



JOSÉ DIOGO TEOTÓNIO PINTO ROMANO

BSc in Computer Science and Engineering

LEARNING TO MAP MICRORNA BIOMARKERS SIGNATURES TO BREAST CANCER SUBTYPES

MASTER IN COMPUTER SCIENCE AND ENGINEERING

NOVA University Lisbon

Draft: July 11, 2025



LEARNING TO MAP MICRORNA BIOMARKERS SIGNATURES TO BREAST CANCER SUBTYPES

JOSÉ DIOGO TEOTÓNIO PINTO ROMANO

BSc in Computer Science and Engineering

Adviser: David Semedo

Assistant Professor, NOVA School of Science and Technology

Co-adviser: Bárbara Mendes

Post-Doctoral Researcher, NOVA Medical School

MASTER IN COMPUTER SCIENCE AND ENGINEERING

NOVA University Lisbon

Draft: July 11, 2025

ABSTRACT

Accurate classification of breast cancer subtypes is essential for enabling personalized and effective treatment strategies. However, the manual stratification of patients based solely on clinical and pathological criteria presents significant challenges, due to the molecular complexity and heterogeneity of breast cancer. In this context, the integration of molecular biomarkers such as microRNAs (miRNAs), small non-coding RNAs with key roles in post-transcriptional regulation, offers a promising avenue for improving diagnostic precision and the possibility for personalized treatments.

This dissertation investigates how Machine Learning and Deep Learning techniques can be used to map miRNA expression profiles to the intrinsic molecular subtypes of breast cancer: Luminal A, Luminal B, HER2-enriched, and Basal-like (Triple-Negative). The main research direction is centred on evaluating the effectiveness of combining miRNA data with clinical variables to develop robust and interpretable subtype classifiers. Particular attention is given to comparing discriminative approaches (e.g., Logistic Regression, XGBoost) with correlation-based or latent space methods (e.g., DIABLO), in order to identify the most appropriate modelling strategies for multi-omics integration and evaluate the impact of incorporating data from different natures.

Preliminary results suggest that ML models can successfully capture miRNA expression signatures associated with certain subtypes, although Luminal B and HER2-enriched remain challenging to distinguish. These findings represent an initial step toward the broader goal of this work: to contribute to the development of reliable, subtype-specific decision support tools that incorporate molecular and clinical data for more precise breast cancer stratification and personalised treatment planning.

Keywords: microRNAs, breast cancer, biomarkers, machine learning applied to Medicine, multi-omics data

RESUMO

A classificação precisa dos subtipos de cancro da mama é essencial para permitir estratégias de tratamento personalizadas e eficazes. No entanto, a estratificação manual de pacientes com base apenas em critérios clínicos e patológicos apresenta desafios significativos, devido à complexidade molecular e heterogeneidade do cancro da mama. Neste contexto, a integração de biomarcadores moleculares, tais como microRNAs (miRNAs), pequenos RNAs não codificantes com papéis fundamentais na regulação pós-transcricional, oferece uma via promissora para melhorar a precisão do diagnóstico e a possibilidade de tratamentos personalizados.

Esta dissertação investiga como as técnicas de Machine Learning e Deep Learning podem ser usadas para mapear perfis de expressão de miRNA para os subtipos moleculares intrínsecos do cancro da mama: Luminal A, Luminal B, HER2-enriquecido e Basal-like (Triplo-Negativo). A principal direção da pesquisa está centrada na avaliação da eficácia da combinação de dados de miRNA com variáveis clínicas para desenvolver classificadores de subtipos robustos e interpretáveis. É dada especial atenção à comparação de abordagens discriminatórias (por exemplo, regressão logística, XGBoost) com métodos baseados em correlação ou espaço latente (por exemplo, DIABLO), a fim de identificar as estratégias de modelação mais adequadas para a integração multi-ómica e avaliar o impacto da incorporação de dados de diferentes naturezas.

Os resultados preliminares sugerem que os modelos de ML podem capturar com sucesso as assinaturas de expressão de miRNA associadas a determinados subtipos, embora o Luminal B e o HER2-enriched continuem a ser difíceis de distinguir. Estas descobertas representam um primeiro passo para o objetivo mais amplo deste trabalho: contribuir para o desenvolvimento de ferramentas de apoio à decisão fiáveis e específicas para cada subtipo, que incorporem dados moleculares e clínicos para uma estratificação mais precisa do cancro da mama e um planeamento personalizado do tratamento.

Palavras-chave: microRNAs, cancro da mama, biomarcadores, aprendizagem automática aplicado à Medicina, dados multi-ômics

CONTENTS

List of Figures	v
Acronyms	viii
1 Introduction	1
1.1 Motivation & Problem Statement	1
1.1.1 Breast Cancer - A Global Health Challenge	1
1.1.2 Can we improve the classification of Breast Cancer subtypes?	2
1.1.3 How to identify the most relevant miRNAs?	3
1.2 Challenges and research hypothesis	3
1.3 Expected Contributions	5
1.4 Document Organization	5
2 Background and Related Work	7
2.1 Biological Context	7
2.1.1 Cancer - A Complex Disease	8
2.1.2 Breast Cancer & its Subtypes	8
2.1.3 MicroRNAs - The Regulators of Gene Expression	10
2.2 Related work	12
2.2.1 Leveraging AI models for Cancer Classification	12
2.2.2 Machine Learning (ML) unravelling microRNAs (miRNAs) as biomarkers	20
2.2.3 Comparison with thesis approach	28
3 Preliminary Work	30
3.1 Dataset and Pre-processing	30
3.2 First tests with baseline models	31
3.3 Second Segment of Experiments	31
4 Work Plan	35

Bibliography	36
Appendices	
A Appendix 1 Covers Showcase	40
Annexes	
I Unveiling microRNA Biomarkers for Breast Cancer Sub-typing	41

LIST OF FIGURES

1.1	Age-standardized incidence rate (ASR, per 100,000 inhabitants) of breast cancer in both sexes in 2022. The data represent global estimates based on International Agency for Research on Cancer [21].	2
1.2	Visual comparison of the estimated number of new cases and deaths caused by Breast Cancer in 2020 (in blue) and 2040 (in red). [21]	2
2.1	In (A) we have a normal breast tissue, while in (B) we can see the presence of a malignant tumor [6].	9
2.2	The figure shows the process of gene expression: DNA is transcribed into mRNA, which is then translated into protein by the ribosome. miRNAs are shown as regulators acting on the mRNA before translation.	11
2.3	Structure of an deep neural network (DNNs). It shows the input, hidden, and output layers, with connections between neurons responsible for processing information [4].	13
2.4	Schematic of the CNN used by Esteva et al. [13] with Inception v3 architecture, adapted to classify skin lesions based on clinical images. The network generates a probability distribution over clinical classes, based on a structured medical taxonomy.	13
2.5	A subset of the hierarchical taxonomy developed in the study by Esteva et al. [13], with diseases organized by clinical and visual similarity into three major groups: benign, malignant, and non-neoplastic.	14
2.6	Performance evaluation of the convolutional neural network (CNN) in skin lesion classification. (a) Confusion matrices of CNN and two dermatologists. The concentration on the diagonal indicates correct classifications; CNN shows less dispersion and better overall performance [13]. (b) Reliability of the CNN demonstrated by <i>AUC</i> curves on a larger, independent dataset [13].	15
2.7	t-SNE projection of the internal representations of the last hidden layer of the CNN [13]. The different classes of lesions are grouped into distinct clouds, revealing the model's ability to extract relevant discriminative features.	15

2.8 Examples of underfitting, proper fitting, and overfitting. From left to right: the model underfits the data, fits it appropriately, and overfits by capturing noise instead of the underlying pattern [2].	16
2.9 Illustration of the k-Nearest Neighbors (k-NN) classification process. The top-left panel shows the initial labeled data (Class A in yellow, Class B in purple) and a new unlabeled sample (?). The top-right panel demonstrates the calculation of distances from the new sample to all existing points. The bottom panel shows the selection of the k=3 nearest neighbors and class assignment based on majority voting, resulting in the classification of the new sample.	17
2.10 Example of a Decision Tree for Classification Based on Attributes: Income, Age, Student Status, and Credit Rating (CR). The tree predicts a binary decision outcome (Yes/No) using hierarchical decision rules [23].	17
2.11 Visualization of a Support Vector Machine (SVM) classifier. The red line represents the optimal hyperplane that separates two classes (blue and green points), while the dashed lines indicate the margins.	18
2.12 Fluxogram of the methodology used in this study, including all the results obtained on each step. [8]	21
2.13 Differential regulation levels of the 29 miRNAs selected by the SVM model in the study by Azari et al. [8]. Positive logFC (log fold-change) values indicate overexpressed miRNAs (magenta) and negative values indicate underexpressed miRNAs (green) in gastric cancer samples compared to healthy tissue.	23
2.14 (a) ROC curve comparing the diagnostic performance of individual microRNAs (hsa-miR-29c and hsa-miR-93) versus their combination. The combined model (blue) shows superior sensitivity and specificity, indicating improved discriminative power for gastric cancer classification. (b) Venn diagram illustrating the overlap of predicted gene targets among five microRNAs (miR-21, miR-133, miR-204, miR-146, and miR-29c) associated with gastric cancer. The central overlap of 426 genes indicates shared regulatory targets across all five miRNAs, suggesting involvement in common biological pathways. Unique and partially overlapping regions highlight miRNA-specific and combinatorial regulatory potential, supporting their relevance as diagnostic and prognostic biomarkers.	24
2.15 Workflow chart used in the following study.	25
2.16 (a) Entropy formula. (b) Information Gain formula.	26
3.1 (a) PCA projection of the training and testing sets. The colors represent the different sets, and the points represent the samples. (b) Same PCA projection, but with the classes represented by different colors.	32
3.2 Class proportions in the training and testing sets. The colors represent the different sets, and the bars represent the proportion of each class.	32

4.1 Gantt chart illustrating the work plan and expected timeline for the project. 35

ACRONYMS

AI	Artificial Intelligence (<i>pp. 3, 12, 15, 16</i>)
BC	Breast Cancer (<i>pp. 1–5, 7, 9, 10, 12, 15, 16, 21, 24, 25</i>)
DL	Deep Learning (<i>pp. 3, 5, 12, 19</i>)
DNA	Deoxyribonucleic Acid (<i>p. 7</i>)
miRNAs	microRNAs (<i>pp. iii, v, vi, 2–5, 7, 10–12, 20–25, 27, 28</i>)
ML	Machine Learning (<i>pp. iii, 3, 5, 12, 16, 19–21, 24, 25, 27–30</i>)
RNA	Ribonucleic Acid (<i>pp. 7, 10</i>)

INTRODUCTION

In this chapter, we will present the motivation and problem statement that led to the development of this dissertation, highlighting the global health challenge posed by breast cancer and the potential of microRNAs as biomarkers for the development of personalized treatments. We will also outline the challenges and research hypothesis that guide this work, as well as the expected contributions and organization of the document.

1.1 Motivation & Problem Statement

1.1.1 Breast Cancer - A Global Health Challenge

Breast Cancer (BC) is currently one of the biggest public health challenges worldwide. In 2022, an article from Bray et al. [10] showed that more than 2.3 million new cases of BC were diagnosed, resulting in around 665,000 global deaths . Other studies estimate that BC will continue to not only be the most commonly diagnosed cancer but also to increase in incidence, with projections indicating that by 2040, the number of deaths will almost double and the number of new cases will be around 3.2 million [7]. Figure 1.1 and 1.2 underline the high incidence and mortality associated with the disease, highlighting the geographical variations in disease burden and the ongoing need to develop more effective strategies for its diagnosis and treatment.

BC is characterized by marked biological heterogeneity, manifested in multiple molecular subtypes that exhibit distinct clinical behaviors. Each subtype exhibits substantial differences in terms of tumor aggressiveness, metastatic potential, and behavior to specific therapies as demonstrated by Prat et al. [33] and Perou et al. [31]. However, in the work of Testa, Castelli, and Pelosi [38], we get a comprehensive explanation of why accurate classification of these subtypes is essential to enable personalized therapeutic approaches, with a direct impact on treatment efficacy and disease prognosis.

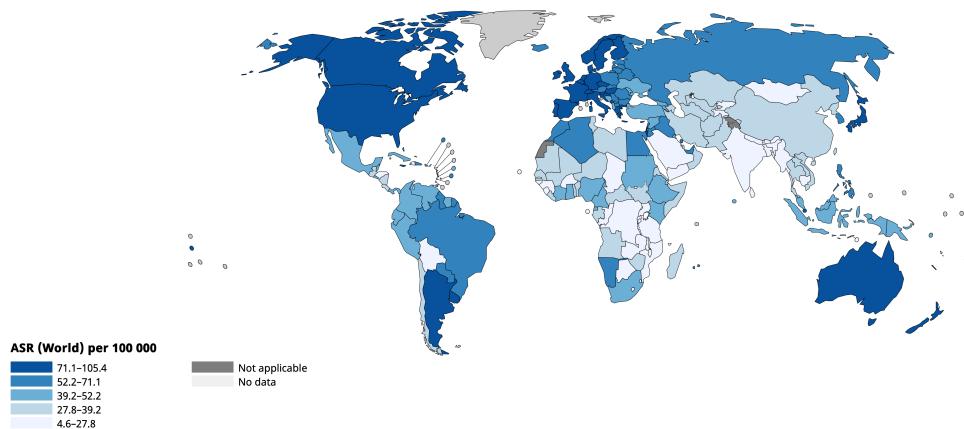


Figure 1.1: Age-standardized incidence rate (ASR, per 100,000 inhabitants) of breast cancer in both sexes in 2022. The data represent global estimates based on International Agency for Research on Cancer [21].

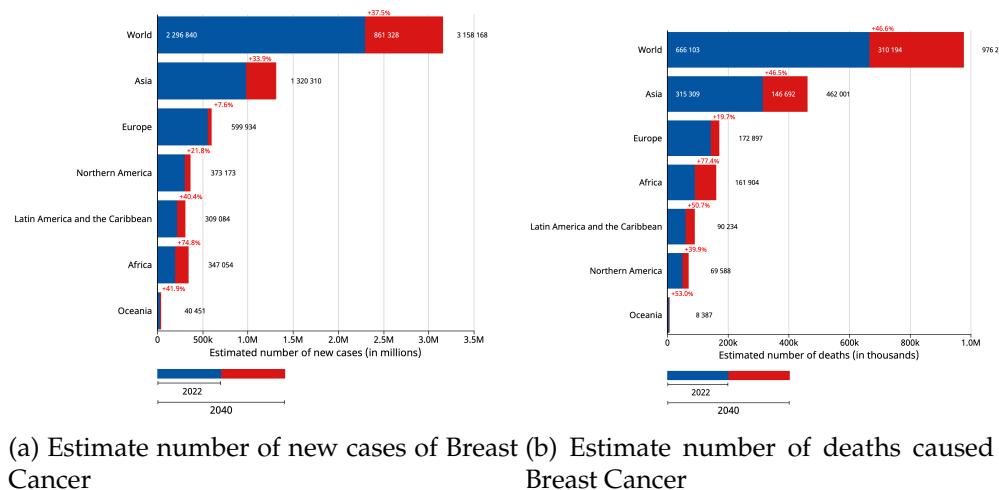


Figure 1.2: Visual comparison of the estimated number of new cases and deaths caused by Breast Cancer in 2020 (in blue) and 2040 (in red). [21]

1.1.2 Can we improve the classification of Breast Cancer subtypes?

Among the emerging candidates for robust biomarkers of BC subtypes are miRNAs, small non-coding RNA molecules that play a crucial regulatory role in gene expression. They are estimated to modulate the expression of about one-third of the genes in the human genome [17] and are implicated in the regulation of multiple physiological and pathological processes, including various human diseases [18].

Given their regulatory nature, several studies have demonstrated a significant association between miRNAs expression profiles and relevant clinical characteristics in the context of BC, including processes such as tumor progression and metastasis development, as seen in [18, 27, 28]. In addition to these aspects, a seminal study by Blenkiron et

al. [9] demonstrated that miRNAs expression profiles can effectively distinguish between different molecular subtypes of BC, highlighting their potential as a precise subtyping tool. This ability to discriminate between subtypes reinforces the value of miRNAs as promising clinical biomarkers.

1.1.3 How to identify the most relevant miRNAs?

The identification of the most relevant miRNAs for BC profiling represents a major analytical challenge due to the complexity of these high dimensional regulatory molecules, non-linearity interactions between and clinical phenotypes that require advanced computational approaches to be effectively modelled. Recent advances in Artificial Intelligence (AI), some of which are explored in the work of Luo et al. [25], particularly in ML and Deep Learning (DL), have demonstrated remarkable potential in extracting meaningful patterns from high-dimensional and heterogeneous (data from distinct nature) biomedical data. These approaches enable not only the accurate classification of BC subtypes but also the identification of discriminative miRNAs signatures, supporting their integration as actionable biomarkers in clinical workflows.

In this context, ML and DL models are particularly well suited for the task of robustly characterizing and explaining the profiles of miRNAs - should such biomarkers exist - with the potential to effectively discriminate between different BC subtypes, as already seen in a study done by Azari et al. [8] where ML algorithms identified potential diagnostic and prognostic miRNAs in gastric cancer, showing high accuracy in the identification of reliable biomarkers for this disease.

This reality reinforces the urgency of developing advanced computational tools that can enable more precise molecular characterization and guide personalized therapeutic decisions, ultimately improving clinical outcomes for patients with aggressive and hard-to-treat BC subtypes.

1.2 Challenges and research hypothesis

Based on the assumption that it is possible to use microRNA expression values and clinical data to map BC subtypes, as shown by Ho, Clark, and Le [18] and Muñoz et al. [28], this dissertation proposes to explore several complementary directions for this pathology where the application of AI techniques is still growing.

First, we intend to assess whether discriminative linear models perform better than latent representation models (in a context where there are two different data sources and many dimensions) - such as DIABLO [37], a widely used model in multi-omics problems. At the same time, we will investigate the impact of patient clinical information (such as age, presence or absence of metastases, hormone levels, among others) on the classification performance of the models, where we will be able to gain valuable insights into possible relationships between these features and BC subtypes.

If substantial results are obtained by any of the models, we will be able to conduct a more extensive study on our main point: whether or not there are miRNAs that are potential biomarkers for BC subtypes. In a more advanced approach, I will be able to explore the applicability of *Conformal Prediction* [5], which provides statistically based confidence intervals for each prediction (which is widely used in areas where risk must be justified and well-founded, such as finance and healthcare). The latter is an approach that is still gaining ground in healthcare and, considering our problem, it makes sense to be able to give a prediction based on a confidence interval, giving our model greater transparency and reliability, something particularly relevant in clinical contexts where error must be minimized and uncertainty well characterized.

Even though the base seems promising, there are several challenges to overcome in order to achieve the desired results such as:

1. Biological heterogeneity: Biological heterogeneity is characterized by the diversity of living organisms, including species, genotypes, and populations, which exhibit a variety of biological characteristics, such as morphology, physiology, genetics, and biogeography. The human body is a highly complex system in which the behavior of each component depends on its interaction with countless other parts. Exposure to the same treatment by two bodies can result in completely different reactions.

2. Functional complexity of microRNAs:

The role of miRNAs in biological regulation and cancer progression is extremely complex and still relatively new from a scientific point of view. The action of a single microRNA is not isolated, but rather part of a network of interactions with dozens (or hundreds) of other miRNAs and contextual factors. This highly interdependent behavior raises questions about the effectiveness of overly simplistic or linear models. The application of non-linear models allows for the discovery of complex relationships and cross-interactions between different miRNAs or between them and clinical variables. These relationships and interactions would be invisible to more traditional approaches.

3. No control set: Another relevant challenge is the absence of a control set that includes data from healthy individuals. Since this type of analysis (miRNAs expression profiling) is not routinely performed in individuals without cancer, it is difficult to define what would be a “normal level” of expression. The implementation of a control group would not only broaden the scope of the model’s task (e.g., distinguishing between the presence and absence of cancer before predicting the subtype), but also optimize the robustness of biomarker identification. An illustrative example of this phenomenon is presented in the study Azari et al. [8] where the implementation of a control set was fundamental to the identification of discriminative markers.

The morphology of this dataset limits the choice of approaches to be used and requires extra caution in how we work with certain models, such as nonlinear ones, given their high adaptability in high dimensions, which, in a context of limited data, can easily lead to overfitting. This requires careful selection of algorithms and attention to the pipeline that is set up to ensure that the results obtained are robust and relevant.

We will eventually find a control set because it would be an important step toward increasing the generalizability of our model. This data set would not only allow us to expand the problem to include the distinction between healthy and sick individuals, but also improve the identification of discriminative biomarkers. The latter has already been successfully tested in other types of cancer, and a key step in the pipeline used is precisely the comparison with a control set [8] to isolate clinically relevant markers.

1.3 Expected Contributions

The main contribution of this dissertation is the development of a computational framework for the classification of BC subtypes based on miRNAs expression and patient clinical data. This framework integrates and compares different ML and DL approaches (still underexplored for this disease) applied to a problem of high biological and statistical complexity. A preliminary version of this framework, focused on discriminative models and subtype-specific biomarker identification, was presented in the article accepted at the EPIA 2025 conference under the track AI in Medicine I.

In addition to the classification process, this framework includes a statistical analysis component aimed at validating the predictions made, in order to confer robustness to the decisions generated by the model. This robustness is particularly relevant in a clinical setting, where transparency and reliability of predictions are essential for potential future translation into medical practice. The findings of this work will be tested using both *in-vivo* and *in-vitro* methods by the team of Dr. Bárbara Mendes NOVA Medical School.

Throughout the work, a critical and comparative analysis of the explored approaches will be promoted, focusing on their applicability to complex and heterogeneous biomedical data. It is thus hoped to contribute to the development of more robust, explainable computational solutions adapted to the reality of biological systems, reinforcing the potential of miRNAs as relevant molecular markers in the stratification of BC patients.

1.4 Document Organization

This dissertation is organized into four main chapters, structured to guide the reader from the conceptual framework to the definition of the work plan. Each chapter contributes to building a foundation for the topic addressed in this thesis: the classification of breast cancer subtypes using miRNA expression profiles and machine learning techniques.

- Chapter 1 – Introduction: Presents the motivation for this work, the main objectives of the research, and the context of the problem in the context of precision medicine assisted by computational algorithms.
- Chapter 2 – Background and State of the Art: Begins with an exposition of the fundamental concepts of molecular biology and genomics necessary for understanding the topic addressed. This is followed by a critical analysis of the relevant literature,

highlighting the main works that apply machine learning models to the identification of biomarkers and the classification of cancer subtypes. This analysis provides a **methodological basis for the work developed**.

- Chapter 3 – Experiments and Preliminary Results: This chapter presents the first practical phase of the work and the preliminary results obtained.
- Chapter 4 – Work Planning: Proposes a detailed schedule of future stages of the research, using a Gantt chart to establish time goals and organize the activities planned.

BACKGROUND AND RELATED WORK

This chapter will provide the necessary background to understand the biological context related to this work, as well as the work that has been already made in the field of miRNAs research and its applications in BC as a biomarker and a subtype classifier. The chapter is divided into two main sections: the first one will present the biological background, including the central dogma of molecular biology, the role of miRNAs in gene expression regulation, and the characteristics of BC and its subtypes. The second section will review the related work in the field of miRNAs research, focusing on the use of miRNAs as biomarkers and subtype classifiers in BC, as well as the challenges and limitations of current approaches.

2.1 Biological Context

The modern understanding of how genetic information is stored, interpreted, and regulated in cells is based on a fundamental principle known as the Central Dogma of Molecular Biology. This concept, first formulated in Pray [34] and Watson and Crick [40], describes the unidirectional flow of genetic information in cells: from Deoxyribonucleic Acid (DNA) to Ribonucleic Acid (RNA) and from there to protein synthesis. According to this model, genes encoded in DNA are transcribed into messenger RNA (or mRNA), which in turn is translated into proteins - the functional molecules responsible for most essential biological processes. This dogma has served as the basis for much of the research in molecular biology and biotechnology.

However, in recent decades, it has become clear that this flow of information is regulated in a much more complex way than initially thought. In particular, it has been discovered that a substantial part of the genome is transcribed into non-coding RNA, i.e., RNA that does not give rise to proteins but plays fundamental regulatory roles. It is in this context that small RNA molecules with central functions in the regulation of gene expression. Their discovery has broadened the classical view of the central dogma, introducing new layers of post-transcriptional control that decisively influence normal and pathological biological phenomena.

2.1.1 Cancer - A Complex Disease

Cancer is a disease characterized by the uncontrolled proliferation of transformed cells, which can invade neighboring tissues and spread to other parts of the body through processes such as metastasis. This definition, based on Brown et al. [11] and National Cancer Institute [29], has recently been expanded to recognize the role of natural selection in the evolution of cancer: it is a cellular system that continuously evolves, adapting to internal and external pressures to ensure its survival.

Under normal conditions, the body's cells divide only when necessary, die when damaged or obsolete, and are replaced by new ones. However, in cancer, this biological balance is disrupted: abnormal cells gain the ability to multiply independently of the body's signals and to resist programmed cell death (apoptosis). These transformed cells become autonomous units that not only ignore normal growth controls but also interact with the tumor microenvironment to promote their own survival, using angiogenesis, immune evasion, and other adaptive mechanisms [11, 29].

The result is a heterogeneous cell population, subject to natural selection within the human body. Cells that acquire adaptive advantages (e.g., higher proliferation rate, drug resistance, or migration ability) tend to prevail, making cancer a constantly evolving disease [11].

Although cancer can arise in virtually any tissue, not all cellular changes are malignant. There are precancerous conditions, such as *hyperplasia* or *dysplasia*, which represent an increase in the number of cells or changes in their morphology, but which do not yet invade surrounding tissues.

The advance of cancer stages is a complex process that involves the **acquisition of invasive and metastatic capacity - properties that distinguish malignant tumors from benign ones**. This process can be silent for years, until more severe symptoms arise, often related to the invasion of vital organs.

2.1.2 Breast Cancer & its Subtypes

Breast cancer is the most commonly diagnosed cancer in women worldwide and is one of the leading causes of cancer death in developed and developing countries as seen in Hong and Xu [19] and Romanowicz, Smolarz, and Nowak [36]. It is estimated that **one in eight women** will be diagnosed with this disease during their lifetime, although it can also affect men - albeit with a much lower incidence.

Most breast tumors are originated in the epithelial cells of the ducts or lobules of the breast, which acquire malignant properties after the accumulation of genetic and epigenetic changes. These events alter the normal control of cell proliferation, differentiation, and apoptosis, allowing for unregulated tumor growth [32].

The development of the disease is associated with a set of well-established risk factors, which include:

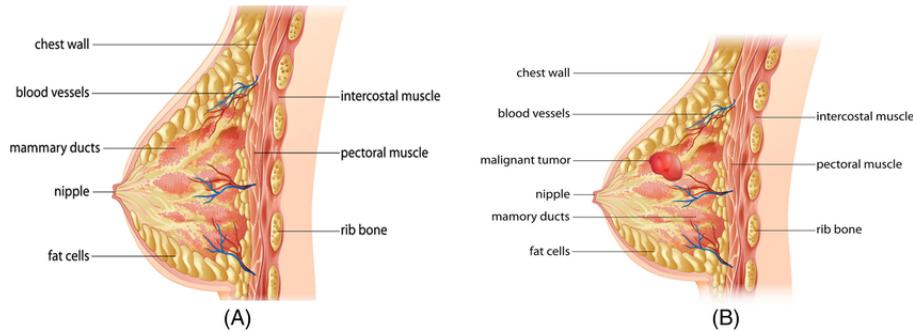


Figure 2.1: In (A) we have a normal breast tissue, while in (B) we can see the presence of a malignant tumor [6].

- Age and family history of the disease;
- Hereditary genetic mutations, especially in the *BRCA1* and *BRCA2* genes;
- Prolonged exposure to endogenous or exogenous hormones (e.g., early menarche, late menopause, hormone therapy);
- Environmental and behavioral factors, such as obesity, physical inactivity, alcohol consumption, and a diet rich in saturated fats [1, 36].

From a molecular and clinical point of view, **BC is highly heterogeneous**. Each tumor may have unique combinations of genetic alterations, signaling pathways, and gene expression profiles, which are reflected in different clinical behaviors, degrees of aggressiveness, and response to treatment [20, 32]. Given this, early detection is crucial for prognosis. When diagnosed in its early stages, BC has survival rates of over 90%. However, in more advanced stages, especially when metastases appear, controlling the disease becomes substantially more difficult and the therapeutic goal shifts from curative to palliative [1, 19].

The therapeutic approach is typically multimodal, combining surgery, radiotherapy, chemotherapy, hormone therapy, and targeted or biological therapies, depending on the characteristics of the tumor and the patient's general condition. The most significant advance in the last decade has been the transition from a uniform model to a personalized treatment approach, tailored to the molecular subtype and individual risk as studied by Romanowicz, Smolarz, and Nowak [36].

In addition, a study made by Polyak [32] recognized that breast tumors are not static entities. Due to phenomena of intra-tumor heterogeneity and clonal evolution, tumors adapt to the selective pressure of treatments, often leading to the development of therapeutic resistance and disease progression.

Given the molecular complexity and clinical diversity of breast tumors, it was established in the papers Adamo et al. [1] and Prat et al. [33] that BC is not a single disease but rather a collection of biologically distinct entities that arise from a common anatomical

Table 2.1: Summary of intrinsic BC subtypes and typical characteristics of each one.
References: [1], [20], [33], [19], [26].

Subtype	Receptors / HER2	Prolif.	Prognosis	Treatment
Luminal A	ER+/PR+, HER2-	Low	Favorable	Endocrine only
Luminal B	ER+, PR↓, HER2±	High	Intermediate	Hormone ± Chemo ± Anti-HER2
HER2-enriched	HER2+, ER-, PR-	High	Improved	Anti-HER2 + Chemo
Basal-like / TNBC	ER-, PR-, HER2-	High	Poor	Chemo; ± PARP/IO (selected)

Acronyms: Chemo = Chemotherapy , IO = Immunotherapy .

site. This heterogeneity is reflected in major differences in tumor progression, metastatic behavior, response to therapy, and long-term prognosis.

To better capture this complexity and inform clinical decision-making, researchers from various studies, such as Perou et al. [31] and Prat et al. [33], have developed a molecular classification system that subdivides breast tumors into intrinsic subtypes which are: *Luminal-A*, *Luminal-B*, *Basal-like* (or Triple-negative) and *HER2-enriched* . These subtypes are defined based on the expression status of three key hormonal biomarkers: **estrogen receptor (ER)**, **progesterone receptor (PR)**, and **human epidermal growth factor receptor 2 (HER2)** as well as other proliferation indices (e.g., Ki-67) and gene expression patterns. The relationship between BC subtypes and these hormone receptors can be seen in Table 2.1. This classification underpins modern precision oncology approaches and has profound implications for therapy and prognosis.

Recent evidence by Polyak [32] and Yeo and Guan [45] suggests that multiple subtypes can coexist within the same tumor (a phenomenon called intra-tumor heterogeneity). This complexity contributes to therapeutic resistance and disease progression.

2.1.3 MicroRNAs - The Regulators of Gene Expression

miRNAs are small non-coding RNA molecules, approximately 20 to 25 nucleotides in length, that play a key role in regulating gene expression at the post-transcriptional level [16, 24, 42]. Instead of encoding proteins, they control the production of proteins from genes.

In simple terms, **miRNAs function as molecular switches that bind to messenger RNA (mRNA) molecules**, blocking their translation into protein or promoting their degradation. This mechanism depends on the degree of complementarity between the miRNAs sequence and that of the target mRNA:

- When there is high complementarity, the mRNA tends to be degraded;

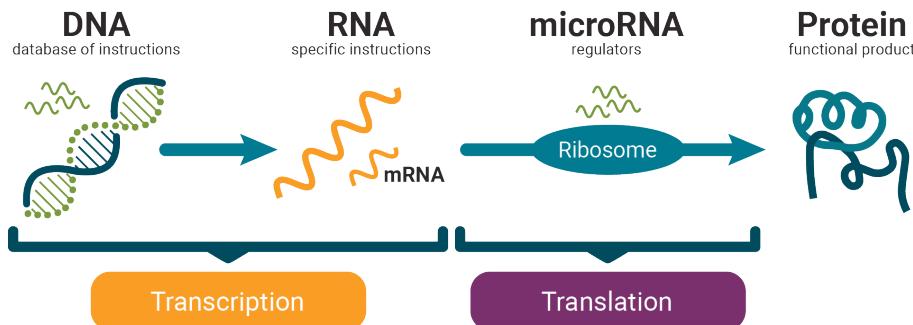


Figure 2.2: The figure shows the process of gene expression: DNA is transcribed into mRNA, which is then translated into protein by the ribosome. miRNAs are shown as regulators acting on the mRNA before translation.

- When complementarity is partial, the miRNAs generally acts by inhibiting translation without destroying the mRNA.

A study made by Calaf et al. [12] demonstrates the high efficiency of this mechanism of regulation: **a single miRNAs can control dozens to hundreds of different genes**, and it is estimated that more than 60% of human coding genes are targeted for regulation by miRNAs.

Given this broad regulatory capacity, miRNAs play a central role in multiple cellular processes such as proliferation, differentiation, apoptosis, and stress response. Consequently, changes in miRNAs expression profiles are associated with several diseases, including cancer, neurodegenerative and cardiovascular diseases. In an oncological context, miRNAs can act as oncogenes (promoting tumor growth) or as tumor suppressors, depending on the biological context and cell type as shown by Gulyaeva and Kushlinskiy [16].

Due to their specificity, stability, and direct involvement in relevant molecular mechanisms, miRNAs have been extensively investigated as promising biomarkers for diagnosis, prognosis, and classification in various diseases - including cancer - because if we understand how these regulators are present in our body, we can ease the clinical decision-making process and improve the treatment of patients.

2.2 Related work

This section will present a critical review of computational approaches developed to date to explore the potential of miRNAs as biomarkers in the context of oncology, covering both BC and other malignant neoplasms. The contributions of ML and DL models applied to the task of classifying different cancer subtypes will also be analyzed, with a special focus on methodologies that integrate molecular data with data from other nature, like clinical characteristics for example.

ML, a branch of AI, involves developing computational models that learn from data to make predictions or decisions. These models are typically trained using either supervised learning, where the target outcomes are known and used during training, or unsupervised learning, in which no explicit labels or outcome variables are provided. In both paradigms, the goal is to uncover meaningful patterns in the data that can be used to generate predictive insights, such as detecting the presence of cancer, estimating survival probabilities, or stratifying patients into risk categories. ML techniques are particularly valuable when dealing with unstructured or complex clinical datasets, as is often the case in oncology.

In recent years, the application of ML algorithms to the field of biomedicine has led to significant advances in the analysis of complex and high-dimensional data, including the expression of miRNAs in cancer [15]. In this context, several studies have explored the use of computational models for the classification of tumor subtypes and/or the identification of discriminative biomarkers, with promising results but also with important limitations.

In this context, we will review and analyze scientific works that have leveraged ML algorithms in contexts similar to the stratification of BC subtypes based on miRNAs expression profiles, complementing them with data of other types (multi-omics data). For each study, it will be important to define the specific work in question so that we can analyze the methodology, algorithms used, and results obtained, all based on the specific context of the research in question, in order to capture and consolidate a ground on which we can work. At the end of the review, we will discuss how these contributions inform and substantiate the methodological choices made in the present work, justifying, whenever possible, the algorithmic and experimental choices based on the available scientific evidence.

2.2.1 Leveraging AI models for Cancer Classification

The classification of different types of cancer using computational models has been one of the most explored areas within the application of AI to medicine. Let's take a look at the work of Esteva et al. [13], a remarkable advance in this “new” relationship between computers and dermatology, where deep neural networks (Figure 2.3) have demonstrated capabilities comparable to those of human experts in the diagnosis of malignant skin lesions.

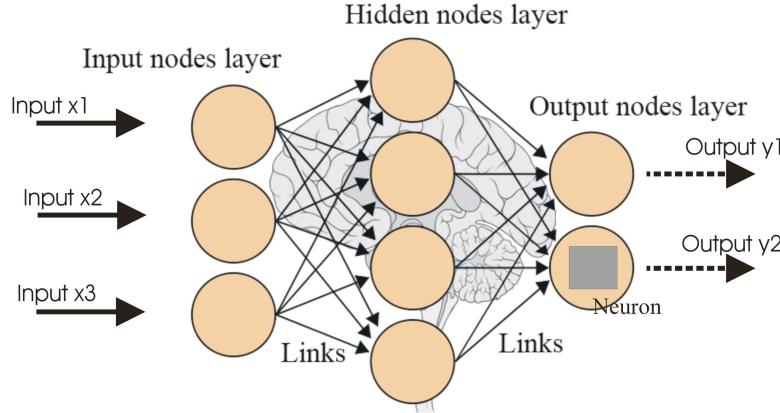


Figure 2.3: Structure of an deep neural network (DNNs). It shows the input, hidden, and output layers, with connections between neurons responsible for processing information [4].

This work was made possible by the use of an architecture based on convolutional neural networks (CNNs) - a type of DNN that is particularly effective in image processing (Figure 2.4). CNNs work by applying convolutional filters that extract visual patterns at different levels of complexity, allowing the model to identify relevant features directly from the image pixels, without the need for specialized preprocessing [3].

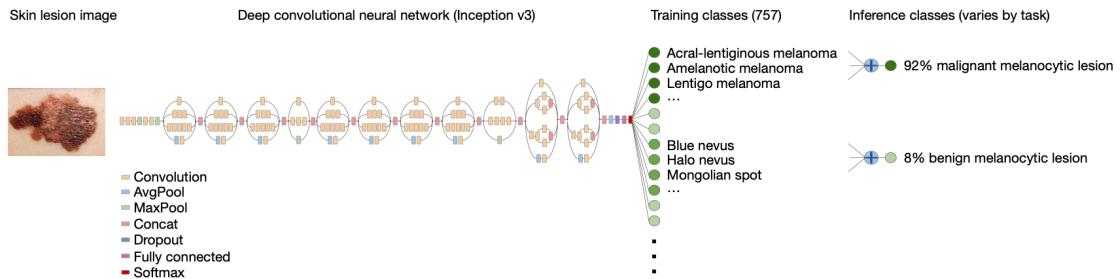


Figure 2.4: Schematic of the CNN used by Esteva et al. [13] with Inception v3 architecture, adapted to classify skin lesions based on clinical images. The network generates a probability distribution over clinical classes, based on a structured medical taxonomy.

The biggest problem with this methodology is that these architectures require a large number of cases (positive and negative) in order to learn the necessary patterns. In this case, 129,450 clinical images covering more than 2,000 different diseases were used. Some factors that determined the good results of this model were:

1. The photographic variability of the samples on which it was trained, since they covered not only images taken with mobile phones, but also dermoscopy images;
2. The manipulation of images during training, enlarging and inverting them to increase the adaptability and robustness of the model;

3. The use of a structured medical taxonomy (Figure 2.5), built on clinical and visual criteria, which allowed for the organization of more than 2,000 diseases into a hierarchy of 757 fine-grained training classes, such as *acrolentiginous melanoma* and *amelanotic melanoma*.

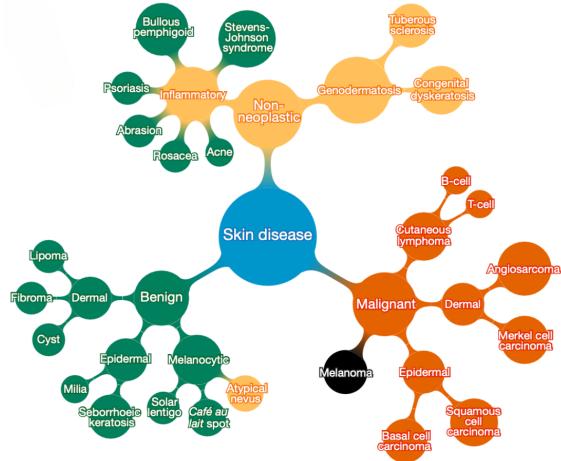


Figure 2.5: A subset of the hierarchical taxonomy developed in the study by Esteva et al. [13], with diseases organized by clinical and visual similarity into three major groups: benign, malignant, and non-neoplastic.

The result? A computational model that not only achieved performance comparable to that of certified dermatologists, but in several scenarios even demonstrated superiority over average human performance, verifiably by this confusion matrix (Figure ??). The trained convolutional neural network was able to classify two critical clinical cases with high accuracy: keratinocytic carcinomas versus benign seborrheic keratoses, and malignant melanomas versus benign nevi. In these binary scenarios, it obtained areas under the curve (*AUC*) of 0.96 and 0.94, respectively (Figure ??) - values higher than those obtained by dermatologists in the same tasks. *AUC* is a metric that quantifies a model's ability to distinguish between classes, with values close to 1 indicating excellent performance.

Furthermore, in more complex scenarios with multiple classes (three and nine disease categories), the model maintained **remarkable levels of accuracy** (72.1% and 55.4%), surpassing or equaling human experts. The robustness of the methodology, that is, its ability to maintain performance under different conditions or test data, was also confirmed in larger test sets, where the network's performance remained stable, with **minimal variations in evaluation metrics**. From a technical standpoint, it was an efficient, scalable system with **potential for application in mobile devices**, which gives it relevant clinical applicability, especially in contexts with limited access to specialists.

Internal analyses further reinforced confidence in the model, showing that it learned consistent clinical representations: the network tended to **group diseases with similar visual characteristics** (Figure 2.7) and focused its attention on the damaged areas of the

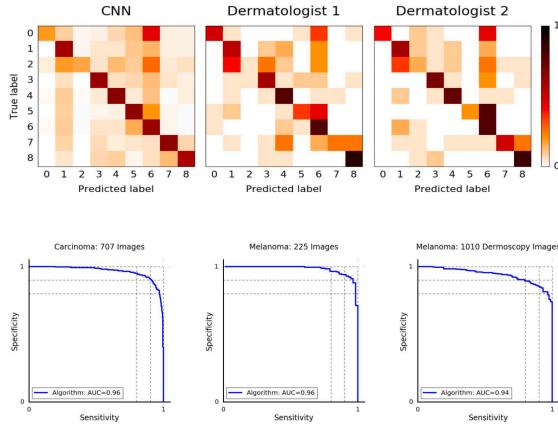


Figure 2.6: Performance evaluation of the convolutional neural network (CNN) in skin lesion classification. (a) Confusion matrices of CNN and two dermatologists. The concentration on the diagonal indicates correct classifications; CNN shows less dispersion and better overall performance [13]. (b) Reliability of the CNN demonstrated by AUC curves on a larger, independent dataset [13].

images, ignoring irrelevant regions such as background or healthy skin - promising evidence of *automated clinical focus* with real practical utility.

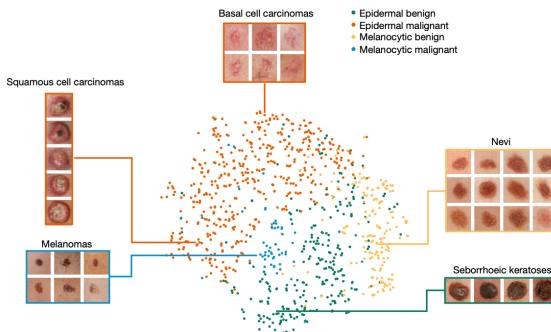


Figure 2.7: t-SNE projection of the internal representations of the last hidden layer of the CNN [13]. The different classes of lesions are grouped into distinct clouds, revealing the model's ability to extract relevant discriminative features.

The effectiveness demonstrated in the previous project shows the magnitude of the benefits that AI can bring to the world of medicine, helping doctors diagnose and stratify diseases with an accuracy that, in some cases, surpasses that of human specialists. This capability is not limited to imaging data: it also extends to the field of molecular data and dermatology, as demonstrated by the study by Wu et al. [43], which applies machine learning algorithms to the task of classifying BC subtypes (in this case, the goal was to distinguish between Triple Negative, or basal-like, and non-Triple Negative tumors, since TNBC is the most deadly cancer with the most difficult prognosis, as we saw in the table).

In the study by Wu et al. [43], when working with gene expression data from thousands of patients, additional challenges arise related to the high dimensionality of the

data, requiring robust feature selection methods and predictive models capable of dealing with complex and often non-linear correlations. This type of study brings us closer not only to the context of this thesis, but also to the type of challenges we will encounter and how we can take advantage of the methodologies used in this paper to achieve good results. From all the algorithms that were tested, Support Vector Machine (SVM) stood out for its performance. This approach allowed the authors to achieve high levels of accuracy, sensitivity, and specificity, demonstrating the potential of AI in classifying BC subtypes based on genomic information. The type of information for this study was the RNA-Sequence (RNA-Seq) profiles made available by The Cancer Genome Atlas (TCGA), a public database containing thousands of tumor samples characterized at the genomic level. This dataset, after pre processing, had 934 tumor samples and over 57,000 genes per sample, a typical high-dimensionality scenario where the number of variables far exceeds the number of observations, and considering that not all of them are necessary or have a great impact, **differential expression analysis was applied** - a bioinformatics technique used to detect which genes are significantly more or less expressed between different conditions - resulting in the selection of 5,502 differentially expressed genes, which served as input for the predictive models (a large reduction of over 50,000 genes). This step corresponds to feature selection, which is essential in problems where there is a high risk of *overfitting* - that is, when the model memorizes the training data but fails to generalize to new examples.

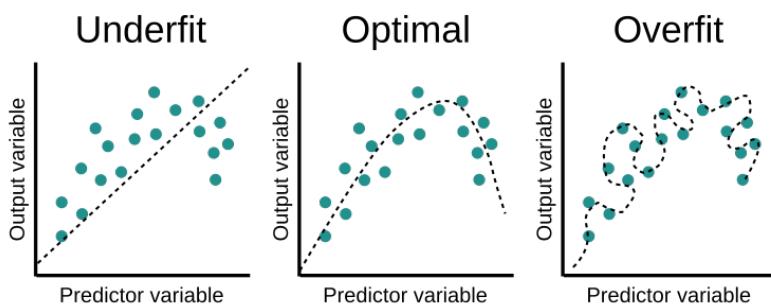


Figure 2.8: Examples of underfitting, proper fitting, and overfitting. From left to right: the model underfits the data, fits it appropriately, and overfits by capturing noise instead of the underlying pattern [2].

Now that the dataset had been reduced to a more informative and manageable subset of features, the authors moved on to the predictive modeling phase. This is a crucial moment in the ML pipeline, where the ability of different algorithms to learn discriminative patterns present in the data is tested - in this case, distinguishing between TNBC and non-TNBC tumors based on the expression levels of selected genes. Several classic supervised learning algorithms were then evaluated, representing different approaches to the classification task:

- **K-nearest Neighbors (kNN):** classifies new data based on the K nearest neighbors in the feature space [47].

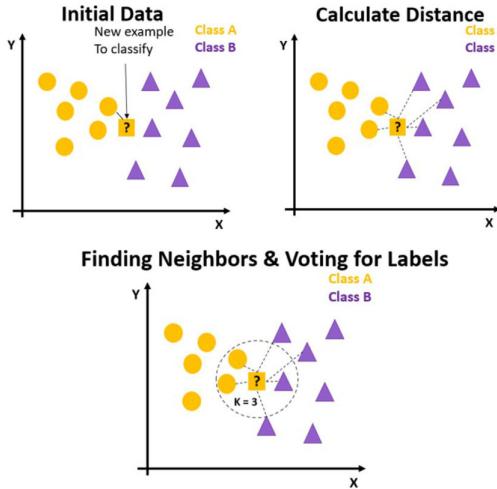


Figure 2.9: Illustration of the k-Nearest Neighbors (k-NN) classification process. The top-left panel shows the initial labeled data (Class A in yellow, Class B in purple) and a new unlabeled sample (?). The top-right panel demonstrates the calculation of distances from the new sample to all existing points. The bottom panel shows the selection of the $k=3$ nearest neighbors and class assignment based on majority voting, resulting in the classification of the new sample.

- **Naïve Bayes (NB):** uses Bayes' Theorem to estimate the most likely class of a sample, assuming that the input variables are independent of each other [41].

$$P(C_k | \mathbf{x}) = \frac{P(C_k) \prod_{i=1}^n P(x_i | C_k)}{P(\mathbf{x})}$$

- **Decision Tree (DT):** a model that makes decisions through a hierarchical tree-shaped structure, where each internal node represents a condition on a variable, and each branch represents a possible outcome of that condition. The process continues until it reaches a leaf node, which indicates the final class or value [22].

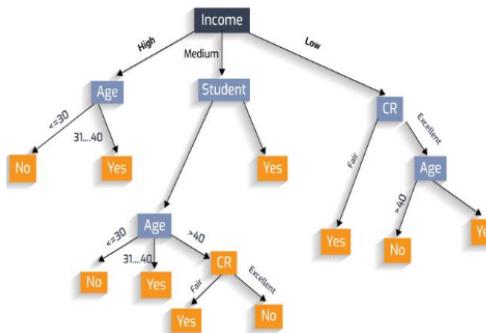


Figure 2.10: Example of a Decision Tree for Classification Based on Attributes: Income, Age, Student Status, and Credit Rating (CR). The tree predicts a binary decision outcome (Yes/No) using hierarchical decision rules [23].

- **Support Vector Machine (SVM):** This algorithm constructs an optimal hyperplane that best separates samples from different classes. The goal of SVM is to maximize

the margin between the two classes for better generalization. Support vectors are the data points that lie closest to the hyperplane and define its position.

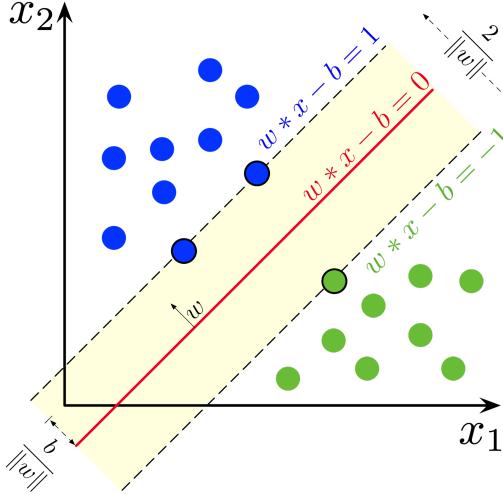


Figure 2.11: Visualization of a Support Vector Machine (SVM) classifier. The red line represents the optimal hyperplane that separates two classes (blue and green points), while the dashed lines indicate the margins.

In order to evaluate a model, performance metrics are used to assess how well it performs across different dimensions, such as accuracy, relevance, and sensitivity to different classes. These metrics provide quantitative insight into the strengths and limitations of a classification algorithm, allowing researchers and practitioners to make informed decisions when comparing models or tuning parameters. Evaluating models using multiple metrics is especially important in scenarios involving imbalanced datasets, where a single metric (such as accuracy) may not provide a complete picture. The most commonly used evaluation metrics in classification tasks, according to Yaseen and Abdulazeez [44], are:

1. Accuracy

Accuracy is the ratio of correctly predicted instances to the total number of instances. It measures the overall effectiveness of a classification model.

2. Precision

Precision measures the proportion of correctly predicted positive observations to the total predicted positive observations. It reflects the model's ability to return only relevant results.

3. Recall (Sensitivity or True Positive Rate)

Recall is the ratio of correctly predicted positive observations to all actual positives. It indicates the model's ability to identify all relevant cases.

4. F1-Score

The F1-Score is the harmonic mean of Precision and Recall, providing a balance between the two. It is particularly useful when the class distribution is imbalanced.

5. Support

Support refers to the number of actual occurrences of each class in the dataset. While not a performance metric per se, it helps in understanding how many examples a classifier is making predictions on for each class.

Now that we are familiar with the main concepts of classification models and the metrics used to evaluate them, let us return to the study by Wu and Hicks (2021) in light of these indicators. The analysis of the results, presented in the following table, allows us to see more clearly how each algorithm performed in the task of distinguishing between TNBC and non-TNBC, highlighting the performance of SVM in virtually all scenarios evaluated.

In the complete gene set, SVM achieved 90% accuracy, 87% recall, and 90% specificity - metrics that indicate, respectively, the proportion of correct classifications, the ability to correctly identify TNBC cases, and the ability to avoid false positives. The analysis of these metrics is essential to correctly interpret the results in a clinical context, where the consequences of classification errors can be significant.

To validate the robustness of the differentially expressed gene selection approach, the authors compared it with other classic feature selection methods, such as *SVM-RFE* (an iterative technique that removes the least relevant features based on *SVM* weights), *Relief* (which weights features based on their correlation with the class), *ARCO*, and *mRMR* (which maximizes relevance and minimizes redundancy between variables). Even so, the model based on differential expression and *SVM* demonstrated better performance in most cases, confirming the soundness of the strategy adopted.

The study further deepened the analysis of feature importance by evaluating the performance of models with different sizes of gene subsets, from the initial 5,502 to only 16 genes. Interestingly, *SVM* performance remained high even with reduced sets, achieving the best results with 256 genes. This stability suggests that discriminative information is concentrated in a small subset of features, which is relevant for scenarios with a low number of samples and a large number of features, as we have in this dissertation, and quite positive because:

- it allows molecular tests to be cheaper and faster since fewer genes are needed;
- it makes models easier to validate in a clinical setting since it is simpler to obtain quality samples with few targets;
- greater interpretability for physicians.

The analysis of the two studies presented here allows us to consolidate a fundamental idea for this dissertation: ML and DL algorithms demonstrate a remarkable ability to **classify different types of cancer based on complex data**, whether imaging or molecular. From deep neural networks applied to dermatological imaging to discriminative algorithms used to analyze gene expression, these models have proven to be effective

tools for supporting clinical decisions, sometimes proving superior to human experts. More than just efficient classifiers, these systems have also proven to be interpretable, robust, and applicable in real clinical scenarios. Both the image-based approach [13] and the gene expression-based approach [43] have faced and overcome challenges typical of medical practice and biomedical research: sample scarcity, high dimensionality, and the need for models with good overall performance, but also **confidence in the prediction of clinically critical cases.**

These findings provide the **conceptual support needed to explore a more specific direction:** the use of ML models to discover and validate **miRNAs as biomarkers** in oncology. Like coding genes, miRNAs carry rich and discriminative information about the biological state of cells and have shown promise in the stratification of tumor subtypes. The next section addresses precisely this line of research, focusing on how ML has been used to reveal miRNAs signatures with diagnostic and prognostic value - a foundation point for the objectives of this dissertation.

2.2.2 ML unravelling miRNAs as biomarkers

As previously discussed, artificial intelligence models have demonstrated proficiency in the task of classifying tumor type, or subtype, by either analysing clinical images, high-dimensional genomic data and other forms of data. In this thesis, our focus is into a more specific and pertinent domain: the identification of miRNAs as biomarkers in oncological contexts through ML techniques. We know for a fact that miRNAs are a class of small molecules that have been demonstrated to possess a substantial regulatory capacity over gene expression 2.1.3, but even though this capacity distinguishes them as optimal candidates for utilization as molecular biomarkers, the number of expressed in human tissues, in conjunction with their variability across individuals, demands the implementation of robust computational methodologies.

In this context, ML models have gained prominence as powerful tools for revealing latent patterns in miRNAs expression data, allowing the identification of subsets with diagnostic, prognostic, or tumor subtype stratification value. This research trajectory is particularly auspicious, as it proffers more readily implementable, non-invasive methodologies that can be substantiated within a clinical environment. That's what we are going to explore in this section, where we will present three papers that illustrate different stages of this scientific effort: the first being a general reference of the type of work that we will be analyzing giving us a glimpse of how ML can work in this context; the second one being a more detailed pipeline applied to the identification of miRNAs as biomarkers for gastric cancer, and lastly, the intersection of this kind of approach with which is the central carcinoma focus of this thesis.

The first example of work that we will address in this subsection focused on a major challenge in modern oncology: “Is there a reliable and clinically relevant molecular

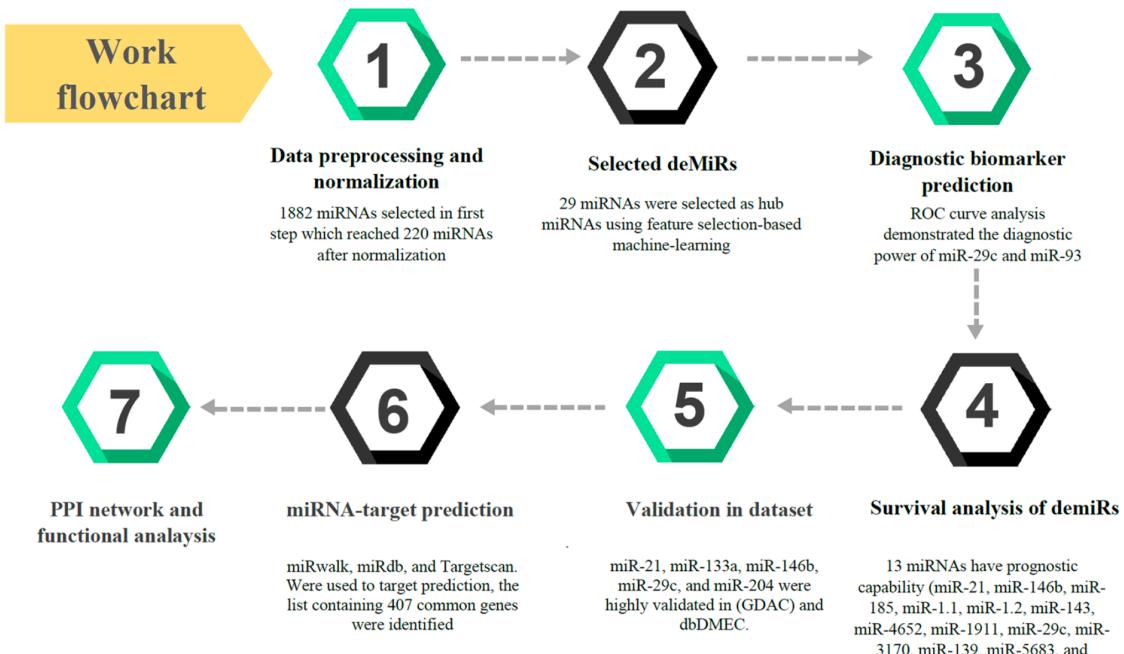


Figure 2.12: Fluxogram of the methodology used in this study, including all the results obtained on each step. [8]

biomarker for the diagnosis and prognosis of gastric cancer?" Despite growing evidence of the regulatory role of miRNAs in tumor progression and aggressiveness, their biological complexity (as already mentioned in section 1.2) makes it difficult to select those with true clinical value from among the thousands that exist, and despite these barriers, considering that gastric cancer has one of the lowest survival rates in a 5 year period (between 20% and 30% survival rate) with a big share due to late detection, usually when the tumor is already in a advanced stage and metastasized. Given this scenario, the authors Azari et al. [8] propose a methodical and well-structured ML-based approach, as we can see in the workflow chart, to automate the process of discovering these miRNAs that may be biomarkers, where the result also has added robustness so that it can be replicated in a real clinical environment.

The research was based on data from 576 samples from gastric cancer patients, collected from the TCGA repository (previously used on other studies), including miRNAs expression profiles and associated clinical information - this clinical information includes patient characteristics such as age, gender, among others - adding 1,882 features whose expression in a real context is not uniform. Therefore, similar to the study on BC subtyping with ML, a differential expression analysis was performed, resulting in a reduction to only 220 of these molecular regulators (a much more acceptable and interpretable value than the 1882 previously). Following this extensive data preprocessing, the authors set up a classification pipeline consisting of five classic ML algorithms: Support Vector Machine, Random Forest, Decision Tree, Logistic Regression, and K-Nearest Neighbors 2.2.1 (all of which have been discussed previously except for Random Forest, which consists

of a set of several DTs), all of which were evaluated using metrics such as F1-score, *AUC*, and confusion matrices.

Table 2.2: Performance of the evaluated classification algorithms.

Algorithm	Accuracy (%)	AUC (%)
DTS	88	47.0
Random Forest	93	39.5
SVM	93	88.5
KNN	93	41.7
Logistic	93	88.0

With regard to the performance of the predictive models (as presented in table 2.2), the results show substantial differences between the classifiers. Although four of the five models - all except Decision Trees - achieved an accuracy of 93%, the analysis of the area under the curve (*AUC*) revealed a more complex and informative picture. SVM stood out as the most balanced model, achieving an *AUC* of 88.5%, which indicates a strong discriminatory capacity (i.e. the model not only gets it right often, but also assigns probabilities reliably). In contrast, models such as KNN and RF, despite their high accuracy, obtained very low *AUCs* (41.7% and 39.5%, respectively), suggesting poor calibration of probabilistic predictions and possible overfitting to the training set. The Decision Tree (DT), with slightly lower accuracy (88%) and an *AUC* of only 47%, showed a more modest performance compared to the others, probably due to its vulnerability to variance in data sets with more noise. Finally, logistic regression also showed robust performance, with an *AUC* of 88%, very close to that obtained by SVM. However, SVM offered better generalization or sensitivity, justifying its selection as the final model (as already discussed in the work by Wu et al. [43]). This analysis highlights the importance of considering multiple metrics in model evaluation, especially in clinical contexts, where the reliability of the assigned probabilities - and not just the hit rate - can be decisive for a safe medical decision.

Considering the chosen model, a processing step was performed where the most important features were selected using a heatmap analysis, which helps to identify patterns and the relevance of miRNAs in this disease. This step resulted in the reduction of 220 candidates to only 29 (5 of which are significantly up-regulated and 24 considerably down-regulated - in Figure 2.13), which is a very important result as it makes the interpretation of the relationships between these potential biomarkers much more human-friendly.

After identifying 29 candidate miRNAs from the reduced set of 220 miRNAs, the authors proceeded to a validation and refinement stage. This phase aimed to ensure that the

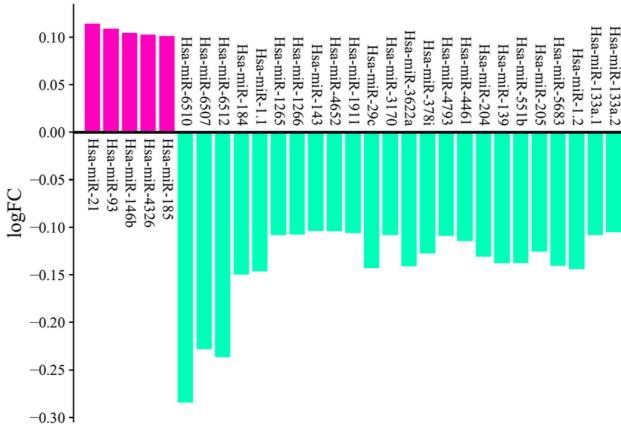


Figure 2.13: Differential regulation levels of the 29 miRNAs selected by the SVM model in the study by Azari et al. [8]. Positive logFC (log fold-change) values indicate overexpressed miRNAs (magenta) and negative values indicate underexpressed miRNAs (green) in gastric cancer samples compared to healthy tissue.

selected miRNAs not only stood out statistically in the training set but also maintained biological relevance and predictive robustness in independent scenarios. To this end, cross-validation was performed in external databases, such as the Global Data Assembly Centres (GDAC) and dbDMEC, which aggregate information from public repositories such as GEO, SRA, ArrayExpress, and TCGA. This process allowed us to identify five miRNAs with consistent differential expression in multiple cancer contexts: hsa-miR-21, hsa-miR-133a, hsa-miR-146b, hsa-miR-29c, and hsa-miR-204, all classified as “highly validated”. These markers stood out for their high ability to discriminate between healthy and pathological states and, at the same time, predict patient prognosis, positioning themselves as a potential tool of clinical value for the diagnosis and monitoring of gastric cancer. Complementarily, ROC curve analyses were conducted to estimate the diagnostic potential of each miRNAs, as well as survival analyses to assess their prognostic value, demonstrating a better indicator performance when combining two miRNAs together (Figure 2.14a).

However, the study’s major revelation is not limited to predictive capacity. The final set of four/five selected miRNAs (hsa-miR-21, hsa-miR-133a, hsa-miR-146b, and hsa-miR-29c + hsa-miR-204) underwent additional functional validation, which allowed the identification of the genes targeted by the selected microRNA. This identification was achieved by constructing a protein interaction network (PPI) (Figure 2.14b). This phase confirmed that the genes regulated by these miRNAs are involved in processes essential to gastric carcinogenesis, such as Wnt signaling or epigenetic regulation. hsa-miR-29c stood out above all for its simultaneous predictive value for diagnosis and prognosis, reinforcing its clinical potential as a dual biomarker.

It’s now a fact that identifying miRNAs that are promising biomarkers requires a study involving several different stages, as we saw in the study by Azari et al. [8], but the impact of the results obtained is of such magnitude that it fully justifies this complexity. From

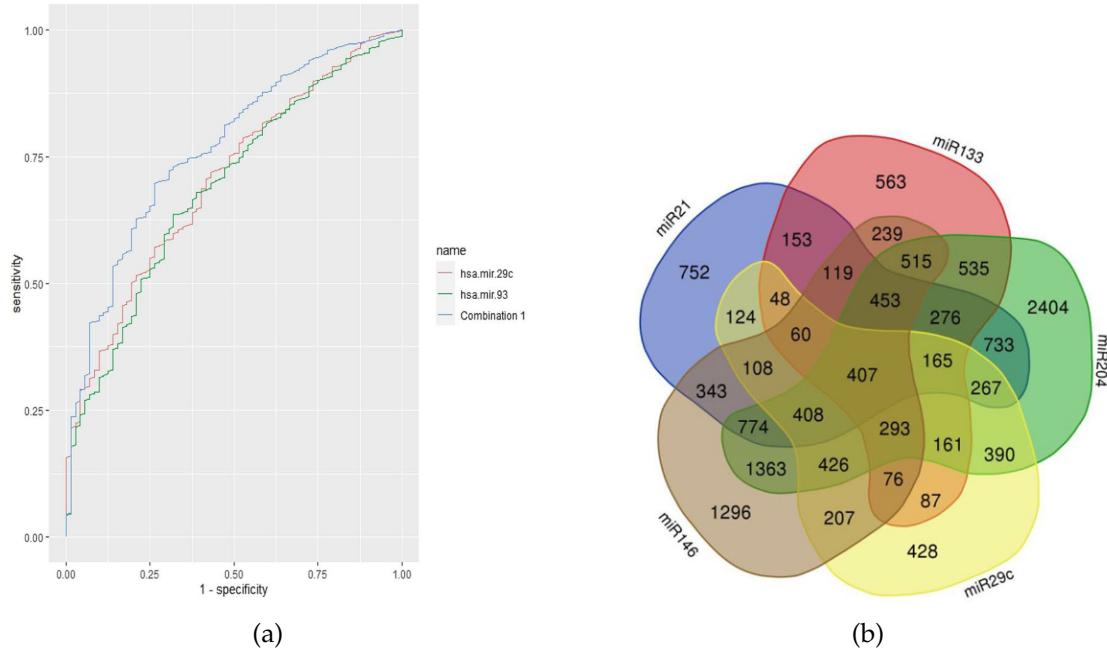


Figure 2.14: (a) ROC curve comparing the diagnostic performance of individual microRNAs (hsa-miR-29c and hsa-miR-93) versus their combination. The combined model (blue) shows superior sensitivity and specificity, indicating improved discriminative power for gastric cancer classification.

(b) Venn diagram illustrating the overlap of predicted gene targets among five microRNAs (miR-21, miR-133, miR-204, miR-146, and miR-29c) associated with gastric cancer. The central overlap of 426 genes indicates shared regulatory targets across all five miRNAs, suggesting involvement in common biological pathways. Unique and partially overlapping regions highlight miRNA-specific and combinatorial regulatory potential, supporting their relevance as diagnostic and prognostic biomarkers.

statistical pre-selection to functional and biological validation, each phase contributed to ensuring that the identified miRNAs were not only statistically relevant but also biologically plausible and clinically actionable. The final panel of four miRNAs, validated at multiple levels, represents a concrete contribution to the construction of more accessible, interpretable, and potentially applicable diagnostic and prognostic tools in clinical practice. This work not only illustrates the power of ML in biomarker discovery, but also establishes a solid methodological foundation that can be replicated and adapted to other diseases - including BC, the focus of this dissertation.

If the previous study demonstrated how it is possible, through a robust machine learning pipeline, to reduce thousands of candidate miRNAs to a small panel with proven clinical potential for gastric cancer, it is natural to ask: is it possible to apply the same principles of statistical selection, biological validation, and predictive modeling in a context closer to the focus of this dissertation, namely This is precisely the proposal of the work of Rehman et al. [35], which starts from a clinically identified set of miRNAs associated with

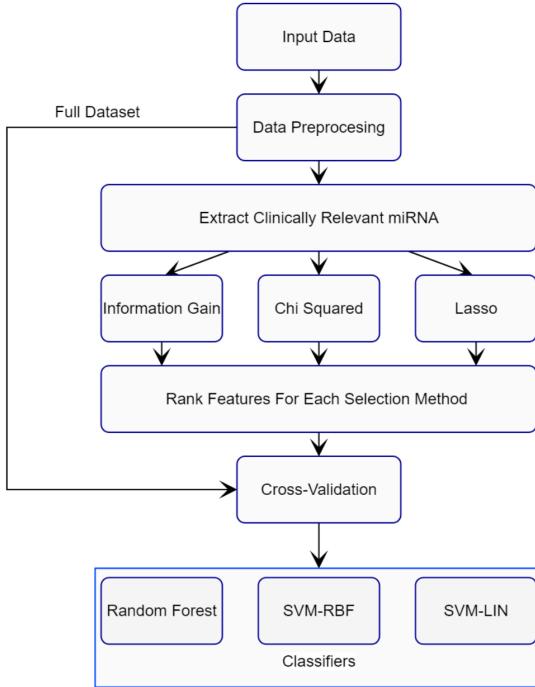


Figure 2.15: Workflow chart used in the following study.

BC and seeks, with the aid of ML algorithms, not only to validate their relevance, but also to build classifiers capable of distinguishing between healthy and cancerous tissue with high accuracy. This study introduces a perspective that complements that of the previous work: while Azari et al. [8] take an exploratory approach to biomarker discovery, this new work focuses on the next step: the verification and validation of already known miRNAs 2.3, assessing the extent to which they are, in fact, discriminative and clinically actionable. It is this transition, from discovery to validation and based on the workflow chart used in this study 2.15, that we will analyze in detail in the following paragraphs.

The study by Rehman et al. [35] focuses on validating a set of miRNAs previously identified as clinically relevant for BC, using a ML-based approach. Unlike exploratory methods, this work starts from an established list of candidates and evaluates their discriminatory power in distinguishing healthy from tumor tissue. The dataset comprises 1207 breast cancer samples with 1881 miRNAs expression profiles from the TCGA-BRCA repository, already normalized and labeled. Notably, the class distribution is balanced, reducing bias in classification and strengthening the reliability of evaluation metrics. Given the high dimensionality, an initial feature cleaning step — including removal of null values — was performed, reducing the set to 1626 features. Based on prior studies, the authors manually selected the 36 most promising miRNAs for classification. Three complementary feature selection methods were then applied:

- (a) *Information Gain* - a concept from information theory, consists of reducing uncertainty (entropy) by dividing the data based on a given attribute. The greater the

Table 2.3: List of miRNAs clinically verified.

miRNA [14]			
hsa-mir-10b	hsa-let-7d	hsa-mir-206	hsa-mir-34a
hsa-mir-125b-1	hsa-let-7f-1	hsa-mir-17	hsa-mir-27b
hsa-mir-145	hsa-let-7f-2	hsa-mir-335	hsa-mir-126
hsa-mir-21	hsa-mir-206	hsa-mir-373	hsa-mir-101-1
hsa-mir-125a	hsa-mir-30a	hsa-mir-520c	hsa-mir-101-2
hsa-mir-17	hsa-mir-30b	hsa-mir-27a	hsa-mir-146a
hsa-mir-125b-2	hsa-mir-203a	hsa-mir-221	hsa-mir-146b
hsa-let-7a-2	hsa-mir-203b	hsa-mir-222	hsa-mir-205
hsa-let-7a-3	hsa-mir-213	hsa-mir-200c	
hsa-let-7c	hsa-mir-155	hsa-mir-31	

information gain, the better that attribute is for classifying the data.

$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

(a)

$$\text{Information Gain} = \text{Entropy}(\text{before}) - \sum_{j=1}^K \text{Entropy}(j, \text{after})$$

(b)

Figure 2.16: (a) Entropy formula. (b) Information Gain formula.

- (b) *Chi-Squared* - the objective of this statistical method is to assess whether there is a relationship between a feature and the target variable. If this is verified, then it means that this feature is potentially useful for the model.

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

- O : observed frequency.
- E : expected frequency

- (c) *LASSO* - this technique starts with a simple linear regression and during training automatically chooses the most important variables, eliminating irrelevant ones by forcing their coefficients to zero. All of this works with a mechanism that penalizes irrelevant features during training [46].

Now, with these 36 miRNAs, the authors moved on to the predictive modeling phase, testing the effectiveness of different classifiers with progressively smaller subsets of miRNAs. Three main classification algorithms were used (Support Vector Machine with linear kernel, SVM with RBF kernel, and Random Forest) evaluated according to several relevant metrics: accuracy, F1-score, sensitivity, specificity, and AUC. Validation was conducted using 10-fold cross-validation, ensuring the reliability of the results.

One of the most relevant results of this study was the finding that the use of only three highly informative miRNAs could support high-quality classifications, with performance comparable (or superior) to that obtained with the complete set of 1881 miRNAs. This has significant practical implications, as it makes the models simpler, more interpretable, and potentially easier to translate into a clinical context. Although the overall accuracy is very high in all scenarios (often above 99%), the authors point out that this metric can be misleading in an unbalanced dataset. For this reason, they placed particular emphasis on the analysis of sensitivity (recall) and specificity, essential metrics in the healthcare context, where the correct identification of patients (avoiding false negatives) and non-patients (avoiding false positives) is critical.

The analysis of the results shows that sensitivity remained consistently high in all models (in many cases above 0.99), reflecting the strong ability of the classifiers to correctly identify tumor samples. Specificity, in turn, showed clear improvements after feature selection, rising from more modest values (such as 0.875 in SVM-RBF without selection) to levels above 0.98 in scenarios with only three selected miRNAs (e.g., IG-3 and CHI2-3 with SVM-RBF). This improvement reveals that reducing dimensionality helped eliminate noise and avoid misclassifying healthy samples as tumorous. Additionally, the AUC metric remained consistently high, with values close to or even equal to 1.00, reinforcing the robustness of the classifiers even in scenarios with a reduced number of variables. The results also show that the three feature selection methods (as seen in (a), (b) and (c)) led to very similar performances, with no clear advantage of one over the others.

The results obtained in the two studies analyzed in this subsection clearly demonstrate the central role that miRNAs can play in the future of cancer diagnosis. From the discovery of new candidates through exploratory pipelines based on ML, as illustrated in the work of Azari et al. [8], to the rigorous validation of miRNAs previously identified as clinically relevant, as presented in Rehman et al. [35], these small molecular regulators consistently prove to be strong indicators of tumor presence and progression. The ability to reduce thousands of candidates to a small panel of reliable biomarkers with clinical applicability, without compromising predictive performance, represents a significant step toward more accessible, non-invasive, and interpretable diagnostic tools. These findings reinforce the fundamental premise of this dissertation: ML models, when properly designed and validated, can effectively exploit the complexity of miRNAs expression profiles, not only to detect cancer, but also to support more specific tasks such as tumor subtype stratification, with real potential for clinical translation.

2.2.3 Comparison with thesis approach

The analysis of the studies presented throughout this chapter allows us not only to understand the state of the art in the use of ML algorithms for biomarker identification, but also how these computational resources can become key elements in the classification of tumor masses, directly outlining some possible paths for this thesis.

One of the first aspects to highlight is the **importance of differential expression analysis as an initial step in feature selection**. This approach was common to several of the studies analyzed, namely Wu et al. [43] and Azari et al. [8], and proved effective in reducing dimensionality in scenarios with high feature density, such as miRNA expression data. This strategy allows the analysis to focus on more informative subsets, reducing the risk of overfitting and making the models more robust. In addition, the studies mentioned demonstrate that this initial selection can (and should) be complemented with additional computational methods. Techniques such as LASSO regression, Information Gain, Chi-Squared, and Recursive Feature Elimination (RFE) have been explored in different papers and show that combining statistical approaches with model-based methods increases the probability of finding truly discriminative features. This idea will be incorporated into the pipeline developed in this thesis, with the aim of building more reliable and interpretable models.

Another valuable methodological element drawn from the analyzed works is the **rigorous validation of the identified models and biomarkers**, either through internal cross-validation or through the use of external data. The study by Azari et al. [8], for example, stands out for having validated miRNA candidates in several public repositories such as GEO, dbDMEC, and GDAC. This type of cross-validation with other sources, even if partial, **gives greater solidity and generalization to the results obtained**, a key aspect that will be taken into account in the structure of this dissertation, with the consideration of robust validation strategies. Regarding the **origin of the data**, the reviewed studies reinforce the relevance of using **The Cancer Genome Atlas (TCGA)** repository as a reliable, comprehensive, and widely used data source in similar studies. This will also be the base repository for this study, ensuring comparability with previous studies and data quality.

It is also important to mention that both studies showed that **the performance of the models does not necessarily depend on all available features**. On the contrary, several models have shown excellent performance with only a fraction of the input data (such as subsets of 3 to 10 miRNAs) highlighting the feasibility of building simpler, more interpretable, and economically applicable models in a clinical context. The possibility of finding a small group of miRNAs that would have the same performance as using all of them is the real breakthrough of a thesis like this, since it would open doors to more affordable testing, leading to more prevention tests and the possibility of earlier detection of these malignant cells.

Finally, the workflow diagrams proposed in the reference studies - with special emphasis on that of Azari et al. [8] - offer a clear, sequential, and methodologically rigorous

visual representation of all the steps involved in the discovery and validation of molecular biomarkers using ML. These flowcharts not only facilitate understanding of the process, but also help to ensure the reproducibility of studies and effective communication of methods between researchers. In the particular case of this dissertation, the structure presented in these works serves as direct inspiration for the design of the pipeline that will be implemented here, covering all critical phases: from data acquisition and pre-processing, through careful feature selection and predictive modeling, to internal statistical validation and model performance analysis. This systematic approach ensures alignment with best practices in the literature and maximizes the translational potential of the results obtained.

In short, the literature review allows us to consolidate a solid methodological basis, anchored in studies that have demonstrated the effectiveness of machine learning in the discovery and validation of molecular biomarkers. From this evidence, it is clear that well-structured approaches such as combining differential expression analysis, careful feature selection, and rigorous validation offer a promising path to building robust and clinically relevant predictive models. This accumulated knowledge will now be adapted and applied to the specific context of this dissertation, which seeks to identify miRNA signatures with discriminatory power for the classification of breast cancer subtypes. The next chapter describes in detail how this approach will be implemented.

PRELIMINARY WORK

This chapter presents the initial experiments using machine learning for breast cancer molecular subtype classification based on the patients' microRNAs expression levels. Early work focused on testing baseline ML models (Logistic Regression, SVM, Random Forest) to assess their classification performance using datasets such as TCGA-BRCA [39].

3.1 Dataset and Pre-processing

The dataset utilized in this study is derived from the Cancer Atlas Program (TCGA), a widely recognized and reliable source of this particular type of data. The TCGA-BRCA project provides researchers with access to a comprehensive set of data and resources, including microRNA expression profiles associated with breast cancer patients, gene expression data, clinical data, and additional data from various sources. As demonstrated in Table 3.1, this dataset poses significant challenges, primarily concerning its dimensionality, class imbalance, redundancy of microRNA expression values, and feature names. The absence of standardized naming guidelines for miRNAs further complicates this issue. For instance, a miRNA may be designated as "*hsa-mir-145*" or "*MIR-145*," necessitating the incorporation of a "translation" layer during the pre-processing stage to ensure compatibility with data from external projects or repositories.

Table 3.1: Dataset Description wtih the number of samples, microRNAs expression features, and clinical attributes before and after pre-processing.

Statistics	Before pre-processing	After pre-processing
# Samples	463	256
# miRNAs	888	442 ¹
# Clinical Features	22	20
# Classes	5	4 ¹

¹Only applied in the second part of the tests, more details in 3.3

During the data preparation phase, we eliminated features that had more over 60% of missing values and, additionally, the "Normal-like" class due to its limited sample size of only five instances, which is insufficient for training a reliable model. In terms of samples, there was a few samples with missing values in some miRNAs and clinical features, which were removed from the dataset to prevent the use of SMOTE or other imputation methods that could introduce bias or noise into the data. Even though we have 888 different miRNA features, there are some columns that are duplicates of others - why? well, let's consider the base *MIR-16* as an example, this miRNA has 6 different forms of expression: *MIR-16/3P*, *MIR-16/16*, *MIR-16/5P*, *MIR-16-1/16*, *MIR-16-13/16*, *MIR-16-2/16* -, all of these have the same values because the extraction technique that was used to obtain the miRNA expression values is not sensitive enough to detect the differences. Given that, to remove redundancy, we decided to remove all miRNAs that share the same base, ending up with 442 miRNAs. Regarding the contextual importance of features, there are some features (i.e. study Id) that were neither important, nor their values had much variance, which would not benefit the model. Adding to all of this, I mapped the categorical features to numerical values to ensure that the models could process them correctly. The table provides a summary of the dataset after pre-processing, including the number of samples, features, and clinical attributes.

3.2 First tests with baseline models

Given the reduced data set, we did an initially split of 80% for training and 20% for testing, with no validation set. The models were trained using the *scikit-learn* library, and the results were evaluated using accuracy, precision, recall, and F1-score metrics. First, I evaluated the quality of the split using a 2D Projection of the data using Principal Component Analysis (PCA) wiht k = 2 PCA is a dimensionality reduction technique that transforms the data into a lower dimensional space while preserving the variance of the data. The PCA projection is shown in Figure 3.1. The PCA projection shows that the training and testing sets are well mixed, indicating that the split is representative of the data. For this evaluation, I also plotted the proportion of each class in the training and testing sets, as shown in Figure 3.2, which indicates that the split is balanced and representative of the data.

After the split, I trained three baseline models: Logistic Regression, XGBoost, and Naive Bayes. The results of the models are shown in Table 3.2.

3.3 Second Segment of Experiments

In this section, I re-processed the dataset from scratch to ensure a clean and leakage-free setup. Given the concerns raised about the normalization procedures in the first segment, I implemented a more robust pre-processing pipeline to address these issues. The new pipeline includes the following steps:

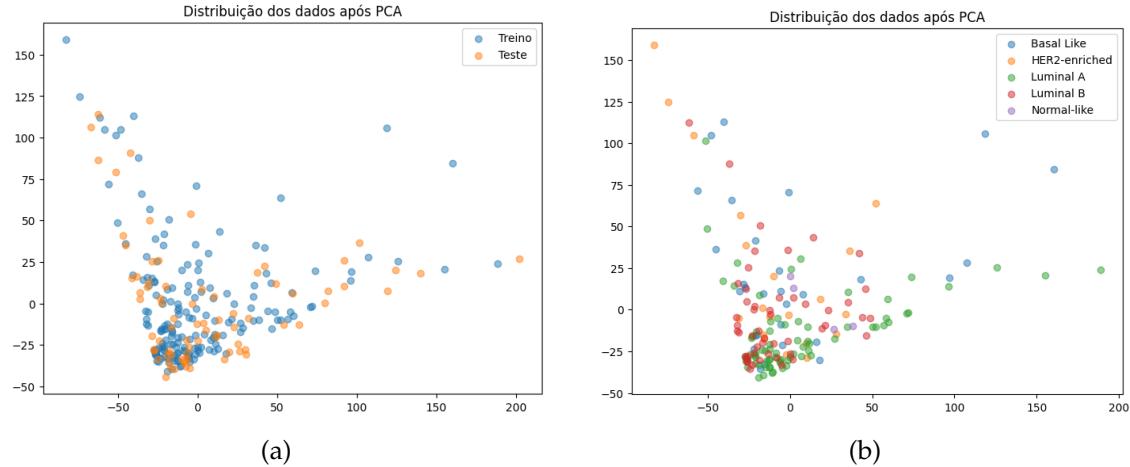


Figure 3.1: (a) PCA projection of the training and testing sets. The colors represent the different sets, and the points represent the samples. (b) Same PCA projection, but with the classes represented by different colors.

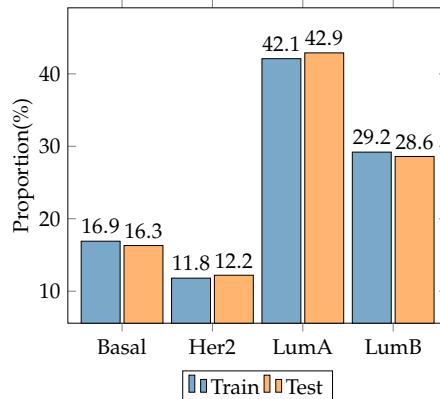


Figure 3.2: Class proportions in the training and testing sets. The colors represent the different sets, and the bars represent the proportion of each class.

1. Removal of features with more than 60% missing values.
2. Removal of the "Normal-like" class due to its limited sample size.
3. Removal of redundant miRNAs based on their base names.
4. Removal of features that were neither important nor had much variance.
5. Mapping of categorical features to numerical values.

Then, I split the dataset into training and testing sets using an 80/20 split, ensuring that the split is representative of the data. Only after the split, I applied normalization to the training set using the *StandardScaler* from the *scikit-learn* library. The normalization was then applied to the testing set using the same scaler fitted on the training set. This approach ensures that the normalization is consistent across the training and testing sets, preventing data leakage. A scaler is the mean and standard deviation of the training set,

Table 3.2: Results of the baseline models on the testing set. The table shows the precision, recall, F1-score, and support for each class, as well as the overall accuracy of the models.

Model	Class	Precision	Recall	F1-score	Support
Random Forest	0	0.89	1.00	0.94	8
	1	1.00	0.17	0.29	6
	2	0.64	0.76	0.70	21
	3	0.50	0.50	0.50	14
<i>Overall Accuracy</i>					0.65
XGBoost	0	1.00	1.00	1.00	8
	1	0.75	0.50	0.60	6
	2	0.86	0.90	0.88	21
	3	0.73	0.79	0.76	14
<i>Overall Accuracy</i>					0.84
Naive Bayes	0	0.80	1.00	0.89	8
	1	0.40	0.33	0.36	6
	2	0.79	0.52	0.63	21
	3	0.45	0.64	0.53	14
<i>Overall Accuracy</i>					0.61

which is then used to normalize the testing set. This approach ensures that the testing set is not used to fit the scaler, preventing data leakage. The formula used for normalization, from Pedregosa et al. [30].

With this new pre-processing pipeline, I re-ran the same baseline models and leveraged the use of greater computational resources for hyperparameter tuning using a grid search approach. These efforts showed that XGBoost outperformed any other model tested, including latent space models such as *DIABLO*.

The comparison between baseline models (trained on potentially leaked data) and heterogeneous models (trained on properly preprocessed data with clinical features) shows a clear improvement in model generalizability and robustness.

In the first setup, XGBoost achieved high performance (accuracy = 84%), but likely benefited from data leakage. Once the pipeline was corrected, all models showed more realistic and balanced results. XGBoost remained the top performer (accuracy = 84%, F1 = 1.00 for Basal-like, 0.85 for Luminal A, and 0.77 for Luminal B), while Logistic Regression also performed well (accuracy = 80%).

However, HER2-enriched and Luminal B subtypes consistently showed lower F1-scores across models. This bottleneck is likely due to their low representation in the dataset and the biological similarity of their miRNA expression profiles, which reduces class separability and hampers learning.

Overall, the results confirm the importance of careful preprocessing and the benefit of

Table 3.3: Results of the heterogeneous models (miRNAs + Clinical Data) on the testing set. The table shows overall performance and F1-score per class.

Model	Class	Precision	Recall	F1-score	Ovr. Accuracy
DIABLO All Blocks	Basal	1.00	1.00	1.00	0.80
	HER2	0.50	0.33	0.40	
	Luminal A	0.75	1.00	0.86	
	Luminal B	0.89	0.57	0.70	
Decision Tree	Basal	0.80	1.00	0.89	0.61
	HER2	0.60	0.50	0.55	
	Luminal A	0.62	0.62	0.62	
	Luminal B	0.46	0.43	0.44	
Logistic Regression	Basal	1.00	0.88	0.94	0.80
	HER2	0.67	0.67	0.55	
	Luminal A	0.90	0.86	0.87	
	Luminal B	0.63	0.71	0.67	
XGBoost	Basal	1.00	1.00	1.00	0.84
	HER2	1.00	0.67	0.67	
	Luminal A	0.83	0.95	0.85	
	Luminal B	0.85	0.79	0.77	

integrating clinical data, especially for improving performance on harder-to-distinguish subtypes.

WORK PLAN

Figure presents a Gantt chart to illustrate an outline of the timeline and key milestones from this point forward. The chart is divided into several phases, each with specific tasks and deadlines.

Task	July	August	September	October	November	December	January	February	March
1. Assessing Statistical Robustness of Preliminary Results									
2. Pipeline Consolidation with Non-linear Models									
3. Feature Selection from microRNA and Clinical Attributes									
4. Improving Statistical Robustness of microRNA Biomarkers Identification									
5. Evaluation Cross-Checking with In-vitro/In-vivo Analysis									
6. Thesis writing									

Figure 4.1: Gantt chart illustrating the work plan and expected timeline for the project.

Although the Gantt chart provides a visual overview of the timeline, some details may not be fully legible. For clarity, the main tasks and milestones are listed below.

1. Assessing Statistical Robustness of Preliminary Results
2. Pipeline Consolidation with Non-linear Models
3. Feature Selection from microRNA and Clinical Attributes
4. Improving Statistical Robustness of microRNA Biomarkers Identification
5. Evaluation Cross-Checking with In-vitro/In-vivo Analysis
6. Thesis writing

BIBLIOGRAPHY

- [1] B. Adamo et al. "Clinical implications of the intrinsic molecular subtypes of breast cancer." In: *Breast* 24 Suppl 2 (2015), S26–35. doi: [10.1016/j.breast.2015.07.008](https://doi.org/10.1016/j.breast.2015.07.008).
- [2] AI Lab MTI Vietnam. *Imagen ilustrativa de análise de dados*. <https://ailab.mti-vietnam.vn/wp-content/uploads/2020/12/image-2.png>. Imagem retirada do site. 2020. (Visited on 2025-07-03).
- [3] Saad Albawi, T. A. Mohammed, and Saad Al-Zawi. "Understanding of a convolutional neural network". In: *2017 International Conference on Engineering and Technology (ICET)* (2017), pp. 1–6. doi: [10.1109/ICENGTECHNOL.2017.8308186](https://doi.org/10.1109/ICENGTECHNOL.2017.8308186).
- [4] Analytics Vidhya. *Processo de análise de dados*. <https://www.analyticsvidhya.com/wp-content/uploads/2016/08/Artificial-Intelligence-Neural-Network-Nodes.jpg>. Imagem retirada do site. s.d. (Visited on 2025-07-03).
- [5] Anastasios N. Angelopoulos and Stephen Bates. "Conformal Prediction: A Gentle Introduction". In: *Found. Trends Mach. Learn.* 16.4 (2023-03), pp. 494–591. issn: 1935-8237. doi: [10.1561/2200000101](https://doi.org/10.1561/2200000101). URL: <https://doi.org/10.1561/2200000101>.
- [6] Arul Edwin Raj Anthony Muthu, Sundaram Muniasamy, and Thirassama Jaya. "Thermography based breast cancer detection using self-adaptive gray level histogram equalization color enhancement method". In: *International Journal of Imaging Systems and Technology* 31 (2020-10). doi: [10.1002/ima.22488](https://doi.org/10.1002/ima.22488).
- [7] M. Arnold et al. "Current and future burden of breast cancer: Global statistics for 2020 and 2040". In: *The Breast : Official Journal of the European Society of Mastology* 66 (2022), pp. 15–23. doi: [10.1016/j.breast.2022.08.010](https://doi.org/10.1016/j.breast.2022.08.010).
- [8] H. Azari et al. "Machine learning algorithms reveal potential miRNAs biomarkers in gastric cancer". In: *Scientific Reports* 13 (2023). doi: [10.1038/s41598-023-32332-x](https://doi.org/10.1038/s41598-023-32332-x).
- [9] C. Blenkiron et al. "MicroRNA expression profiling of human breast cancer identifies new markers of tumor subtype". In: *Genome Biology* 8 (2007), R214–R214. doi: [10.1186/gb-2007-8-10-r214](https://doi.org/10.1186/gb-2007-8-10-r214).

- [10] Freddie Bray et al. "Global cancer statistics 2022: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries". In: *CA: A Cancer Journal for Clinicians* 74.3 (2024), pp. 229–263. doi: <https://doi.org/10.3322/caac.21834>. eprint: <https://acsjournals.onlinelibrary.wiley.com/doi/pdf/10.3322/caac.21834>. URL: <https://acsjournals.onlinelibrary.wiley.com/doi/abs/10.3322/caac.21834>.
- [11] Joel S Brown et al. "Updating the Definition of Cancer". In: *Molecular Cancer Research* 21 (2023), pp. 1142–1147. doi: [10.1158/1541-7786.MCR-23-0411](https://doi.org/10.1158/1541-7786.MCR-23-0411).
- [12] Gloria M. Calaf et al. "The Role of MicroRNAs in Breast Cancer and the Challenges of Their Clinical Application". In: *Diagnostics* 13 (2023). doi: [10.3390/diagnostics13193072](https://doi.org/10.3390/diagnostics13193072).
- [13] Andre Esteva et al. "Dermatologist-level classification of skin cancer with deep neural networks". In: *Nature* 542.7639 (2017), pp. 115–118. doi: [10.1038/nature21056](https://doi.org/10.1038/nature21056). URL: <https://doi.org/10.1038/nature21056>.
- [14] S. W. Fu, L. Chen, and Y. G. Man. "miRNA Biomarkers in Breast Cancer Detection and Management". In: *Journal of Cancer* 2 (2011), pp. 116–122. doi: [10.7150/jca.2.116](https://doi.org/10.7150/jca.2.116). URL: <https://doi.org/10.7150/jca.2.116>.
- [15] Maryellen L. Giger. "The Role of Artificial Intelligence in Early Cancer Diagnosis". In: *Radiologic Clinics of North America* 59.6 (2021), pp. 933–946. doi: [10.1016/j.rcl.2021.06.006](https://doi.org/10.1016/j.rcl.2021.06.006). URL: <https://doi.org/10.1016/j.rcl.2021.06.006>.
- [16] L. Gulyaeva and N. E. Kushlinskiy. "Regulatory mechanisms of microRNA expression". In: *Journal of Translational Medicine* 14 (2016). doi: [10.1186/s12967-016-0893-x](https://doi.org/10.1186/s12967-016-0893-x).
- [17] S. Hammond. "An overview of microRNAs." In: *Advanced drug delivery reviews* 87 (2015), pp. 3–14. doi: [10.1016/j.addr.2015.05.001](https://doi.org/10.1016/j.addr.2015.05.001).
- [18] P. T. Ho, I. Clark, and L. Le. "MicroRNA-Based Diagnosis and Therapy". In: *International Journal of Molecular Sciences* 23 (2022). doi: [10.3390/ijms23137167](https://doi.org/10.3390/ijms23137167).
- [19] R. Hong and Bing-he Xu. "Breast cancer: an up-to-date review and future perspectives". In: *Cancer Communications* 42 (2022), pp. 913–936. doi: [10.1002/cac2.12358](https://doi.org/10.1002/cac2.12358).
- [20] N. Howlader et al. "Differences in Breast Cancer Survival by Molecular Subtypes in the United States". In: *Cancer Epidemiology, Biomarkers and Prevention* 27 (2018), pp. 619–626. doi: [10.1158/1055-9965.EPI-17-0627](https://doi.org/10.1158/1055-9965.EPI-17-0627).
- [21] International Agency for Research on Cancer. *Global Cancer Observatory - Cancer Today: Breast cancer incidence heatmap* (2022). <https://gco.iarc.fr/today/en/dataviz/maps-heatmap?mode=population&types=0&cancers=20>. Accessed: 2025-06-22. 2022.

BIBLIOGRAPHY

- [22] Bahzad Taha Jijo and Adnan Mohsin Abdulazeez. "Classification Based on Decision Tree Algorithm for Machine Learning". In: *Journal of Applied Science and Technology Trends* 2.1 (2021), pp. 20–28. ISSN: 2708-0757. doi: [10.38094/jastt20165](https://doi.org/10.38094/jastt20165). URL: <http://www.jastt.org/index.php/index>.
- [23] Bahzad Taha Jijo and Adnan Mohsin Abdulazeez. "Classification Based on Decision Tree Algorithm for Machine Learning". In: *Journal of Applied Science and Technology Trends* 02.01 (2021), pp. 20–28. ISSN: 2708-0757. doi: [10.38094/jastt20165](https://doi.org/10.38094/jastt20165). URL: <http://www.jastt.org/index.php/index>.
- [24] Rosalind C. Lee, Rhonda L. Feinbaum, and Victor Ambros. "The *C. elegans* heterochronic gene lin-4 encodes small RNAs with antisense complementarity to lin-14". In: *Cell* 75 (1993), pp. 843–854. doi: [10.1016/0092-8674\(93\)90317-3](https://doi.org/10.1016/0092-8674(93)90317-3).
- [25] Yuxun Luo et al. "Machine learning in the development of targeting microRNAs in human disease". In: *Frontiers in Genetics* 13 (2023), p. 1088189. doi: [10.3389/fgene.2022.1088189](https://doi.org/10.3389/fgene.2022.1088189). URL: <https://doi.org/10.3389/fgene.2022.1088189>.
- [26] D. Mahalingam, E. Vagia, and M. Cristofanilli. "The Landscape of Targeted Therapies in TNBC". In: *Cancers* 12 (2020). doi: [10.3390/cancers12040916](https://doi.org/10.3390/cancers12040916).
- [27] Bárbara B. Mendes et al. "Nanodelivery of nucleic acids". In: *Nature Reviews Methods Primers* 2.24 (2022). Article citation ID: (2022) 2:24. doi: [10.1038/s43586-022-00104-y](https://doi.org/10.1038/s43586-022-00104-y). URL: <https://doi.org/10.1038/s43586-022-00104-y>.
- [28] Juan P Muñoz et al. "The Role of MicroRNAs in Breast Cancer and the Challenges of Their Clinical Application". In: *Diagnostics* 13 (2023). doi: [10.3390/diagnostics13193072](https://doi.org/10.3390/diagnostics13193072).
- [29] National Cancer Institute. *What Is Cancer?* Accessed: 2025-06-15. 2021. URL: <https://www.cancer.gov/about-cancer/understanding/what-is-cancer>.
- [30] Fabian Pedregosa et al. "Scikit-learn: Machine Learning in Python". In: *Journal of Machine Learning Research* 12 (2011), pp. 2825–2830.
- [31] Charles M. Perou et al. "Molecular portraits of human breast tumours". In: *Nature* 406.6797 (2000), pp. 747–752. doi: [10.1038/35021093](https://doi.org/10.1038/35021093). URL: <https://doi.org/10.1038/35021093>.
- [32] K. Polyak. "Breast cancer: origins and evolution." In: *The Journal of clinical investigation* 117 11 (2007), pp. 3155–63. doi: [10.1172/JCI33295](https://doi.org/10.1172/JCI33295).
- [33] A. Prat et al. "Clinical implications of the intrinsic molecular subtypes of breast cancer." In: *Breast* 24 Suppl 2 (2015), S26–35. doi: [10.1016/j.breast.2015.07.008](https://doi.org/10.1016/j.breast.2015.07.008).
- [34] Leslie A. Pray. "Discovery of DNA structure and function: Watson and Crick". In: *Nature Education* 1.1 (2008), p. 100. URL: <https://www.nature.com/scitable/topicpage/discovery-of-dna-structure-and-function-watson-397/>.
- [35] Oneeb Rehman et al. "Validation of miRNAs as Breast Cancer Biomarkers with a Machine Learning Approach." In: *Cancers* (2019). doi: [10.3390/cancers11030431](https://doi.org/10.3390/cancers11030431).

- [36] H. Romanowicz, B. Smolarz, and Anna Zadrożna Nowak. "Breast Cancer—Epidemiology, Classification, Pathogenesis and Treatment (Review of Literature)". In: *Cancers* 14 (2022). doi: [10.3390/cancers14102569](https://doi.org/10.3390/cancers14102569).
- [37] Amrit Singh et al. "DIABLO: an integrative approach for identifying key molecular drivers from multi-omics assays". In: *Bioinformatics* 35.17 (2019), pp. 3055–3062. doi: [10.1093/bioinformatics/bty1054](https://doi.org/10.1093/bioinformatics/bty1054). URL: <https://doi.org/10.1093/bioinformatics/bty1054>.
- [38] U. Testa, G. Castelli, and E. Pelosi. "Breast Cancer: A Molecularly Heterogenous Disease Needing Subtype-Specific Treatments". In: *Medical Sciences* 8 (2020). doi: [10.3390/medsci8010018](https://doi.org/10.3390/medsci8010018).
- [39] The Cancer Genome Atlas (TCGA) Research Network. *TCGA-BRCA Project: Breast Invasive Carcinoma*. <https://portal.gdc.cancer.gov/projects/TCGA-BRCA>. Accessed: 2025-07-11. 2011.
- [40] J. Watson and F. Crick. "The structure of DNA." In: *Cold Spring Harbor symposia on quantitative biology* 18 (1953), pp. 123–31. doi: [10.1101/SQB.1953.018.01.020](https://doi.org/10.1101/SQB.1953.018.01.020).
- [41] Thomas J. Watson. "An empirical study of the naive Bayes classifier". In: 2001. URL: <https://api.semanticscholar.org/CorpusID:14891965>.
- [42] Bruce Wightman, Iva Ha, and Gary Ruvkun. "Posttranscriptional regulation of the heterochronic gene lin-14 by lin-4 mediates temporal pattern formation in C. elegans". In: *Cell* 75 (1993), pp. 855–862. doi: [10.1016/0092-8674\(93\)90318-4](https://doi.org/10.1016/0092-8674(93)90318-4).
- [43] Jiande Wu et al. "Breast Cancer Type Classification Using Machine Learning." In: *Journal of Personalized Medicine* (2021). doi: [10.3390/jpm11020061](https://doi.org/10.3390/jpm11020061).
- [44] Mohammed Yaseen and Adnan Mohsin Abdulazeez. "Performance Evaluation Metrics in Machine Learning Models: A Comparative Study". In: *Journal of Soft Computing and Data Mining* 2.2 (2021), pp. 21–31. ISSN: 2710-139X. doi: [10.30880/jscdm.2021.02.02.003](https://doi.org/10.30880/jscdm.2021.02.02.003).
- [45] S. Yeo and J. Guan. "Breast Cancer: Multiple Subtypes within a Tumor?" In: *Trends in cancer* 3.11 (2017), pp. 753–760. doi: [10.1016/j.trecan.2017.09.001](https://doi.org/10.1016/j.trecan.2017.09.001).
- [46] Huaqing Zhang et al. "Feature Selection for Neural Networks Using Group Lasso Regularization". In: *IEEE Transactions on Knowledge and Data Engineering* 32.4 (2020), pp. 659–673. doi: [10.1109/TKDE.2019.2893266](https://doi.org/10.1109/TKDE.2019.2893266).
- [47] Lili Zhu and P. Spachos. "Support Vector Machine and YOLO for a Mobile Food Grading System". In: *Internet of Things* 13 (2021-01), p. 100359. doi: [10.1016/j.iot.2021.100359](https://doi.org/10.1016/j.iot.2021.100359).

A

APPENDIX 1 COVERS SHOWCASE

UNVEILING MICRORNA BIOMARKERS FOR BREAST CANCER SUB-TYPING

This annex includes the full version of the paper titled "**Unveiling microRNA Biomarkers for Breast Cancer Sub-typing using Discriminative Models**", submitted and accepted to EPIA 2025 under the category AI for Medicine.

Unveiling microRNA biomarkers for Breast Cancer Sub-typing using Discriminative Models

Anonymous Submission

No Institute Given

Abstract. Breast cancer is a heterogeneous disease comprising multiple molecular subtypes, each with distinct clinical outcomes and therapeutic challenges. MicroRNAs (miRNAs), as key regulators of gene expression, hold great promise as biomarkers for cancer subtyping and developing personalized treatments. In this paper, we propose a machine-learning discriminative modeling framework to uncover subtype-specific miRNA biomarkers for breast cancer. Our approach jointly integrates miRNA expression data with patient clinical data, to identify miRNA signatures that differentiate between luminal A, luminal B, HER2-enriched, and Basal-like subtypes. We conduct extensive validation across multiple discriminative models and provide evidence that microRNAs are strong discriminators within breast cancer subtypes, achieving up to 90% F1. Furthermore, we contrast our findings against state-of-the-art multi-omics integration and biomarker discovery, and provide a reassessment of its predicted miRNA signatures. To this end, a model explainability approach is employed to analyze and pinpoint subtype-specific miRNA profiles. These potentially highlight subtype-specific biologically meaningful and functionally relevant miRNAs, that can now be therapeutically validated through in vitro and in vivo experiments.

Keywords: breast cancer · microRNA · biomarker discovery · machine learning.

1 Introduction

Breast cancer is one of the most prevalent and challenging malignancies in women worldwide, with complex molecular mechanisms driving its progression and treatment resistance. It affects millions of women and causes thousands of deaths annually. Its natural heterogeneity presents a clinical challenge, by being categorized according to the expression of hormone receptors, namely estrogen receptor (ER) and progesterone receptor (PR), and human epidermal growth factor receptor 2 (HER2). These molecular subtypes, Luminal A, Luminal B, HER2-enriched, and Basal-like (also referred as Triple Negative), differ significantly in prognosis, response to treatment and chance of metastasis. Identifying new biomarkers and personalized therapeutics is urgently necessary to reduce mortality and address the critical lack of targeted treatment options.

Recent studies have demonstrated that microRNAs (miRNAs), non-coding RNAs involved in gene regulation, play a significant role in breast cancer development [15,16]. With miRNAs expression levels demonstrating unique signatures across breast cancer molecular subtypes, the necessity of systematically pinpointing a restricted set of relevant microRNAs for each subtype, that can posteriorly be analysed through in-vitro and in-vivo analyses, is paramount, towards the design of gene-based nanotherapies. Machine learning-based approaches have demonstrated great potential in cancer profiling [6,4,10]. These have demonstrated to be capable of leveraging high-dimensional multi-omics data, and either uncover latent microRNA, genes and proteins interactions [14] or cancer diagnosis [4].

In this work we introduce a multi-source breast cancer biomarker identification framework, combining individuals' microRNAs expression levels and clinical attributes (e.g., age, disease progression indicators, among others). We follow related work, applied to gastric cancer biomarker identification [4], and follow an approach where discriminative classification tasks are used as a proxy to assess multi-source feature relevance. In particular, we leverage well-established but explainable discriminative classification models to robustly conduct breast cancer subtyping in two fundamental settings: 1) miRNA-only expression-levels, and 2) hybrid miRNA and patient clinical data.

Our systematic evaluation demonstrates that not only breast cancer subtypes can be robustly predicted with high effectiveness solely from miRNA expression levels (up to 87% precision), but that integrating patient clinical features positively improves performance. From these high-performing models, we further conduct a model explainability study where we systematically identify key responsible miRNA biomarkers per subtype, drawing key insights to both subtype-specific and shared microRNA biomarkers.

2 Related Work

Recent advances have positioned miRNAs as highly promising biomarkers in breast cancer research. Their stability in biological fluids and specificity for tumor subtypes make them particularly suitable for non-invasive diagnostics such as liquid biopsies [15,16]. Panels of miRNAs have demonstrated strong discriminatory power across breast cancer subtypes, aiding early detection and patient stratification for tailored therapies [11,17]. Simultaneously, machine learning (ML) techniques have proven to be powerful tools for cancer subtype classification. Supervised classifiers outperform traditional methods like immunohistochemistry by effectively leveraging high-dimensional datasets, including miRNA expression profiles [7]. Moreover, these facilitate the identification of the most predictive features through embedded techniques like regularization and recursive feature elimination, enhancing both classification accuracy and biological interpretability.

However, key challenges remain. Class imbalance—particularly for rare subtypes like Basal-like—continues to hinder generalizability, requiring strategies such as resampling and synthetic data generation (e.g., SMOTE) [2]. Moreover, while complex models often offer improved performance, their opacity limits

clinical trust. This has spurred the adoption of explainable AI (XAI) methods like SHAP [8] and LIME [13], which attribute model predictions to input features, increasing transparency and biological credibility [5]. Research in adjacent domains reinforces these trends. Azari et al. extended this methodology using multiple ML models to detect 29 miRNAs linked to gastric cancer prognosis, yielding 93% accuracy and 88,5% AUC with SVM [1]. These studies illustrate the potential of ML-driven miRNA biomarker discovery, which is now being translated into breast cancer research with increasing success.

Therapeutically, miRNAs are gaining traction beyond diagnostics. Approaches such as AntagomiRs and miRNA mimics are under clinical investigation for restoring tumor-suppressive miRNA activity or inhibiting oncogenic miRNAs [12,9]. Additionally, nanocarrier systems are being developed to deliver these molecules with high specificity, offering promising routes for RNA-based personalized therapies [3].

Together, these advances suggest that ML-enabled miRNA analysis holds transformative potential for breast cancer subtype classification, early detection, and individualized treatment design.

3 A Multi-source Breast Cancer Biomarker Identification Framework

In this work, we introduce a general-purpose framework for identifying microRNA biomarkers relevant to breast cancer subtyping using discriminative machine learning models. Rather than relying on latent variable methods or integration through component extraction, our approach focuses on a direct supervised learning strategy, leveraging both miRNA and clinical data as sources of predictive signal.

Our pipeline is composed of three core stages: data preparation, discriminative modeling, and biomarker evaluation. As illustrated in Figure 1, we process miRNA expression and clinical data independently, cleaning each source before concatenating the features into a unified input space for posterior split and normalization. This integration enables the models to access multi-source information directly, without the need for explicit projection into a shared latent space.

Once the data is prepared, we train classification models to predict breast cancer subtypes using two configurations: homogeneous (miRNA only) and heterogeneous (miRNA + clinical). The goal is not only to optimize subtype classification performance, but also to uncover which features — especially miRNAs — are most indicative of each subtype. We later employ explainability techniques, such as SHAP and coefficient analysis, to extract and rank relevant microRNAs from our top-performing model.

The entire workflow can be summarized in five key steps:

1. **Data curation and preprocessing**, including filtering of redundant miRNAs and cleaning of clinical records.
2. **Feature integration** by concatenating miRNA and clinical data into a single matrix.

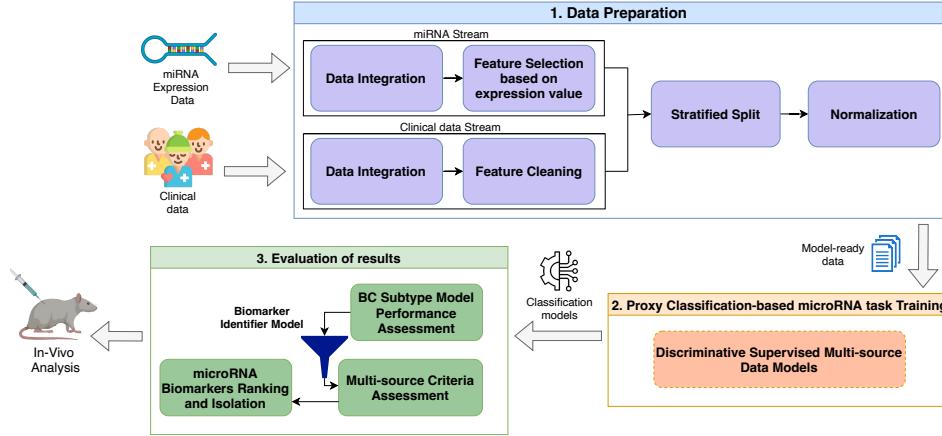


Fig. 1. Unified framework for miRNA biomarker identification and subtyping in Breast Cancer, using discriminative models.

3. **Model training** under homogeneous (miRNA-only) and heterogeneous (miRNA + clinical) settings.
4. **Model evaluation** using accuracy, precision, recall and F1-score.
5. **Feature importance analysis** to identify subtype-specific biomarkers.

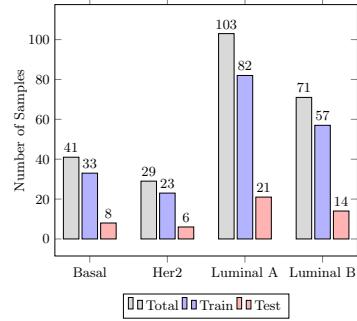
Discriminative Models Considering the limited size of our dataset (244 individuals/samples with an aggregate set of 462 features — more details in Section 4.1), we opt for classical discriminative models that are well-suited to high-dimensional, low-sample-size scenarios, offering a more robust and interpretable foundation for biomarker identification. Larger models, with higher number of parameters to tune, would pose a severe risk of overfitting. Therefore, To support our framework, we selected three discriminative models that balance performance and interpretability:

- **Decision Tree:** A decision tree classifier, using the Gini impurity as classification criteria, offering fully interpretable decision rules.
- **Logistic Regression:** A linear multinomial logistic regression classifier, with ℓ_2 regularization.
- **XGBoost:** A gradient-boosted decision tree model, optimized with a multinomial negative likelihood.

This combination of tree-like and regression models, including simpler (Decision Tree) to boosted models (XGBoost), enable us to understand how discriminative models of different complexities, can leverage microRNAs to conduct breast cancer profiling.

Fig. 2. Descriptive statistics of the TCGA full dataset (left table), sample clinical features (right table) and histogram per breast cancer subtype (right figure).

Sample Clinical Attrs.	
Statistic	Total
# Samples	244
# Individuals	244
# microRNAs	442
# Clinical Attributes	20
...	



4 Evaluation

4.1 Dataset

The primary dataset utilized in this study is a curated subset of the publicly available TCGA (The Cancer Genome Atlas) database, specifically comprising miRNA expression profiles from breast cancer patients. Initially, the dataset included expression levels for 888 distinct miRNAs across 251 patients.

Upon detailed examination, it was observed that several miRNAs with the same base identifier—for example, *MIR-758*, *MIR-758/3P*, and *MIR-758/5P*—exhibited identical expression values across all samples. This redundancy is likely due to limitations in the measurement technique, which may lack the sensitivity required to distinguish between these closely related variants. Consequently, to reduce redundancy and eliminate features that do not contribute additional variance or informative signal to the model, all such miRNAs sharing the same base and expression profile were removed. This preprocessing step reduced the number of miRNA features from 888 to 464, eliminating 446 non-informative miRNAs.

In addition to the miRNA expression data, clinical details from the same dataset, containing complementary patient information was incorporated. As shown in Figure 2, it comprises features such as age, cancer stage, presence of metastasis, number of lymph nodes, among others. These clinical variables are particularly valuable as they are often readily accessible without the need for costly procedures and may provide significant predictive power for subtype classification.

Features with more than 90% missing values were excluded to ensure data integrity. Additionally, six patients were discarded due to incomplete clinical attributes. Given the intrinsic biological heterogeneity among individuals—where similar diagnoses can result in divergent outcomes—it was deemed inappropriate to impute missing values. All features were appropriately encoded and normalized. In particular, to mitigate the risk of data leakage, z-score normalization was applied exclusively based on the training set: the scaling parameters were

computed from the training data and then used to transform both the training and testing sets. This normalization was applied to all continuous numerical features, including miRNA expression values and clinical variables.

In overall, we obtain a final dataset of 244 patients (samples), comprising a total of 462 features (microRNAs and clinical attributes), annotated with the corresponding breast cancer subtype.

4.2 Protocol

To establish a solid evaluation benchmark, we trained three classical machine learning models: **Decision Tree**, **Logistic Regression** and **XGBoost**. These were selected due to their proven performance on high-dimensional biological datasets and their interpretability [4,10].

To benchmark the effectiveness of our discriminative pipeline, we compare against **DIABLO** [14]¹, a leading multi-omics integration method based on latent variable modeling. While DIABLO aims to capture correlated structures across modalities, our pipeline aims purely at optimizing discriminative performance. We applied the recommended cross-validation and hyperparameter tuning steps.

The combined dataset of miRNA expression values and clinical features (see Section 4.1) was randomly split into training (80%) and testing (20%) subsets using stratified sampling to preserve class distribution. This resulted in 195 training and 49 testing samples (see Figure 2).

Scaling parameters (mean and standard deviation) were computed from the training data and applied to both training and testing subsets. This normalization was applied to all continuous variables, including miRNA expression levels and clinical features. All models were evaluated using well-established classification metrics: *accuracy*, *precision*, *recall* and *F1-score*.

4.3 Results

To evaluate the performance of discriminative models in breast cancer subtype profiling, we conducted experiments using both homogeneous data (with miRNAs expression values only) and heterogeneous data (miRNA + clinical features from the patient). This section presents a comparative analysis of model performance across different metrics and class distributions, succeeded by a discussion on the various advantages of discriminative models over the latent representation approach used in DIABLO.

Performance of Discriminative Models using miRNAs Only In this first segment, we assessed the predictive performance of three discriminative models - Decision Trees , XGBoost, Logistic Regression - trained solely on microRNA expression values. Table 1 presents their results in comparison with DIABLO using a single omics block (in DIABLO, datasets of different natures are called omics blocks - in this case, the single omics block that was used to classify the subtype of BC was miRNA).

Table 1. Homogeneous classifiers (miRNAs only). Results are shown in percentage.

Model	Overall				Basal HER-2 Lum A Lum B			
	P	R	F1	Acc			F1	
DIABLO Single Block	75.27	67.26	67.69	73.47	100.00	40.00	80.77	50.00
Decision Tree	57.35	56.40	56.79	57.14	87.50	36.36	61.90	41.38
Logistic Regression	80.26	75.60	73.50	77.55	94.12	61.54	85.71	52.63
XGBoost	86.96	74.40	76.73	75.51	100.00	80.00	76.92	50.00

Table 2. Heterogeneous models (miRNAs + Clinical Data). Results are shown in percentage.

Model	Overall				Basal HER-2 Lum A Lum B			
	P	R	F1	Acc			F1	
DIABLO All Blocks	78.47	72.62	73.82	79.59	100.00	40.00	85.71	69.57
Decision Tree	62.01	63.69	62.45	61.22	88.89	54.55	61.90	44.44
Logistic Regression	77.22	75.60	75.57	79.59	94.12	54.55	86.96	66.67
XGBoost	90.06	79.17	82.17	83.67	100.00	66.67	85.11	76.92

XGBoost stood out with the highest overall F1-score (76.73%) and precision (86.96%), followed closely by Logistic Regression (F1 = 73.50%, precision = 80.26%). The findings of this experiment corroborate the initial suppositions that a dataset solely with microRNA (miRNA) values can serve as a resource to distinguish breast cancer subtypes, achieving up to 87% precision. In contrast, the Decision Tree showed lower performance across all metrics, suggesting limited generalization.

When analyzing subtype-specific results reveal perfect classification (F1 = 100%) for the Basal-like subtype, and this results is sustained by multiple models. However, performance drops for less distinct subtypes such as HER2-enriched and Luminal B. For instance, DIABLO's F1-score on HER2 was only 40.00%, where XGBoost doubled that score with 80.00%. Logistic Regression also performed well in these conditions across all classes, achieving the highest F1-score for Luminal A (85.71%) and Luminal B (52.63%).

These results indicates that discriminative models trained solely on miRNAs can produce clinically relevant predictions, reinforcing the value of miRNA expression patterns alone.

Impact of Adding Clinical Data In the second phase, we incorporated patient clinical features along with miRNA features to evaluate whether this multimodal integration could enhance model performance. Table 2 shows a consistent performance boost across all models.

¹ Available in the `mixOmics` R package.

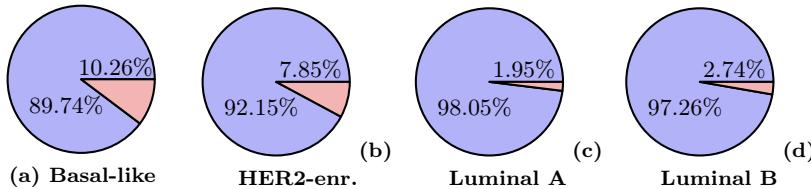


Fig. 3. XGBoost-Based Feature Group Importance per Breast Cancer Subtype.

When we analyzed the performance of the models, XGBoost showed the most substantial improvement, with the F1-score increasing from 76.73% to 82.17% and accuracy from 75.51% to 83.67%. Logistic Regression also benefited from the aggregation of clinical data (F1-score from 73.50% to 75.57% and accuracy from 77.55% to 79.59%), but being dethroned by XGBoost across all general metrics. These improvements suggest that clinical features provided complementary information that helped refine the subtype decision boundaries, particularly for overlapping classes.

At the subtype level, overall, the models exhibited enhanced classification of previously ambiguous groups. XGBoost's F1 for Luminal B rose from 50.00% to 76.92%, and Logistic Regression increased from 52.63% to 66.67%. This came with the cost of reducing XGBoost's F1 on HER2-enriched (which can be attributed to class imbalances or sensitivity to the integration process), but maintaining the top performance versus other models. This cost is seen as worth it, especially when we consider the rise of the other metrics for the other classes and the improvement in the model's classification ability.

To better understand the contribution of the clinical data, we analyzed the relative feature importance of each dataset as a whole using XGBoost model (the model that performed the best in this heterogeneous environment). As shown in Figure 3, miRNAs accounted for the vast majority of total feature importance across all subtypes (often above 90%). Yet, the small contribution from clinical features still had a measurable impact on the model performance, especially for the more ambiguous subtypes (Luminal B and HER2-enriched). This underscores the high informational density of clinical data, which helps resolve cases where miRNA profiles alone are ambiguous.

Discriminative vs. Latent Representation Approaches A direct comparison between DIABLO and discriminative models underscores the advantages of our approach. While DIABLO (All Blocks) improved its performance over its single-block version (F1: 73.82% vs. 67.69%), it still fell behind both XGBoost (F1: 82.17%) and Logistic Regression (F1: 75.57%) under the same multimodal setup.

Confusion matrices (Figure 7) offer additional insights into the strengths of the discriminative models. Both XGBoost and Logistic Regression made fewer misclassifications in difficult subtypes such as HER2-enriched and Luminal B. DIABLO, on the other hand, often confused HER2-enriched with Luminal A or

Luminal B, suggesting that its latent space did not adequately disentangle these phenotypically close subtypes.

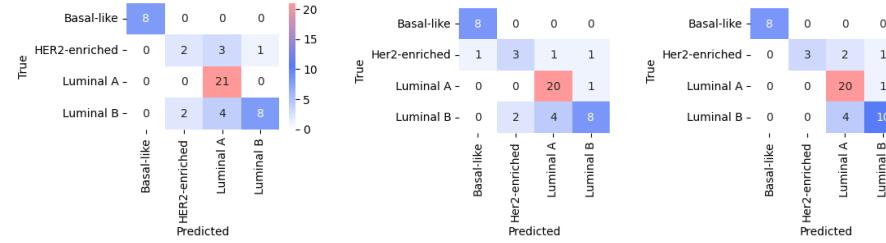


Fig. 4. DIABLO

Fig. 5. Log. Regression

Fig. 6. XGBoost

Fig. 7. Confusion matrices for each model using miRNA + clinical data.

In addition to that, the interpretability of discriminative models such as XGBoost further supports their applicability in clinical settings. Their coefficients can be directly linked to miRNA or clinical features, enabling biomarker identification and validation (more on that in Section 4.4. In contrast, DIABLO's latent components are harder to interpret and require additional steps to trace back to the original features.

In summary, discriminative models not only achieved higher predictive performance but also offered clearer insights into feature importance and subtype boundaries. These attributes are essential in biomedical contexts where transparency and traceability are paramount.

4.4 Identifying microRNA Biomarkers Candidates

To derive the final set of candidate miRNA biomarkers, we applied an intersecting strategy across multiple machine learning models and explainability methods (shown in Figure 3). We computed feature importance using both model-specific techniques (e.g., coefficients) and model-agnostic methods like SHAP. For each breast cancer subtype, we ranked features according to their contribution and retained the top features per method. The final set was obtained by selecting miRNAs that consistently appeared among the top-ranked features across different models and explanation techniques. This cross-validation across interpretability sources ensured that only robust, repeatable signals were considered, increasing confidence in their potential as true subtype-discriminative biomarkers. We take it that the top three miRNAs distinguish themselves highly from the rest. These miRNAs often, in the individual analyses, scored high for different models and methods of importance extraction. The MIR-934/934 ranked first half of the subtypes, it alone has great distinction capabilities.

One of the XAI (eXplainable AI) techniques used in our analysis of understanding which features were most influential is SHAP (SHapley Additive

Table 3. Top 10 most important miRNA features, and corresponding score ($\times 10^{-1}$ scale), per breast cancer subtype, using XGBoost. Common top miRNAs across subtypes are highlighted in bold. The prefix "MIR" is omitted for presentation purposes.

Basal-like		HER2-enriched		Luminal A		Luminal B	
Feature	Score	Feature	Score	Feature	Score	Feature	Score
135B/135B*	1.487	192/192*	0.814	105-1/105	0.456	135B/135B*	0.231
934/934	1.026	187/187	0.523	205/205*	0.423	934/934	0.188
190B/190B	0.974	34C/5P	0.436	25/25	0.342	125A/3P	0.173
455/5P	0.615	30C-2/2*	0.436	19A/19A	0.326	205/205*	0.159
577/577	0.564	148/148A	0.407	1301/1301	0.293	320B-2/320B	0.144
455/3P	0.410	497/497*	0.378	942/942	0.261	339/3P	0.144
18A/18A*	0.359	30C/3P	0.349	3200/3P	0.244	187/187	0.130
99A/99A*	0.205	2115/2115*	0.262	342/5P	0.195	148/148A	0.115
140/3P	0.205	3150B/3P	0.233	1270/1270	0.147	744/744*	0.115
9-1/9*	0.154	210/210	0.203	26A/26A	0.130	125B-2/2*	0.115

exPlanations). The spread and direction of SHAP values offer insight into how consistently each feature contributes across samples. Narrow distributions indicated potential subtype patterns, while large spreads, though still insightful, not as it's own biomarker alone. Overall, the SHAP analysis improved model transparency and provided biological understanding of miRNAs' role in separation supporting its promise as potential biomarkers.

The plots (Fig. 8) display the results for the SHAP analysis on the XGBoost model. It shows the top 10 features influencing predictions for each breast cancer subtype. miRNA expression features were consistently among the most impactful, suggesting their strong discriminative power. Clinical variables such as HER2 Status and ER status also appeared with high importance in relevant subtypes, aligning with established clinical markers. These plots unveil what impact each miRNA has towards the individual predictions per breast cancer subtype. For example, the MIR-190B/190B, when low expressed indicates an increase in the chance of a classification of the Basal-like subtype. Another relevant distinction is how the ER status feature strongly influences the classification of the Luminal B subtype, where a positive label is associated with the subtype.

5 Conclusion

In this paper, we introduce a multi-source breast cancer biomarker identification framework that integrates individuals' microRNA expression profiles with clinical attributes to enhance subtype classification. We provide compelling evidence that microRNA expression profiles alone offer strong discriminative power for accurately classifying breast cancer subtypes. Furthermore, our findings reveal that integrating clinical features further boosts predictive performance, highlighting the value of multi-modal data in subtype assessment. Through a model

Unveiling microRNA biomarkers for Breast Cancer Sub-typing

11

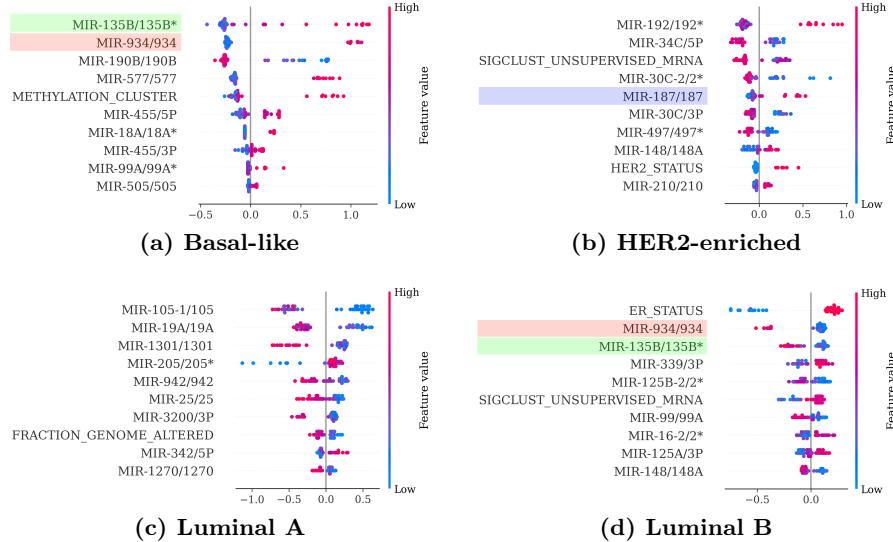


Fig. 8. SHAP summary plots for the top 10 features for each BC subtype (XGBoost model).

explainability approach, we isolate a focused set of subtype-specific potential microRNA biomarkers. The next step consists of conducting in-vitro and in-vivo laboratory validation, towards further validating their potential to develop targeted therapeutics.

References

1. Azari, H., Nazari, E., Mohit, R., Asadnia, A., Maftooh, M., Nassiri, M., Hassanian, S.M., Ghayour-Mobarhan, M., Shahidsales, S., Khazaei, M., Ferns, G.A., Avan, A.: Machine learning algorithms reveal potential miRNAs biomarkers in gastric cancer. *Sci. Rep.* **13**(1), 6147 (Apr 2023)
2. Blagus, R., Lusa, L.: SMOTE for high-dimensional class-imbalanced data. *BMC Bioinformatics* **14**(1), 106 (Mar 2013)
3. Bravo-Vázquez, L.A., Méndez-García, A., Rodríguez, A.L., Sahare, P., Pathak, S., Banerjee, A., Duttaroy, A.K., Paul, S.: Applications of nanotechnologies for miRNA-based cancer therapeutics: current advances and future perspectives. *Front. Bioeng. Biotechnol.* **11**, 1208547 (Jul 2023)
4. Gilani, N., Arabi Belaghi, R., Aftabi, Y., Faramarzi, E., Edgünlü, T., Somi, M.H.: Identifying potential miRNA biomarkers for gastric cancer diagnosis using machine learning variable selection approach. *Front. Genet.* **12**, 779455 (2021)
5. Hrinivich, W.T., Wang, T., Wang, C.: Editorial: Interpretable and explainable machine learning models in oncology. *Front. Oncol.* **13**, 1184428 (Mar 2023)
6. Hu, Z., Lai, C., Liu, H., Man, J., Chen, K., Ouyang, Q., Zhou, Y.: Identification and validation of screening models for breast cancer with 3 serum miRNAs in an 11,349 samples mixed cohort. *Breast Cancer* **31**(6), 1046–1058 (Jul 2024)

12 F. Author et al.

7. Li, J., Zhang, H., Gao, F.: Identification of miRNA biomarkers for breast cancer by combining ensemble regularized multinomial logistic regression and cox regression. *BMC Bioinformatics* **23**(1), 434 (Oct 2022)
8. Lundberg, S.M., Lee, S.I.: A unified approach to interpreting model predictions. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) *Advances in Neural Information Processing Systems*. vol. 30. Curran Associates, Inc. (2017)
9. Mendes, B.B., Connio, J., Avital, A., Yao, D., Jiang, X., Zhou, X., Sharf-Pauker, N., Xiao, Y., Adir, O., Liang, H., Shi, J., Schroeder, A., Conde, J.: Nanodelivery of nucleic acids. *Nat. Rev. Methods Primers* **2**(1) (Apr 2022)
10. Mendes, B.B., Zhang, Z., Connio, J., Sousa, D.P., Ravasco, J.M.J.M., Onweller, L.A., Lorenc, A., Rodrigues, T., Reker, D., Conde, J.: A large-scale machine learning analysis of inorganic nanoparticles in preclinical cancer research. *Nat Nanotechnol* **19**(6), 867–878 (May 2024)
11. Muñoz, J.P., Pérez-Moreno, P., Pérez, Y., Calaf, G.M.: The role of MicroRNAs in breast cancer and the challenges of their clinical application. *Diagnostics (Basel)* **13**(19) (Sep 2023)
12. Pagoni, M., Cava, C., Sideris, D.C., Avgeris, M., Zoumpourlis, V., Michalopoulos, I., Drakoulis, N.: MiRNA-based technologies in cancer therapy. *J. Pers. Med.* **13**(11), 1586 (Nov 2023)
13. Ribeiro, M., Singh, S., Guestrin, C.: “why should I trust you?”: Explaining the predictions of any classifier. In: DeNero, J., Finlayson, M., Reddy, S. (eds.) *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations*. pp. 97–101. Association for Computational Linguistics, San Diego, California (Jun 2016)
14. Singh, A., Shannon, C.P., Gautier, B., Rohart, F., Vacher, M., Tebbutt, S.J., Lê Cao, K.A.: Diablo: an integrative approach for identifying key molecular drivers from multi-omics assays. *Bioinformatics* **35**(17), 3055–3062 (01 2019)
15. Valihrach, L., Androvic, P., Kubista, M.: Circulating mirna analysis for cancer diagnostics and therapy. *Molecular aspects of medicine* (2020)
16. Wang, W.T., Chen, Y.Q.: Circulating miRNAs in cancer: from detection to therapy. *J. Hematol. Oncol.* **7**(1), 86 (Dec 2014)
17. Wang, Z., Liao, H., Deng, Z., Yang, P., Du, N., Zhanng, Y., Ren, H.: miRNA-205 affects infiltration and metastasis of breast cancer. *Biochem. Biophys. Res. Commun.* **441**(1), 139–143 (Nov 2013)

