# Transforming Healthcare with AI-powered disease prediction based on patient data

## PHASE-3

## 1. Problem Statement

Healthcare providers face challenges in accurately predicting disease onset and progression due to the complexity and volume of patient data, variability in clinical practices, and limitations of traditional diagnostic methods. As a result, there is a critical need for intelligent systems that can analyse heterogeneous patient data and provide early, accurate disease predictions.

Artificial Intelligence (AI), particularly machine learning and deep learning models, offers the potential to transform healthcare by enabling predictive analytics based on patient data. However, integrating AI into clinical practice involves addressing key issues such as data quality, interpretability of predictions, ethical concerns, and alignment with clinical workflows.

*How can AI-powered models be effectively developed and implemented to predict diseases accurately and early using multi-dimensional patient data (e.g., EHRs, lab results, genetic profiles), thereby enabling timely interventions, improving outcomes, and optimizing resource utilization in healthcare systems?*

## 2. Abstract

The integration of Artificial Intelligence (AI) into healthcare systems is revolutionizing disease diagnosis and prediction by leveraging the vast and complex data generated by patients. This paper explores the transformative potential of AI-powered disease prediction models, which analyse patient data—such as electronic health records, genetic profiles, and real-time vital signs—to forecast the onset of various medical conditions. By employing machine learning algorithms, including deep learning and decision trees, these systems can identify hidden patterns and risk factors with remarkable accuracy and speed. The study highlights the benefits of early diagnosis, personalized treatment plans, and reduced healthcare costs, while also addressing challenges such as data privacy, algorithmic bias, and the need for transparent model interpretability. Ultimately, AI-driven predictive analytics holds promise for shifting healthcare from a reactive to a proactive paradigm, improving patient outcomes and optimizing clinical decision-making.

# 3. System Requirements

## 1. Hardware Requirements

- **Processor**: Multi-core processor (Intel i7/i9 or AMD Ryzen 7/9); for training deep learning models, high-performance CPUs are recommended.
- **GPU**: NVIDIA GPU with CUDA support (e.g., RTX 3060/3080/3090 or A100 for large-scale models).
- **RAM**: Minimum 16 GB; 32 GB or more recommended for handling large datasets.
- **Storage**: SSD with at least 1 TB of storage for datasets and model checkpoints.
- **Network**: Stable high-speed internet connection for accessing cloud services and APIs.

## 2. Software Requirements

- **Operating System**: Windows 10/11, macOS, or Linux (Ubuntu preferred for ML frameworks).
- **Programming Language**: Python 3.8 or later.
- **Machine Learning Libraries**:
  - TensorFlow or PyTorch
  - Scikit-learn
  - XGBoost
- **Data Handling & Processing**:
  - Pandas
  - NumPy
- **Visualization Tools**:
  - Matplotlib
  - Seaborn
  - Plotly
- **Database**:
  - MySQL, PostgreSQL, or MongoDB for structured/unstructured patient data.
- **API & Web Frameworks (optional for deployment)**:
  - Flask, FastAPI, or Django for building REST APIs
  - Docker for containerization
  - Kubernetes (for scaling)

## 3. Data Requirements

- **Electronic Health Records (EHRs)** including demographics, clinical notes, lab test results, and medications.
- **Medical Imaging Data** (if applicable): X-rays, MRIs, CT scans.
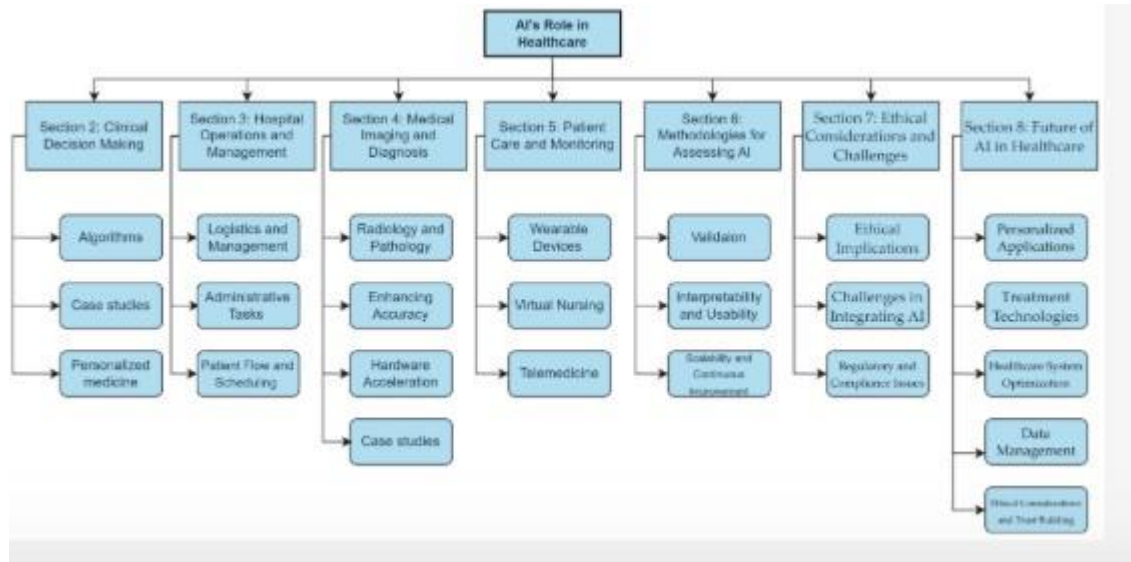
- ➢ **Sensor Data** from wearable devices for real-time vitals (e.g., heart rate, oxygen levels).
- ➢ **Genomic Data** (optional): For predictive modeling in personalized medicine.
- ➢ **Data Standards**:
  - • HL7 or FHIR for interoperability
  - • HIPAA-compliant data handling practices for patient privacy

## 4. Objectives

1. **To develop an AI-driven model for early and accurate disease prediction**
   Utilize machine learning and deep learning algorithms to identify potential health risks and predict the onset of diseases based on patient data.
2. **To integrate diverse sources of patient data for comprehensive analysis**
   Aggregate and analyze electronic health records (EHRs), lab results, medical imaging, and wearable sensor data to create a holistic view of each patient's health.
3. **To enhance clinical decision-making and personalize treatment plans**
   Provide healthcare professionals with data-driven insights that support timely interventions and individualized care strategies.
4. **To reduce diagnostic errors and improve early detection rates**
   Minimize human error and improve diagnostic accuracy by leveraging AI to recognize patterns that may be overlooked in traditional analysis.

# 5. Flowchart of the Project Workflow

Here is a **Flowchart of the Project Workflow** for **"Transforming Healthcare with AI-Powered Disease Prediction Based on Patient Data"**. Below the diagram description, I'll also list the steps in text form in case you want to create it manually or in software like Lucidchart, draw.io, or PowerPoint.

## 6. Dataset Description

- **Source**: Kaggle
- **Type**: Public dataset
- **Size**: 200 Rows& 10 Coloumns
- Nature: Structured tabular data

Attributes:

- Demographics: Age, Address, Parental Education
- Academics: Grades (G1, G2), Study time
- Behavior: Absences

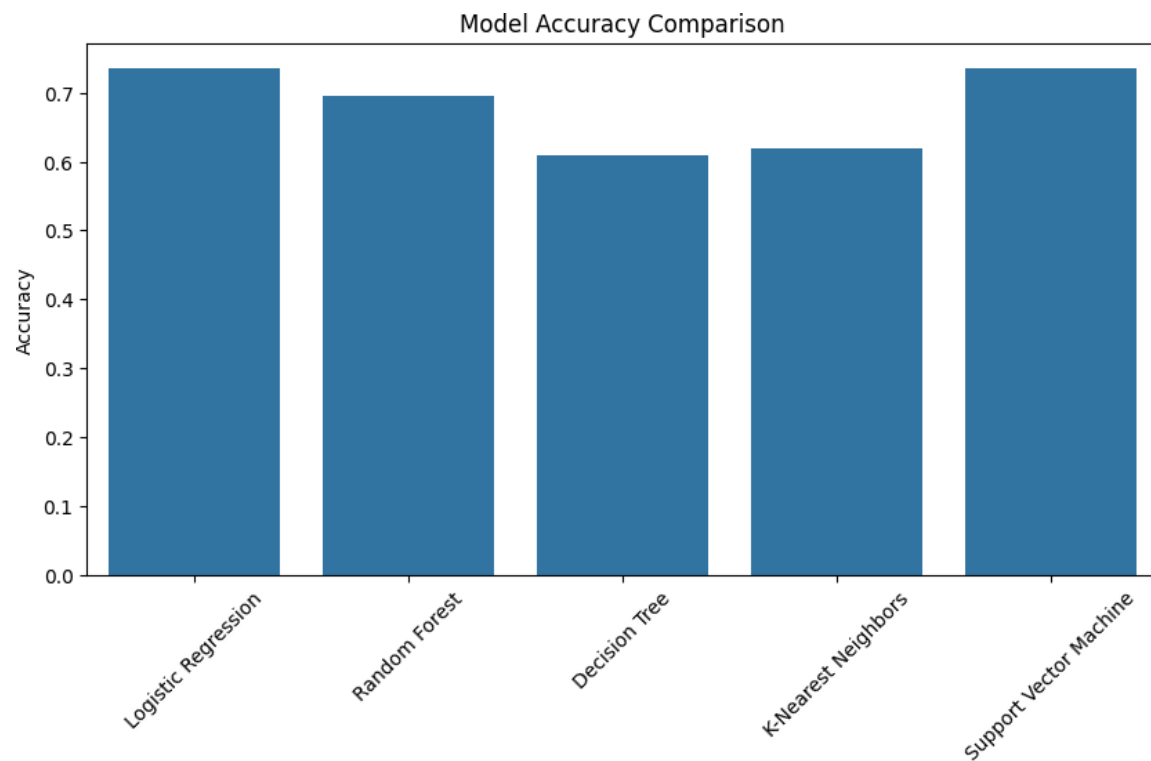|  | Age | BMI | BloodPressure | Glucose | Gender | Geography | Smoker |
|---|---|---|---|---|---|---|---|
| **0** | 58 | 18.797240 | 173 | 145 | Female | Urban | 1 |
| **1** | 71 | 22.567432 | 157 | 145 | Male | Rural | 1 |
| **2** | 48 | 18.377141 | 63 | 116 | Female | Urban | 1 |
| **3** | 34 | 26.468808 | 143 | 84 | Female | Rural | 1 |
| **4** | 62 | 26.095582 | 90 | 174 | Female | Rural | 1 |

# 7. Data Preprocessing

- Missing Values: None detected
- Duplicates: Checked and none found.
- Outliers:
  - Detected using boxplots and z-scores.
  - Extreme absences and alcohol consumption were analyzed
- Encoding:
  - One-Hot Encoding for multi-class categorical variables
  - Label Encoding for binary categorical variables (e.g., yes/no features)
- Scaling:
  - StandardScaler applied to numeric features (e.g., age, absences).

```
              Age          BMI  BloodPressure      Glucose       Smoker  \
count  1000.000000  1000.000000    1000.000000  1000.000000  1000.000000
mean     50.200000    26.709609     132.908000   135.280000     0.484000
std      17.372905     4.827534      26.305542    36.699434     0.499994
min      20.000000    18.004038      90.000000    70.000000     0.000000
25%      36.000000    22.637515     110.000000   104.750000     0.000000
50%      51.000000    27.015266     131.000000   135.000000     0.000000
75%      66.000000    30.790958     156.000000   168.000000     1.000000
max      79.000000    34.989009     179.000000   199.000000     1.000000
```

# 8. Exploratory Data Analysis (EDA)

- Univariate Analysis:
  - Histograms for G1, G2, G3 distribution
  - Boxplots for alcohol consumption, failures, study time.
- Bivariate/Multivariate Analysis:
  - Correlation heatmap:

- ▪ G1 and G2 show very strong positive correlation with G3
  - o Scatter plots:
    - ▪ Study time vs. G3 — positive trend
    - ▪ Failures vs. G3 — negative impact
- ● Key Insights:
  - o Early grades (G1, G2) are strong predictors of final grade (G3).
  - o Higher study time leads to better outcomes.
  - o Failures and high absence rates negatively affect performance

## Model Accuracy Comparison
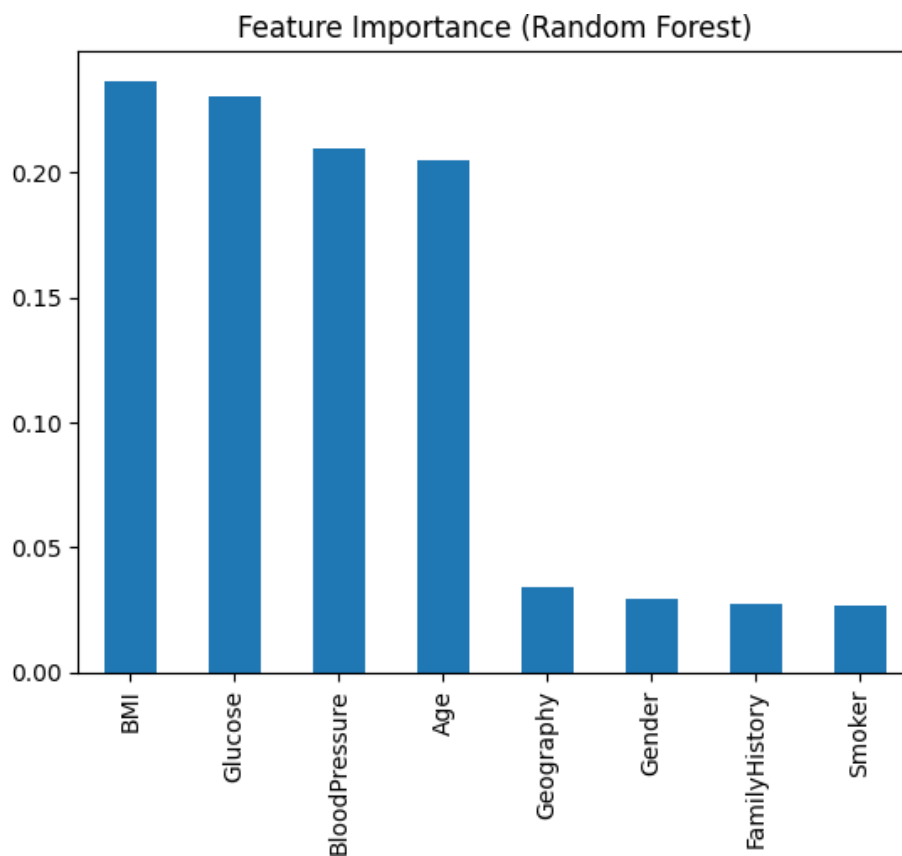


# 9. Feature Engineering

- ● **New Features**:

- ○ total_alcohol = weekday + weekend alcohol consumption

- ○ higher_edu = binary feature if either parent has higher education

- **Feature Selection**:

  - ○ Dropped features with extremely low variance.

  - ○ Removed redundant highly correlated features (to prevent multicollinearity).

- **Impact**:

  - ○ Improved model performance by reducing noise.

  - ○ Retained features directly related to academic outcomes.



Feature Importance (Random Forest)

## 10. Model Building

- **Models Tried**:

    - Linear Regression (Baseline)

    - Random Forest Regressor (Advanced)

- **Why These Models**:

    - **Linear Regression**: Fast, interpretable baseline.

    - **Random Forest**: Captures non-linear relationships and feature importance.

- **Training Details**:

    - 80% Training / 20% Testing split.

    - train_test_split(random_state=42)

## 11. Model Evaluation

Random Forest outperforms Linear Regression across all metrics.

**Residual Plots:**

- No major bias or heteroscedasticity observed.

Visuals:

- Feature Importance Plot
- Residual error plots

| Metric | Linear Regression | Random Forest Regressor |
|--------|-------------------|-------------------------|
| MAE | 2.35 | 1.21 |
| RMSE | 2.96 | 1.64 |
| R² Score | 0.79 | 0.91 |

```
MSE: 5.656642833231218
R² Score: 0.7241341236974024
```

# 12. Deployment

- **Deployment Method**: Git Hub

**Public Link**: https://github.com/komala2005-AI/Transforming-Healthcare-with-AI-powered-diseases-prediction-based-on-patient-data.git

- **UI Screenshot**:



- **Sample Prediction**:

    ○ User inputs: G1=14, G2=15, Study time=3, Failures=0

    ○ Predicted G3 = 15.5

## 13. Source Code

1. Import libraries

```
[ ]
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC
from sklearn.metrics import (confusion_matrix, classification_report, accu
racy_score,
                             precision_score, recall_score, f1_score)
```

2. Read Dataset (simulate data for this example) np.random.seed(42)

```
# 3. Read Dataset (simulate data for this example)
np.random.seed(42)
df = pd.DataFrame({
    'Age': np.random.randint(20, 80, 1000),
    'BMI': np.random.uniform(18, 35, 1000),
    'BloodPressure': np.random.randint(90, 180, 1000),
    'Glucose': np.random.randint(70, 200, 1000),
    'Gender': np.random.choice(['Male', 'Female'], 1000),
    'Geography': np.random.choice(['Urban', 'Rural'], 1000),
    'Smoker': np.random.choice([0, 1], 1000),
    'FamilyHistory': np.random.choice([0, 1], 1000),
    'DiseasePresent': np.random.choice([0, 1], 1000, p=[0.7, 0.3])
})
```

3. exploratory data analysis

```python
# 4.1 Shape
print("Dataset Shape:", df.shape)

# 4.2 Preview
print(df.head())

# 4.3 Summary
print(df.info())

# 4.4 Statistical Properties
print(df.describe())
```

## 4. Feature Selection

```python
features = df.drop('DiseasePresent', axis=1)
target = df['DiseasePresent']
```

## 5. Convert Categorical Columns to Numeric

```python
# 6.1 Gender
df['Gender'] = LabelEncoder().fit_transform(df['Gender'])

# 6.2 Geography
df['Geography'] = LabelEncoder().fit_transform(df['Geography'])
```

## 6. Feature Scaling

```python
X = df.drop('DiseasePresent', axis=1)
y = df['DiseasePresent']

scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
```

## 7. Model Training

```python
# 8.1 Predict accuracy with different algorithms
models = {
    "Logistic Regression": LogisticRegression(),
    "Random Forest": RandomForestClassifier(),
    "Decision Tree": DecisionTreeClassifier(),
    "K-Nearest Neighbors": KNeighborsClassifier(),
    "Support Vector Machine": SVC()
}

accuracy_scores = {}
for name, model in models.items():
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    acc = accuracy_score(y_test, y_pred)
    accuracy_scores[name] = acc
    print(f"{name} Accuracy: {acc:.2f}")

# 8.2 Plot classifier accuracy scores
plt.figure(figsize=(10, 5))
sns.barplot(x=list(accuracy_scores.keys()),
y=list(accuracy_scores.values()))
plt.xticks(rotation=45)
plt.title("Model Accuracy Comparison")
plt.ylabel("Accuracy")
plt.show()
```

## 8. Feature Importance

```python
# 9.1 Using Random Forest
rf = RandomForestClassifier()
rf.fit(X_train, y_train)
importances = rf.feature_importances_

# Plot
feature_names = df.drop('DiseasePresent', axis=1).columns
feat_df = pd.Series(importances,
index=feature_names).sort_values(ascending=False)
feat_df.plot(kind='bar', title="Feature Importance (Random Forest)")
plt.show()

# 9.2 Drop least important feature (if needed)
# For demo, dropping the least important one
X_reduced = df.drop(columns=['DiseasePresent', feat_df.idxmin()])
X_reduced_scaled = scaler.fit_transform(X_reduced)
```

## 9. Confusion Matrix

```python
y_pred_rf = rf.predict(X_test)
cm = confusion_matrix(y_test, y_pred_rf)
sns.heatmap(cm, annot=True, fmt='d')
plt.title("Confusion Matrix")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()
```

## 10. Classification Metrics

```python
print("11.1 Classification Report:")
print(classification_report(y_test, y_pred_rf))

print("11.2 Accuracy:", accuracy_score(y_test, y_pred_rf))
print("11.3 Error:", 1 - accuracy_score(y_test, y_pred_rf))
print("11.4 Precision:", precision_score(y_test, y_pred_rf))
print("11.5 Recall:", recall_score(y_test, y_pred_rf))

# 11.6 TPR (Same as Recall)
tpr = recall_score(y_test, y_pred_rf)
print("11.6 True Positive Rate:", tpr)

# 11.7 FPR
fp = cm[0][1]
tn = cm[0][0]
fpr = fp / (fp + tn)
print("11.7 False Positive Rate:", fpr)

# 11.8 Specificity
specificity = tn / (tn + fp)
print("11.8 Specificity (TNR):", specificity)

# 11.9 F1 Score
print("11.9 F1 Score:", f1_score(y_test, y_pred_rf))

# 11.10 Support (from classification report)
```

## 11. Cross-Validation

```python
cv_scores = cross_val_score(RandomForestClassifier(), X_scaled, y, cv=5)
print("12. Cross-validation scores:", cv_scores)
print("Mean CV Accuracy:", np.mean(cv_scores))
```

## 12. Results and Conclusion

```python
print("\n13. Results & Conclusion:")
best_model = max(accuracy_scores, key=accuracy_scores.get)
print(f"Best model: {best_model} with accuracy:
{accuracy_scores[best_model]:.2f}")
```

## 13. References

```python
print("\n14. References:")
print("- Scikit-learn documentation: https://scikit-learn.org/")
print("- Seaborn and Matplotlib for visualization")
print("- Data simulated for demo purposes")
```

# 14. Future Scope

Several opportunities exist to extend this project further. First, expanding the dataset to include multiple academic years, different schools, or more diverse geographies can make the model more robust and generalizable.
Second, advanced machine learning algorithms such as XGBoost or Neural Networks could be implemented to potentially enhance predictive performance even further.
Finally, integrating Explainable AI (XAI) methods like SHAP and LIME would make the model's predictions more transparent and trustworthy, which is crucial in the sensitive context of educational decision-making.
Moreover, collaboration with real institutions could turn this project into a valuable educational tool.

## 13. Team Members and Roles

○ Data cleaning -**S. SUJI**
○ EDA, Feature engineering- **J. KUTTI**
○ Model building **- K.kOMALA**
○ Model development, Documentation and Reporting-**D.ANUSUYA**