

# Bachelor thesis

Comparative Analysis of Classifiers for Breast Cancer Detection with  
Visualizations



Iván Sotillo del Horno



**UNIVERSIDAD AUTÓNOMA DE MADRID  
ESCUELA POLITÉCNICA SUPERIOR**



**Bachelor as Ingeniería Informática (modalidad bilingüe)**

**BACHELOR THESIS**

**Comparative Analysis of Classifiers for Breast  
Cancer Detection with Visualizations**

**Author: Iván Sotillo del Horno  
Advisor: Alejandro Bellogín Kouki**

**febrero 2024**

**All rights reserved.**

No reproduction in any form of this book, in whole or in part  
(except for brief quotation in critical articles or reviews),  
may be made without written authorization from the publisher.

© February 2024 by UNIVERSIDAD AUTÓNOMA DE MADRID  
Francisco Tomás y Valiente, nº 1  
Madrid, 28049  
Spain

**Iván Sotillo del Horno**

**Comparative Analysis of Classifiers for Breast Cancer Detection with Visualizations**

**Iván Sotillo del Horno**

PRINTED IN SPAIN

*A mi madre y a mi abuela, cuya lucha contra el cáncer de mama me ha inspirado a realizar este trabajo.*

DRAFT



# RESUMEN

---

Esta tesis presenta un análisis comparativo de clasificadores para la detección del cáncer de mama y el uso de Inteligencia Artificial Explicable (XAI) para interpretar los resultados. En la fase inicial se realizará la construcción y optimización de los modelos de clasificación, estos clasificadores analizarán los resultados de las biopsias de aguja fina y clasificarán las muestras como benignas o malignas.

Posteriormente, se realiza una comparación de rendimiento comparando métricas como la puntuación F1 o la *recall*. El objetivo es identificar el mejor clasificador de acuerdo a nuestras métricas. Una vez encontrado el mejor modelo, nos adentramos más en él para entender cómo funciona. Para esto, utilizaremos SHAP (SHapley Additive exPlanations), un método de XAI que nos permite ver la importancia de cada característica y cómo contribuyen a la decisión final del modelo. Esto nos permitirá no solo clasificar las muestras, sino también entender por qué el modelo ha tomado esa decisión, lo que puede ser un avance en la comprensión de los modelos de IA para fines médicos.

## PALABRAS CLAVE

---

Detección de Cáncer de Mama, Clasificadores, Análisis Comparativo, Interpretabilidad, SHAP, IA Explicable, Visualización





# ABSTRACT

---

This thesis presents a comparative analysis of base and ensemble classifiers for breast cancer detection and the use of eXplainable AI (XAI) to interpret the results. The initial phase involves constructing and optimizing the classifier models, these classifiers will analyze the results from fine needle biopsy aspirations and classify the samples as benign or malignant.

Following this, a performance comparison is conducted comparing metrics such as the F1 score or the recall. The aim is to identify the best classifier regarding our metrics. Once the best classifier model is found, we dive deeper into it to understand how it works. For this, we will use SHAP (SHapley Additive exPlanations), a method of XAI (eXplainable AI) that allows us to see the importance of each feature, and how they contribute to the final decision of the model. This will allow us to not only classify the samples but also to understand why the model has made that decision which can be a step forward in understanding AI models for medical purposes.

## KEYWORDS

---

Breast Cancer Detection, Classifiers, Comparative Analysis, Interpretability, SHAP, eXplainable AI, Visualization



# TABLE OF CONTENTS

---

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation .....	1
1.2	Objectives .....	1
1.3	Structure of the document .....	1
<b>2</b>	<b>State of the art</b>	<b>3</b>
2.1	Base and Ensemble Classifiers .....	3
2.1.1	Base Classifiers .....	3
2.1.2	Ensemble Classifiers .....	3
2.2	Classifier Optimization .....	3
2.3	Evaluation of Classifiers .....	3
2.4	Explainable AI .....	3
2.4.1	Explainable AI .....	3
2.4.2	SHAP .....	3
2.5	Web Application Development .....	3
<b>3</b>	<b>Design and Implementation</b>	<b>5</b>
3.1	Project Structure .....	5
3.2	Exploratory Data Analysis .....	5
3.2.1	Descriptive Statistics .....	5
3.2.2	Data Visualization .....	5
3.3	Data Preprocessing .....	5
3.3.1	Scaling and Normalization .....	5
3.3.2	Principal Component Analysis .....	5
3.4	Building and Optimizing Classifiers .....	5
3.5	SHAP Implementation .....	5
3.6	Web Application Development .....	5
<b>4</b>	<b>Experiments and Results</b>	<b>7</b>
4.1	Classifier Comparison .....	7
4.1.1	Base Classifiers Comparison .....	7
4.1.2	Ensemble Classifiers Comparison .....	7
4.1.3	Choosing the Best Classifier .....	7
4.2	SHAP Analysis .....	7
4.2.1	Global Interpretability .....	7

4.2.2 Local Interpretability .....	7
<b>5 Conclusions and Future Work</b>	<b>9</b>
5.1 Conclusions .....	9
5.2 Future Work .....	9
<b>Bibliography</b>	<b>11</b>
<b>Appendices</b>	<b>13</b>

DRAFT

# LISTS

---

**List of algorithms**

**List of codes**

**List of equations**

**List of figures**

**List of tables**

DRAFT



# INTRODUCTION

---

## 1.1. Motivation

Breast cancer is the most common cancer type among women [1]; in 2020, there were more than 2.26 million women diagnosed with breast cancer [1], being the second leading cause of death among women in the United States [2]. Early detection is a crucial step for improving survival rates. With the current analysis techniques of FNAB (Fine Needle Aspiration Biopsy), we have a sensitivity (ability of a test to identify positive cases correctly) of 0.927 [3]. Therefore, there is a need for a more accurate interpretation of those tests.

Machine learning is a branch of artificial intelligence that focuses on developing algorithms that can learn from data and extract patterns from it to be then able to generalize it to unseen data. In this case, we care about classifiers, whose potential is in the ability to learn from a dataset and then on unseen data being able to classify it as one class or another; in this case, we will be able to classify as benign or malign the results of a fine needle aspiration.

The potential of classifiers in breast cancer detection is immense. However, the effectiveness of the different classifiers can vary; this is why it is crucial to understand how each classifier works, how to tweak it, and how to make them as precise and effective as possible, which is the goal of this thesis.

Finding the best possible classifier for this problem would impact cancer detection tasks, facilitating healthcare professionals in their diagnostic responsibilities and, ultimately, improving patient outcomes.

## 1.2. Objectives

## 1.3. Structure of the document





## STATE OF THE ART

---

### 2.1. Base and Ensemble Classifiers

#### 2.1.1. Base Classifiers

#### 2.1.2. Ensemble Classifiers

### 2.2. Classifier Optimization

### 2.3. Evaluation of Classifiers

### 2.4. Explainable AI

#### 2.4.1. Explainable AI

#### 2.4.2. SHAP

### 2.5. Web Application Development



## DESIGN AND IMPLEMENTATION

---

- 3.1. Project Structure**
- 3.2. Exploratory Data Analysis**
  - 3.2.1. Descriptive Statistics**
  - 3.2.2. Data Visualization**
- 3.3. Data Preprocessing**
  - 3.3.1. Scaling and Normalization**
  - 3.3.2. Principal Component Analysis**
- 3.4. Building and Optimizing Classifiers**
- 3.5. SHAP Implementation**
- 3.6. Web Application Development**



# EXPERIMENTS AND RESULTS

---

## 4.1. Classifier Comparison

### 4.1.1. Base Classifiers Comparison

### 4.1.2. Ensemble Classifiers Comparison

### 4.1.3. Choosing the Best Classifier

## 4.2. SHAP Analysis

### 4.2.1. Global Interpretability

### 4.2.2. Local Interpretability



## CONCLUSIONS AND FUTURE WORK

---

### 5.1. Conclusions

### 5.2. Future Work





# BIBLIOGRAPHY

---

- [1] WCRF International, "Breast cancer statistics | World Cancer Research Fund International."
- [2] American Cancer Society, "Breast Cancer Statistics | How Common Is Breast Cancer?."
- [3] Y.-H. Yu, W. Wei, and J.-L. Liu, "Diagnostic value of fine-needle aspiration biopsy for breast mass: a systematic review and meta-analysis," *BMC Cancer*, vol. 12, p. 41, Jan. 2012.



# APPENDICES







Universidad Autónoma  
de Madrid