

DATA SCIENCE PROJECT 1 – DIABETES HEALTHCARE SYSTEM

Personalized Healthcare Recommendations Data Science Project

Project:	Personalized Healthcare Recommendations
Tools:	Jupyter Notebook
Technologies:	Machine Learning
Domain:	Data Analytics, Data Science
Difficulty:	Advanced
Dataset:	https://www.kaggle.com/uciml/pima-indians-diabetes-database
GitHub Link:	https://github.com/Vicjosh07/Predicting-Diabetes/

• ABOUT DATASET:

Predict Diabetes from Medical Records

Diabetes mellitus (DM), commonly referred to as diabetes, is a group of metabolic disorders in which there are high blood sugar levels over a prolonged period. Type 1 diabetes results from the pancreas's failure to produce enough insulin. Type 2 diabetes begins with insulin resistance, a condition in which cells fail to respond to insulin properly. As of 2015, an estimated 415 million people had diabetes worldwide, with type 2 diabetes making up about 90% of the cases. This represents 8.3% of the adult population.

Source: https://en.wikipedia.org/wiki/Diabetes_mellitus

The [Pima Indians Diabetes Database](#) can be used to train machine learning models to predict if a given patient has diabetes. This dataset contains measurements relating to Pregnancies, Glucose, BloodPressure, SkinThickness, Insulin, BMI, DiabetesPedigreeFunction, and Age.

Features: Includes Pregnancies, Glucose, Blood Pressure, SkinThickness, Insulin, BMI, DiabetesPedigreeFunction, and Age.

Label: Outcome (1 for diabetic, 0 for healthy).

A	B	C	D	E	F	G	H	I
Pregnancies	Glucose	BloodPress	SkinThickness	Insulin	BMI	DiabetesPe	Age	Outcome
6	148	72	35	0	33.6	0.627	50	1
1	85	66	29	0	26.6	0.351	31	0
8	183	64	0	0	23.3	0.672	32	1
1	89	66	23	94	28.1	0.167	21	0
0	137	40	35	168	43.1	2.288	33	1
5	116	74	0	0	25.6	0.201	30	0
3	78	50	32	88	31	0.248	26	1
10	115	0	0	0	35.3	0.134	29	0
2	197	70	45	543	30.5	0.158	53	1
8	125	96	0	0	0	0.232	54	1

PROJECT AIM: Diabetes Prediction and Recommendations System

Introduction;

Concept:

This project aims to develop a diabetes prediction system to classify individuals as diabetic or non-diabetic using machine learning models. It provides personalized recommendations for lifestyle and medical guidance based on key health metrics.

Use Cases:

Identification of Diabetic or Non-Diabetic patients.

Early identification of individuals at risk of diabetes.

Guidance for monitoring glucose levels and other related parameters.

Case Scenario:

A healthcare provider uses the tool to identify at-risk patients during regular health check-ups.

Problem Addressed:

Creating and Improving accuracy and interpretability of diabetes prediction while providing actionable insights.

Step-by-Step Description/ Implementation in my Program

1. Understanding the Problem

- The goal is to provide personalized healthcare recommendations to patients based on some of their health metrics, medical history, lifestyle, and other relevant factors.
- Use machine learning techniques to analyze patient data and generate actionable insights.

2. Dataset Preparation

- Data Sources: I sourced for different data from various platforms such as electronic health records (EHRs), kaggle.com, UCI repository, online patient surveys, and publicly available health datasets.
- **Features:** Includes Pregnancies, Glucose, Blood Pressure, SkinThickness, Insulin, BMI, Diabetes Pedigree Function, and Age.
- **Labels:** Outcome (1 for diabetic, 0 for healthy).

3. Initializing library tools and Exploratory Data Analysis (EDA)

- I imported the necessary libraries, tools and ML models to work with e.g (pandas, numpy, Decision tree)

o I Loaded and explored the dataset using descriptive statistics and visualization techniques.

o Visualizing Features:

Histograms and scatter plots were used to examine the relationships between features like Insulin, SkinThickness, and the target label.

Example finding: Individuals with lower insulin levels tend to have a higher risk of diabetes.

Dataset Exploration/Visualization (EDA)

```
[ ] # Calculate and visualize correlation matrix
correlation_matrix = df.corr()
plt.figure(figsize=(12, 8))
sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm", fmt=".2f")
plt.title("Feature Correlation Matrix")
plt.show()
```

```
# lets visualize the feature against the label, especially the insulin and skinthickness
# to determine if there is correlation between a particular outcome and the nulled insulin and skintype
# Loop through each feature except the label (Outcome)
for feature in df.columns[:-1]:
    def plotHistogram(values, label, feature, title):
        sns.set_style("whitegrid")

        # Define a custom palette for the outcomes
        custom_palette = [0: "blue", 1: "red"]

        # Using sns.histplot
        plt.figure(figsize=(10, 6))
        sns.histplot(
            data=values,
```

+ Code + Text

```
data=values,
x=feature,
hue=label,
kde=False,
palette=custom_palette,
bins=20,
alpha=0.5
)
plt.title(title)
plt.xlabel(feature)
plt.ylabel('Proportion')
plt.legend(labels=["diabetic", "healthy"], title=label)
plt.show()

# Call the function
plotHistogram(
    df,
    df.columns[-1],
    feature,
    f"{feature} vs {df.columns[-1]} (red=diabetic; blue=healthy)"
)
```

So i can generalize that the absence of insulin doesnt necessary correlate to diabetic thus our model wouldnt be biased

4. Data Cleaning and Preprocessing

- Replacing zero values in features like Glucose, Blood Pressure, and SkinThickness with NaN.
- Imputed missing values using the median strategy to retain statistical integrity.

5. Feature Engineering

- Dropped less relevant features (SkinThickness, BloodPressure) to optimize model performance.
- Correlation heatmap was used to identify relationships among features and with the label.

FEATURE ENGINEERING/SELECTION

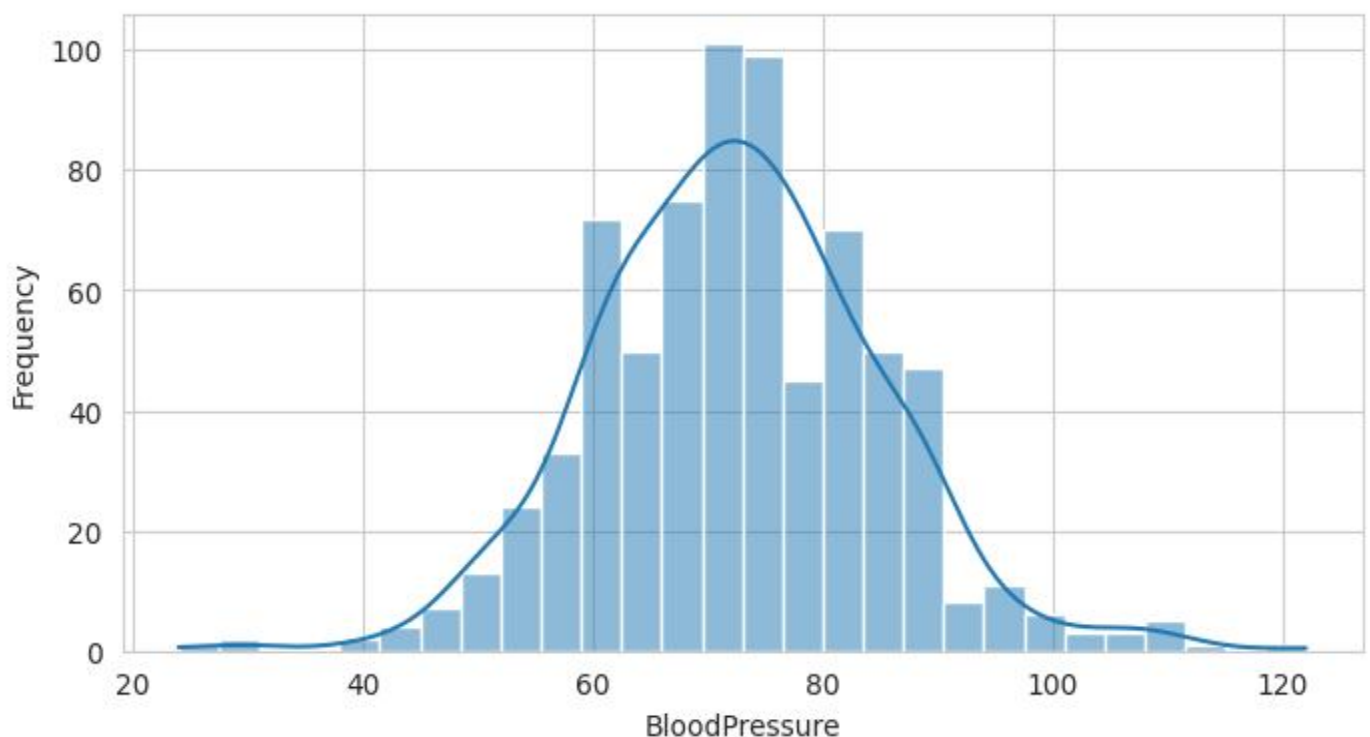
At this point i am going to split the dataset to training and testing. After this i will apply the median imputer because of the outlier and skewed nature to fill up the null value..

I will now use the imputer to fit(learn) and transform the xtrain...then apply this value to the xtest to prevent data leakage, overfitting or re-learning

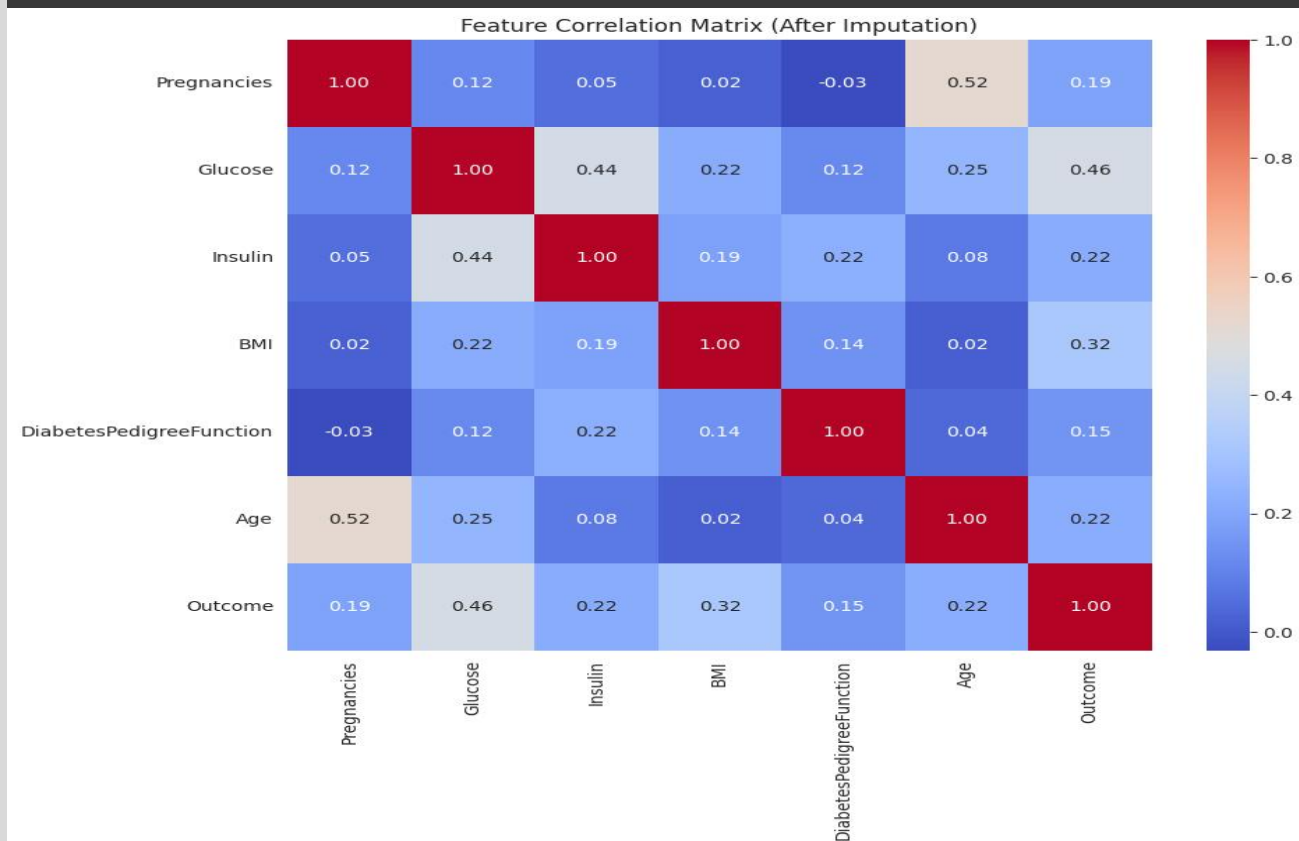
Also this imputation is not applied before splitting so our model would'nt learn bad data..we'll be test on our real-based dataset and also prevent overly optimistic metrics

```
[9] # Iterate over the specified columns
for col in df.columns[1:6]:
    plt.figure(figsize=(8, 4)) # Set the figure size
    sns.histplot(df[col], kde=True) # Plot histogram with KDE
    plt.title(f'Distribution of {col}') # Add a title
    plt.xlabel(col) # Label the x-axis
    plt.ylabel('Frequency') # Label the y-axis
    plt.show() # Display the plot
```

Distribution of BloodPressure



```
# Generating the heatmap
plt.figure(figsize=(10, 8))
correlation_matrix = train_data.corr()
sns.heatmap(correlation_matrix, annot=True, fmt=".2f", cmap="coolwarm")
plt.title("Feature Correlation Matrix (After Imputation)")
plt.show()
```



5. Advanced Data Processing (Data Splitting/Imputation)

- I splitted the dataset into the training and testing with the relevant features and the label as the outcome into 70% to 30%.
- After confirming the splitting the dataset and plotting to pinpoint the imputation technique to be used, I imputed the median values for the NaN occurences in the features to allow proper interpretability.

After the plot, i am going to use median strategy for handlig missing values because of the outliers in most features

```
[ ] #-----> i used this to test for the features to drop <-----
# X = df[df.columns[:-1]] # All columns except the last one
# y = df[df.columns[-1]] # The last column
# X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=1)
# imputer = SimpleImputer(missing_values=np.nan, strategy='median')
# X_train_imputed = imputer.fit_transform(X_train)
# X_test_imputed = imputer.transform(X_test)

[ ] # another aspect of the feature engineering is the selection of the most
# relevant features to be used
X = df.drop(['SkinThickness', 'BloodPressure', 'Outcome'], axis=1) # All columns except the last one
y = df[df.columns[-1]] # The last column
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=1)
imputer = SimpleImputer(missing_values=np.nan, strategy='median')
X_train_imputed = imputer.fit_transform(X_train)
X_test_imputed = imputer.transform(X_test)
```

6. Data Splitting and Scaling

- StandardScaler was also applied to normalize features, ensuring that models could train effectively without bias from magnitude differences.

```
+ Code + Text
# Scale the features (most models typically work better with scaled features)
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train_imputed)
X_test_scaled = scaler.transform(X_test_imputed)

/usr/local/lib/python3.11/dist-packages/sklearn/utils/validation.py:2739: UserWarning: X does not have valid feature names, but StandardScaler was fitted with fe
warnings.warn(
```

7. Model Training and Evaluation

- **Models Used:** I looped through and tested different models to pick the most optimal with different hyper-parameters; The following were Explored: Decision Tree, SVM, Logistic Regression, Random Forest, K-Nearest Neighbors, Gradient Boosting, Naive Bayes, Neural Network.

MODEL	ACCURACY (%)	PRECISION (AVG)	RECALL (AVG)	F1-SCORE (AVG)
Decision Tree	77%	0.76	0.72	0.73
SVM	83%	0.82	0.83	0.82
Logistic Regression	80%	0.79	0.80	0.79
Random Forest	82%	0.81	0.80	0.80
K-Nearest Neighbors	77%	0.75	0.75	0.75
Gradient Boosting	82%	0.81	0.80	0.80

Naive Bayes	77%	0.76	0.73	0.74
Neural Network	~83%	0.81	0.80	0.80

+ Code + Text		File Edit View Insert Runtime Tools Help				
Training Decision Tree...		+ Code + Text				
Decision Tree Model Performance:		Training Logistic Regression...				
precision recall f1-score support		Logistic Regression Model Performance:				
0 0.77 0.89 0.83 146		precision recall f1-score support				
1 0.75 0.55 0.64 85		0 0.88 0.79 0.83 146				
accuracy 0.77 231		1 0.70 0.81 0.75 85				
macro avg 0.76 0.72 0.73 231		accuracy 0.80 231				
weighted avg 0.76 0.77 0.76 231		macro avg 0.79 0.80 0.79 231				
Training SVM...		weighted avg 0.81 0.80 0.80 231				
SVM Model Performance:		Training Random Forest Classifier...				
precision recall f1-score support		Random Forest Classifier Model Performance:				
0 0.90 0.83 0.86 146		precision recall f1-score support				
1 0.74 0.84 0.78 85		0 0.85 0.87 0.86 146				
accuracy 0.83 231		1 0.77 0.73 0.75 85				
macro avg 0.82 0.83 0.82 231		accuracy 0.82 231				
weighted avg 0.84 0.83 0.83 231		macro avg 0.81 0.80 0.80 231				
		weighted avg 0.82 0.82 0.82 231				

diabetes_project.ipynb

File Edit View Insert Runtime Tools Help

+ Code + Text

MODELS

```
# Models to test with hyperparameters
models = {}

"Decision Tree": DecisionTreeClassifier(random_state=1, max_depth=5, min_samples_leaf=5),
"SVM": SVC(random_state=1, class_weight='balanced'),
"Logistic Regression": LogisticRegression(random_state=1, class_weight='balanced'),
"Random Forest Classifier": RandomForestClassifier(random_state=1, min_samples_split=10, class_weight='balanced'),
"K-Nearest Neighbors": KNeighborsClassifier(n_neighbors=7),
"Gradient Boosting": GradientBoostingClassifier(random_state=1),
"Naive Bayes": GaussianNB()

# Train and evaluate each model
for model_name, model in models.items():
    print(f"Training {model_name}...")
    # Train for Medicine prediction
    model.fit(X_train_scaled, y_train)
    y_pred = model.predict(X_test_scaled)
    print(f"{model_name} Model Performance:")
    print(classification_report(y_test, y_pred))
```

Show hidden output

Connecting to Python 3 Google Compute Engine backend

8. Performance Metrics:

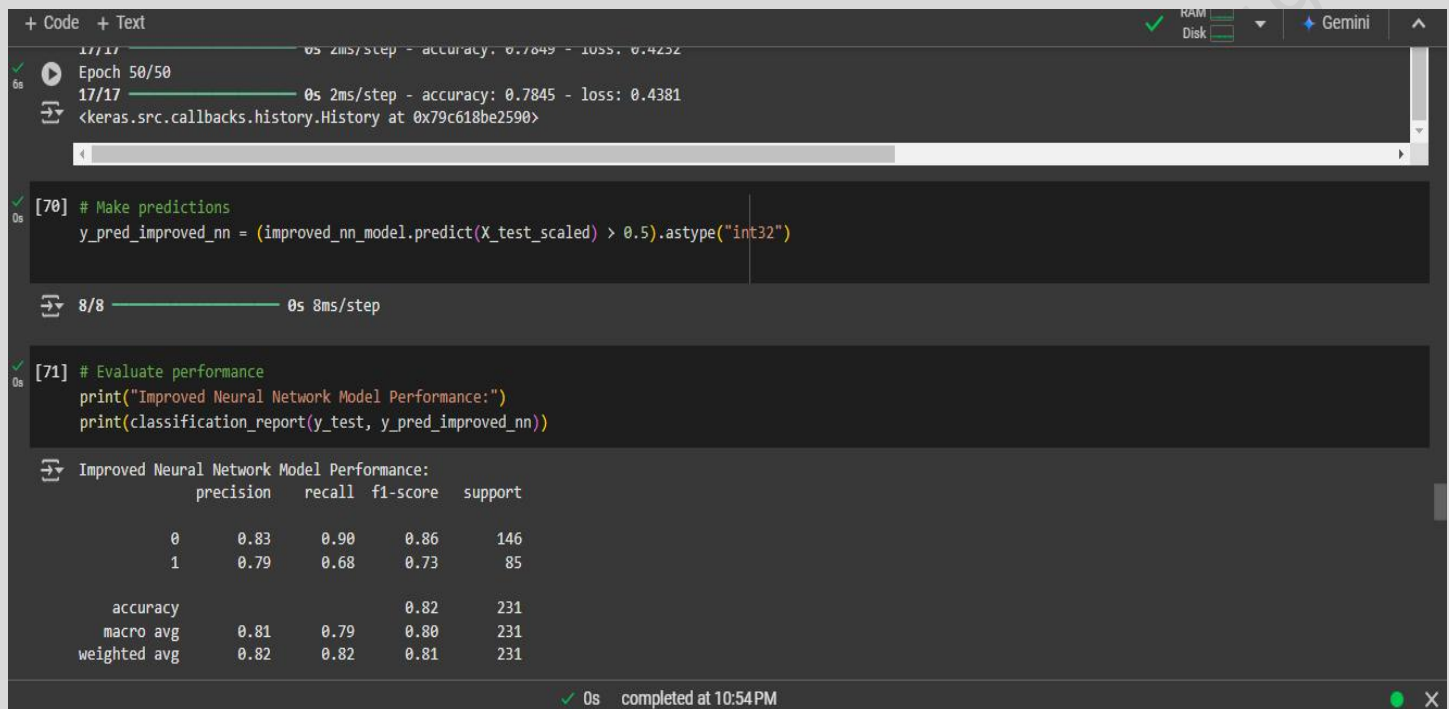
Classification report for each model, including precision, recall, F1-score, and accuracy.

Example result: SVM achieved the highest accuracy of 83%. Which when I applied to real world scenario from diverse dataset did quite well!!

Advanced Model: Neural Network

A neural network was implemented using TensorFlow/Keras with three dense layers and dropout for regularization.

Evaluation: Achieved an accuracy of ~83% after training for 50 epochs.



```
+ Code + Text
Epoch 50/50
17/17 — 0s 2ms/step - accuracy: 0.7849 - loss: 0.4232
17/17 — 0s 2ms/step - accuracy: 0.7845 - loss: 0.4381
<keras.src.callbacks.history.History at 0x79c618be2590>

[70] # Make predictions
y_pred_improved_nn = (improved_nn_model.predict(X_test_scaled) > 0.5).astype("int32")
8/8 — 0s 8ms/step

[71] # Evaluate performance
print("Improved Neural Network Model Performance:")
print(classification_report(y_test, y_pred_improved_nn))

Improved Neural Network Model Performance:
              precision    recall  f1-score   support

     0       0.83         0.90      0.86         146
     1       0.79         0.68      0.73          85

 accuracy          0.81
 macro avg         0.81         0.79      0.80         231
weighted avg         0.82         0.82      0.81         231

0s completed at 10:54 PM
```

Challenges Encountered

Honestly, I encountered many challenges while working on this project. From the start, I had to search for the most appropriate dataset to work with. Initially, I started off with a medicine recommendation system. After many setbacks with, with processing the dataset and models, I moved to search for another dataset when I came across this - diabetes. Below are some of the major setbacks I also encountered while building the predicting diabetes and recommendation system:

- Firstly, I had problem getting documentations to work with from the various official website of each libraries and tool to organize them into my code. Eg Keras for NN model building.

- Handling missing values in critical features like Insulin and Glucose. I was in a fix in handling the missing values. What type of imputation technique? Where should I put in? before or after scaling? Before or after splitting the dataset?
- Determining optimal hyperparameters for diverse models. I wanted to maximize the accuracy, f1 score and other metrics for my models. To this effect, I tried different hyperparameters before getting the best among the range.
- Managing imbalanced data, as fewer individuals were classified as diabetic.
- Another huge issue I encountered was during the building of the recommendation algorithm. The complexities involved in finding the mapping to each health patient metric inputted.
- Also, during the algorithm development, I forgot to inverse_transform the inputted scaled data to generate the correct recommendation per diabetic or non-diabetic

Recommendations System

This algorithm is a complex function that allow for a trustable healthcare prediction whether diabetic or non-diabetic. This system also generates actionable insights based on predictions.

Real-Life Scenario – An example: I want to predict a patient that has had 1 pregnancies, a glucose level of 145, Insulin level of 110, BMI as 29, Diabetes pedigree function(a metric to record the extent to which a person lineage is prone to diabetes) as 0.7 and lastly the age as 45.,

My algorithm takes this feed and predicts this patient as DIABETIC. Also, it generates actionable recommendation “Your glucose levels are under moderate control, but ongoing monitoring and a balanced diet are still crucial. Consider incorporating more physical activity and regularly consulting with your healthcare provider.”

"The recommendation system works by analyzing key health metrics and mapping them to predefined health categories. The system follows these steps:"

1. Predicts whether the patient is diabetic or not using the trained model.
2. Extracts important health factors (e.g., glucose, insulin, BMI, age, and diabetes pedigree function).
3. Matches these values to predefined thresholds to generate personalized recommendations.

Conclusion

The project successfully developed a robust diabetes prediction system, achieving high accuracy and providing meaningful recommendations.

Future improvements:

Though this program is quite a hit in many ways, due to time constraint and other limiting factors there are many provide improvement. Below are the future improvement and integrations I look to implement:

- I look forward to getting the similar dataset to populate the present one and conversely improve model's data interpretability and improve the model accuracy so as to boost the reliability in the healthcare industry.
- Simplifying the health metrics on a deployed system as a web or mobile application to serve anyone one both the lay-man and the health care personnel

Olaleye Joshua - UMIP271961