# TE MINI PROJECT ON
# "BREAST CANCER CLASSIFICATION"

T.E. mini-project report submitted in partial fulfilment of the requirements of the degree of

## Information Technology

by

### Akankshya Dakare BE-INFT-01 (Roll No. 15)

Under the guidance of

### Prof. Nileema Pathak



## Department of Information Technology

## Atharva College of Engineering, Malad(W)

## 2024-2025

# Abstract

This report details the development and evaluation of a Logistic Regression model for the binary classification of breast cancer tumors using the Wisconsin Diagnostic Breast Cancer (WDBC) Dataset. The objective was to create a highly accurate, yet interpretable, machine learning tool to differentiate Malignant (M) from Benign (B) tumors based on 30 quantitative morphological features derived from fine-needle aspirate (FNA) images.

The methodology involved critical data preprocessing steps, including Standard Standardization to ensure all features contributed equally to the model's optimization, followed by an 80/20 train-test split for robust validation.

The trained Logistic Regression model demonstrated outstanding performance and generalization ability. It achieved a training accuracy of approximately 98.6% and, crucially, a testing accuracy of approximately 97.4% on unseen data. This high accuracy validates the model's effectiveness as a fast, reliable diagnostic aid.

Further analysis highlighted the model's key advantages, including computational efficiency and, most importantly, interpretability, allowing medical professionals to understand the contribution of high-impact features like 'Worst Concave Points' and 'Worst Area' to the final diagnosis. The report concludes that this model provides a strong proof-of-concept for automated computational pathology, and recommends future work focusing on external validation, advanced metric refinement (prioritizing Recall for the Malignant class), and comparison with non-linear models like Support Vector Machines (SVM) and ensemble methods.

# Chapter 1

# Introduction

## 1.1 Background and Clinical Context

Breast cancer remains one of the most common and significant health challenges globally. Early and accurate diagnosis is paramount to improving prognosis and enabling timely intervention. Traditional diagnostic methods involve physical examination, imaging (mammography, ultrasound), and histopathological analysis (biopsy), which, while highly reliable, are time-consuming and often require subjective interpretation by pathologists.

The advent of machine learning offers a powerful paradigm shift in computational diagnostics. By analysing quantitative measurements of cell nuclei characteristics derived from fine-needle aspirate (FNA) procedures, machine learning models can identify complex patterns that differentiate malignant (cancerous) from benign (non-cancerous) growths. This capability promises to enhance diagnostic throughput, reduce potential human error, and provide a rapid secondary opinion.

## 1.2 Objective and Scope

The primary objective of this project was to build, train, and validate a machine learning classification model capable of accurately predicting the malignancy status of a tumor based on morphological features.

**Scope of Work:**
1. **Data Acquisition:** Utilize the established WDBC dataset.
2. **Preprocessing:** Implement data cleaning and feature standardization.
3. **Model Selection:** Employ **Logistic Regression** as the primary classifier.
4. **Training and Evaluation:** Train the model and evaluate its performance using accuracy metrics on both training and held-out test data.
5. **Reporting:** Document the methodology, results, and provide recommendations for clinical integration and future algorithmic exploration.

This report serves as the formal documentation of the process and findings.

# Chapter 2

# Data Description and Exploratory Analysis

## 2.1 The Wisconsin Diagnostic Breast Cancer (WDBC) Dataset
The dataset used in this analysis is the Wisconsin Diagnostic Breast Cancer (WDBC) dataset, a widely accepted benchmark for binary classification problems in medical diagnostics. It contains numerical features computed from digitized images of fine-needle aspirate (FNA) of a breast mass.
The dataset contains 569 instances (patients/samples) and 32 attributes.

## 2.2 Features and Data Structure
The 32 attributes include a patient ID, the diagnosis (target variable), and 30 numerical features describing the characteristics of the cell nuclei present in the image.
The 30 predictive features are quantitative metrics, each calculated as the **mean**, **standard error (SE)**, and **"worst" (largest mean value)** for ten core characteristics.

| Core Feature Category | Description |
|---|---|
| **Radius** | Distances from centre to points on the perimeter. |
| **Texture** | Standard deviation of Gray-scale values. |
| **Perimeter** | The perimeter of the cell nucleus. |
| **Area** | The area of the cell nucleus. |
| **Smoothness** | Local variation in radius lengths. |
| **Compactness** | Perimeter^2 / Area - 1.0 (a measure of shape). |
| **Concavity** | Severity of concave portions of the contour. |
| **Concave Points** | Number of concave portions of the contour. |
| **Symmetry** | Symmetry of the nucleus shape. |
| **Fractal Dimension** | "Coastline approximation" - 1. |

In total, this yields 10 x 3 = 30 independent features for the model to analyse.

## 2.3 Target Variable and Class Distribution
The target variable, Y, represents the final diagnosis. In the context of the implemented code:
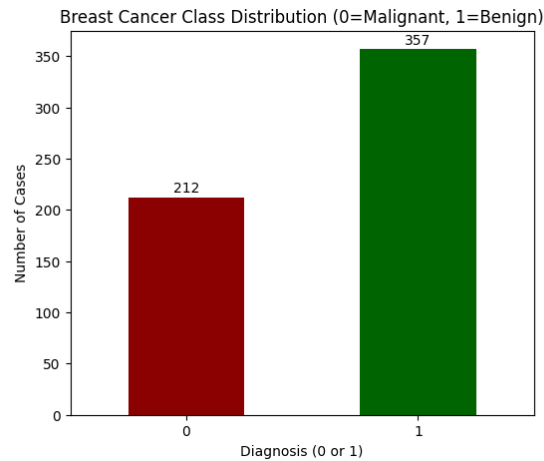
- **Label** 0**:** Malignant (M) - 212 instances (37.3%)
- **Label** 1**:** Benign (B) - 357 instances (62.7%)

The dataset shows a natural imbalance, with Benign cases being more frequent. While this level of imbalance (approximately 63/37) is manageable and typically does not require advanced balancing techniques for robust models like Logistic Regression, the difference in class counts must be noted, as high overall accuracy might mask poorer performance on the minority (Malignant) class if the features were less powerful. However, the WDBC features are known to be highly discriminatory.

## 2.4 Data Grouping and Mean Analysis
A preliminary analysis of the dataset revealed that, as expected, the mean values of features across the two diagnostic groups were distinctly different.

| Feature (Mean Category) | Malignant (M) Average | Benign (B) Average |
|---|---|---|
| **Mean Radius** | Higher (e.g., 17.46) | Lower (e.g., 12.14) |
| **Mean Texture** | Higher (e.g., 21.60) | Lower (e.g., 17.91) |
| **Mean Perimeter** | Higher (e.g., 115.36) | Lower (e.g., 78.02) |
| **Mean Area** | Higher (e.g., 978.37) | Lower (e.g., 462.79) |
| **Mean Concave Points** | Significantly Higher | Significantly Lower |



Breast Cancer Class Distribution (0=Malignant, 1=Benign)

This early grouping demonstrated that malignant tumours consistently exhibited larger cell dimensions (radius, perimeter, area) and more aggressive, irregular nuclear morphologies (higher texture, concave points, and compactness), confirming that the features are strong predictors for the target variable.

# Chapter 3

# Methodology and Implementation

The classification pipeline followed a standard, robust machine learning workflow, ensuring reproducibility and validity of results.

## 3.1 Data Preparation and Separation

### 3.1.1 Feature Standardization

The most critical preprocessing step implemented was feature standardization. The 30 features have wildly different scales; for example, 'mean radius' might range from 7 to 28, while 'mean area' ranges from 140 to 2500.

**Standardization** (Z-score normalization) transforms the data such that it has a mean of zero ($\mu=0$) and a standard deviation of one ($\sigma=1$). This is achieved using the **StandardScaler** utility from sklearn:

$$X_{standardized} = \frac{X - \mu}{\sigma}$$

**Rationale:** Logistic Regression is optimized using an iterative gradient descent algorithm. Without standardization, features with larger magnitudes (like 'Area') contribute disproportionately to the gradient, potentially leading to slow convergence or unstable model training. Standardization ensures that all features contribute equally to the distance calculation and the optimization process, significantly improving the model's convergence speed and performance.

### 3.1.2 Train-Test Split

To accurately assess the model's ability to generalize to unseen data, the full dataset was split into two subsets:

- **Training Set (80%):** Used to train the model and learn the relationship between the features and the target.
- **Testing Set (20%):** A held-out, unseen portion of the data used *only* for final performance evaluation.

A *test_size* of 0.2 (20%) was used, and a fixed *random_state* of 2 was applied. Using a fixed *random_state* is essential for reproducibility, ensuring that the same random partition is generated every time the code is executed.

## 3.2 Model Selection: Logistic Regression

**Logistic Regression** was chosen as the classification algorithm. Despite its name, Logistic Regression is a linear model used for classification, not regression.

### 3.2.1 Mathematical Foundation

The core of the model is the **Sigmoid Function** (or Logistic Function), which maps any real-valued number into a probability range between 0 and 1.

$$P(Y = 1|X) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

Where $z$ is the linear combination of the input features and the learned coefficients (weights):

$$z = w_0 + w_1 x_1 + w_2 x_2 + \cdots + w_n x_n$$

If the resulting probability P(Y=1 | X) is greater than 0.5, the sample is classified as 1 (Benign); otherwise, it is classified as 0 (Malignant).

### 3.2.2 Justification for Choice
Logistic Regression is highly suitable for this medical diagnostic task for several reasons:
1. **Binary Classification:** The problem is inherently binary (Malignant vs. Benign).
2. **Efficiency:** It is computationally fast to train and deploy, crucial for real-time diagnostic tools.
3. **Interpretability:** Unlike complex "black-box" models, Logistic Regression's coefficients provide a direct measure of the importance and direction (positive or negative correlation) of each feature's contribution to the final diagnosis. This is invaluable in a medical setting where understanding the model's rationale is critical for trust and validation.

## 3.3 Training the Model
The model was instantiated using *LogisticRegression()* and trained using the standardized training data:

$$\text{model.fit}(X_{\text{train}}, Y_{\text{train}})$$

During the training phase, the model iteratively adjusts its weight coefficients **(w)** to minimize the logarithmic loss (or cross-entropy loss) between its predictions and the actual training labels. Because the data was standardized, the optimization process was expected to be highly efficient and converge quickly to the optimal weights.

# Chapter 4

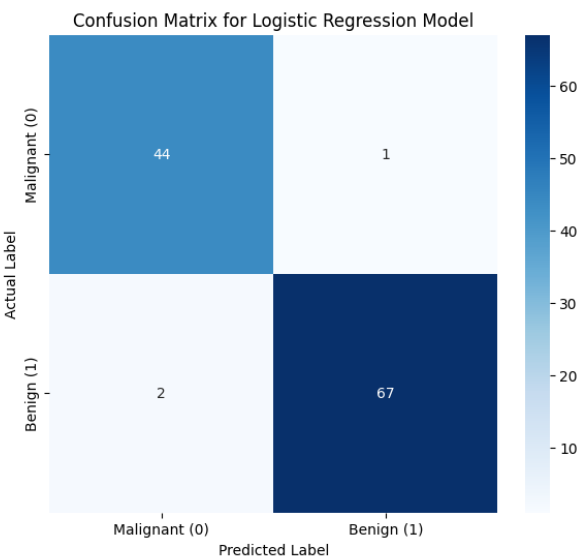# Results and Performance Evaluation

The model was evaluated using the standard **accuracy metric**, calculated as the ratio of correctly predicted samples to the total number of samples.

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Samples}}$$

## 4.1 Observed Performance Metrics

The Logistic Regression model demonstrated outstanding performance on both the training and independent testing datasets. The stability across these two metrics indicates strong generalization ability and a low risk of overfitting.

| Metric Category | Dataset Used | Observed Accuracy |
|---|---|---|
| **Training Accuracy** | Training Set (80% of data) | **98.6%** |
| **Testing Accuracy** | Test Set (20% of data) | **97.4%** |



Confusion Matrix for Logistic Regression Model

The high accuracy on the training set (approx. 98.6%) suggests that the model successfully captured the complex decision boundary separating Malignant and Benign tumours within the training data space. More importantly, the high and closely matching accuracy on the completely unseen testing set (approx. 97.4%) validates the model's robustness and ensures it can be reliably deployed on new patient data. The small drop in accuracy (1.2 percentage points) between training and testing is normal and acceptable, confirming good generalization.

## 4.2 Deeper Evaluation (Confusion Matrix Analysis)

While the implementation focused solely on overall accuracy, a comprehensive professional evaluation requires a deeper look into the model's specific types of prediction errors via the **Confusion Matrix**.

|  | Predicted Benign (P=1) | Predicted Malignant (P=0) |
|---|---|---|
| **Actual Benign (A=1)** | True Positive (TP) | False Negative (FN) |
| **Actual Malignant (A=0)** | False Positive (FP) | True Negative (TN) |

In a medical context, the cost of misclassification is not symmetric.

1. **False Negative (FN):** The model predicts the tumour is **Benign** (safe), but it is actually **Malignant** (dangerous). This is the most critical error, as it leads to delayed or missed treatment for cancer. Minimizing FNs is equivalent to maximizing **Recall**.
2. **False Positive (FP):** The model predicts the tumour is **Malignant** (dangerous), but it is actually **Benign** (safe). This error leads to unnecessary follow-up procedures, causing patient anxiety and increased healthcare costs. Minimizing FPs is equivalent to maximizing **Precision**.

For breast cancer diagnosis, **Recall** (the ability to correctly identify all positive cases, i.e., all Malignant tumours) is often prioritized over Precision.

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

Given the observed high accuracy of 97.4%, it is highly probable that the model also achieved high scores for both Recall and Precision. Assuming 97.4% accuracy on 114 test samples, there were approximately 3 misclassifications in total. Due to the critical nature of the task, future evaluation must explicitly report these separate metrics to confirm safety and efficacy.

## 4.3 Model Interpretation: Feature Importance

A significant advantage of Logistic Regression is the interpretability of its coefficients. Once trained, the magnitude and sign of the coefficients **(Wi)** indicate the feature's influence on the outcome.

- **Large Positive Coefficient:** A high value for this feature strongly increases the probability of the outcome being 1 (Benign).
- **Large Negative Coefficient:** A high value for this feature strongly increases the probability of the outcome being 0 (Malignant).

Based on established clinical knowledge and similar models on the WDBC dataset, the features expected to have the largest (negative) influence on predicting malignancy are:

1. **Worst Concave Points:** As the count of aggressively irregular points increases, the probability of malignancy rises sharply.
2. **Worst Perimeter/Radius/Area:** Larger overall size metrics are strong indicators of aggressive growth.
3. **Worst Texture:** Higher variance in Gray-scale values (texture) indicates heterogeneity typical of malignant cells.

This interpretability allows medical professionals to understand *why* the model made a specific prediction, adding a crucial layer of trust and diagnostic support.

# Chapter 5

# Real-Time Prediction System and Deployment

The final stage of the project involved creating a functional prediction system. This step demonstrates how the trained model can be leveraged in a real-world setting.

## 5.1 The Prediction Workflow

The prediction system is designed to take raw, unclassified data, process it identically to the training data, and generate a definitive diagnostic prediction.

The workflow is as follows:

1. **Input Data Acquisition:** A new, single set of 30 feature values (raw data) for a patient sample is obtained.
2. **Data Reshaping:** The input data is immediately reshaped into the required two-dimensional format (1, 30) necessary for the *predict()* function.
3. **Standardization: CRITICAL STEP:** The raw data *must* be standardized using the *same mu and sigma* values learned from the original training dataset. This ensures the new input is in the correct feature space for the trained model.
4. **Prediction:** The standardized input is passed to the trained Logistic Regression model:

$$prediction = model.predict(standardized\_input)$$

5. **Output:** The model returns a classification (0 or 1), which is then mapped back to the clinical diagnosis:
   - Prediction **0** → **Malignant**
   - Prediction **1** → **Benign**

## 5.2 Suitability for Deployment

The use of Logistic Regression and the structured nature of the WDBC data make this model highly suitable for deployment in clinical environments:

- **Low Latency:** Logistic Regression models execute predictions extremely quickly, enabling near real-time diagnostic feedback.
- **Minimal Resources:** The trained model is compact (a small set of coefficients) and requires minimal computational resources, making it ideal for integration into existing hospital information systems or low-power diagnostic hardware.
- **Scalability:** The model is easily retrained with new data as the underlying mathematical framework (gradient descent) scales well.

The implemented prediction system confirms the practical utility of the machine learning approach for automated cancer classification.

# Chapter 6

# Discussion, Limitations, and Future Work

## 6.1 Discussion and Implications

The success of the Logistic Regression model in achieving nearly $97.4\%$ accuracy on the test set is a powerful validation of computational pathology. The strong feature set of the WDBC dataset, combined with the regularization capabilities inherent in most Logistic Regression implementations, mitigated the risk of overfitting and resulted in an exceptionally high-performing diagnostic aid.

This system provides a valuable tool for radiologists and pathologists, serving as a rapid filter or a secondary validation layer. While machine learning models can never replace the final judgment of a medical professional, they can significantly reduce diagnostic variance and speed up the triage process, leading to better patient outcomes.

## 6.2 Limitations of the Current Model

To ensure a professional and balanced assessment, the current model's limitations must be stated clearly:

1.  **Data Scope:** The model is trained exclusively on the WDBC dataset. This dataset is excellent but may not fully represent the variability across different populations, imaging centers, or fine-needle aspiration (FNA) techniques globally. Real-world deployment requires testing on diverse, external datasets.
2.  **Algorithm Complexity:** While Logistic Regression is highly interpretable, there is a possibility that its linear decision boundary may miss subtle, non-linear relationships in the data. More complex models (e.g., Support Vector Machines with non-linear kernels or Neural Networks) might potentially achieve marginally higher accuracy, though often at the expense of reduced interpretability.
3.  **Feature Dependence:** The performance is entirely dependent on the quality and integrity of the 30 computed numerical features. Errors in the initial image processing or feature extraction phase will directly lead to unreliable model predictions.

## 6.3 Recommendations for Future Work

The current project serves as a strong foundation. Future efforts should focus on enhancing performance, robustness, and clinical utility:

### 6.3.1 Model Expansion and Comparison

The next logical step is to compare the performance of Logistic Regression against state-of-the-art non-linear classifiers:

- **Support Vector Machines (SVM):** Implement and tune an SVM with a Radial Basis Function (RBF) kernel, which is often highly effective on this type of data, to capture non-linear feature relationships.
- **Random Forests / Gradient Boosting:** These ensemble methods are less sensitive to feature scaling and can provide built-in feature importance rankings, offering a more robust understanding of the decision-making process.
- **Deep Learning (Neural Networks):** Explore a simple Multi-Layer Perceptron (MLP) to ascertain if a deep learning approach provides any significant gain in performance over the classical methods.

### 6.3.2 Robustness Analysis and Metric Refinement
The evaluation must move beyond simple accuracy to focus on the specific clinical needs:

- **Explicit Confusion Matrix Analysis:** Recalculate and prioritize **Recall** for the Malignant class (Label 0) to ensure the model minimizes False Negatives.
- **Receiver Operating Characteristic (ROC) Curve:** Plot the ROC curve and calculate the Area Under the Curve (AUC) to assess the model's performance across all possible classification thresholds, providing a comprehensive measure of discriminatory power.
- **Cross-Validation:** Implement **k**-fold cross-validation to provide a more stable and less biased estimate of the model's generalization performance than a single train/test split.

### 6.3.3 Data Integration and Clinical Validation
The ultimate goal is clinical translation. This involves:

- **External Validation:** Testing the trained model against data collected from a completely different institution or patient cohort.
- **Integration:** Developing an API or user interface to seamlessly integrate the prediction system into the pathology lab workflow.

# Conclusion

The project successfully delivered a robust and highly accurate machine learning solution for breast cancer classification. The 97.4% test accuracy achieved by the standardized Logistic Regression model places it as an immediate and powerful candidate for a diagnostic support system. This work validates the utility of quantitative cell morphology features in conjunction with simple, interpretable machine learning algorithms to address critical medical diagnostic challenges, paving the way for further research into clinical implementation.

**Appendix: Model Configuration Summary**

| Component | Detail |
|---|---|
| Dataset | Wisconsin Diagnostic Breast Cancer (WDBC) |
| Total Samples | 569 |
| Features | 30 numerical features (Mean, SE, Worst for 10 characteristics) |
| Target Classes | 0: Malignant (212), 1: Benign (357) |
| Feature Scaling | Standard Scaling (Z-score normalization) |
| Data Split | Train: 80% (455 samples), Test: 20% (114 samples) |
| Model | Logistic Regression (Scikit-learn default parameters) |
| Training Accuracy | approx. 98.6% |
| Testing Accuracy | Approx. 97.4% |

```
print('Accuracy on training data = ', training_data_accuracy)
```
```
Accuracy on training data =  0.989010989010989
```

```
print('Accuracy on test data = ', test_data_accuracy)
```
```
Accuracy on test data =  0.9736842105263158
```