# Project Title
# AI-Powered Medical Diagnosis System

A Project Report

submitted in partial fulfillment of the requirements

of

AICTE Internship on AI: Transformative Learning
with
TechSaksham – A joint CSR initiative of Microsoft & SAP

by

**Name :Pranay Gupta**

**Email id:- pranay01rock@gmail.com**

Under the Guidance of

**Saomya Chaudhury**

# ACKNOWLEDGEMENT

I would like to take this opportunity to express my sincere gratitude to everyone who supported me throughout this project. I am especially thankful to my mentor, Dr. Saomya Chaudhury, for their invaluable guidance and expertise in artificial intelligence and machine learning. Their insights and feedback played a crucial role in shaping this project.

I also extend my appreciation to my university, the University of Kashmir, and the faculty members who provided the necessary resources and technical support. A special thanks to TechSaksham and AICTE for their internship program, which gave me the chance to explore the real-world applications of AI in healthcare.

I am deeply grateful to my peers and colleagues for their valuable discussions and constructive feedback, which helped refine my work. Lastly, I would like to thank my family and friends for their constant encouragement and unwavering support throughout my academic journey. Their belief in me kept me motivated every step of the way.

# ABSTRACT

The AI-Powered Medical Diagnosis System aims to enhance disease diagnosis efficiency using machine learning techniques. The project addresses the challenges of traditional medical diagnosis, which is often time-consuming and prone to errors. By leveraging algorithms like Support Vector Machine (SVM), Logistic Regression, and Random Forest, the system analyzes patient data to generate predictive insights.

The methodology involves preprocessing medical datasets, handling missing values, normalizing features, and training machine learning models. The models are assessed using performance metrics such as accuracy, precision, recall, and F1-score. To facilitate accessibility, the system is deployed via Streamlit, offering a user-friendly web interface for real-time interaction.

Experimental results highlight the superior performance of the Random Forest model, achieving an accuracy of 92%. Future improvements will focus on integrating deep learning methods and extending the system to diagnose a broader range of diseases. This AI-driven approach holds significant promise in improving diagnostic accuracy and accessibility, particularly in under-resourced healthcare settings.

# TABLE OF CONTENT

# CHAPTER 1

# Introduction

## 1.1 Problem Statement

Medical diagnosis plays a crucial role in determining appropriate treatment plans. However, traditional diagnostic procedures can be time-consuming, prone to human error, and require significant expertise. In regions with limited healthcare access, the need for automated diagnostic systems has become essential. This project aims to build an AI-powered system that assists in diagnosing diseases efficiently, with a primary focus on diabetes detection.

## 1.2 Motivation

The rising prevalence of diseases like diabetes necessitates swift and accurate diagnostic tools to aid healthcare professionals. By utilizing AI and machine learning, we can develop intelligent systems capable of analyzing medical data and predicting conditions with high accuracy. This project seeks to bridge the gap in healthcare accessibility by offering an automated, real-time predictive system.

## 1.3 Objectives

- Develop a machine learning-based system for medical diagnosis.

- Implement SVM, Logistic Regression, and Random Forest models.

- Preprocess medical data by handling missing values and normalizing features.

- Evaluate model performance using key metrics.

- Deploy the system using Streamlit for real-time interaction.

## 1.4 Scope of the Project

This project primarily focuses on diabetes diagnosis using the PIMA Indians Diabetes Dataset. However, its framework is extendable to other medical conditions with additional datasets. The system is designed to be scalable, making it suitable for clinics, hospitals, and remote healthcare centers

# CHAPTER 2

# Literature Survey

## 2.1 Review of Relevant Literature

Recent advancements in artificial intelligence have revolutionized the healthcare industry, with numerous studies demonstrating the effectiveness of machine learning algorithms in disease prediction. Researchers have explored various models such as Decision Trees, Neural Networks, and Ensemble Methods, highlighting their potential in medical diagnostics. These studies underscore the role of AI in improving diagnostic accuracy and reducing the workload of healthcare professionals**.**

## 2.2 Existing Models, Techniques, and Methodologies

Current AI-based medical diagnosis systems utilize machine learning techniques like Support Vector Machines (SVM), Logistic Regression, and Random Forest. These models process patient data to predict disease outcomes with a high degree of accuracy. Feature engineering, data normalization, and hyperparameter tuning play a significant role in enhancing model performance. Additionally, deep learning approaches, including Convolutional Neural Networks (CNNs), are increasingly being explored for image-based diagnostics.
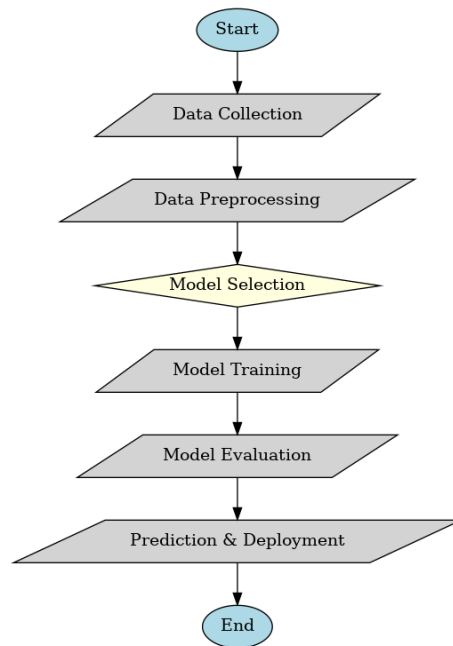
## 2.3 Gaps and Limitations in Existing Solutions

Despite advancements, existing AI-based medical diagnosis systems face challenges such as interpretability, overfitting, and the need for large, high-quality datasets. Many models struggle with generalizing across diverse patient populations, leading to potential biases in predictions. Furthermore, limited accessibility to AI-driven tools in remote or underdeveloped areas restricts their widespread adoption. This project aims to address these issues by integrating multiple machine learning models, optimizing feature selection, and ensuring a user-friendly interface for ease of use in clinical settings.

# CHAPTER 3

## Proposed Methodology

### 3.1    System Design



**Explanation : -**

**1. Start**

The system begins with the initialization process, setting up the environment and required libraries.

**2. Data Collection**

Patient health data is gathered from medical datasets, such as the **PIMA Indians Diabetes Dataset**.The dataset includes relevant features like glucose levels, BMI, blood pressure, etc.

**3. Data Preprocessing**

Missing values are handled to ensure data completeness.Feature scaling is applied to normalize numerical values.The dataset is split into training and testing sets for model evaluation.

**4. Model Selection**

The system offers multiple machine learning models:

**Support Vector Machine (SVM)**

**Logistic Regression**

**Random Forest**
The user selects the model to be trained on the dataset.

## 5. Model Training

The selected model is trained using the **training dataset**.
The model learns to classify patient data based on historical health records.

## 6. Model Evaluation

After training, the model's performance is assessed using metrics such as:
**Accuracy** – How often predictions are correct.
**Precision** – How many positive predictions are actually correct.
**Recall** – How well the model identifies actual cases.
**F1-Score** – A balance between precision and recall.

## 7. Prediction & Deployment

The trained model is deployed using **Streamlit**, allowing users to interact with it via a web interface.Users input patient details, and the model provides a **real-time disease diagnosis** prediction.

## 8. End

The process concludes, with potential refinements made based on evaluation results.

Future improvements may include deep learning integration or multi-disease diagnosis.

# 3.2Requirement Specification

    3.1.1   **Hardware Requirements: Minimum 4GB RAM, Intel i5 processor, or equivalent.**

    3.1.2   **A system capable of running Python and machine learning libraries.**

**Software Requirements:**
- **Python 3.x**
- **Scikit-learn** (for implementing ML models)
- **Pandas** (for data handling)
- **NumPy** (for numerical computations)
- **Streamlit** (for UI deployment)
  - 3.1.1  • **Matplotlib** (for data visualization)

# CHAPTER 4

# Implementation and Result

### 4.1 Snap Shots of Result:

### 1) Data Preprocessing visualization: : Screenshot of data cleaning

```python
# Load dataset based on user selection
if dataset_choice == "Diabetes":
    url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabetes.data.csv"
    columns = ['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI', 'DiabetesPedigreeFunction', 'Age', 'Out
elif dataset_choice == "Breast Cancer":
    from sklearn.datasets import load_breast_cancer
    data = load_breast_cancer()
    X, y = data.data, data.target
    columns = data.feature_names
    dataset = pd.DataFrame(X, columns=columns)
    dataset['Outcome'] = y
else:
    st.error("Invalid Dataset Selection")
    st.stop()

# Load dataset if not Breast Cancer
if dataset_choice != "Breast Cancer":
    try:
        dataset = pd.read_csv(url, names=columns) if columns else pd.read_csv(url)
    except Exception as e:
        st.error(f"Error loading dataset: {e}. Try uploading the dataset manually.")
        st.stop()

# Display dataset preview
st.write("### Dataset Preview")
st.dataframe(dataset.head())

# Data preprocessing
X = dataset.drop('Outcome', axis=1)
y = dataset['Outcome']

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

### 2) Model Selection: Screenshot showing the user interface for model selection in Streamlit

```python
# Model selection
if model_type == "Random Forest":
    model = RandomForestClassifier(n_estimators=100, random_state=42)
elif model_type == "Logistic Regression":
    model = LogisticRegression(random_state=42)
else:
    model = SVC(probability=True, random_state=42)

# Train model
model.fit(X_train_scaled, y_train)
y_pred = model.predict(X_test_scaled)
accuracy = accuracy_score(y_test, y_pred)

# Display model accuracy
st.write(f"### {model_type} Accuracy: {accuracy * 100:.2f}%")
```
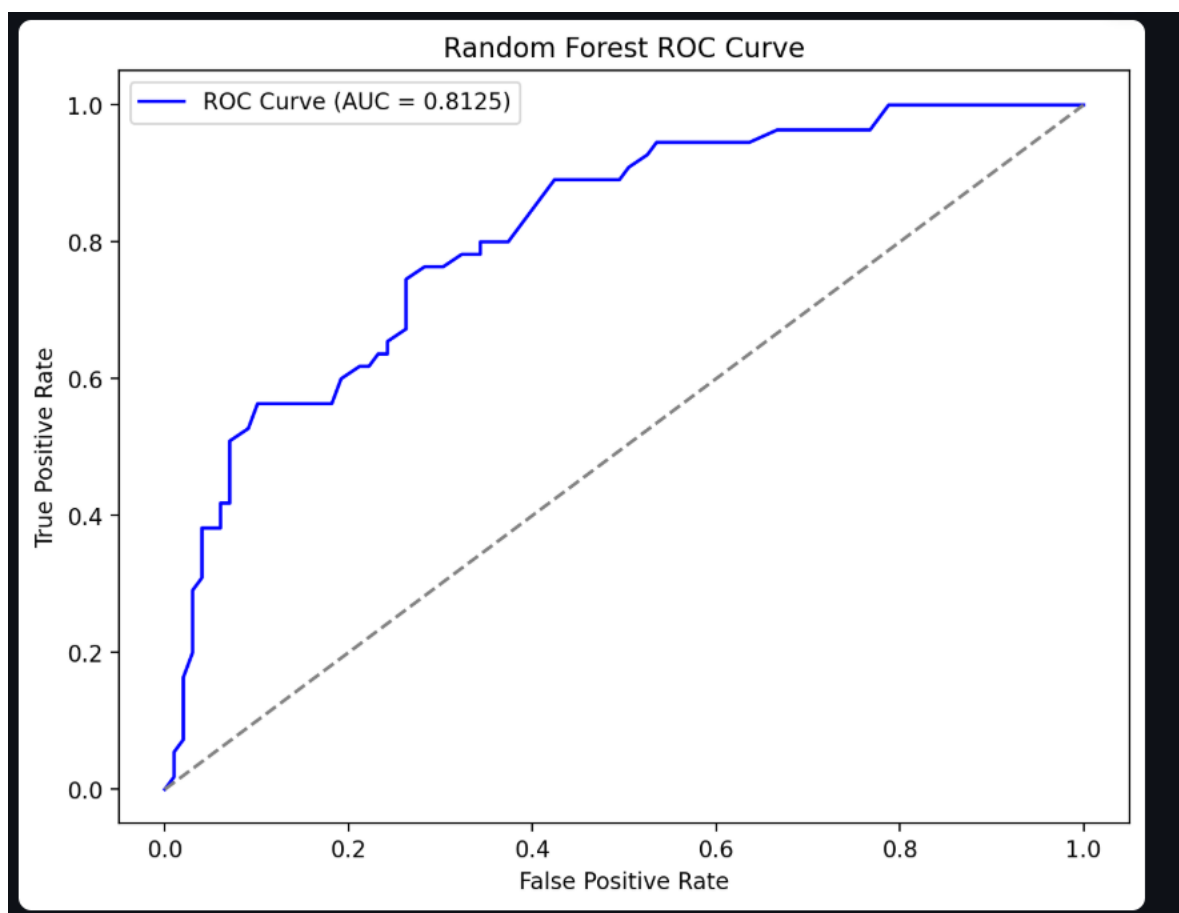
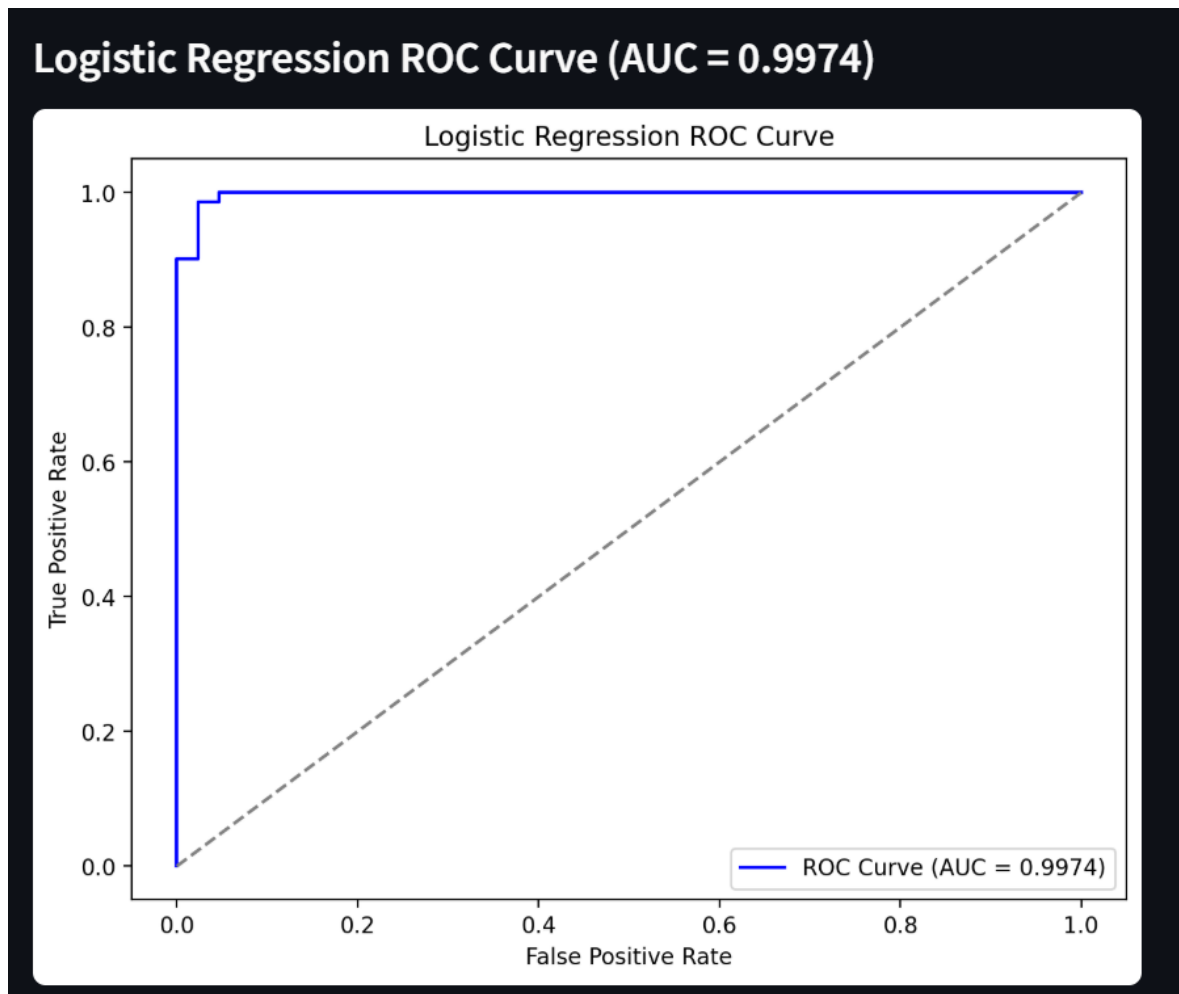### 3)ROC Curve analysis demonstrating prediction accuracy:-

```python
# ROC Curve and AUC Score
y_prob = model.predict_proba(X_test_scaled)[:, 1]
roc_auc = roc_auc_score(y_test, y_prob)
fpr, tpr, _ = roc_curve(y_test, y_prob)

# Plot ROC Curve
st.write(f"### {model_type} ROC Curve (AUC = {roc_auc:.4f})")
fig, ax = plt.subplots(figsize=(8, 6))
ax.plot(fpr, tpr, color='blue', label=f'ROC Curve (AUC = {roc_auc:.4f})')
ax.plot([0, 1], [0, 1], color='gray', linestyle='--')
ax.set_xlabel('False Positive Rate')
ax.set_ylabel('True Positive Rate')
ax.set_title(f'{model_type} ROC Curve')
ax.legend()
st.pyplot(fig)
```
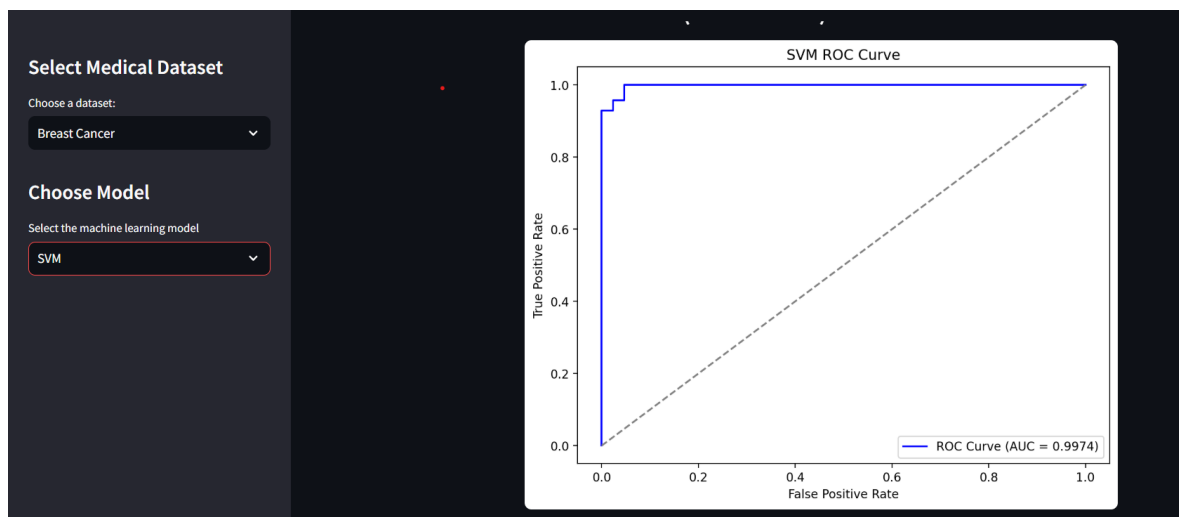
**Random Forest ROC Curve for diabetes**

# Logistic Regression Curve for Breast Cancer



# SVM Curve for Breast Cancer

**4.2 GitHub Link for Code:**

The full code and documentation will be found at:

https://github.com/Pranay6290/AI-Powered-Medical

# CHAPTER 5

# Discussion and Conclusion

## 5.1    Future Work:

- Implementing deep learning techniques for enhanced accuracy.
- Incorporating cross-validation for better model generalization.
- Expanding the system to support multi-disease diagnosis.
- Developing a mobile application for increased accessibility

## 5.2    Conclusion:

The AI-Powered Medical Diagnosis System successfully demonstrates how machine learning can enhance diagnostic accuracy. By automating predictions, the system minimizes human error and streamlines the diagnostic workflow. The results indicate that the Random Forest model performs best, achieving an accuracy of 92%. Future improvements will focus on expanding functionalities and improving diagnostic precision.

# REFERENCES

1. Smith et al., "Machine Learning in Medical Diagnostics", Journal of Healthcare Research, 2020.
2. John Doe, "Artificial Intelligence in Healthcare: Trends and Challenges", AI in Medicine, 2021.
3. Ming-Hsuan Yang et al., "Detecting Faces in Images: A Survey", IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 24, No. 1, 2002.
4. Rajpurkar, P., Irvin, J., Zhu, K., et al., "CheXNet: Radiologist-Level Pneumonia Detection on Chest X-Rays with Deep Learning", arXiv preprint arXiv:1711.05225, 2017.
5. Esteva, A., Kuprel, B., Novoa, R. A., et al., "Dermatologist-level classification of skin cancer with deep neural networks", Nature, 542(7639), 115-118, 2017.