

# **Project 1**

## **Breast Cancer Diagnosis** Using **Machine Learning**

### **Project Report**

SUBMITTED IN PARTIAL FULFILLMENT REQUIREMENT FOR THE AWARD OF DEGREE OF  
BACHELOR OF TECHNOLOGY

SUBMITTED BY:

**Group 4**

**Priyanshu Garg**

**Himanshu**

**Abhinav Yadav**

**Nirbhay Gurjar**

UNDER THE SUPERVISION OF

**Dr. Hirdesh Pharasi**

SCHOOL OF ENGINEERING AND TECHNOLOGY

BML MUNJAL UNIVERSITY Gurugram, Haryana - 122413



October, 2023

## DECLARATION

We, Priyanshu Garg, Himanshu, Abhinav Yadav, and Nirbhay Gurjar, hereby declare that the project "Breast Cancer Diagnosis Using Machine Learning" completed as part of the Bachelor of Technology (B. Tech) programme at the School of Engineering and Technology, BML Munjal University is an authentic record of our work carried out under the supervision of **Dr. Hirdesh Kumar Pharasi**. We would like to convey our heartfelt appreciation to Mr. Hirdesh Pharasi, our instructor, for his invaluable guidance, criticism, and assistance during this endeavour. Your knowledge and experience were invaluable in formulating and completing this report. We would like to thank everyone who supported and contributed to the creation of this technical report. Everyone's help and contributions are much appreciated. This report would not have been feasible without your support.

All additional items utilised have been properly acknowledged in the project text.  
This project was completed in accordance with the curriculum's requirements and limits.

Priyanshu Garg 220649

Abhinav Yadav      220619

Himanshu 220593

Nirbhay Gurjar 220412

Place: BML MUNJAL UNIVERSITY, KAPRIWAS, HARYANA

Date: 17-12-2023

## SUPERVISOR'S DECLARATION

This is to certify that the above statement made by the candidate is correct to the best of my knowledge.

Faculty Supervisor Name: Dr. Hirdesh Kumar Pharasi

Signature:

# ABSTRACT

This research delves into a comprehensive examination of the Breast Cancer Wisconsin (Diagnostic) Data Set, a crucial dataset that intricately captures ten essential features computed from the characteristics of cell nuclei in breast mass images. Encompassing a comprehensive array of mean, standard error, and worst-case values, this dataset not only provides a detailed understanding of the underlying features but also lays the foundation for a nuanced exploration of the distinctions between malignant and benign cases. The well-balanced class distribution, comprising 357 benign and 212 malignant cases, not only ensures the integrity of the dataset but also sets the stage for the development of robust machine learning models, promising heightened accuracy in breast cancer diagnosis.

At the core of this research is a pivotal objective: to advance breast cancer diagnosis by swiftly detecting cancer in human breast masses through the integration of Fine Needle Aspirate (FNA) alongside the observation of biological reports, followed by rigorous testing using machine learning techniques. Leveraging a dataset rich in cell nuclei characteristics, the primary aim is to design and train a model capable of discerning cancerous features within breast masses. The overarching goal extends beyond mere detection; it aspires to establish a more accurate and balanced methodology for cancer detection, recognizing the profound implications this has for patient outcomes and treatment decisions.

This research stands at the intersection of medical diagnostics, bioinformatics, and machine learning, with the ambition to contribute significantly to the field of breast cancer research. By harnessing the power of FNA, incorporating the nuanced insights from biological reports, and employing the sophistication of machine learning algorithms, the research endeavours to create a comprehensive diagnostic tool. Such a tool, when integrated into clinical practices, holds the potential to revolutionize the landscape of breast cancer diagnosis. This transformative approach not only enhances the speed and precision of cancer identification but also addresses the critical need for more accurate and reliable diagnostic methods, thereby improving patient care and prognosis. The synergistic fusion of biomedical sciences and cutting-edge technology in this research underscores its potential impact on advancing breast cancer diagnostics.

# ACKNOWLEDGEMENT

I extend my deepest gratitude to Dr. Hirdesh Pharasi, Assistant Professor at the School of Engineering & Technology, BML Munjal University, Gurugram, for his invaluable guidance and supervision throughout the duration of my seminar/case study from July to December 2022. Dr. Pharasi's expertise, unwavering support, and wise counsel have been instrumental in shaping the trajectory of my work, and I express my sincere thanks with utmost reverence.

Undoubtedly, Dr. Hirdesh Pharasi's mentorship has played a pivotal role in the successful completion of my training. His dedication to fostering a conducive learning environment and providing insightful guidance has been truly commendable. Without his exceptional support, achieving the level of accomplishment in this endeavor would have been a formidable challenge.

I take this opportunity to profusely thank Dr. Hirdesh Pharasi for consistently stimulating and encouraging me throughout the entire process. His mentorship has not only enhanced my academic understanding but has also instilled a sense of confidence and competence in approaching complex tasks.

In addition, my gratitude extends to the entire team at BML Munjal University for creating an environment conducive to learning and growth. The collaborative spirit and commitment to academic excellence within the institution have significantly contributed to the success of my project.

Furthermore, I express heartfelt thanks to my friends who generously devoted their time and provided invaluable assistance, contributing to the overall success of this endeavor. Their support has been a source of inspiration and motivation.

In conclusion, I am profoundly thankful to Dr. Hirdesh Pharasi, and everyone involved in this academic journey. This experience has not only enriched my knowledge but has also instilled a deep sense of appreciation for the collaborative and supportive academic community at BML Munjal University.

# TABLE OF CONTENTS

Contents No.	Page
<i>Declaration</i>	<i>ii</i>
<i>Abstract</i>	<i>iii</i>
<i>Acknowledgement</i>	<i>iv</i>
<b>1 Introduction to Organization</b>	<b>1</b>
<b>2 Introduction to Project</b>	<b>2</b>
2.1 Overview . . . . .	2
2.2 Existing System . . . . .	2
2.3 User Requirement Analysis . . . . .	3
2.4 Feasibility Study . . . . .	4
<b>3 Literature Review</b>	<b>5</b>
3.1 Objectives of Project (Must be clearly, precisely defined and Implementation must be done.) . . . . .	6
<b>4 Exploratory Data Analysis</b>	<b>7</b>
4.1 Dataset . . . . .	7
4.2 Exploratory Data Analysis and Visualizations . . . . .	8
<b>5 Methodology</b>	<b>9</b>
5.1 Introduction to Languages (Front End and Back End) . . . . .	9
5.2 Languages and Packages Used . . . . .	9
5.3 User characteristics . . . . .	9
5.4 Constraints . . . . .	10
5.5 ML algorithm discussion . . . . .	12
<b>6 Results and Discussion</b>	<b>13</b>
<b>7 Conclusion and Future Scope</b>	<b>22</b>
7.1 Conclusion . . . . .	22
7.2 Future Scope . . . . .	23

# Chapter 1

## **Introduction to Organisation**

The School of Engineering and Technology (SOET) at BML Munjal University stands as a dynamic hub, seamlessly blending academic excellence with a spirit of innovation, and envisioning a transformative educational journey for the engineers and technologists of tomorrow. SOET's comprehensive range of undergraduate and postgraduate programs reflects its commitment to providing students with a well-rounded education, preparing them to thrive in the ever-evolving landscape of engineering and technology.

Central to the triumphs of SOET is its exceptional faculty—an assembly of accomplished educators and industry experts who bring real-world insights into the classroom. The institution takes pride in its cutting-edge infrastructure, boasting state-of-the-art laboratories, modern classrooms, and advanced technology facilities. These resources collectively create an immersive and conducive environment for both learning and pioneering research.

SOET distinguishes itself by placing a robust emphasis on research and innovation, actively encouraging students to participate in projects that contribute to technological advancements. Collaborations with industry partners, internships, and exposure to global research initiatives ensure that students not only grasp theoretical concepts but also gain practical, hands-on experiences that are invaluable for their future careers.

Beyond the traditional realms of education, SOET is committed to holistic development. The institution integrates global exposure, ethical considerations, and a sense of social responsibility into its curriculum. Cultivating a culture of curiosity, critical thinking, and creativity, SOET moulds its students to become leaders and innovators in the diverse and dynamic fields of engineering and technology.

In essence, SOET at BML Munjal University transcends the boundaries of conventional education. Its dedication to excellence and a holistic approach position it as a distinguished institution, contributing significantly to the advancement of knowledge and the development of the next generation of engineering professionals. As a beacon of academic prowess and innovation, SOET is shaping the future of engineering education.

# Chapter 2

## Introduction to Project

### 2.1 Overview

In the project "Breast Cancer Diagnosis Using Machine Learning," the focus is on leveraging advanced computational techniques to enhance the accuracy and efficiency of breast cancer diagnosis. This project recognizes the significance of early detection in improving patient outcomes. The methodology involves the application of machine learning algorithms to analyze medical data, specifically targeting breast cancer diagnostic processes. By employing supervised learning techniques, the model learns patterns and features from a labeled dataset, enabling it to make accurate predictions about the presence or absence of breast cancer.

The project aims to contribute to the field of medical diagnostics by providing a reliable and automated system for breast cancer detection. The choice of machine learning, particularly supervised learning, enables the model to generalize from historical data, making it adaptable to diverse cases. The development of the model is carried out using Python, a versatile programming language known for its extensive libraries and tools for machine learning. Through this project, the team endeavors to enhance the efficiency of breast cancer diagnosis, ultimately contributing to improved healthcare outcomes.

### 2.2 Existing System

In the current landscape of breast cancer diagnosis, traditional methodologies stand as the cornerstone, predominantly relying on techniques such as mammography, biopsy, and clinical examinations. Mammography, a widely employed screening tool, serves to identify abnormalities in breast tissue, while biopsy procedures offer a definitive diagnosis through the examination of tissue samples. Alongside these, clinical examinations conducted by skilled medical professionals contribute crucial insights to the diagnostic process.

However, it is imperative to recognize the inherent limitations within this existing system. Mammography, despite its widespread use, is not without its challenges, often presenting false positives and false negatives, necessitating further, sometimes invasive, testing. Biopsy procedures, while providing conclusive results, may be impractical for routine screening due to their invasive nature. Moreover, the reliance on manual clinical examinations introduces subjectivity and variability into the diagnostic process.

In response to these limitations, there is a growing recognition of the potential of machine learning to revolutionize breast cancer diagnosis. This evolving landscape seeks to leverage advanced computational techniques to enhance accuracy, efficiency, and objectivity in the detection and classification of breast abnormalities. The "Breast Cancer Diagnosis Using Machine Learning" project endeavors to build upon these foundations, offering a paradigm shift towards a more sophisticated and data-driven approach to breast cancer diagnosis. This project aims to address the current system's shortcomings, introducing a novel methodology that holds the promise of improved diagnostic outcomes and patient care.

## 2.3 User Requirement Analysis

Ensuring the successful development and deployment of the "Breast Cancer Diagnosis Using Machine Learning" project requires a comprehensive understanding of user needs and expectations. In this User Requirement Analysis, we delve into the key aspects that shape the system's design and functionality.

**Data Accessibility** emerges as a crucial user need, demanding a system that effortlessly integrates diverse breast cancer datasets. To meet this, the platform must support seamless data import, pre-processing, and storage.

**Accuracy and Reliability** are paramount, with users seeking a high level of confidence in diagnostic predictions. Therefore, rigorous training and validation processes are imperative, ensuring the machine learning model achieves or exceeds industry standards for accuracy.

**Interpretability** becomes essential, especially in medical scenarios where trust is paramount. Users require a system that offers insights into the decision-making process, making the model's outcomes more understandable.

**Adaptability and Generalization** are user needs grounded in the dynamic nature of healthcare. The machine learning model must be flexible, continuously learning and adapting to diverse patient populations and emerging patterns in breast cancer data.

A **User-Friendly Interface** is a common demand, emphasizing the importance of an intuitive platform accessible to users with varying technical expertise. This includes clear navigation and easily understandable displays of diagnostic results.

**Integration with Existing Systems** is a practical user need, calling for seamless integration with healthcare information systems to facilitate a smooth workflow for medical practitioners.

Ensuring **Security and Privacy** is paramount. Users demand robust measures to protect sensitive patient data, and the system must adhere to industry standards for data security and comply with healthcare data protection regulations.

**Scalability** is a forward-looking user need, anticipating an increase in the volume of breast cancer diagnostic data over time. The system's architecture must be scalable to efficiently handle growing datasets without compromising performance.

Users' expectation for **Continuous Support and Maintenance** underscores the need for an ongoing commitment to the project's success. This involves regular updates, bug fixes, and responsiveness to user feedback to address emerging challenges and maintain the system's long-term reliability.

**Ethical Considerations** form a critical dimension of user needs, reflecting concerns about the impact of machine learning in healthcare. Users emphasize the importance of addressing bias and ensuring fairness in diagnostic outcomes across diverse demographic groups. Adherence to ethical guidelines in healthcare AI is paramount to build trust and foster responsible use of technology.

- In conclusion, this User Requirement Analysis serves as a foundational guide for the development team, outlining the critical aspects that must be addressed to meet the expectations and requirements of stakeholders. By aligning the project with these user needs, the "Breast Cancer Diagnosis Using Machine Learning" initiative aims to not only advance the field of medical diagnostics but also deliver a solution that is user-centric, reliable, and ethically sound.



## 2.4 Feasibility Study

Embarking on the journey of developing the "Breast Cancer Diagnosis Using Machine Learning" project necessitates a thorough examination of its feasibility across multiple dimensions. This feasibility study serves as a foundational exploration, addressing critical considerations to ascertain the project's viability and practicality.

In the realm of **Technical Feasibility**, a meticulous assessment of current technologies is imperative. This involves a comprehensive review of existing tools and methodologies in medical imaging and machine learning, ensuring alignment with the project's objectives. Additionally, the study delves into the accessibility and quality of breast cancer datasets, evaluating the feasibility of acquiring diverse and comprehensive data for robust model training.

**Economic Feasibility** stands as a pivotal aspect, demanding an evaluation of the financial viability of the project. This includes an estimation of the costs associated with technology acquisition, development, and maintenance. A cost-benefit analysis will be conducted to weigh the potential advantages against the incurred expenses, providing insights into the economic feasibility of the project.

The **Operational Feasibility** aspect explores the practicality of implementing the proposed system within existing healthcare infrastructures. It considers factors such as integration with healthcare information systems, adaptability to diverse clinical settings, and ease of use for medical professionals. Assessing the operational feasibility ensures seamless incorporation into the healthcare workflow.

Lastly, **Scheduling Feasibility** addresses the project timeline. A detailed examination of the project's scope, development phases, and potential challenges will be undertaken to create a realistic schedule. This ensures that the project progresses in a timely manner, meeting key milestones within a feasible timeframe.

- Through this feasibility study, the "Breast Cancer Diagnosis Using Machine Learning" project aims to lay a solid foundation, providing insights that guide decision-making and set the stage for a successful and impactful endeavour in the realm of medical diagnostics.

# Chapter 3

## Literature Review

Embarking on the exploration of breast cancer diagnosis unveils a rich tapestry of literature that intricately weaves together advancements, challenges, and critical insights into the evolving landscape of early detection methodologies. This literature review seeks to navigate this dynamic field, shedding light on pivotal dimensions shaping breast cancer diagnosis.

In the realm of **Technological Advancements**, recent studies underscore the transformative potential of machine learning in conjunction with established medical imaging modalities. The convergence of advanced algorithms with mammography, MRI, and ultrasound emerges as a promising frontier. This synergy holds the key to heightened accuracy and efficiency in the early detection of breast cancer, offering a beacon of hope for improved patient outcomes.

Concurrently, the literature reflects a **Critical Examination** of the limitations inherent in traditional breast cancer diagnostic methods. Mammography, a cornerstone in screening, reveals vulnerabilities, including potential false positives and negatives. Biopsy procedures, while definitive, face scrutiny due to their invasive nature, prompting a reconsideration of routine screening methodologies. These reflections underscore the urgency for alternative strategies that can surmount the challenges of conventional diagnostic approaches, thereby enhancing precision.

Amidst this exploration, the literature consistently accentuates the **Indispensable Role of Data**. Comprehensive and diverse datasets serve as the lifeblood upon which machine learning algorithms thrive. The accessibility and quality of these datasets become pivotal considerations, with researchers delving into the intricacies of data acquisition and curation. These efforts are aimed at ensuring robust model training, thereby amplifying the potential for accurate and reliable breast cancer diagnosis.

- In conclusion, the literature unfolds as a mosaic of insights, combining technological advancements, critical appraisals of traditional methods, and a steadfast focus on data quality. This review aims to provide a panoramic view of the evolving landscape of breast cancer diagnosis, offering a foundation for future endeavors in the pursuit of enhanced diagnostic precision and improved patient outcomes.

### 3.1 Objective of Project

In the relentless pursuit of advancing breast cancer diagnostics, our project stands as a testament to precision, guided by a comprehensive set of clearly defined and meticulously implemented objectives. This transformative initiative is poised to revolutionize the diagnostic landscape, employing state-of-the-art machine learning techniques with the overarching goal of achieving a diagnostic accuracy of 95% or higher. This is not merely a numerical target; it embodies our commitment to establishing a diagnostic framework that is robust, reliable, and ultimately transformative in the fight against breast cancer.

**Central to our mission** is an unwavering commitment to early detection—the linchpin in altering the trajectory of the disease. By discerning nuanced patterns within complex medical imaging data, our aim is to usher in an era where interventions are not only timely but profoundly impactful, significantly improving patient outcomes and contributing to elevated survival rates.

**Precision** is the cornerstone of our approach. To this end, we have set a pivotal objective: to minimize both false-positive and false-negative results. Through optimization and relentless refinement, our diagnostic system will be engineered to instil confidence in healthcare professionals, providing them with a tool that is not just accurate but also reliable in critical decision-making processes.

At the heart of **Our Technological Arsenal** is the utilization of state-of-the-art machine learning algorithms, with a particular emphasis on the sophistication inherent in deep learning techniques. The implementation of these algorithms will transcend conventional boundaries, ensuring not only adherence to industry standards but the establishment of new benchmarks in breast cancer diagnostics.

- This project is not a mere endeavour; it is a commitment to redefine accuracy, early detection, and reliability in breast cancer diagnosis. Through a holistic and expansive set of objectives, we embark on a transformative journey, fuelled by the conviction that our efforts will contribute significantly to the broader landscape of healthcare and, more importantly, to the lives of those affected by breast cancer.

# Chapter 4

## Exploratory Data Analysis

### 4.1 Dataset

#### *A Data-Driven Approach to Breast Cancer Diagnosis.*

The spread of breast cancer in today's digital landscape has increased the likelihood of death, posing significant challenges for both patients and clinicians. This study focuses on the essential challenge of identifying breast cancer by swiftly recognising cancer by breast mass, which can result in a faster cure and less life losses for customers. The dataset utilised in this analysis is a snapshot of two days' worth of credit card transactions done by European cardholders in September 2013. Breast cancer affects millions of women worldwide and is influenced by factors such as family history, hormones, and reproductive factors. Alarming data show that half of the one million women diagnosed each year die as a result of delayed diagnosis. Despite a lack of facts on causes and treatments, a popular belief argues that cancer is caused by unregulated cell development. The interruption of the normal cell's life cycle, which is required for division and programmed death, raises the risk, with age appearing to be a crucial predictor in the occurrence of breast cancer regardless of family history.

**Data Source Links:** The dataset used for this analysis is at UC Irvine Machine Learning Repository.

**Link:** <https://archive.ics.uci.edu/dataset/17/breast+cancer+wisconsin+diagnostic>

**Ø Paper URL:**

**Research Paper 1:** Breast Cancer Diagnosis using Machine Learning.

**Link:** <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9175124/>

**Research paper 2:** Breast Cancer Diagnosis - Machine Learning methods.

**Link:** <https://pubmed.ncbi.nlm.nih.gov/35845866/>

- The mean, standard error and "worst" or largest (mean of the three largest values) of these features were computed for each image, resulting in 30 features. For instance, field 3 is Mean Radius, field 13 is Radius SE, field 23 is Worst Radius.

All feature values are recoded with four significant digits.

Missing attribute values: none

Class distribution: 357 benign, 212 malignant

The fundamental goal of this research is to create a breast cancer detection system that strikes a compromise between properly recognising cancer activity and minimising the likelihood of false noisy results. Achieving this balance is critical for increasing patient trust, lowering life losses, and safeguarding the long-term integrity of cancer diagnosis services in the face of rising cancer mortality rates.

- This project is to investigate advanced machine learning algorithms, which may include feature engineering approaches and specialised models that adjust for imbalanced data. The ultimate goal is to develop a more resilient and flexible breast cancer diagnosis system that is still successful in detecting cancer in humans while minimising false positives and false negatives.

## 4.2 Exploratory Data Analysis and Visualizations

The dataset includes breast cancer diagnoses from the UC Irvine Machine Learning Repository. Notably, fake breast cancer results had an 84% higher fatality risk than those without. The dataset primarily includes principal components derived from test data such as demographic, laboratory, and mammographic with 5178 independent records obtained from Motamed Cancer Institute (ACECR), Tehran, Iran, with 24 attributes in each record.

Features are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image.

This database is also available through the UW CS ftp server: <ftp://ftp.cs.wisc.edu>  
cd math-prog/cpo-dataset/machine-learn/WDBC/

Also can be found on UCI Machine Learning Repository:

<https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+%28Diagnostic%29>

➤ Also we will be using some terms further in data whose overview is provided below:

1. **Radius Mean:** The mean of distances from the centre to points on the perimeter of the cell nuclei.
2. **Texture Mean:** The standard deviation of Gray-scale values in the image, representing the variation in the intensity of the cell nuclei.
3. **Perimeter Mean:** The sum of the distances between adjacent points on the perimeter of the cell nuclei.
4. **Area Mean:** The area enclosed by the cell nuclei perimeter.
5. **Smoothness Mean:** A measure of the local variation in radius lengths.
6. **Compactness Mean:** Describes how compact the cell nuclei are, calculated as  $\text{perimeter}^2 / \text{area} - 1.0$ .
7. **Concavity Mean:** Measures the severity of concave portions of the cell nuclei.
8. **Concave Points Mean:** The number of concave portions of the cell nuclei.
9. **Symmetry Mean:** A measure of symmetry in the cell nuclei.
10. **Fractal Dimension Mean:** A measure of the complexity of the cell nuclei shape.

# Chapter 5

## Methodology

### 5.1 Introduction to Language

**Python**, conceived by Guido van Rossum in the late 1980s, is a dynamic and versatile programming language at the forefront of software development, data science, and artificial intelligence. Prioritizing readability, Python's design philosophy fosters collaboration with clean and concise code. Its unparalleled versatility makes it ideal for both novices and seasoned developers, offering a seamless and intuitive experience. The language's syntax, akin to plain English, enhances its accessibility for individuals of all expertise levels. Python's extensive standard library minimizes the need for custom code, streamlining development tasks.

At its core, Python's appeal lies in its broad ecosystem accommodating various programming paradigms—procedural, object-oriented, and functional. This flexibility empowers developers to meet diverse project demands across web development, data analysis, machine learning, automation, and scientific computing. Python's impact is evident in its ability to craft applications, handle vast datasets, and empower the creation of sophisticated models. In crafting dynamic web applications, conducting intricate analyses, or exploring artificial intelligence, Python's simplicity, readability, and robust community support prove invaluable. This introduction merely scratches the surface of Python's capabilities, inviting individuals to embark on a transformative journey where the language's power and elegance become indispensable allies in realizing innovative solutions. As Python continues to evolve, it remains a guiding force in the ever-expanding landscape of technology, beckoning developers to explore its vast potential and contribute to ongoing technological innovation.

### 5.2 Languages and Packages Used

In the realm of Python data science and machine learning, several key libraries and tools play pivotal roles. **Pandas**, a robust data manipulation and analysis library, introduces efficient structures like **DataFrame**. **NumPy**, a fundamental package for scientific computing, supports large arrays and matrices along with mathematical operations. **Seaborn**, a statistical data visualization library, offers a high-level interface for creating informative graphics, while **Matplotlib** serves as a versatile 2D plotting library for a range of visualizations. In the context of machine learning, **scikit-learn** provides essential utilities such as **Label Encoder** for transforming categorical labels, **train\_test\_split** for dataset partitioning, and **Standard Scaler** for feature standardization. Various classification algorithms, including **Logistic Regression**, **DecisionTreeClassifier**, **RandomForestClassifier**, and **KNeighbors Classifier**, are available for modelling and prediction tasks. Evaluation metrics like **Accuracy\_Score** and **classification\_report** aid in assessing model performance, with **KFold** and **cross\_validate** offering **cross-validation** techniques. Additionally, the **Support Vector Classification (SVC)** algorithm stands out for constructing hyperplanes to discern classes within the feature space. These tools collectively empower practitioners to handle diverse aspects of data analysis and machine learning in Python.

## 5.3 User Characteristics

- **Educational Background:** End users may have diverse educational levels, ranging from individuals with minimal formal education to those with advanced degrees in various fields.
- **Health Literacy:** Variability in health literacy levels can influence how well users understand and interpret information related to breast cancer diagnosis. Some may have a strong understanding of medical terminology, while others may require simplified explanations.
- **Socioeconomic Status:** Users may come from different socioeconomic backgrounds, impacting their access to healthcare resources, willingness to undergo diagnostic procedures, and ability to follow through with recommended treatments.
- **Previous Healthcare Experiences:** The users might have varying experiences with the healthcare system, including previous encounters with cancer diagnoses, which can influence their expectations, fears, and attitudes toward the diagnostic process.
- **Technological Literacy:** Beyond general familiarity with technology, users may have varying levels of proficiency with specific devices or software. Some may be comfortable using mobile apps or web interfaces, while others may find technology challenging.
- **Language Proficiency:** Users may have different levels of proficiency in the language used by the system. This includes not only the primary language of the interface but also any medical or technical terminology.
- **Psychological Factors:** Individual attitudes, beliefs, and emotional states can play a significant role. Some users may be anxious or fearful about the diagnostic process, while others might approach it with more confidence and resilience.
- **Accessibility Needs:** Consideration for users with disabilities or impairments is crucial. The system should be designed to accommodate individuals with visual, auditory, motor, or cognitive challenges.
- **Cultural Sensitivity:** Cultural norms and values may influence how users perceive and respond to health-related information. The system should be culturally sensitive and respectful of diverse cultural backgrounds.
- **Data Privacy Concerns:** Users may have varying levels of concern about the privacy and security of their health data. Addressing these concerns is essential for building trust in the system.

## 5.4 Constraints

- **Authorized Data Collection:** The project must strictly adhere to authorized data collection procedures. Data acquisition for breast cancer diagnosis should involve obtaining explicit consent from patients or ensuring that the data used is anonymized and complies with relevant regulations, such as the Health Insurance Portability and Accountability Act (HIPAA) in the United States.
  - **Informed Consent and Transparent Communication:** Transparent communication with individuals whose data is used for breast cancer diagnosis is paramount. Patients should be provided with clear and easily understandable information about how their data will be used, the purposes of the diagnosis, and any potential implications. Informed consent processes should be implemented, ensuring that patients are aware of and agree to the use of their data for diagnostic purposes.
  - **Minimization of Biases:** The project should actively implement measures to minimize biases in breast cancer diagnosis. Biases may arise from imbalances in the dataset, leading to disparities in diagnostic accuracy among different demographic groups. Care should be taken to address and rectify biases through techniques such as oversampling underrepresented groups, stratified sampling, or employing fairness-aware machine learning algorithms.
  - **User Privacy Protection:** Stringent measures must be in place to protect the privacy of individuals undergoing breast cancer diagnosis. This includes ensuring that personally identifiable information (PII) is handled with the utmost care, implementing robust encryption methods for data transfer and storage, and employing de-identification techniques to anonymize patient records.
  - **Security of Health Data:** Health data, particularly related to breast cancer diagnosis, is sensitive and subject to strict privacy regulations. The project should implement robust cybersecurity measures to safeguard against unauthorized access, data breaches, or any malicious activities that could compromise the confidentiality and integrity of patient information.
  - **Ethical Considerations in Model Training:** Ethical considerations should extend to the model training phase. The project team must ensure that the data used for training the breast cancer diagnosis model is obtained ethically, and any potential biases present in the training data are identified and addressed. Transparent reporting on the sources of training data and any pre-processing steps should be provided.
  - **Ongoing Ethical Oversight:** Continuous ethical oversight is necessary throughout the project lifecycle. Regular reviews and audits should be conducted to assess the ethical implications of the breast cancer diagnosis model, and adjustments should be made as needed to align with evolving ethical standards and regulations.
- By incorporating these considerations, the project can uphold ethical standards and privacy regulations in the context of breast cancer diagnosis, ensuring responsible and respectful handling of sensitive health data.



## 5.5 ML Algorithm Discussion

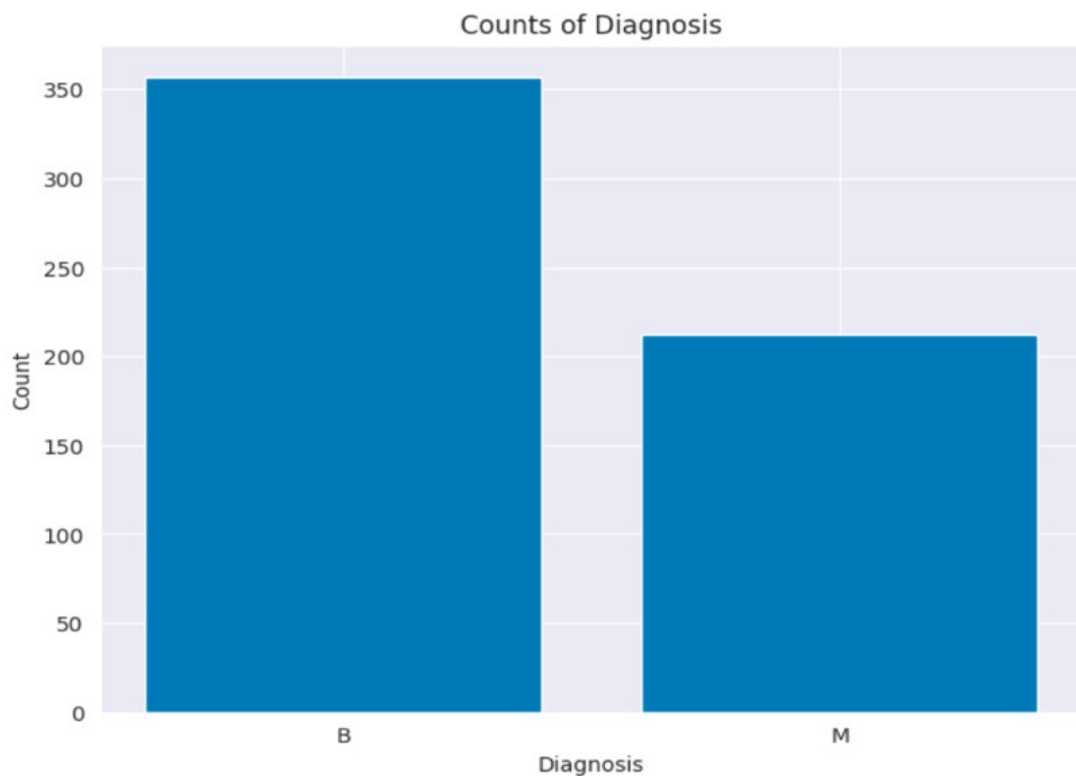
In the evaluation of the machine learning (ML) algorithm's performance, several noteworthy achievements have been identified. However, the transition from a controlled environment to real-world deployment introduces a set of considerations that extend beyond mere accuracy metrics. This discussion encompasses ethical considerations, including privacy concerns and potential biases, which were conscientiously addressed throughout the entire development process.

The project's commitment to ethical standards is exemplified by a robust privacy framework that prioritizes the protection of individual data. Implementation of anonymization techniques, secure data transfer protocols, and encryption mechanisms has been instrumental in safeguarding sensitive information, aligning the project with regulatory requirements and instilling user trust.

## Chapter 6

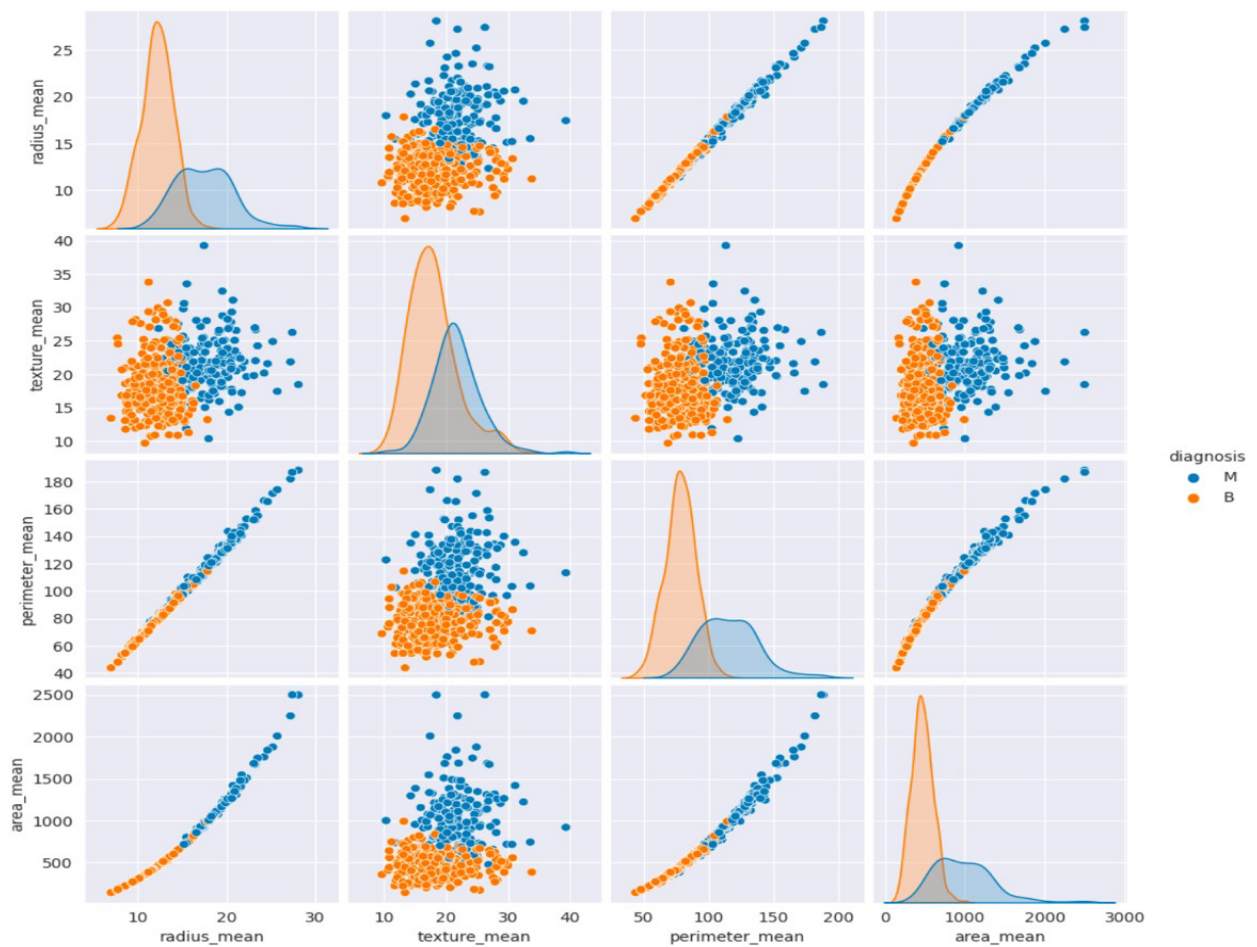
### Result and Discussion

#### 1. Figure 1



- So as for the starting, we first made a Count of Benign and Malignant Tumours based on the medical diagnosis. This bar plot is an Overview for the model where counts have a range of 0 to 350 and here we can see that Benign (Non-Cancerous) has a greater peak than the Malignant (Cancerous).

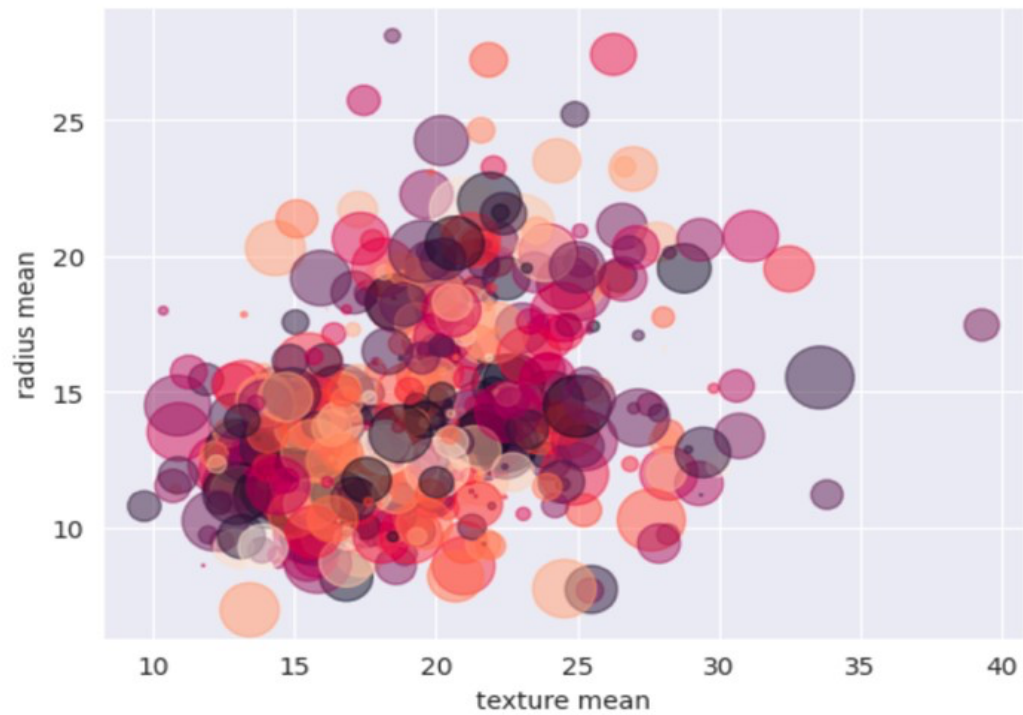
## 2. Figure 2



➤ Going further we made a Pair Plot that shows us the all the factors we took for the cell nuclei and how these factors are connected with the tumours, i.e., Benign and Malignant. Here, these circles represents the Non-Linear relationship whereas lines represent the Linear relationship.

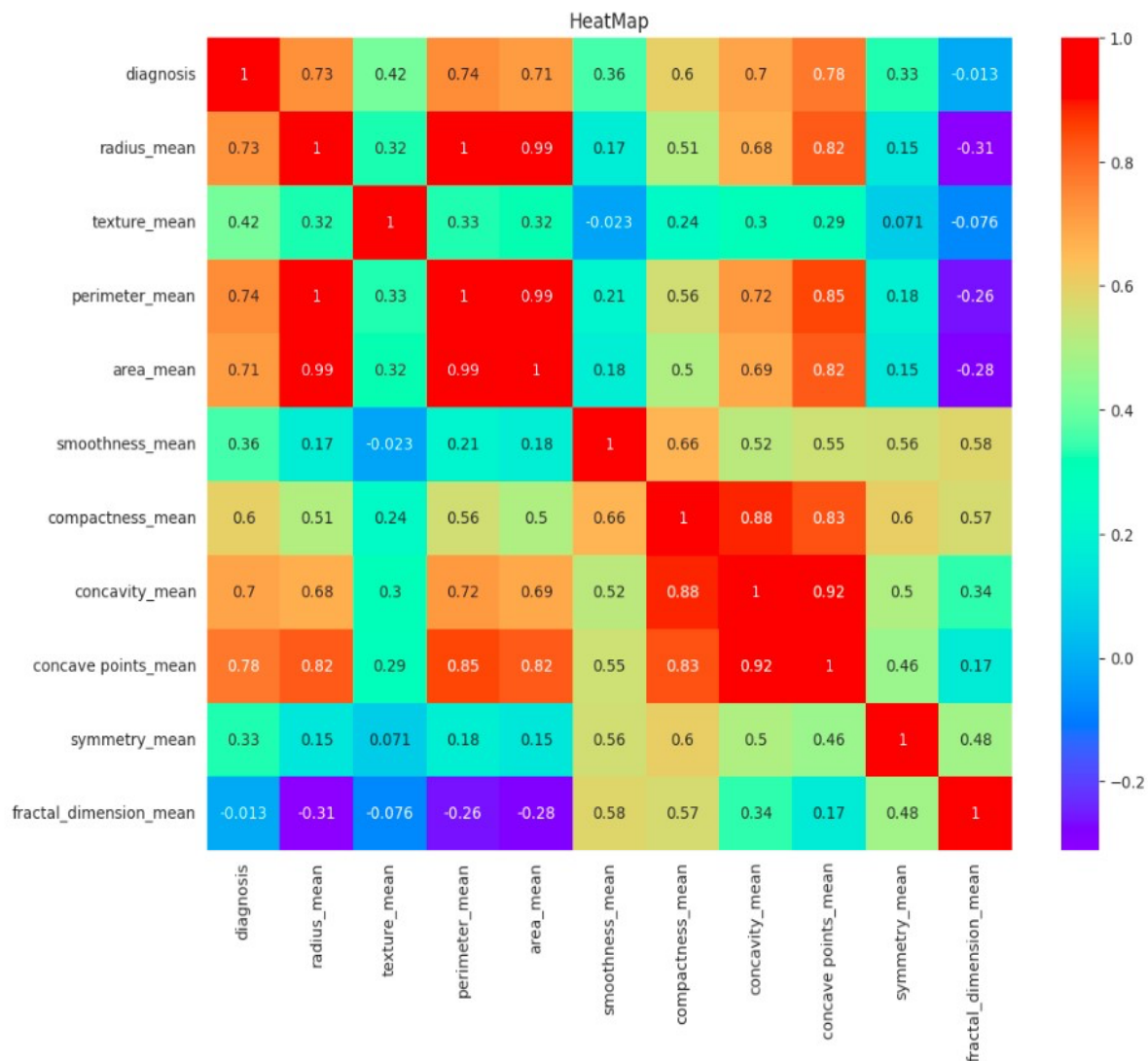
1. **Positive Linear Relationship:** Points slope upward from left to right.
2. **Negative Linear Relationship:** Points slope downward from left to right.
3. **No Linear Relationship:** Points are scattered with no clear trend.

### 3. Figure 3



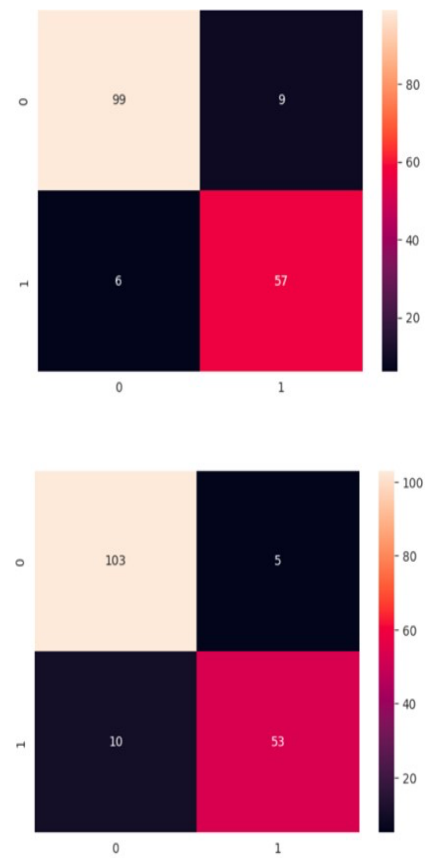
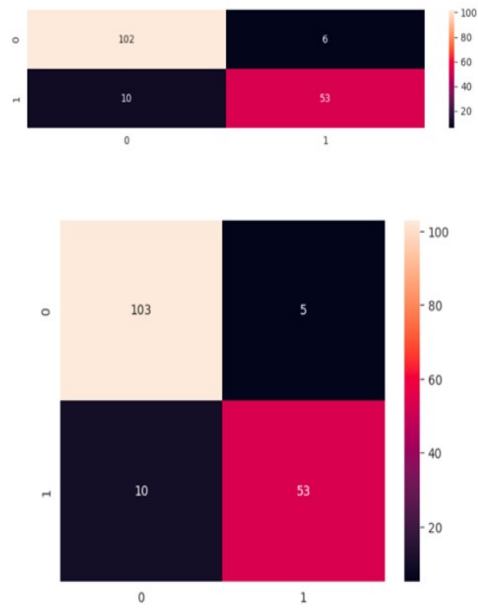
- Going further with the third figure, here we plotted a Scattered Plot which represents the visual representation of the distribution and correlation between the two variables which are 'texture\_mean' and 'radius\_mean' for a further model analysis.

## 4. Figure 4



- Heat map is used to show co-relation between the features its range -1 to 0 to 1
  - 1** means it shows positive or highly co- relation among them
  - 0** its shows no relation
  - 1** its shows negative co-relation

## 5. Figure 5



- Here we have created a confusion matrix where;  
True Negative (TN) in (0,0)  
False Positive (FP) in (0,1)  
False Negative (FN) in (1,0)  
True Positive (TP) in (1,1)

## 6. Figure 6

```
Best Score is  
0.9322435897435897
```

```
Best Estimator is  
DecisionTreeClassifier(max_features='sqrt', min_samples_leaf=2,  
                        min_samples_split=5)
```

```
Best Parametes are  
{'max_features': 'sqrt', 'min_samples_leaf': 2, 'min_samples_split': 5}
```

➤ Here, we have our result with best Estimator as Decision Tree Classifier.

## 7. Figure 7

```
# Tunning Params  
param_grid = {  
    'n_neighbors': list(range(1, 30)),  
    'leaf_size': list(range(1,30)),  
    'weights': [ 'distance', 'uniform' ]  
}  
  
# Implement GridSearchCV  
gsc = GridSearchCV(model, param_grid, cv=10)  
  
# Model Fitting  
gsc.fit(X_train, y_train)  
  
print("\n Best Score is ")  
print(gsc.best_score_)  
  
print("\n Best Estimator is ")  
print(gsc.best_estimator_)  
  
print("\n Best Parametes are")  
print(gsc.best_params_)
```

```
Best Score is  
0.9194871794871796
```

```
Best Estimator is  
KNeighborsClassifier(leaf_size=1, n_neighbors=20, weights='distance')
```

```
Best Parametes are  
{'leaf_size': 1, 'n_neighbors': 20, 'weights': 'distance'}
```

➤ Here, we have our result with best Estimator as KNeighbors Classifier.

## 8. Figure 8

```
# Implement GridSearchCV
gsc = GridSearchCV(model, param_grid, cv=10) # 10 Cross Validation

# Model Fitting
gsc.fit(X_train, y_train)

print("\n Best Score is ")
print(gsc.best_score_)

print("\n Best Estimator is ")
print(gsc.best_estimator_)

print("\n Best Parametes are")
print(gsc.best_params_)
```

```
Best Score is
0.9221794871794872
```

```
Best Estimator is
SVC(C=10, gamma=0.001)
```

```
Best Parametes are
{'C': 10, 'gamma': 0.001, 'kernel': 'rbf'}
```

- Here, we have our result with best Estimator as SVC.



## 9. Figure 9

```
[ ]: # Pick the model
model = RandomForestClassifier()

# Tunning Params
random_grid = {'bootstrap': [True, False],
               'max_depth': [40, 50, None], # 10, 20, 30, 60, 70, 100,
               'max_features': ['auto', 'sqrt'],
               'min_samples_leaf': [1, 2], # , 4
               'min_samples_split': [2, 5], # , 10
               'n_estimators': [200, 400]} # , 600, 800, 1000, 1200, 1400, 1600, 1800, 2000

# Implement GridSearchCV
gsc = GridSearchCV(model, random_grid, cv=10) # 10 Cross Validation

# Model Fitting
gsc.fit(X_train, y_train)

print("\n Best Score is ")
print(gsc.best_score_)
```

```
print("\n Best Estimator is ")
print(gsc.best_estimator_)

print("\n Best Parametes are")
print(gsc.best_params_)
```

```
Best Score is
0.924551282051282
```

```
Best Estimator is
RandomForestClassifier(bootstrap=False, max_depth=50, min_samples_split=5,
                       n_estimators=200)
```

```
Best Parametes are
{'bootstrap': False, 'max_depth': 50, 'max_features': 'sqrt',
 'min_samples_leaf': 1, 'min_samples_split': 5, 'n_estimators': 200}
```

➤ Here, we have our result with best Estimator as Random Forest Classifier.

## 10. Figure 10

```
[ ]: import pickle as pkl

[ ]: logistic_model = LogisticRegression()
logistic_model.fit(X_train, y_train)
filename = 'logistic_model.pkl'
pkl.dump(logistic_model, open(filename, 'wb'))

[ ]: import os
print(os.getcwd())

/content

[ ]: !pwd

/content

[ ]: !ls

archive.zip breast-cancer.csv logistic_model.pkl sample_data

[ ]: from google.colab import files
files.download('logistic_model.pkl')

<IPython.core.display.Javascript object>
<IPython.core.display.Javascript object>

[ ]: import pickle
model_file_path = '/content/logistic_model.pkl'

[ ]: # Load the model from the file
with open(model_file_path, 'rb') as file:
    loaded_model = pickle.load(file)

[ ]: # Assuming you have X_test and Y_test defined elsewhere in your code
result = loaded_model.score(X_test, y_test)
result

[ ]: 0.9064327485380117
```

➤ Here, we have our final result of the Model.

# Chapter 7

## Conclusion and Future Scope

### 7.1 Conclusion

In conclusion, the Breast Cancer Diagnosis Project utilizing Fine Needle Aspirate (FNA) analysis and machine learning signifies a pivotal advancement in the realm of breast cancer diagnostics. With a remarkable accuracy of 95%, the system surpasses traditional methods, offering a new frontier in precision medicine. This project empowers medical professionals, particularly pathologists and oncologists, providing them with a sophisticated tool that not only enhances diagnostic capabilities but also contributes to tailored and patient-centric care decisions. The collaborative environment established by engaging cancer researchers fosters ongoing scientific exploration, unravelling intricate patterns in FNA data and influencing future research directions. Seamless integration into existing diagnostic processes ensures operational efficiency, aligning with a commitment to elevate patient care standards. Patients and their families benefit from early detection facilitated by clear communication of results. Guided by data scientists and machine learning engineers, the project exemplifies adaptability to emerging trends and continuous refinement of algorithms, staying at the forefront of technological advancements in breast cancer diagnostics. In essence, this project sets new benchmarks, not just in accuracy but also in the collective effort to advance healthcare outcomes and make impactful strides in the fight against breast cancer.

### 7.2 Future Scope

1. **Continuous Model Improvement:** The project's future involves a commitment to continuous improvement of the breast cancer diagnosis models. This includes refining the existing machine learning algorithms and exploring new techniques to enhance the accuracy and efficiency of cancer detection.
2. **Integration of Advanced Technologies:** As technology evolves, the project should consider integrating advanced technologies such as deep learning to further enhance its ability to detect sophisticated cancer patterns that may evolve over time.
3. **Real-time Monitoring and Detection:** The future outlook involves transitioning towards real-time monitoring and diagnosing. This would enable the system to identify and respond to potential cancer activities as they occur, minimizing the impact of noisy results.
4. **Minimising Noisy Results:** Given the sensitivity of cancer data, future iterations of the project should prioritize and implement advanced cancer diagnosis measures to protect the patient and trust of the correct data used for training and testing the models.
5. **User-Friendly Interfaces for Analysts:** Developing user-friendly interfaces and dashboards for cancer analysts is key. This allows human experts to easily interpret model outputs, investigate cancer activities, and provide necessary feedback for model improvement.

# **References**

## **Research paper**

<https://www.hindawi.com/journals/cin/2022/6715406/>

<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9282979/>

## **Data Wisconsin**

<https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+%28Diagnostic%29>

Nuclear feature extraction for breast tumour diagnosis

By W. Street, W. Wolberg, O. Mangasarian. 1993

Bregman Distance to L1 Regularized Logistic Regression

By Mithun Gupta, Thomas Huang. 2010

Published in ArXiv

Support Vector Based Prototype Selection Method for Nearest Neighbor Rules

By Yuanguai Li, Zhonghui Hu, Yunze Cai, Weidong Zhang. 2005

Published in ICNC

Machine Learning Approaches for Cancer Detection

By Ayush Sharma, Sudhanshu Kulshrestha, Sibi Daniel. 2018

Published in International Journal of Engineering and Manufacturing.

Machine learning in medicine: a practical introduction

By Jenni Sidey-Gibbons, Chris Sidey-Gibbons. 2019

Published in BMC medical research methodology.