue areas were manually annotated and tessellated [36], we use 625 non-overlapping tissue tiles of dimension  $150px \times 150px$  ( $74\mu m \times 74\mu m$ ). In summary, the specifically required images correspond to previously labeled histopathology images of hematoxylin-eosin (H&E) from colorectal cancer. That data is covered by an MIT license [35].

Class	Name	Diagnostic	Description
1	Adipose	No Cancer	Adipose tissue
2	Complex	<b>Cancer</b>	Containing single tumour cells and/or few immune cells
3	Debris	No Cancer	Including necrosis, hemorrhage and mucus
4	Empty	No Cancer	No tissue
5	Lympho	<b>Cancer</b>	Immune-cell conglomerates and sub-mucosal lymphoid follicles
6	Mucosa	No Cancer	Normal mucosal glands
7	Stroma	<b>Cancer</b>	Homogeneous composition, includes tumour stroma, extra-tumoural stroma and smooth muscle
8	Tumor	<b>Cancer</b>	Tumour epithelium

TABLE 3.1: Description of selected types of cancer for the study.

### 3.2.2 Image Features

The second phase contains the extraction of both morphological and color features. We present in Table 3.2, the corresponding features to analyze and build the model.

Color Features	Morphological Features*
Mean	Area
Median	
Standard deviation	Perimeter
Skewness	
Kurtosis	Circularity
Energy	
Entropy	Eccentricity
Color Histogram (8 bins)	

TABLE 3.2: Features to extract. \* For each of the morphological characteristics, a mean and standard deviation is obtained

- **Color Features**

The obtention of the color characteristics will be done with first-level statistics, which are the mean, median, standard deviation, kurtosis, skewness, and the color histogram (which obtains an even quantity of distribution and normalizes). The anterior process is done for different color spaces such as HSV, RGB, YUV, and grayscale. Figure 3.5 displays the process for the RGB space, Figure 3.6 shows the process for the HSV components and Figure 3.7 exhibits the process for the YUV components.

- **Morphological Features**

On the other hand, the morphological features are extracted through a nucleation segmentation. This nucleation segmentation process is carried out by creating a hematoxylin and eosin color space; then, the color component in hematoxylin is taken out applying a watershed algorithm to detect nuclei as shown in Figure 3.8. Besides, morphological characteristics have to undergo a process of core segmentation. After having the segmented image we proceed to obtain the statistical measurements mentioned above. The extraction of the main features, such as circularity, eccentricity, solidity, area, and perimeter, is performed — features from which the mean and standard deviation are also extracted.

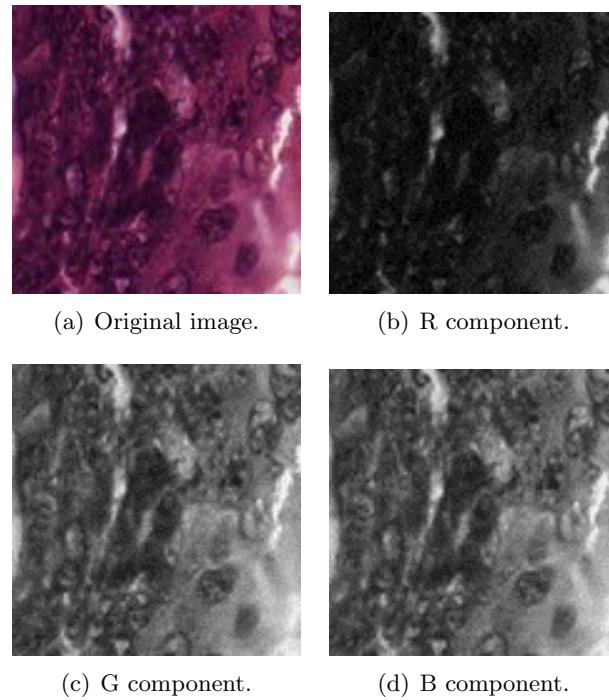


FIGURE 3.5: Decomposition of an image in the RGB components.

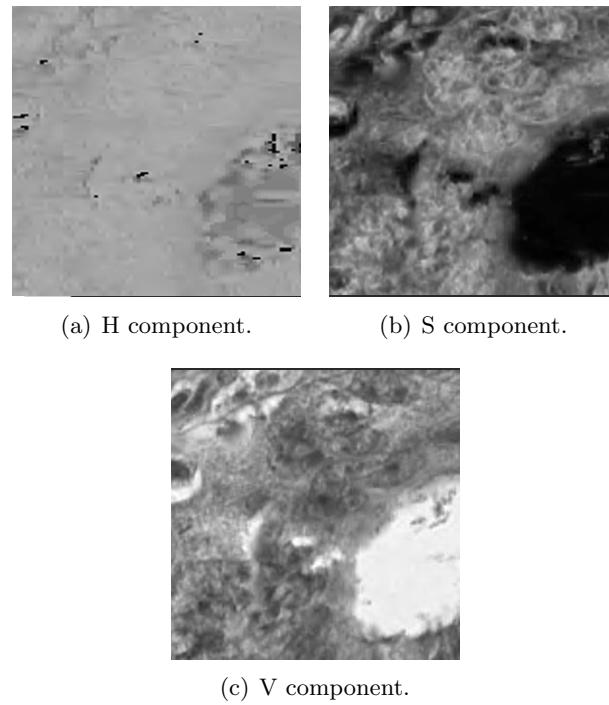


FIGURE 3.6: Decomposition of a image in the HSV components.

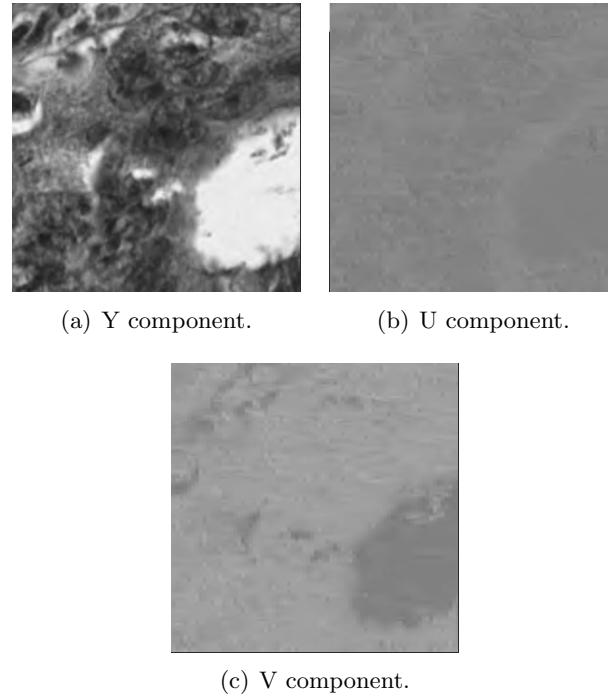


FIGURE 3.7: Decomposition of a image in the YUV components.

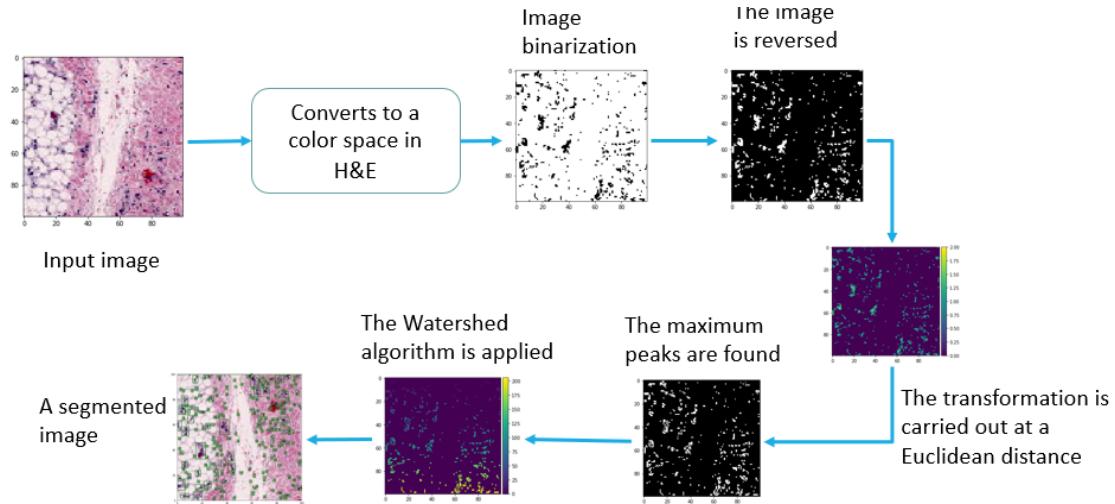


FIGURE 3.8: Nuclei segmentation.

After obtaining the characteristics of histopathology images, we need to mine useful data. Thus, we create a clean dataframe on python, using *pandas* library to prepare features to being spare for next classification algorithms. In order to perform machine

learning, the categorical variables in our dataset are usually regarded as discrete entities and coded as feature vectors. "Dirty" unorganized data will produce categorical variables with redundancy: multiple categories reflect the same entity [38]. In this database, this problem is solved by converting categorical data into numbers. We apply a **one hot encoding**. It is a representation of categorical variables as binary vectors, where each integer value is represented as a binary vector that is all zero values except the index of the integer, which is marked with a 1. A one hot encoding allows the representation of categorical data to be more expressive, and causes that redundancy to learning algorithms will be deleted, bringing significant benefits to our model performance.

Additionally, we reduce size of dataset by dropping variables which do not represent relevant data to build our model. Along with this, we evaluate if obtained dataset is balanced, it means that each class have equal number of samples and allow to build a most reliable model.

Finally, we separate this clean dataset into 80% for *training set*, and 20% for *test set*. Both are filled with random values of the clean dataset. This methodology allows to find the appropriate performance of the next classification algorithms, where, in the end, we seek to have a generalization of the detection problem and obtain a classifier model as output of each algorithm.

### 3.2.3 Cancer Classification Algorithms

As a third stage, we build a classification of a colorectal cancer diagnosis model based on machine-learning algorithms. To begin, we need to define an input and output set to feed our classification algorithms. First set is conformed by the features dataset (already prepared), and second set consists of the labels of eight classes, previously described on section 2.1.

Four machine learning algorithms of *supervised type* are considered for applying at this image dataset. The algorithms are the Naive Bayes, Random Forest, Support Vector Machine (SVM) and Multilayer Perceptron (MLP).

- **Naive Bayes Algorithm**

It is a classification technique based on Bayes' Theorem with an assumption of independence among predictors. The reason why it is called 'Naive' because it requires rigid independence assumption between input variables. In simple terms, a Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature. Naive Bayes Classifier can be trained easily and fast and can be used as benchmark model. Also, Naive Bayes requires a strong assumption of independent predictors, so when the model has a bad performance, the reason leading to that may be the dependence between predictors.

We deploy this algorithm on python using **Sklearn** library, constructing a classifier model with the command "*GaussianNB*" adjusted with default values (Ignoring prior probabilities of the classes and a portion of the largest variance of all features equal to  $1e - 9$ ).

- **Random Forest Algorithm**

Random forest is a supervised classification algorithm, which consists of a large number of individual decision trees that operate as an ensemble. Each individual tree in the random forest spits out a class prediction and the class with the most votes becomes our model's prediction.

In this case, we also consider **Sklearn** library to set up the classifier model. The command "*RandomForestClassifier*" is adjusted with default values, except the maximum depth of the tree determined 50, and a value of zero (0) to the randomness of the bootstrapping<sup>1</sup>.

- **Support Vector Machine Algorithm (SVM)**

Support Vector Machine (SVM) analyzes data for regression and classification analysis. The objective of SVM is to find a surface in the n-space that separates the space in regions. The shape of the surface depends on the type used to make the separation. The surface is usually known as kernel. In the simplest form, the kernel is lineal and represents a hyperplane the space of dimension  $n$ .

At this part, we continue using **Sklearn** library. In order to ensemble a SVM model

---

<sup>1</sup><https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>

we make use of command `"make_pipeline"` adjusted to default values and `"gamma SVC"` value set on `auto`.

- **Multilayer Perceptron Algorithm (MLP)**

It is a artificial neural network as explained in section 1.3, consisting of many layers such that it allows to solve problems that are not linearly separables. It consists of an input layer, intermediate layers and an output layer. The set of the layers represent non linear functions

To set up this model, we call **TensorFlow** library and use a series of commands to define the architecture of our MLP. We establish a 3-layer-sequential model configured with parameters to training our data exposed on Table 3.3.

Parameter	Value
Dense 1	100
Dense 2	50
Dense 3	8
Dense 1,2 - Activation Function	ReLU
Dense 3 - Activation Function	Softmax
Epochs	30
Training Set	90%
Validation Set	10%

TABLE 3.3: Parameters of MLP model proposed (Training set).

### 3.2.4 Assessment Metrics

To measure their performance of presented machine-learning algorithms applied, we use the concepts of accuracy, Precision, Recall and F1-Score were used as defined as follows:

- **Accuracy** In general, the accuracy metric measures the ratio of correct predictions

over the total number of instances evaluated. Its values are determined by Equation 3.1.

$$\text{Accuracy} = \frac{\text{CorrectPredictions}}{\text{TotalPredictions}} \quad (3.1)$$

- **Precision**

*Precision* is used to measure the positive patterns that are correctly predicted from the total predicted patterns in a positive class, defined by Equation 3.2.

$$\text{Precision} = \frac{\text{TruePositive}}{\text{PredictedPositive}} \quad (3.2)$$

- **Recall**

*Recall* is used to measure the fraction of positive patterns that are correctly classified. Its values can be obtained using Equation 3.3.

$$\text{Recall} = \frac{\text{TruePositive}}{\text{ActualPositive}} \quad (3.3)$$

- **F1-Score**

*F1-Score* is the harmonic mean between *Precision* and *Recall*, where the average is calculated per label (class), then, averaged across all labels  $Q$  (classes). If  $p_j$  and  $r_j$  are the *Precision* and *Recall* for all  $\lambda_j \in h(x_i)$  from  $\lambda_j \in y_i$ , the *F1-Score* is represented by Equation 3.4.

$$\text{F1Score} = \frac{1}{Q} \sum_{j=1}^Q \frac{2p_j r_j}{p_j + r_j} \quad (3.4)$$

## Chapter 4

# RESULTS AND DISCUSSION

In this chapter, models applied presents the result of cancer classification approach in section 4.1, but also we describes an experimental algorithm to reduce dimension of the dataset before classification. We apply the same classification model on two different datasets to discuss and compare its performance and measure robustness of the built model. Since, we can obtain a better approach after comparing its results. Finally, in section 4.2, we compare classification performances on these two different datasets based on metrics approaches previously commented.

Before presenting results, it is important to introduce the two datasets used.

- **First Dataset - 8 Classes:** This dataset consists of 8 classes previously described on 3.2.1, containing 5000 rows and 161 of visual features: 15 features correspond to color features, and around 150 to morphological. These classes were selected from the most representative images from our dataset.
- **Second Dataset - 2 Classes:** The dataset used consists of 165 images derived from 16 *H&E* stained histological sections of stage *T3* or *T4* colorectal adenocarcinoma. Each section belongs to a different patient, and the sections were processed in the laboratory on different occasions. Digitization of these histological sections into whole slide images (WSI) was performed using a Zeiss MIRAX MIDI slide scanner with  $0.465\mu m$  pixel resolution. The WSIs were subsequently rescaled to a pixel resolution of  $0.620\mu m$  (equivalent to a  $20 \times$  objective magnification).

A total of 52 visual fields of benign and malignant areas across the entire set of WSIs were selected to cover the widest possible tissue variety. An expert pathologist (DRJS) then rated each visual field as "benign" or "malignant", according to the general glandular architecture. The pathologist also delineated the boundary of each individual glandular object in that visual field.

## 4.1 Cancer Classification Results

### 4.1.1 Experimental dimensionality reduction algorithm

This **Experimental stage** consists in applying a dimensionality reduction algorithm to reduce the size of the data without significantly reducing the useful information based on the variance analysis and finding the direction in which the variance changes in a maximum way. In our case, we use the Principal Component Analysis (PCA), which is one of the most used algorithms on CAD field [40] [39]. We pretend to decrease computational load searching an optimal dimension that contains most useful information. However, after dimensionality reduction, every classification model designed provide less than 0.6 values for every performance measurements. Then, we assume that reduced dataset does present abnormalities on used features, which affect classification and suffer a high level of information losses in the distribution of tissue components, especially on morphological features.

Because of low performance using the PCA - dimensionality reduction algorithm, we use both whole datasets on next stage.

### 4.1.2 8-Classes Dataset - Clasification Results

The **four machine learning algorithms** are applied to first dataset (8 classes). The aim here is to compare performance through SVM, MLP, Random Forest and Naive Bayes algorithms.

The following Table 4.1 presents results of performance for the algorithms. It should be noted that Random Forest algorithm attains a better result than other three algorithms when compared by Accuracy, Precision and Recall. When compared by F1-Score the best

performance is carried out by the Support Vector Machine (SVM). Nevertheless, every algorithm reach a performance above 0.8, leaving the Accuracy value of MLP with 0.89 as the lowest value.

Algorithm	Accuracy	Precision	Recall	F1-Score
Support Vector Machine (SVM)	0,98	0,98	0,98	0,98
Random Forest	0,94	0,95	0,94	0,95
Naive-Bayes	0,93	0,93	0,93	0,93
Multilayer Perceptron (MLP)	0,89	0,92	0,89	0,89

TABLE 4.1: Result of the performance for the 4 algorithms - 8-Classes Dataset

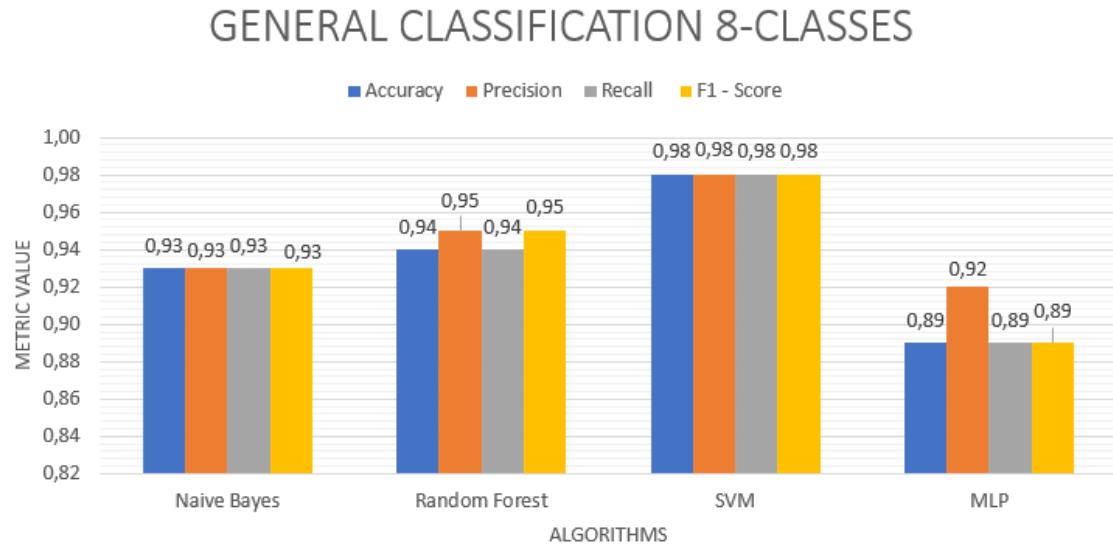


FIGURE 4.1: General Performance of the 4 algorithms - 8-Classes Dataset.

The performance of the four algorithms is also computed taking each class into account. The result of the four metrics of general classification of each algorithm is shown in Table 4.1. From this table we concluded that the best classification is carried out by the SVM and Random Forest above 0.93. In contrast, results of MLP and Naive Bayes have a mean value of every measurement below 0.93.

Subsequently, we show in Table 4.2 the performance of F1-score for each class classified by each algorithm. It is important to include precision and recall, to appreciate proportions of F1-score and detail its relation between classes.

Furthermore, in Table 4.2 there are some details about which class represents a challenge to our models. Histopathology images have visual patterns with particularities that make difficult its analysis. Some of these challenges may be found in sixth class for MLP algorithm, where precision is 0.65. This means MLP is not precise to predict a true positive value. However, there are two values in the Table 4.2 under 0.7 (Third and sixth class).

Likewise, SVM represents a reliable computational diagnostic tool by evaluating its quantitative measures. In most cases, proposed metrics achieve a near value to 100% of F1-score. In short, SVM performs a model capable of classifying correctly each class of dataset.

Following, it is possible to say that Random Forest algorithm is, a partially reliable computational diagnosis tool. Although this algorithm performs a value below 0.9 only in 3 cases, it is capable of classifying this large set of features and demonstrate that requires less computational load to achieve a 100% value of precision or recall on classes apparently easy to distinguish.

Nevertheless, Naive Bayes algorithm also has a reliable performance with metrics above of 0.8. Similar to MLP, Naive Bayes presents some weakness to classify sixth class. Finally, MLP presents an acceptable performance, but not enough to make it a reliable CAD tool.

We present too, Figure 4.2 where it can be seen the performance measured by precision for each class to clarify results between algorithms. Same as previously chart, Figure 4.3 shows performance measured by recall and finally Figure 4.4 by F1-Score. At this time, we notice in F1-Score the accurate at the class 4 that were classified properly with the algorithms of Random Forest, Naive-Bayes, SVM and MLP with very high values of Precision, Recall and F1-Score which are around 0.99.

Additionally, we seek to evaluate confusion matrix of all the algorithms to demonstrate classification challenges and similarities through classes.

Algorithm	Class	Precision	Recall	F1-Score
SVM	1	0.98	0.98	0.98
	2	0.95	0.99	0.97
	3	0.98	0.94	0.96
	4	0.97	0.99	0.98
	5	0.97	0.98	0.97
	6	0.98	0.97	0.98
	7	0.99	0.98	0.99
	8	0.98	0.98	0.98
Random Forest	1	0.71	0.95	<b>0.82</b>
	2	0.98	0.91	0.95
	3	0.99	0.90	0.94
	4	0.96	1.00	0.98
	5	1.00	0.96	0.98
	6	0.99	0.93	0.96
	7	1.00	0.96	0.98
	8	1.00	0.93	0.97
Naive-Bayes	1	0.99	0.94	0.96
	2	0.92	0.91	0.91
	3	0.87	0.91	0.89
	4	0.96	1.00	0.98
	5	0.93	0.91	0.92
	6	0.84	0.90	0.87
	7	0.98	0.94	0.96
	8	0.98	0.97	0.98
MLP	1	0.99	1.00	1.00
	2	<b>0.80</b>	1.00	0.89
	3	1.00	<b>0.62</b>	<b>0.77</b>
	4	1.00	1.00	1.00
	5	0.93	0.88	0.90
	6	<b>0.65</b>	1.00	<b>0.79</b>
	7	1.00	<b>0.71</b>	0.83
	8	0.97	0.96	0.97

TABLE 4.2: Result of the performance for the four algorithms at the level of Class.

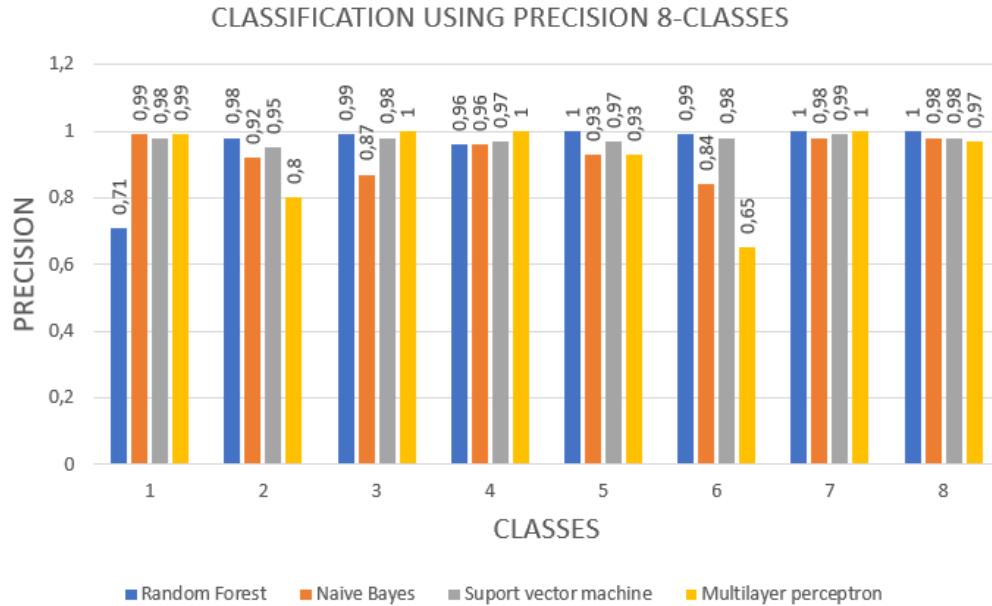


FIGURE 4.2: Performance of the 4 algorithms using Precision - 8-Classes Dataset.

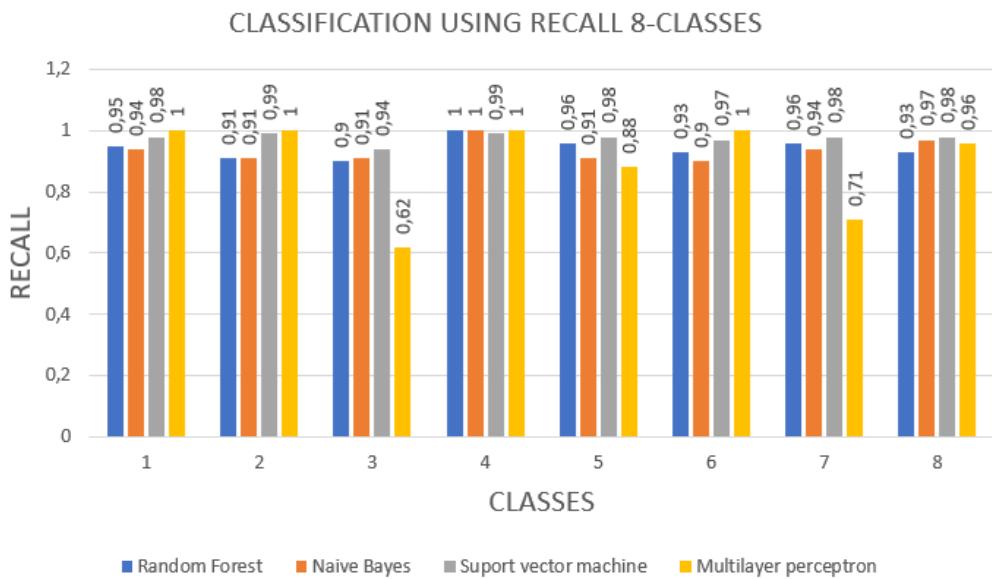


FIGURE 4.3: Performance of the 4 algorithms using Recall - 8-Classes Dataset.

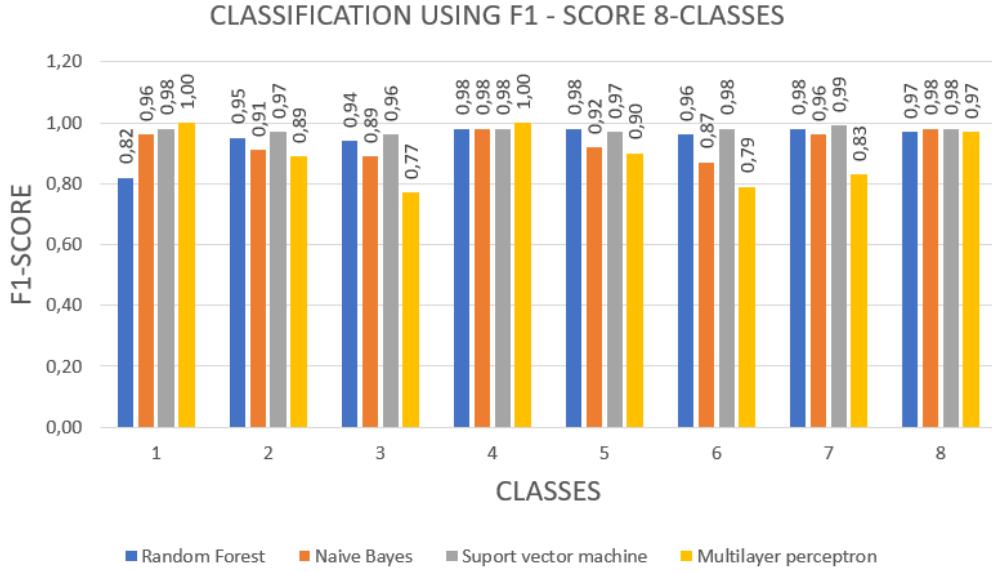


FIGURE 4.4: Performance of the 4 algorithms using F1-Score - 8-Classes Dataset.

First, Figure 4.5(a) describes the confusion matrix of MLP, it can be seen that *class 6* is easily confused with *class 7*. Also, *class 3* has an unexpected value, it could mean that Debris(3) and Stroma (7)

while other classes achieves different values, it means model can classify acceptable through 8-classes.

In the same way, the confusion matrix of Naive Bayes algorithm, on Figure 4.5(b), shows a strong classification where we can see clearly differences between classes. For example, *class 2* and *class 3* are nearly confused, but it does not represent a convincing confusing between classes.

On Figure 4.5(d) Support vector machines tend to get more confused when classifying the Adipose class, classifying it as Debris. Unlike other classes, SVM algorithm represents a robustness on classify each class. Subsequently, Figure 4.5(c) shows alike results obtained with the random forest algorithm, but *class 1* is slightly confused with the same class (*class 3*). This can happen because there is little color variability between the Adipose and Debris classes.

After comparison using general classification and confusion matrix, it is evident that the best algorithm to classify has been the support vector machines having a F1-score of 98%.

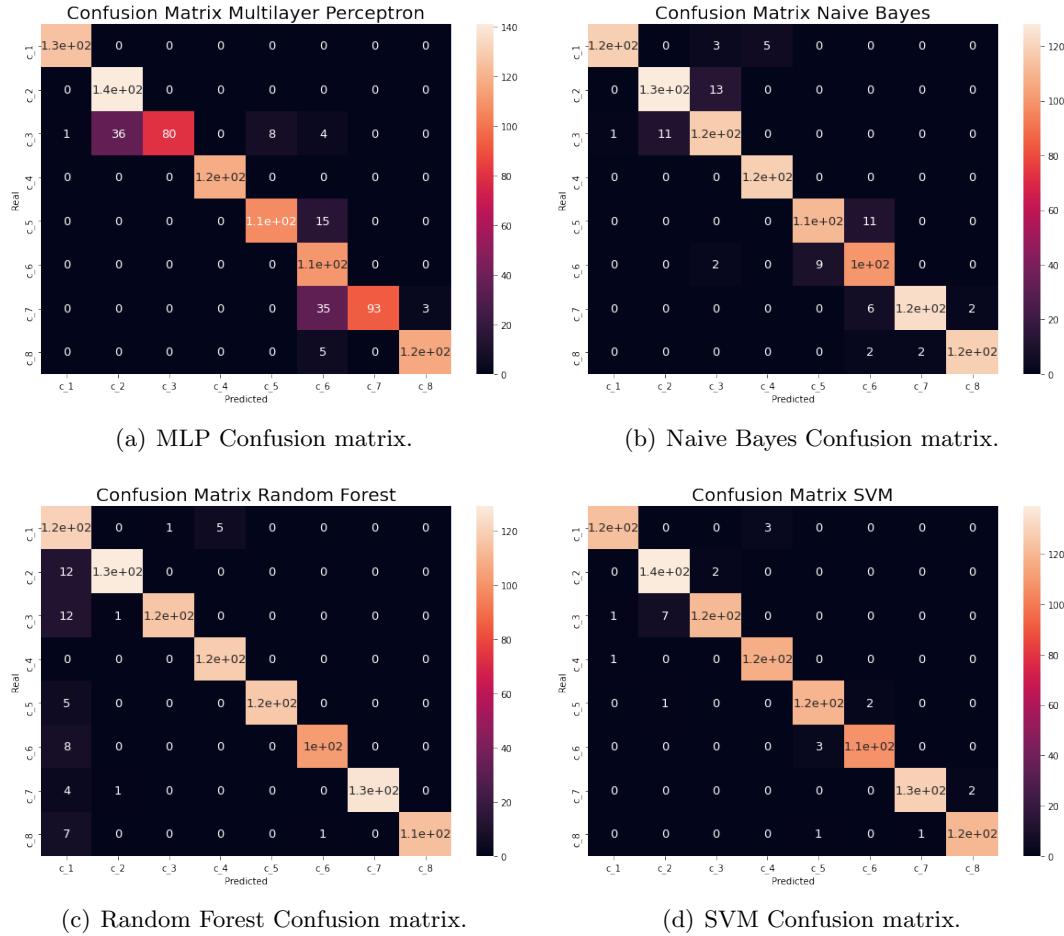


FIGURE 4.5: Confusion matrix of four algorithms - 8-Classes Dataset.

Both Random Forest algorithm and Support Vector Machine are better at classifying these types of features on tissue 8-classes.

#### 4.1.3 Binary Dataset - Classification Results

In this case, performance of SVM, MLP, Random Forest and Naive Bayes algorithms trained on **second dataset (2 classes)** is compared and evaluated.

The performance of the four algorithms is computed to classify into "cancer" and "non-cancer". Overall, results show a great performance for every algorithm with metrics above 0.8. From Figure 4.6 we concluded that the best classification is carried out by the SVM and Random Forest above 0.93. In contrast, results of MLP and Naive Bayes have a mean

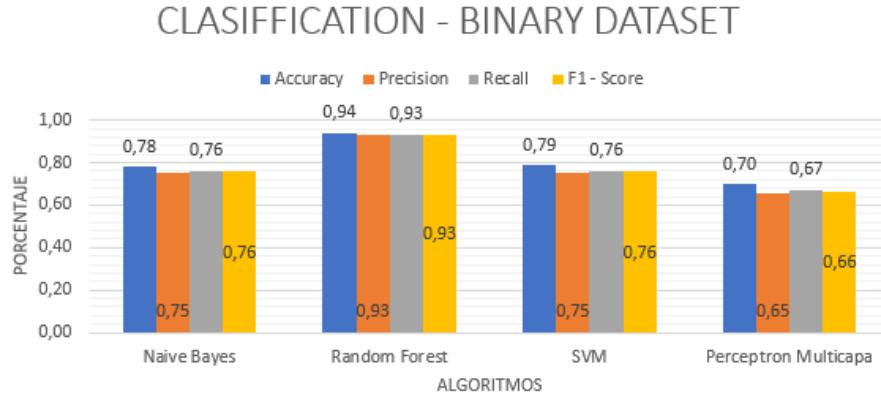


FIGURE 4.6: Performance of the 4 algorithms - Binary Dataset

value of every measurement below 0.93.

In feature selection, a subset of features like color and morphology data can be selected to demonstrate that represents relevant characteristics to being classify into a cancer diagnosis. However, a low size of dataset carries weight training models and giving a properly output, making such an exhaustive process to search a practical and big tissue dataset.

Following, in case of the random forest algorithm where all its performances metrics are above 92%, while the other algorithms are below 79% of accuracy and 76%. Then, as a result we may say SVM performances high values of every metric, near to achieve a 100% of correct classifications. Next algorithms sorted by metrics is Random Forest, Naibe bayes and at last, MLP.

Also, we evaluate confusion matrix of algorithms to see its performance to distinguish between 'cancer' or 'not cancer'.

Differing from Subsection 4.1.2, Binary dataset represents 2-classes: "cancer" or "not cancer" as shown in Figure 4.7(a) which describes the confusion matrix of MLP, it can be seen that is easily confused between classes. Alternatively, On Figure 4.7(d) Support vector machines tends to get a better approach to classify classes. Also, Figure 4.7(b) Naive Bayes shows a similar result as the random forest algorithm on Figure 4.7(c), both are confused at same level. After comparison using general classification and confusion matrix, it is evident that the best algorithm to classify has been the Random Forest having a score of 92%.

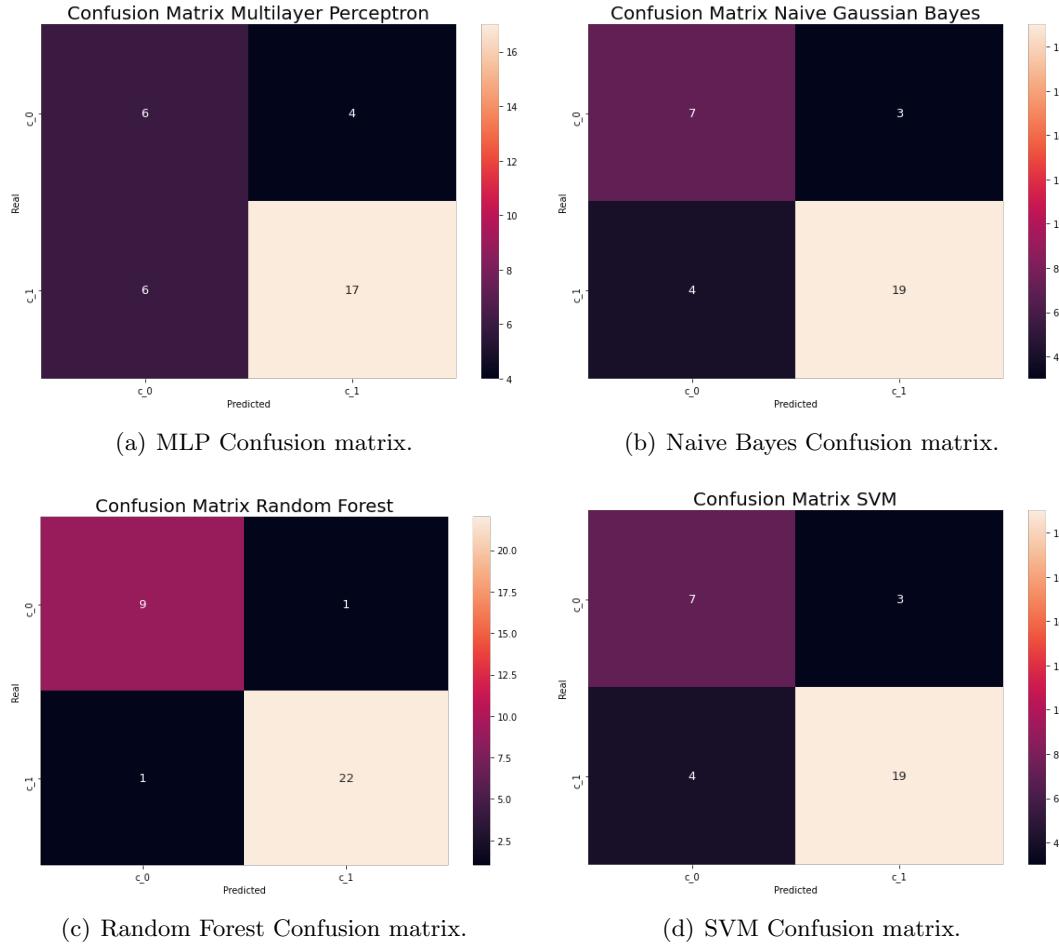


FIGURE 4.7: Confusion matrix of four algorithms - Binary Dataset.

## 4.2 Comparison

In this section, we present differences between the classification performance on two datasets.

Firstly, we compare Naive Bayes and MLP algorithm results from Figure 4.1 and Figure 4.6, its metric values are similar which are not above 0.8. Also, it means that these models could get a better adjust, for example, we could adjust the MLP model adding a new layer or changing activation function of final layer. Similarly, on Naive Bayes algorithm, it could be useful replacing "Gaussian" model and replacing with "Multinomial" configuration. It is important to emphasize that second dataset is relatively smaller, which represents the reason of the difficulty to have better training (misclassification).

Subsequently, from the same Figures 4.1 and 4.6, when comparing Random Forest and

SVM, it is easy to distinguish classes from each other on 8-class dataset using SVM. However, SVM algorithm faces a low level of classification on binary dataset in contrast to 8-class dataset, taking its metric values from 0.9 with its 8-classes classification performance to 0.79 in second dataset. On the other hand, Random Forest deploys a good performance with both datasets. In each case, its metric values stand above 0.9 and introduce a clear difference through classes on confusion matrix.

After analyzing metrics, we propose a way to prevent misclassification of color-morphological features, especially on second dataset. It is necessary to increase training samples. This will force the model to learn more number of samples and increase the number of models trained on the path.

At last, current results shows the Random Forest Algorithm performed the best, when applied to large or small dataset, considering the number of features detected properly. The MLP approach indicates its worst-case performance in 8-classes and 2 classes dataset.

## Chapter 5

# CONCLUSIONS AND FUTURE WORK

The Random Forest algorithm has best performance with respect to Accuracy, Precision and Recall with the values 0.6457, 0.6479 and 0.6429 respectively. For F1-Score the best performance is carried out by Naive-Bayes algorithm followed by Random Forest with values 0.5752 and 0.5235 respectively.

The proposed model is tested on two different datasets to demonstrate the computational load and classification capacity is more effective to keep the accuracy high when less training data are used for learning. This is due to the ability of the random forest algorithm to increase the learner's generalization ability by increasing the size and variation of the training data.

Different types of experiments can be carried out where an improvement of the configuration is done, thus making the algorithms more precise and accurate. Also, it is planned to carry out experiments with more images of this type by applying a binary classification. We also plan to compare other machine learning methods in the near future. We may explore the use of other adjust methods that makes a stronger model such as increasing depth in Random forest algorithm. Additionally, SVM can be improved using a sequence of more appropriate adjustment on its scaler capable of selecting points and such feature vectors more exactly, even with a reduced dataset.

Traditionally, the same methodology has been applied to classify images with cancer but with characteristics of texture, in this work it was possible to determine that the color and morphological characteristics are also useful for this purpose, opening another possibility for the study of this type of images. Including color and morphological features to colorectal cancer classification allows to have high algorithms performance for certain classes.

This thesis successfully addresses the issue of having limited labeled training data in the domain of histopathological tissue image classification. To this end, it presents a nuclei classification algorithm that performs a F1-score higher than 0.9 on Random Forest approach.

Finally, Random forest method performs a precise classification of cancer tissues. We can use its results to apply on a higher dataset with undesirable changes like noise, and we could improve model design to be capable of create an overall structure of potential results. This may benefit the CAD in automated tasks with the use of less features to extract and obtain a truly diagnosis type on posterior colorectal researches.

Although it is specifically designed for histopathological images of colon tissue, the proposed method may be used in different types of images and different types of organizations. This can also be considered as the future research direction of this article.

# Bibliography

- [1] Marr, D. (1982). Vision: A computational investigation into the human representation and processing of visual information. San Francisco, CA: W.H. Freeman.
- [2] J. Thevenot, M. B. López and A. Hadid, "A Survey on Computer Vision for Assistive Medical Diagnosis From Faces," in IEEE Journal of Biomedical and Health Informatics, vol. 22, no. 5, pp. 1497-1511, Sept. 2018, doi: 10.1109/JBHI.2017.2754861.
- [3] Richard Szeliski. 2010. "Computer Vision: Algorithms and Applications" (1st. ed.). Springer-Verlag, Berlin, Heidelberg.
- [4] Kunio Doi, Computer-aided diagnosis in medical imaging: Historical review, current status and future potential, Computerized Medical Imaging and Graphics, Volume 31, Issues 4–5, 2007, Pages 198-211, ISSN 0895-6111, <https://doi.org/10.1016/j.compmedimag.2007.02.002>.
- [5] Bluteau, Ryan, "Obstacle and Change Detection Using Monocular Vision" (2019). Electronic Theses and Dissertations. 7766. <https://scholar.uwindsor.ca/etd/7766>
- [6] Hidefumi Kobatake, Future CAD in multi-dimensional medical images: – Project on multi-organ, multi-disease CAD system –, Computerized Medical Imaging and Graphics, Volume 31, Issues 4–5, 2007, Pages 258-266.
- [7] A. Mhalla, T. Chateau, S. Gazzah and N. E. B. Amara, "An Embedded Computer-Vision System for Multi-Object Detection in Traffic Surveillance," in IEEE Transactions on Intelligent Transportation Systems, vol. 20, no. 11, pp. 4006-4018, Nov. 2019, doi: 10.1109/TITS.2018.2876614.

- [8] Sarah Behrens, Hendrik Laue, Matthias Althaus, Tobias Boehler, Bernd Kuemmerlen, Horst K. Hahn, Heinz-Otto Peitgen, "Computer assistance for MR based diagnosis of breast cancer: Present and future challenges, Computerized Medical Imaging and Graphics", Volume 31, Issues 4–5, 2007, Pages 236-247,ISSN 0895-6111,<https://doi.org/10.1016/j.compmedimag.2007.02.007>.
- [9] Sivic, Josef, C. L. Zitnick and R. Szeliski. "Finding People in Repeated Shots of the Same Scene." BMVC (2006).
- [10] Yang, Ming-Hsuan Kriegman, David Ahuja, Narendra. (2002). Detecting Faces in Images: A Survey. Pattern Analysis and Machine Intelligence, IEEE Transactions on. 24. 34 - 58. [10.1109/34.982883](https://doi.org/10.1109/34.982883).
- [11] Davies, E. (2017). Computer Vision: Principles, Algorithms, Applications, Learning.
- [12] Gu, Jiuxiang Wang, Zhenhua Kuen, Jason Ma, Liyang Shahroudy, Amir Shuai, Bing Liu, Ting Wang, Xingxing Wang, Gang. (2015). Recent Advances in Convolutional Neural Networks. Pattern Recognition. [10.1016/j.patcog.2017.10.013](https://doi.org/10.1016/j.patcog.2017.10.013).
- [13] Alexander Selvikvåg Lundervold, Arvid Lundervold, An overview of deep learning in medical imaging focusing on MRI, Zeitschrift für Medizinische Physik, Volume 29, Issue 2, 2019, Pages 102-127, ISSN 0939-3889, <https://doi.org/10.1016/j.zemedi.2018.11.002>.
- [14] Shen D, Wu G, Suk HI. Deep Learning in Medical Image Analysis. Annual Review of Biomedical Engineering. 2017 Jun;19:221-248. DOI: [10.1146/annurev-bioeng-071516-044442](https://doi.org/10.1146/annurev-bioeng-071516-044442).
- [15] Ian Goodfellow and Yoshua Bengio and Aaron Courville. "Deep Learning". MIT Press. 2016. <http://www.deeplearningbook.org>.
- [16] J. Ker, L. Wang, J. Rao and T. Lim, "Deep Learning Applications in Medical Image Analysis," in IEEE Access, vol. 6, pp. 9375-9389, 2018, doi: [10.1109/ACCESS.2017.2788044](https://doi.org/10.1109/ACCESS.2017.2788044).
- [17] Maier, Andreas Syben, Christopher Lasser, Tobias Riess, Christian. (2019). A gentle introduction to deep learning in medical image processing. Zeitschrift für Medizinische Physik. 29. [10.1016/j.zemedi.2018.12.003](https://doi.org/10.1016/j.zemedi.2018.12.003).

- [18] Michael A. Nielsen, "Neural Networks and Deep Learning", Determination Press, 2015.  
<http://neuralnetworksanddeeplearning.com/>
- [19] G. Litjens et al. "A survey on deep learning in medical image analysis" [Online].2017.  
<https://arxiv.org/abs/1702.05747>
- [20] Cooper LA, Carter AB, Farris AB, et al. Digital Pathology: Data-Intensive Frontier in Medical Imaging: Health-information sharing, specifically of digital pathology, is the subject of this paper which discusses how sharing the rich images in pathology can stretch the capabilities of all otherwise well-practiced disciplines. Proc IEEE Inst Electr Electron Eng. 2012;100(4):991-1003. doi:10.1109/JPROC.2011.2182074
- [21] Gurcan MN, Boucheron LE, Can A, Madabhushi A, Rajpoot NM, Yener B. Histopathological image analysis: a review. IEEE Rev Biomed Eng. 2009;2:147-171. doi:10.1109/RBME.2009.2034865
- [22] Aeffner F, Adissu HA, Boyle MC, et al. Digital Microscopy, Image Analysis, and Virtual Slide Repository. ILAR J. 2018;59(1):66-79. doi:10.1093/ilar/ily007
- [23] Irshad, Humayun Veillard, Antoine Roux, Ludovic Racoceanu, Daniel. (2014). Methods for Nuclei Detection, Segmentation and Classification in Digital Histopathology: A Review Current Status and Future Potential. IEEE reviews in biomedical engineering. 7. 97-114. 10.1109/RBME.2013.2295804.
- [24] Alok Kumar Jain and Shyam Lal. Feature Extraction of Normalized Colorectal Cancer Histopathology Images. Ambient Communications and Computer Systems. Springer Singapore. 2019.
- [25] Ethem Alpaydin. 2010. Introduction to Machine Learning (2nd. ed.). The MIT Press.
- [26] Konstantina Kourou, Themis P. Exarchos, Konstantinos P. Exarchos, Michalis V. Karamouzis, Dimitrios I. Fotiadis, Machine learning applications in cancer prognosis and prediction, Computational and Structural Biotechnology Journal, Volume 13,2015, Pages 8-17, ISSN 2001-0370, <https://doi.org/10.1016/j.csbj.2014.11.005>.
- [27] M. Gupta and B. Gupta, "A Comparative Study of Breast Cancer Diagnosis Using Supervised Machine Learning Techniques," 2018 Second International Conference on

- Computing Methodologies and Communication (ICCMC), Erode, 2018, pp. 997-1002, doi:10.1109/ICCMC.2018.8487537.
- [28] Okun O., Priisalu H. 2007. Random Forest for Gene Expression Based Cancer Classification: Overlooked Issues. In: Martí J., Benedí J.M., Mendonça A.M., Serrat J. (eds) Pattern Recognition and Image Analysis. IbPRIA 2007. Lecture Notes in Computer Science, vol 4478. Springer, Berlin, Heidelberg.
- [29] Kharya, Shweta; Soni, Sunita. Weighted naive bayes classifier: A predictive model for breast cancer detection. International Journal of Computer Applications, 2016, vol. 133, no 9, p. 32-37.
- [30] C. Ferri, J. Hernández-Orallo, R. Modroiu, An experimental comparison of performance measures for classification, Pattern Recognition Letters, Volume 30, Issue 1, 2009, Pages 27-38, ISSN 0167-8655, <https://doi.org/10.1016/j.patrec.2008.08.010>.
- [31] Hosny, K.M., Kassem, M.A. Foaud, M.M. Skin melanoma classification using ROI and data augmentation with deep convolutional neural networks. Multimed Tools Appl 79, 24029–24055. 2020. <https://doi.org/10.1007/s11042-020-09067-2>.
- [32] Demir, Cigdem and B. Yener. “Automated cancer diagnosis based on histopathological images : a systematic survey.” 2005.
- [33] He, Lei Long, L. Antani, Sameer Thoma, George. 2012. Histology image analysis for carcinoma detection and grading. Computer methods and programs in biomedicine. 107. 538-56. 10.1016/j.cmpb.2011.12.007.
- [34] J. Gil, H. Wu, B.Y. Wang, Image analysis and morphometry in the diagnosis of breast cancer, Microsc. Res. Techniq. 2002. pp 59. 109-118.
- [35] Yan Xu, Jun-Yan Zhu, Eric I-Chao Chang, Maode Lai, Zhuowen Tu. Weakly supervised histopathology cancer image segmentation and classification. Medical Image Analysis. Volume 18, Issue 3. 2014. Pages 591-604, ISSN 1361-8415, <https://doi.org/10.1016/j.media.2014.01.010>.
- [36] Kather, J. N., Weis, C. A., Bianconi, F., Melchers, S. M., Schad, L. R., Gaiser, T., Marx, A., Zöllner, F. G. 2016. Multi-class texture analysis in colorectal cancer histology. Scientific reports, 6, 27988. <https://doi.org/10.1038/srep27988>

- [37] Kather, Jakob Nikolas, et al. "Continuous representation of tumor microvessel density and detection of angiogenic hotspots in histological whole-slide images." *Oncotarget* 6.22. 2015.
- [38] Cerdà, Patricio, Gaël Varoquaux, and Balázs Kégl. "Similarity encoding for learning with dirty categorical variables." *Machine Learning* 107.8-10: 1477-1494. 2018.
- [39] J.J Dai, L. Lieu, D. Rocke, Dimension reduction for classification with gene expression microarray data, *Statistical Applications in Genetics and Molecular Biology* 5 (1) 1–15. 2006.
- [40] C. Truntzer, C. Mercier, J. Este've, C. Gautier, P. Roy, Importance of data structure in comparing two dimension reduction methods for classification of microarray gene expression data, *BMC Bioinformatics* 8 (90). 2007.
- [41] Miller, K.D., Nogueira, L., Mariotto, A.B., Rowland, J.H., Yabroff, K.R., Alfano, C.M., Jemal, A., Kramer, J.L. and Siegel, R.L. (2019), Cancer treatment and survivorship statistics, 2019. CA A Cancer J Clin, 69: 363-385. <https://doi.org/10.3322/caac.21565>
- [42] Alan C. Bovik, Chapter 4 - Basic Binary Image Processing, Editor(s): Al Bovik, The Essential Guide to Image Processing, Academic Press, 2009, Pages 69-96, ISBN 9780123744579, <https://doi.org/10.1016/B978-0-12-374457-9.00004-4>.
- [43] Gevers, T., Gijsenij, A., van de Weijer, J., Geusebroek, J.-M., Smeulders, W.c.b.A.W.M. and Bagdanov, A.D. (2012). Color Feature Detection. In *Color in Computer Vision* (eds T. Gevers, A. Gijsenij, J. van de Weijer and J.-M. Geusebroek). <https://doi.org/10.1002/9781118350089.ch13>
- [44] Akbar, B., Gopi, V. P., & Babu, V. S. (2015). "Colon cancer detection based on structural and statistical pattern recognition". 2015 2nd International Conference on Electronics and Communication Systems (ICECS).
- [45] Altunbay, D., Cigir, C., Sokmensuer, C., Gunduz-Demir, C. (2010). "Color Graphs for Automated Cancer Diagnosis and Grading", *IEEE Transactions on Biomedical Engineering*, 57 (3), 665-674.
- [46] Arévalo, J., Cruz-Roa, A. & González, F. (2014). "Histopathology Image Representation for Automatic Analysis: A State of the Art Review", *Revista Med*, 22 (2).

- [47] Dos Santos, L., Alves, L., Botazzo, G., Gonçalves, M., do Nascimento, M. & Azevedo, T. (2018). “Multidimensional and Fuzzy Sample Entropy (SampEnMF) for Quantifying H&E Histological Images of Colorectal Cancer”, Computers in Biology and Medicine, 103, 148-160.
- [48] Dundar, M., Badve, S., Bilgin, G., Raykar, V., Jain, R., Sertel, O. & Gurcan, M. (2011). “Computerized Classification of Intraductal Breast Lesions Using Histopathological Images”, IEEE Transactions on Biomedical Engineering, 58 (7), 1977-1984.
- [49] Hadid, A., Pietikainen, M. & Ahonen, T. (2004). “A Discriminative Feature Space for Detecting and Recognizing Faces”. IEEE Conference on Computer Vision and Pattern Recognition, 797-804.
- [50] Iftikhar, M., Hassan, M. & Alquhayz, H. (2016). “A Colon Cancer Grade Prediction Model using Texture and Statistical Features, SMOTE and mRMR”, 19th International Multi-Topic Conference.
- [51] Jørgensen, A., Rasmussen, A., Mäkinen, A., Andersen, S., Emborg, J., Røge, R. & Østergaard, L. (2017). “Using Cell Nuclei Features to Detect Colon Cancer Tissue in Hematoxylin and Eosin Stained Slides”, Cytometry Part A, 91, 785-793.
- [52] Kather, J., Weis, C., Bianconi, F., Melchers, S., Schad, L., Gaiser, T., Marx, A. & Zöllner, F. (2016). “Multi-class Texture Analysis in Colorectal Cancer Histology”, Nature Scientific Reports, 6, 27988.
- [53] Kothari, S., Phan, J., Moffitt, R., Stokes, T., Hassberger, S., Chaudry, Q., Young, A. & Wang, M. (2011). “Automatic Batch-Invariant Color Segmentation of Histological Cancer Images”, IEEE International Symposium on Biomedical Imaging: from Nano to Macro, 657-660.
- [54] Linder, N., Konsti, J., Turkki, R., Rahtu, E., Lundin, M., Nordling, F., Haglund C., Ahonen, T., Pietikainen, M. & Lundin, J. (2012). “Identification of tumor epithelium and stroma in tissue microarrays using texture analysis”, Diagnostic Pathology, 7 (22), 1-11.
- [55] Masood, K., & Rajpoot, N. (2009). “Texture based classification of hyperspectral colon

- biopsy samples using CLBP". 2009 IEEE International Symposium on Biomedical Imaging: From Nano to Macro.
- [56] Malik, (2019). "Colorectal Cancer Diagnosis from Histology Images: A Comparative Study". ArXiv, Computer Vision and Pattern Recognition. <https://arxiv.org/abs/1903.11210>.
- [57] C. Wang, J. Shi, Q. Zhang, and S. Ying, "Histopathological image classification with bilinear convolutional neural networks," in 2017 39th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC). IEEE, 2017, pp. 4050–4053
- [58] J. N. Kather, J. Krisam, P. Charoentong, T. Luedde, E. Herpel, C.-A. Weis, T. Gaiser, A. Marx, N. A. Valous, D. Ferber et al., "Predicting survival from colorectal cancer histology slides using deep learning: A retrospective multicenter study," PLoS medicine, vol. 16, no. 1, 2019
- [59] Abdelsamea, Mohammed M and Pitiot, Alain and Grineviciute, Ruta Barbora and Besusparis, Justinas and Laurinavicius, Arvydas and Ilyas, Mohammad (2019). "A cascade-learning approach for automated segmentation of tumour epithelium in colorectal cancer". Expert Systems with Applications, vol 118, pp.539–552.
- [60] Ozdemir, Erdem and Sokmensuer, Cenk and Gunduz-Demir, Cigdem. 2012. A resampling-based Markovian model for automated colon cancer diagnosis. IEEE transactions on biomedical engineering, vol 59, number1, pp.281–289.
- [61] Xu, Jun and Cai, Chengfei and Zhou, Yangshu and Yao, Bo and Xu, Geyang and Wang, Xiangxue and Zhao, Ke and Madabhushi, Anant and Liu, Zaiyi and Liang (2019). Multi-tissue Partitioning for Whole Slide Images of Colorectal Cancer Histopathology Images with Deeptissue Net. European Congress on Digital Pathology, pp100–108.
- [62] BackPropagation.tex - Computer Science and Engineering
- [63] Seth Eckhouse, Grant Lewison and Richard Sullivan (2008). Trends in the global funding and activity of cancer research. US National Library of Medicine, National Institutes of Health.

- [64] <https://www.cbsnews.com/news/global-cancer-rates-expected-to-hit-22-million-per-year-in-by-2025-who/#:text=WHO%20%2D%20CBS%20News-,Global%20cancer%20rates%20expected%20to%20hit%2022%20million,per%20year%20by%202030%3A%20WHO&text=By%202030%2C%20the%20number%20of,the%20World%20Health%20Organization%20reported.>