

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**

marzo 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, n^o 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.

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

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| 1.1 | Motivation | 1 |
| 1.2 | Objectives | 1 |
| 1.3 | Structure of the document | 1 |
| 2 | State of the art | 3 |
| 2.1 | Base and Ensemble Classifiers | 3 |
| 2.1.1 | Base Classifiers | 3 |
| 2.1.2 | Ensemble Classifiers | 5 |
| 2.2 | Evaluation of Classifiers | 7 |
| 2.3 | Classifier Optimization | 10 |
| 2.4 | Explainable AI | 11 |
| 2.4.1 | Explainable AI | 11 |
| 2.4.2 | SHAP | 11 |
| 2.5 | Web Application Development | 11 |
| 3 | Design and Implementation | 13 |
| 3.1 | Project Structure | 13 |
| 3.2 | Requirements | 15 |
| 3.2.1 | Functional Requirements | 15 |
| 3.2.2 | Non-Functional Requirements | 15 |
| 3.3 | Exploratory Data Analysis | 16 |
| 3.3.1 | Descriptive Statistics | 16 |
| 3.3.2 | Data Visualization | 17 |
| 3.4 | Data Preprocessing | 17 |
| 3.4.1 | Scaling and Normalization | 17 |
| 3.4.2 | Principal Component Analysis | 17 |
| 3.5 | Building and Optimizing Classifiers | 18 |
| 3.5.1 | Comparison of Classifiers | 19 |
| 3.5.2 | Optimization of Classifiers | 20 |
| 3.6 | SHAP Implementation | 21 |
| 3.7 | Web Application Development | 22 |
| 4 | Experiments and Results | 25 |
| 4.1 | Exploratory Data Analysis | 25 |

| | |
|---|-----------|
| 4.2 Classifier Comparison | 25 |
| 4.2.1 Base Classifiers Comparison | 25 |
| 4.2.2 Ensemble Classifiers Comparison | 25 |
| 4.2.3 Choosing the Best Classifier | 25 |
| 4.3 SHAP Analysis | 25 |
| 4.3.1 Global Interpretability | 25 |
| 4.3.2 Local Interpretability | 25 |
| 5 Conclusions and Future Work | 27 |
| 5.1 Conclusions | 27 |
| 5.2 Future Work | 27 |
| Bibliography | 30 |
| Acronyms | 31 |
| Appendices | 33 |

LISTS

List of algorithms

List of codes

| | | |
|-----|---|----|
| 3.1 | Building the pipeline | 19 |
| 3.2 | Obtaining the predicted probability of the pipeline and calculating the SHAP values ... | 22 |
| 3.3 | Plotting SHAP | 22 |

List of equations

| | | |
|-----|--|----|
| 2.1 | Balanced Accuracy | 9 |
| 2.2 | F Beta Score | 9 |
| 2.3 | Matthews Correlation Coefficient | 10 |
| 3.1 | Standardization | 17 |

List of figures

| | | |
|-----|-------------------------------------|----|
| 2.1 | Gradient Boosting | 5 |
| 2.2 | Random Forest | 6 |
| 2.3 | Bagging Classifier | 7 |
| 2.4 | Stacking Classifier | 8 |
| 2.5 | ROC Curve Example | 10 |
| 3.1 | Structure of the project | 13 |
| 3.2 | Fine Needle Aspiration Biopsy | 14 |
| 3.3 | Standardization | 18 |
| 3.4 | PCA | 18 |
| 3.5 | Web Application | 23 |
| 3.6 | Web Application | 24 |

List of tables

2.1 Confusion Matrix 8

2.2 Confusion Matrix Example 9

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

Classification is the process of taking input and assigning it to a class, it is a fundamental task in machine learning. There are different types of classification, but for this work, we are going to focus on binary supervised classification. Binary means that the classification will only be carried out between two classes, in our case, benign and malignant; supervised means that the training process will be performed with labeled data, that is, we will train that classifier with **Fine Needle Aspiration Biopsy (FNAB)** samples that have been labeled as benign or malignant by a pathologist. Our classifiers will not use the image directly, but rather the features extracted from the images.

A classifier is an algorithm that performs the aforementioned task, it receives an input and assigns it to a class. To train a classifier we must first have a dataset, this dataset will be divided into two parts, the training set, and the test set. The training set will be used to train the classifier, which means that the classifier will learn from this data trying to find patterns that allow it to classify the samples correctly. The test set will be used to evaluate the classifier, this set will be used to see how well the classifier performs with data that it has not seen before, this is done to ensure that the classifier is not overfitting (memorizing the training data) and that it can generalize to new data.

Classifiers can be split into two main groups, base classifiers and ensemble classifiers. Their main difference is that base classifiers are single classifiers, while ensemble classifiers are composed of multiple base classifiers. In this section, we will give an overview of the most common base and ensemble classifiers.

2.1.1. Base Classifiers

Base classifiers as just mentioned are single classifiers, they are the simplest form of classifiers, and they are the building blocks of ensemble classifiers. There are many base classifiers, but for this work, we are going to focus on the most common ones.

Logistic Regression

Logistic Regression (LR) is a probabilistic model that uses a logistic function to predict the probability of a binary outcome. The central idea is to find the best-fitting model to describe the relationship between the binary dependent variable and one or more independent variables.

Multilayer Perceptron

A Multilayer Perceptron (MLP) is a neural network composed of perceptrons, a perceptron is a simple model of a biological neuron, it takes several binary inputs and produces a single binary output using a weighted sum of the inputs and an activation function.

Stochastic Gradient Descent

Stochastic Gradient Descent (SGD) is not a classifier itself, but an optimization algorithm used instead of the standard gradient descent algorithm, it is used to minimize the loss function of a classifier. It is in this list because in Scikit-learn it can be used as a classifier by indicating the linear model to be used [4].

Support Vector Machine

Support Vector Machine (SVM) is a statistical model that separates data into classes by finding a hyperplane that separates the classes, since there can be multiple hyperplanes that do so, the model tries to find the one that maximizes the margin between the classes.

K-Nearest Neighbors

K-Nearest Neighbors (KNN) is a model that classifies the sample based on the majority of its neighbors, the number of neighbors, k , is a parameter that must be set by the user.

Linear and Quadratic Discriminant Analysis

Two different models, Linear Discriminant Analysis (LDA) and Quadratic Discriminant Analysis (QDA). They are used to find a linear or quadratic decision boundary that separates the classes.

Decision Tree

Decision Tree (DT) is a model that separates the data into classes by asking a series of questions, each question is based on the value of a feature, and the answer to the question determines the next question. It is a tree-like model, where the leaves are the classes.

Gaussian Naive Bayes

Gaussian Naive Bayes (GNB) is a probabilistic model that uses Naive Bayes' theorem with the assumption that the features are independent, given the class. It is called Gaussian because it assumes that the features are normally distributed.

2.1.2. Ensemble Classifiers

Ensemble classifiers are a type of classifier that combines the predictions of multiple base classifiers to improve the accuracy of the predictions. Some of the ensemble classifiers allow the base classifiers to be chosen, which means that we can use the best base classifiers to build the ensemble classifier. In this section, we are going to give an overview of the ensemble classifiers that we are going to use in this project.

Gradient Boosting

Gradient Boosting (GB) is an ensemble machine learning technique that combines the predictions from several models to improve the overall predictive accuracy. It works by first building a weak model (a model that is slightly better than random guessing) and then building a second model that corrects the errors made by the first model and so on [5], we can see a graphical representation of the GB in the Figure 2.1.

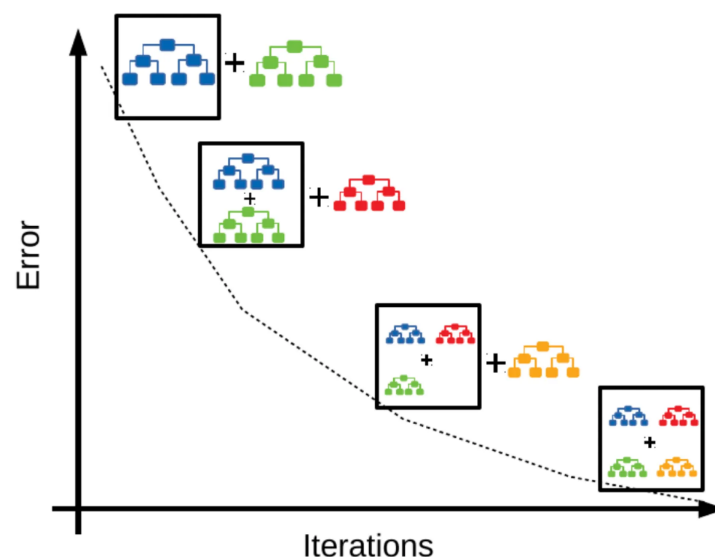


Figure 2.1: Graphical representation of Gradient Boosting [6]

AdaBoost

AdaBoost (AB) is similar to GB in the sense that it combines the predictions of multiple models but it does so by giving more weight to the misclassified data points in each iteration.

Random Forest

Random Forest (RF) is based on using DT as base classifiers. It builds multiple DT and combines their predictions, it is called random because it uses a random subset of the features to build the DT [7], this randomness is called bagging. We can see a graphical representation of the RF in the Figure 2.2.

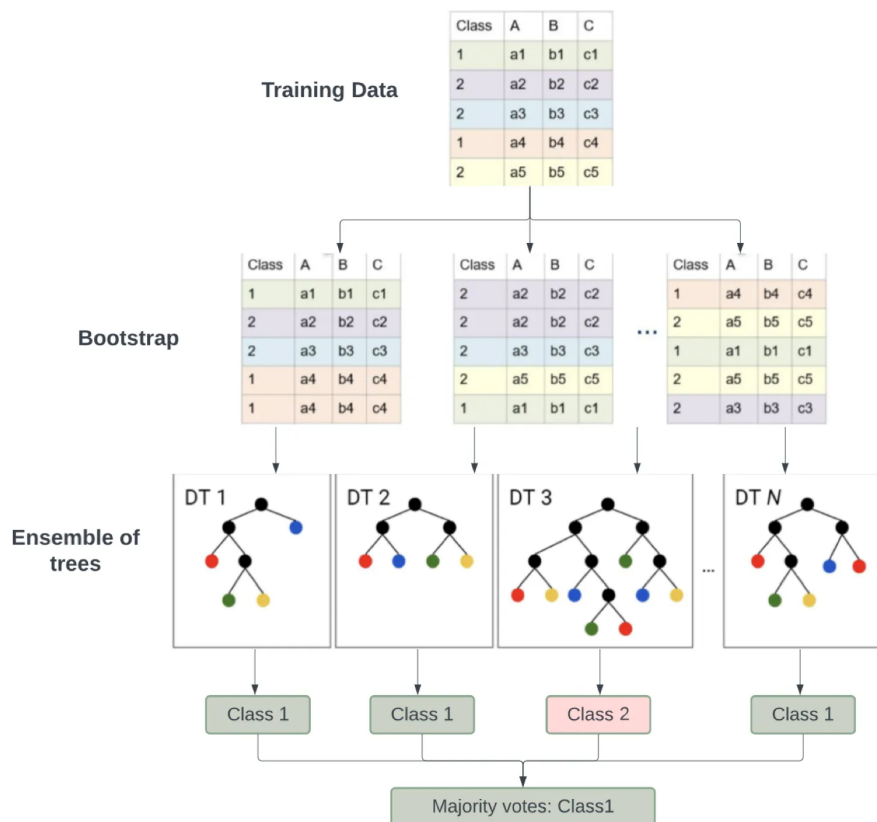


Figure 2.2: Graphical representation of Random Forest [8]

Voting Classifier

A Vote Classifier (VC) as its name suggests, is an ensemble classifier that combines the predictions of multiple classifiers by voting. It can be hard or soft voting, hard voting is the usual voting where the class that gets the most votes is the final prediction, and soft voting is when the average of the probabilities of the classes is used to make the final prediction.

Bagging Classifier

A **Bagging Classifier (BC)** is an ensemble classifier that builds multiple base classifiers and combines their predictions by averaging or voting, as it was mentioned before, in bagging the base classifiers are built using a random subset of the features [9], we can see a graphical representation of the BC in the Figure 2.3.

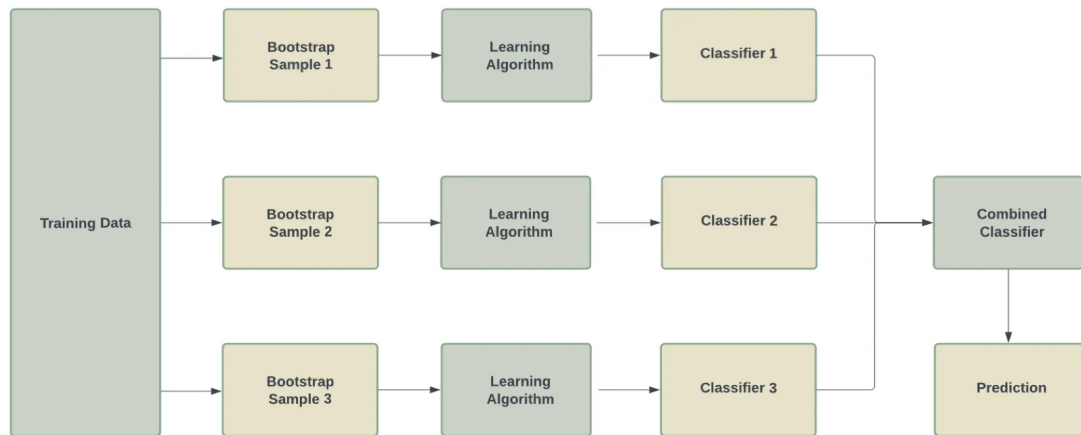


Figure 2.3: Graphical representation of Bagging [10]

Stacking Classifier

A **Stacking Classifier (SC)** is an ensemble that combines the base classifiers into a meta-classifier, that is, the sample to be classified goes through the base classifiers and then the predictions of the base classifiers are used as features to build the meta-classifier [11]. We can see a graphical representation of the SC in the Figure 2.4.

2.2. Evaluation of Classifiers

Now that we know what a classifier is, how it works, and the different types of classifiers, we need to know which classifier is the best for our problem, for this, we need to use some metric to evaluate the performance of the classifiers. Before we describe the metrics we need to understand the values that we are going to use to calculate these metrics, the values are:

- **True Positive (TP):** The number of instances that are positive and the classifier correctly classified as positive.
- **True Negative (TN):** The number of instances that are negative and the classifier correctly classified as negative.
- **False Positive (FP):** The number of instances that are negative and the classifier incorrectly classified as positive.
- **False Negative (FN):** The number of instances that are positive and the classifier incorrectly classified as negative.

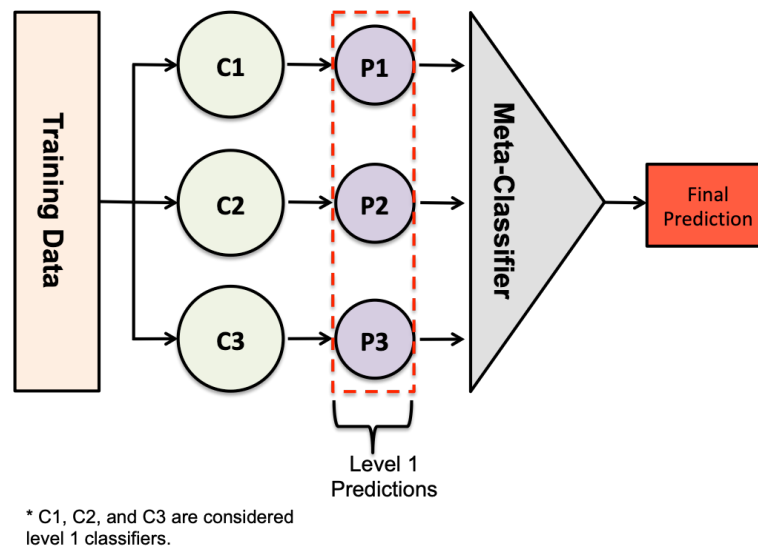


Figure 2.4: Graphical representation of Stacking [12]

Now that we have our building blocks for the metrics, we can describe the main four metrics:

- **Accuracy:** Proportion of correct predictions over the total number of instances evaluated. $Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$
- **Recall:** Proportion of **TP** results among the number of all positive instances. $Recall = \frac{TP}{TP+FN}$
- **Precision:** The proportion of **TP** results among the number of cases that were predicted to be positive. $Precision = \frac{TP}{TP+FP}$
- **Specificity:** The proportion of **TN** results among the number of cases that were actually negative. $Specificity = \frac{TN}{TN+FP}$

A confusion matrix is a table, that is used to see the performance of the classifier in a more detailed way, it shows the number of **TP**, **TN**, **FP**, and **FN**, this is shown in Table 2.1. We can use this table to calculate all the metrics we have seen before. The table is split into four quadrants, the top left quadrant is the **TP**, the top right quadrant is the **FP**, the bottom left quadrant is the **FN**, and the bottom right quadrant is the **TN**, what we want is to have as many **TP** and **TN** as possible since these are the ones that have been correctly classified.

| | Predicted Positive | Predicted Negative |
|-----------------|--------------------|--------------------|
| Actual Positive | True Positive | False Negative |
| Actual Negative | False Positive | True Negative |

Table 2.1: Confusion Matrix

The problem with the accuracy metric is that it is not a good metric when the classes are imbalanced, for example, if we have 95 % of the instances of class 1 and 5 % of the instances of class 2, a classifier that always predicts the class 1 will have a 95 % accuracy, but it is not a good classifier. For this reason, we need to use other metrics, but we can also fix the accuracy metric by using the *Balanced Accuracy*, which is the average of the recall of each class, this metric is better than the accuracy when the classes

are imbalanced, the formula for the balanced accuracy is shown in Equation 2.1.

$$\text{BalancedAccuracy} = \frac{\text{Recall} + \text{Specificity}}{2} \quad (2.1)$$

With this new accuracy metric, we can see a better performance of the classifier when the classes are imbalanced. We can clearly see the difference between the accuracy and the balanced accuracy with a quick example, if we have a dataset with 100 instances, 97 of them are positive and 3 of them are negative, and we have a classifier that predicted 96 samples as positive and 4 as negative, this is expressed in the Table 2.2, lets calculate the accuracy and the balanced accuracy.

| | Predicted Positive | Predicted Negative |
|-----------------|--------------------|--------------------|
| Actual Positive | 94 | 3 |
| Actual Negative | 2 | 1 |

Table 2.2: Confusion Matrix Example

- **Accuracy:** $\frac{TP+TN}{TP+TN+FP+FN} = \frac{94+1}{94+1+2+1} = \frac{95}{100} = 95,00\%$
- **Recall:** $\frac{TP}{TP+FN} = \frac{94}{94+3} = 96,91\%$
- **Specificity:** $\frac{TN}{TN+FP} = \frac{1}{1+3} = 25,00\%$
- **Balanced Accuracy:** $\frac{\text{Recall}+\text{Specificity}}{2} = \frac{0,9691+0,25}{2} = 60,95\%$

As we can see, we obtained a 95% accuracy, but the balanced accuracy is 60.95%, this is a big difference and represents better the performance of the classifier which misclassified a high percentage of the negative instances.

Apart from these basic metrics we have other metrics that are based on these, for example, we have the *F-Beta Score*, which is the weighted harmonic mean of the precision and recall, depending on the value of β we can give more importance to the precision or the recall, for example, if we want to give more importance to the recall we can use the *F2 Score*, which means that the β is 2, and if we want to give more importance to the precision we can use the *F0.5 Score*, which means that the β is 0.5. The general formula for the F-Beta Score is shown in Equation 2.2. This metric can be really useful when we want to give more importance to one of the metrics, for example, if we are predicting a disease, we want to give more importance to the recall since we want to detect as many cases as possible, even if we have some FP, but if we are predicting if a person is going to buy a product, we want to give more importance to the precision, since we want to avoid FP.

$$F_{\beta} = (1 + \beta^2) \cdot \frac{\text{precision} \cdot \text{recall}}{(\beta^2 \cdot \text{precision}) + \text{recall}} \quad (2.2)$$

We also have the *Matthews Correlation Coefficient*, this metric tries to summarize the confusion matrix into a single value, it doesn't take into account the class imbalance, so when the classes are imbalanced F β Score is better [13], but when the classes are balanced the Matthews Correlation Coef-

ficient can be useful, it measures the correlation between the true and predicted values [14]. The formula for the Matthews Correlation Coefficient is shown in Equation 2.3.

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (2.3)$$

The last metric we are going to mention is the *ROC AUC*, this metric is used to measure the performance of the classifier, it is the area under the curve of the Receiver Operating Characteristic (ROC) curve, this curve is a plot of the true positive rate (Recall) against the False Positive rate ($FPR = \frac{FP}{FP+TN}$), the ROC curve is a good way to see the performance of the classifier, the closer the curve is to the top left corner, the better the performance of the classifier. The area under the curve is the ROC AUC, and it is a value between 0 and 1, the closer to 1 the better the performance of the classifier, if the area is 0.5 or the curve is close to the diagonal, the classifier is not better than random guessing. An example of an ROC curve is shown in Figure 2.5, we can see that the curve is somewhat close to the top left corner so it is a good classifier but there is still room for improvement. This metric is really useful when the classes are balanced, if the classes are imbalanced the ROC AUC can be misleading, and thus we will want to use the aforementioned metric $F\beta$ Score.

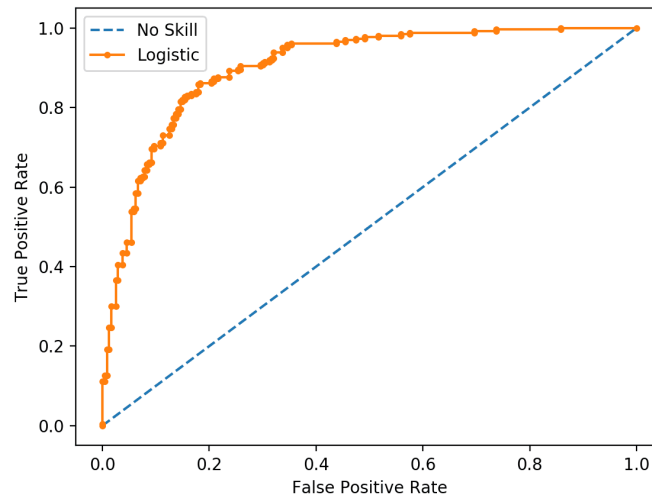


Figure 2.5: Example of a ROC Curve, the y-axis is the True Positive Rate or Recall, and the x-axis is the FP Rate

2.3. Classifier Optimization

Now that we know what a classifier is, how it works and the different types of classifiers, we can start talking about how we are going to optim

2.4. Explainable AI

2.4.1. Explainable AI

2.4.2. SHAP

SHapley Additive exPlanations (SHAP) [15] is a model-agnostic approach to Explainable Artificial Intelligence (XAI), that is, it can be used with any model since it only uses the input and output. In short, how SHAP works, is that it uses the Shapley values from cooperative game theory to explain the output of any machine learning model. The Shapley value is a concept from cooperative game theory, it is a way to fairly distribute the loot obtained by a coalition of players. The difference is that in SHAP instead of players we have features and instead of loot, we have the predicted output as a probability. The Shapley value is the average contribution of a feature to the prediction, and it is calculated by averaging the marginal contributions of a feature to the prediction over all possible orderings of the features. The marginal contribution of a feature to the prediction is the difference between the prediction with the feature and the prediction without the feature. This is done for all possible subsets of features, and then the Shapley value is the average of all these marginal contributions. This is done for all the features, and then we have the Shapley values, which are the importance of the features for the prediction.

2.5. Web Application Development

DESIGN AND IMPLEMENTATION

In this section the design of the project will be shown, showing the structure of the project, its requirements, and the implementation of the different parts of the project. Additionally, the implementation of each part will be explained.

3.1. Project Structure

The structure of the project is divided into four main modules corresponding to the main stages of the project: dataset, classifiers, XAI, and web application. The dataset module is responsible for loading the dataset and splitting it into a 70 % training set and a 30 % test set. Then in the classifiers module, we do the **Exploratory Data Analysis (EDA)**, data preprocessing, and building and optimizing the classifiers. The XAI module is responsible for the implementation of the SHAP algorithm and the analysis of the results. Finally, the web application module is responsible for the development of the web application. The structure of the project is shown in Figure 3.1. Now we will describe each module in more detail.

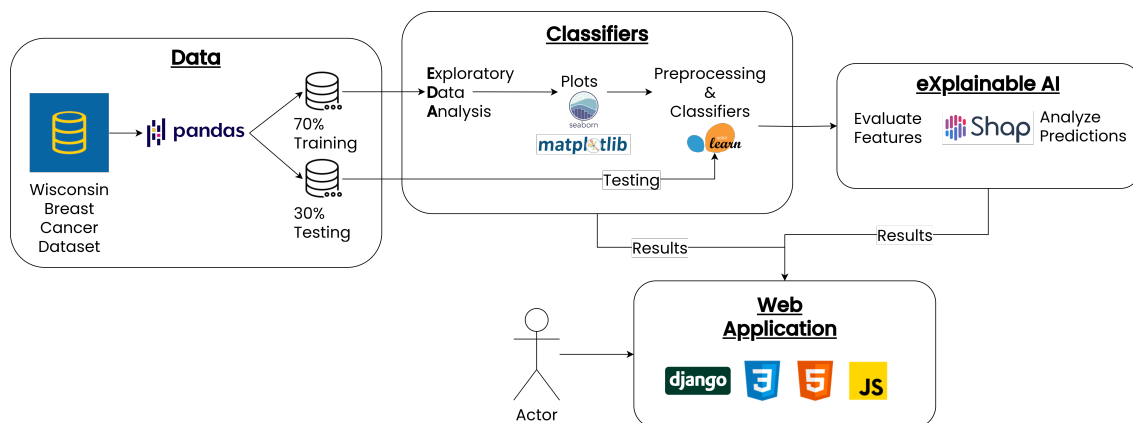


Figure 3.1: Structure of the project

1. **Dataset module:** This module is responsible for loading the dataset and splitting it into a 70 % training set and a 30 % test set. The dataset used in this project is the Breast Cancer Wisconsin dataset [16], which contains the parameters of the cell nuclei of the sample obtained by a **FNAB** of breast

masses. The dataset is loaded using the *pandas* library. The dataset is then split into a 70 % training set and a 30 % test set using *scikit-learn*. From now on, only the training will be used for the **Exploratory Data Analysis**, data preprocessing, and building and optimizing the classifiers. The test set will be used to evaluate the classifiers. This ensures that we avoid data leakage. An example of a **FNAB** is shown in Figure 3.2.

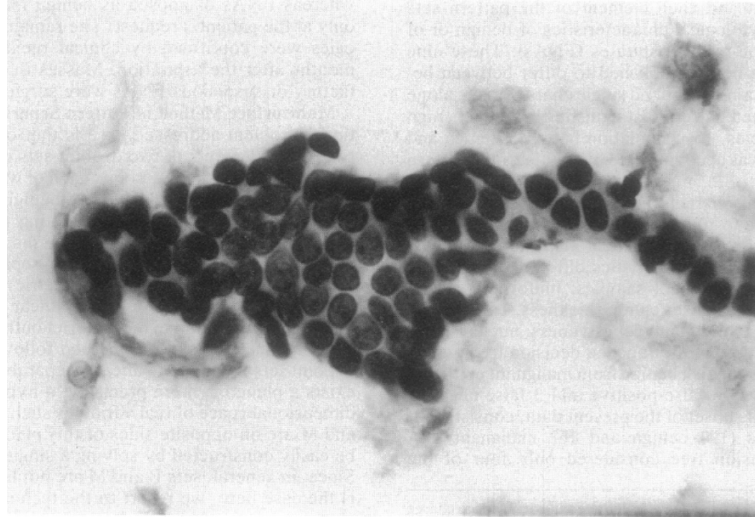


Figure 3.2: Fine Needle Aspiration Biopsy [17]

2. Classifiers module: This module is responsible for the **EDA**, data preprocessing, and building and optimizing the classifiers. The **EDA** is carried out using *matplotlib* and *seaborn* to obtain descriptive and visual statistics of the dataset. The data preprocessing is carried out using *scikit-learn* to scale and normalize the data and to apply **Principal Component Analysis (PCA)**, this is made using pipelines that facilitate the process of building and optimizing the classifiers and their reproducibility. The building and optimization of the classifiers are carried out using *scikit-learn* using the training set to build and optimize the classifiers. Once the classifiers are built and optimized, they are evaluated using the test set.

3. XAI module: This module is responsible for the implementation of the **SHAP** algorithm and the analysis of the results. This is implemented using the *shap* library, which is a Python library that allows us to calculate the **SHAP** values of the classifiers, this values will be used to understand the importance of the features for each classifier and to understand how a prediction was made. This is what makes this project more interpretable than just using classifiers since this will allow a doctor to understand why a prediction was made.

4. Web application module: This module is responsible for the development of the web application. This website was made using *Django* as the backend, *HTML*, *CSS*, and *JavaScript* as the frontend. The website is mainly a way of visualizing the results of the project, it allows the user to see a list of all classifiers and their metrics, which features are the most important for that classifier, and how some predictions were made. It also allows the user to compare two classifiers and see the differences in

their metrics and the importance of the features.

3.2. Requirements

In this section, we will define the requirements of the project. We will divide them into functional (what the system should do) and non-functional (how the system should do it) requirements.

3.2.1. Functional Requirements

Dataset

- **FR-1.-** The application should allow the use of the Wisconsin breast cancer dataset or one with the same features.
- **FR-2.-** The dataset should be split into a training and testing dataset using a 70-30 ratio.

Classifiers

- **FR-3.-** The application should employ the features provided in the dataset for classifier comparison and XAI.
- **FR-4.-** The application should work with multiple classifiers for breast cancer detection.
- **FR-5.-** The application should preprocess the dataset before feeding it to the classifiers. This may include steps such as normalization, and feature selection, focusing on reproducibility and repeatability.
- **FR-6.-** The application should use a test dataset for classifier comparison.
- **FR-7.-** The application should use evaluation metrics to determine the effectiveness of the classifiers taking into account the nature of the problem.

Explainable Artificial Intelligence

- **FR-8.-** The application should show which features are the most important for a classifier.
- **FR-9.-** The application should show how sure the classifier was of its prediction.
- **FR-10.-** The application should show which features determined if the sample was benign or malignant.

Web Application

- **FR-11.-** The web application should have a home page with a slideshow with information.
- **FR-12.-** The web application should display a list with all the classifiers.
- **FR-13.-** The web application should show a details page about each classifier, showing its metrics, most important features, and examples of decisions made by the classifier.
- **FR-14.-** The web application should compare any two classifiers showing their metrics and most important features.

3.2.2. Non-Functional Requirements

Classifiers

- **NFR-1.-** The programming language should be Python.
- **NFR-2.-** The classifiers should be created and trained in a Jupiter Notebook.
- **NFR-3.-** The Pandas library should be used for data manipulation.
- **NFR-4.-** The Scikit-learn library should be used for the classifiers.
- **NFR-5.-** The Matplotlib and Seaborn libraries should be used for the plots.
- **NFR-6.-** The Shap library should be used to understand the feature importance of the models and how a sample decision was made.

Web Application

- **NFR-7.-** The web application should be designed for desktop usage.
- **NFR-8.-** The web application should have a pink color palette and a modern design.
- **NFR-9.-** The web application should be implemented in Django.
- **NFR-10.-** The web application should have an intuitive interface.

3.3. Exploratory Data Analysis

As it has been previously mentioned the dataset used in this project is the Breast Cancer Wisconsin dataset [16]. This dataset is composed of 569 samples of breast masses obtained by **FNAB**, then these samples were parametrized, meaning that now we will work with numbers instead of images. As a result of this process, we have 30 features that describe the cell nuclei of the sample, these features are the mean, standard error, and worst of the ten different parameters of the cell nuclei.

Before using these values in the classifier we first need to understand the dataset, this is done using **EDA** techniques. This involves, first obtaining a general understanding of the dataset, then visualizing the dataset, and finally understanding the relationships between the features. All of this is done only using the training set, this is to avoid data leakage.

3.3.1. Descriptive Statistics

To have a general understanding of the dataset we can use the *pandas* library to load the dataset and then use the *.describe()* method to see an overview of the features in the dataset, their mean, standard deviation, minimum, and maximum values; and then use the *.skew()* method to see the skewness of the features. With this, we would have the needed descriptive statistics of the dataset.

3.3.2. Data Visualization

After seeing a statistical overview of the dataset, we can now dive into visualizing the dataset. This is done using the *matplotlib* and *seaborn* libraries. With these libraries, we will be able to see more visually the distribution of the target variable, the distribution of the features, and the relationships between the features. To achieve this, we will use histograms, correlation matrices, pair, density, and violin plots among others. This will give us a better understanding of the dataset and the relationships between the features which will allow us to make more informed decisions when preprocessing the dataset and building the classifiers.

3.4. Data Preprocessing

Once we have done the **EDA** and we know how the dataset is structured, how the different features are distributed, and the relationships between the features, we can start preprocessing the dataset. To achieve this we will standardize the features and then reduce the dimensionality of the features by using **PCA**, this will result in a dataset that is easier to work with and that will allow us to build better classifiers [18]. To facilitate this process we will use pipelines, which will allow us to easily build and optimize the classifiers and to ensure reproducibility and prevent data leakage [19], this means that we won't directly have to modify the dataset. The Code Snippet 3.1 shows how the pipeline is built.

3.4.1. Scaling and Normalization

This is done using the *scikit-learn* library. The first step is to handle the missing values, this is done using the *SimpleImputer* class, which replaces the missing values with the mean of the feature. Then we scale the features using the *StandardScaler* class, which scales the features to have a mean of 0 and a standard deviation of 1, this is done like the Equation 3.1, where μ is the mean of all the samples of the feature, σ is the standard deviation, and x is the value we want to standardize. We can see in Figure 3.3 how the distribution of the feature stays the same but the mean is 0 and the standard deviation is 1 after standardization.

$$x' = \frac{x - \mu}{\sigma} \quad (3.1)$$

3.4.2. Principal Component Analysis

Lastly, we select the most important variables using the **PCA**, what this does is that it takes the features and transforms them into a new set of features that are linear combinations of the original features trying to maximize the variance of the new features while the correlation between the new

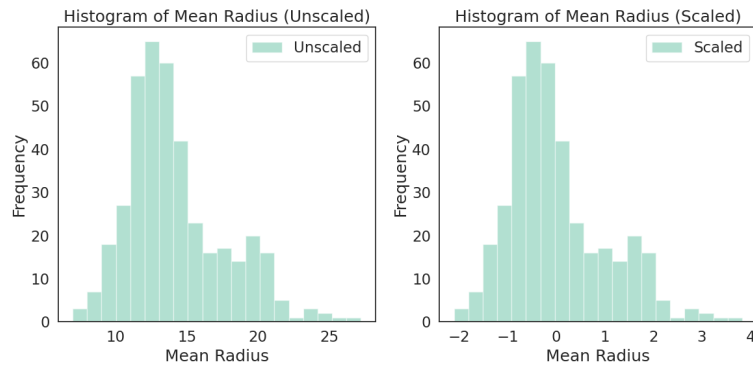


Figure 3.3: Standardization before and after of mean radius

features is 0, this is done using orthogonal transformations [20], in our case instead of specifying the number of components we want to keep, we specify the percentage of variance we want to keep, in this case, we want to keep 95 % of the variance, this is done so that we can reduce the number of features while keeping most of the information. In our case, when we plot the cumulative explained variance we can see that with 10 features we can keep 95 % of the variance, this is shown in the Figure 3.4. The Code Snippet 3.1 shows how the pipeline is built.

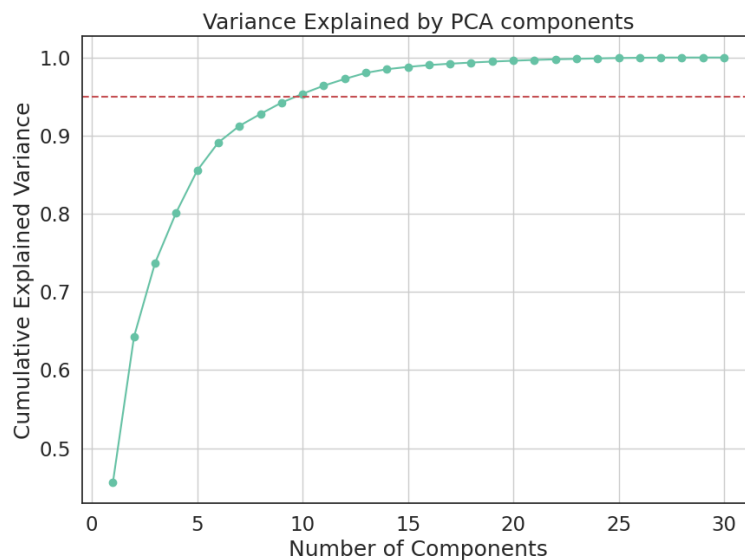


Figure 3.4: Cumulative explained variance when using PCA

3.5. Building and Optimizing Classifiers

To build the classifiers, *Scikit-learn* was used. This library was chosen because it is the most popular library for building machine learning models in *Python* and it has a wide variety of classifiers and ensemble methods. The classifiers were built in a *Jupyter Notebook* to allow reproducibility and to facilitate the process of building and optimizing the classifiers. The classifiers can be divided into two

Code 3.1: Building the pipeline with the StandardScaler and PCA

```

1 pipeline = Pipeline([
2     ('scaler', StandardScaler()),
3     ('pca', PCA(n_components=0.95)),
4     ('clf', Classifier())
5 ])

```

types, base and ensemble, base classifiers are the simplest classifiers, and ensemble classifiers are a combination of base classifiers.

The base classifiers used are:

- **Logistic Regression**
- **Multilayer Perceptron**
- **Support Vector Machine**
- **Decision Tree**
- **Stochastic Gradient Descent**
- **K-Nearest Neighbors**
- **Linear Discriminant Analysis**
- **Gaussian Naive Bayes**
- **Quadratic Discriminant Analysis**

Based on the best results of the base classifiers, the ensemble classifiers were built. Some of these ensemble classifiers allow one to choose the base classifier to use, for these cases, the best base classifiers were used. The ensemble classifiers used are:

- **Random Forest**
- **AdaBoost**
- **Gradient Boosting**
- **Bagging**
- **Voting**
- **Stacking**

3.5.1. Comparison of Classifiers

To decide if a classifier is better than another, evaluation metrics were used, and to choose the appropriate metrics for this problem, we first need to understand the basic metrics used in classification problems [21]:

- **Accuracy:** Proportion of correct predictions over the total number of instances evaluated. $Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$
- **Recall:** Proportion of true positive results among the number of all positive instances. $Recall = \frac{TP}{TP+FN}$
- **Precision:** The proportion of true positive results among the number of cases that were predicted to be positive.

$$Precision = \frac{TP}{TP+FP}$$

- **Specificity:** The proportion of true negative results among the number of cases that were actually negative.

$$Specificity = \frac{TN}{TN+FP}$$

Now we have to choose the appropriate metrics for this problem, the most important thing for this problem is to minimize the number of false negatives, since a false negative means that a patient with cancer was classified as healthy, and this could have fatal consequences. To minimize the number of false negatives, we will use the *Recall*, but the problem of only taking into account the recall is that to maximize it, the model could classify all samples as positive, to avoid this we will use the *Precision* since this will penalize the model for classifying a sample as positive if it was not. To combine these two metrics we will use the *F1 Score*, which is the harmonic mean of the precision and recall. We will also use the *F2 Score*, which is the weighted harmonic mean of the precision and recall, this will give more importance to the recall since we want to minimize the number of false negatives [22]. So, the evaluation metrics sorted by the importance that we will use are:

- **F1 Score:**

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

- **Recall:**

$$Recall = \frac{TP}{TP + FN}$$

- **F2 Score:**

$$F2 = 5 \times \frac{Precision \times Recall}{\frac{4}{Precision} + \frac{1}{Recall}}$$

- **Precision:**

$$Precision = \frac{TP}{TP + FP}$$

3.5.2. Optimization of Classifiers

To maximize these metrics in the classifiers, we have to optimize the hyperparameters of the classifiers. The hyperparameters are the parameters that are not learned by the model, and they are set before the training process, these can, for example, be the number of neighbors in the K-Nearest Neighbors classifier, the number of neurons in the Multilayer Perceptron classifier, or the maximum depth of the Decision Tree classifier. To optimize these hyperparameters, we used the *RandomizedSearchCV* followed by the *GridSearchCV*. The reason for using the *RandomizedSearchCV* first is that can achieve similar results to the *GridSearchCV* but with a lower computational cost [23], this will allow us to explore a wider range of hyperparameters and have a range of values to use in the *GridSearchCV*. Then we used the *GridSearchCV* to find the best hyperparameters in the range of values found by the *RandomizedSearchCV*.

These searches had a 5-fold **Cross-Validation (CV)**, which means that the dataset was split into 5 parts, 4 parts were used for training and 1 part was used for testing, this process was repeated 5 times, each time with a different part for testing, and the results were averaged, using the *F1 Score* as the

metric to optimize. All the *RandomizedSearchCV* and *GridSearchCV* were done using the same seed, this means that the split of the dataset was the same for all the classifiers, this ends up in a more reliable result and also allows the reproducibility of the results. This process was done to avoid overfitting and to have a more reliable result, of course, this process was only done in the training set, to avoid data leakage.

3.6. SHAP Implementation

Building and optimizing the classifiers is a crucial step and can be helpful in a clinical situation, but it can only get us so far. We also need to know why the model made a certain decision, for this, there is a branch of Artificial Intelligence (AI) called Explainable Artificial Intelligence, this branch of AI is focused on making Machine Learning (ML) models more than a black box which is fed with an input and provides an output, it is focused on making these models transparent, allowing clinical practitioners to understand why a certain decision was made [24]. To achieve this we will use the SHapley Additive exPlanations library. The reason for using SHAP is that it has been proven to achieve great explanations for ML models in the medical field [25].

We will use SHAP to have a global understanding of the importance of the features in the models and a local understanding of how a certain decision was made. This will allow us to understand which features are the most important for a classifier, how sure the classifier was of its prediction, and which features determined if the sample was benign or malignant. This will also be done in the *Jupyter Notebook* to ensure reproducibility.

To be able to calculate the SHAP values, we need to create an *Explainer* object, this object is created with the classifier as a parameter, but in our case, we are not using just the classifier but pipelines, this means, that we will have to find a way to obtain the predicted probability of the pipeline. To do this, we will use the *predict_proba* method of the pipeline, this method will return the predicted probability of the sample being malignant no matter which classifier is used, with this, we will be able to create the *Explainer* object and calculate the SHAP values, all this process is show in the Code Snippet 3.2.

After we have the SHAP values, we will be able to use them to understand the importance of the features for the classifiers by plotting a graph with the SHAP values of all the samples and also one for each single feature which will show the behavior of that feature, meaning that we will be able to see if the higher the value of the feature the higher the probability of the sample being malignant or benign. We will also be able to use them to understand how a certain prediction was made by plotting the SHAP values of a certain sample. This will allow us to understand why a certain prediction was made and which features were the most important for that prediction, to do this we will plot the prediction of three cases where the classifier was sure the sample was benign, three where it was sure it benign, the three closest to the decision boundary and the misclassified ones. In the Code Snippet 3.3 we can see how

to code the global feature importance, the behavior of a single feature, and the local explanation of a prediction for a certain classifier, this will be done for all the classifiers, features, and the aforementioned samples

Code 3.2: Obtaining the predicted probability of the pipeline and calculating the SHAP values

```
1 def model_predict_proba(data_asarray):
2     data_asframe = pd.DataFrame(data_asarray, columns=X_train.columns)
3     proba = lr_pipeline.predict_proba(data_asframe)
4     return proba[:, 1]
5
6
7 # explainer object
8 explainer = shap.Explainer(model_predict_proba, X_train)
9
10 shap_values = explainer.shap_values(X_test)
11
12 explainer_x_test = explainer(X_test)
```

Code 3.3: Plotting the SHAP global features, single feature, and local explanation

```
1 shap.initjs()
2
3 # global feature importance
4 shap.plots.bar(explainer_x_test, max_display=25)
5 shap.plots.beeswarm(explainer_x_test, max_display=30, alpha=0.75)
6
7 # single feature importance
8 shap.dependence_plot(feature, shap_values, X_test, alpha=0.65, show=False)
9 ax = plt.gca()
10 x = ax.collections[0].get_offsets()[:, 0]
11 y = ax.collections[0].get_offsets()[:, 1]
12 regression_model = LinearRegression().fit(x.reshape(-1, 1), y)
13 x_range = np.linspace(x.min(), x.max(), 100)
14 y_range = regression_model.predict(x_range.reshape(-1, 1))
15 plt.plot(x_range, y_range, color='red')
16
17 # prediction explanation
18 shap.plots.waterfall(explainers_x_test[prob][idx], show=False)
```

3.7. Web Application Development

To make the results more accessible and visual, we have built a web application using *Django*. This framework was used because *Python* has been the main programming language used in the project and *Django* was the most familiar framework. For the front-end, we used **HyperText Markup Language (HTML)**, **Cascading Style Sheets (CSS)**, and **JavaScript (JS)**. The web application is hosted in *Render*

with the database in *Neon* this makes it accessible using this [link](#), since the free hosting of *Render* is not always available, we have also made a video showing the web application in action, this video can be found in this [link](#) or if not the project can be run locally cloning the [repository](#). The home page of the web application is shown in Figure 3.5, an example of the comparison of two classifiers is shown in Figure 3.6.

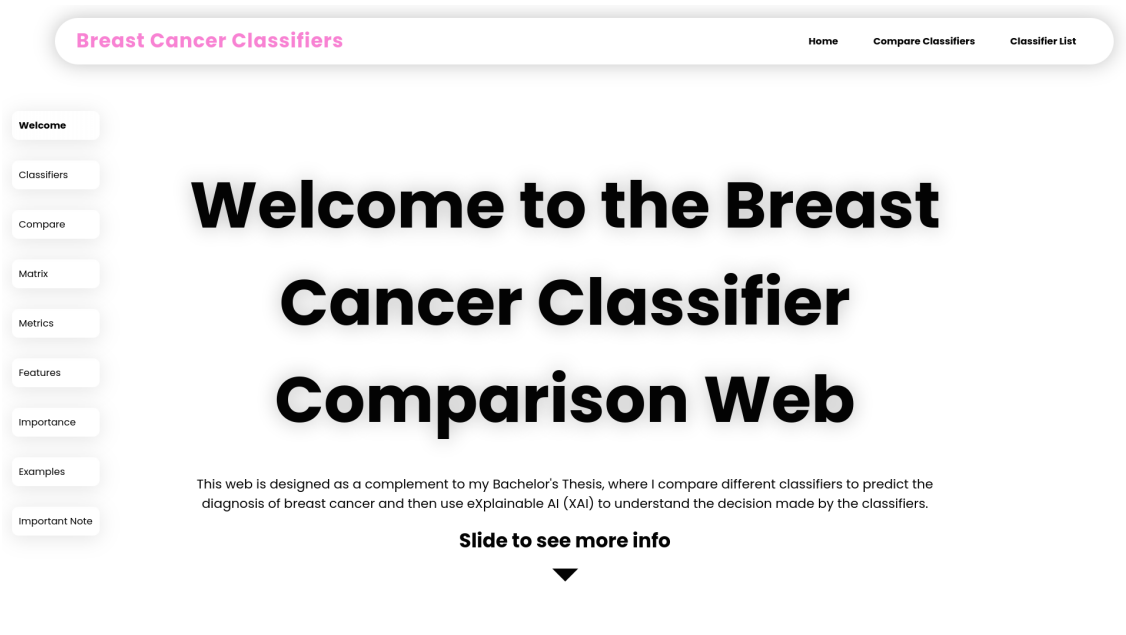


Figure 3.5: Home page of the web application

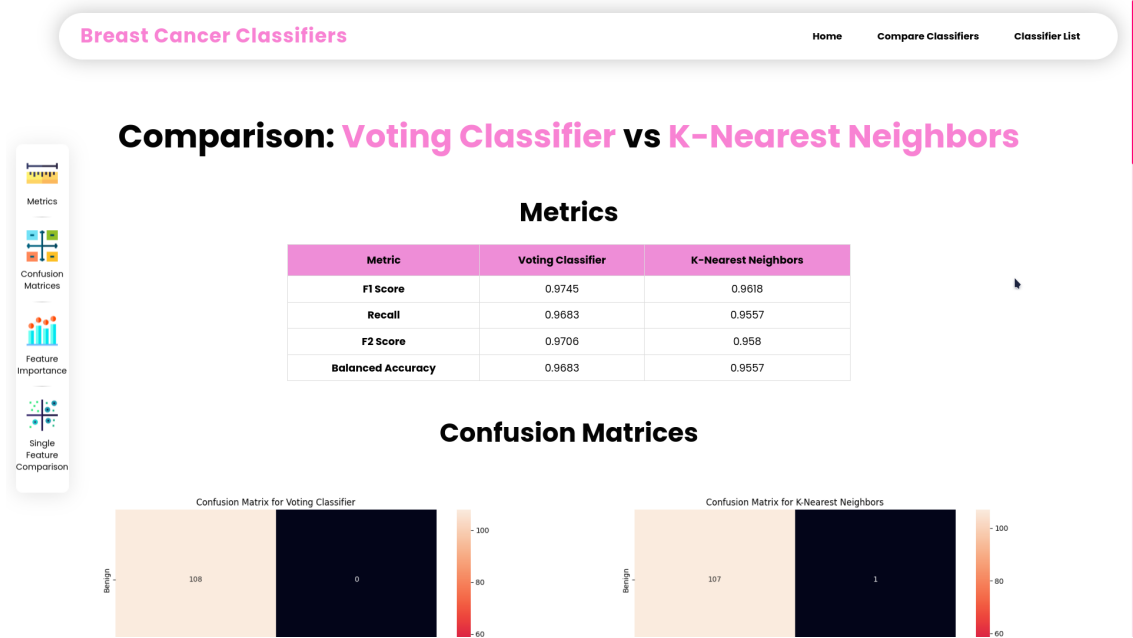


Figure 3.6: Comparison of two classifiers

EXPERIMENTS AND RESULTS

4.1. Exploratory Data Analysis

4.2. Classifier Comparison

4.2.1. Base Classifiers Comparison

4.2.2. Ensemble Classifiers Comparison

4.2.3. Choosing the Best Classifier

4.3. SHAP Analysis

4.3.1. Global Interpretability

4.3.2. 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.
- [4] "sklearn.linear_model.SGDClassifier."
- [5] A. Natekin and A. Knoll, "Gradient boosting machines, a tutorial," *Frontiers in Neurorobotics*, vol. 7, Dec. 2013. Publisher: Frontiers.
- [6] Hemashreekilari, "Understanding Gradient Boosting," Sept. 2023.
- [7] L. Breiman, "Random Forests," *Machine Learning*, vol. 45, pp. 5–32, Oct. 2001.
- [8] Hemashreekilari, "Understanding Random Forest," Aug. 2023.
- [9] Scikit Learn, "Bagging Classifier."
- [10] Hemashreekilari, "Understanding Bagging," Aug. 2023.
- [11] D. H. Wolpert, "STACKED GENERALIZATION,"
- [12] F. Ceballos, "Stacking Classifiers for Higher Predictive Performance," Sept. 2019.
- [13] L. Sisters, "Matthews Correlation Coefficient: when to use it and when to avoid it," May 2020.
- [14] B. Shmueli, "Matthews Correlation Coefficient is The Best Classification Metric You've Never Heard Of," May 2020.
- [15] S. M. Lundberg and S.-I. Lee, "A Unified Approach to Interpreting Model Predictions," in *Advances in Neural Information Processing Systems* (I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, eds.), vol. 30, Curran Associates, Inc., 2017.
- [16] O. M. William Wolberg, "Breast Cancer Wisconsin (Diagnostic)," 1993.
- [17] J. A. M. Sidey-Gibbons and C. J. Sidey-Gibbons, "Machine learning in medicine: a practical introduction," *BMC Medical Research Methodology*, vol. 19, pp. 1–18, Dec. 2019. Number: 1 Publisher: BioMed Central.
- [18] I. Dinç, M. Sigdel, S. Dinç, M. S. Sigdel, M. L. Pusey, and R. S. Aygün, "Evaluation of Normalization and PCA on the Performance of Classifiers for Protein Crystallization Images," *Proceedings of IEEE Southeastcon / IEEE Southeastcon. IEEE Southeastcon*, vol. 2014, p. 10.1109/SECON.2014.6950744, Mar. 2014.
- [19] K. Zhao, "Pre-Process Data with Pipeline to Prevent Data Leakage during Cross-Validation.," Sept. 2019.
- [20] A. Daffertshofer, C. J. Lamoth, O. G. Meijer, and P. J. Beek, "PCA in studying coordination and variability: a tutorial," *Clinical Biomechanics*, vol. 19, pp. 415–428, May 2004.
- [21] H. M and S. M.N, "A Review on Evaluation Metrics for Data Classification Evaluations," *International Journal of Data Mining & Knowledge Management Process*, vol. 5, pp. 01–11, Mar. 2015.

- [22] M. Rutecki, "Best techniques and metrics for Imbalanced Dataset."
- [23] J. Bergstra and Y. Bengio, "Random Search for Hyper-Parameter Optimization,"
- [24] K. Borys, Y. A. Schmitt, M. Nauta, C. Seifert, N. Krämer, C. M. Friedrich, and F. Nensa, "Explainable AI in medical imaging: An overview for clinical practitioners – Beyond saliency-based XAI approaches," *European Journal of Radiology*, vol. 162, p. 110786, May 2023.
- [25] R. Massafra, A. Fanizzi, N. Amoroso, S. Bove, M. C. Comes, D. Pomarico, V. Didonna, S. Diotaiuti, L. Galati, F. Giotta, D. La Forgia, A. Latorre, A. Lombardi, A. Nardone, M. I. Pastena, C. M. Res-sa, L. Rinaldi, P. Tamborra, A. Zito, A. V. Paradiso, R. Bellotti, and V. Lorusso, "Analyzing breast cancer invasive disease event classification through explainable artificial intelligence," *Frontiers in Medicine*, vol. 10, 2023.

ACRONYMS

| | |
|-------------|----------------------------------|
| AB | AdaBoost. |
| AI | Artificial Intelligence. |
| BC | Bagging Classifier. |
| CSS | Cascading Style Sheets. |
| CV | Cross-Validation. |
| DT | Decision Tree. |
| EDA | Exploratory Data Analysis. |
| FN | False Negative. |
| FNAB | Fine Needle Aspiration Biopsy. |
| FP | False Positive. |
| GB | Gradient Boosting. |
| GNB | Gaussian Naive Bayes. |
| HTML | HyperText Markup Language. |
| JS | JavaScript. |
| KNN | K-Nearest Neighbors. |
| LDA | Linear Discriminant Analysis. |
| LR | Logistic Regression. |
| ML | Machine Learning. |
| MLP | Multilayer Perceptron. |
| PCA | Principal Component Analysis. |
| QDA | Quadratic Discriminant Analysis. |
| RF | Random Forest. |
| SC | Stacking Classifier. |
| SGD | Stochastic Gradient Descent. |
| SHAP | SHapley Additive exPlanations. |
| SVM | Support Vector Machine. |
| TN | True Negative. |
| TP | True Positive. |
| VC | Vote Classifier. |

XAI Explainable Artificial Intelligence.

APPENDICES



Universidad Autónoma
de Madrid