



DIABETES PREDICTION

Project submission

Presented by: Group 6

MACHINE LEARNING TUNING

3 CRITICAL ASPECTS TO SOLVE



MODEL INTERPRETABILITY AND EMPERICAL TUNING

BLOOD PRESSURE IMPUTATION - FEMALE ONLY

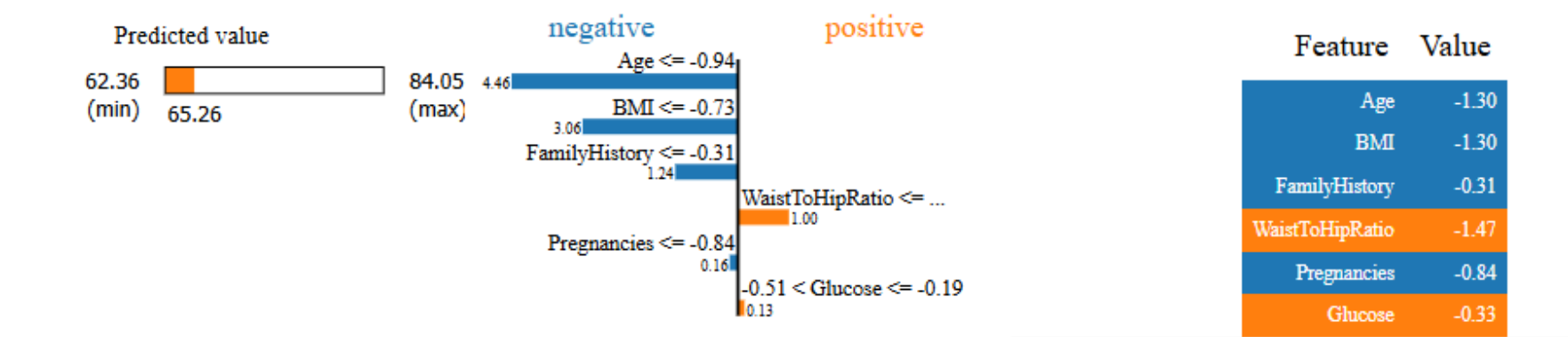
VISUALIZE RELATIONSHIPS

Features and their description

1. Pregnancies - Number of times the person conceived
2. Fasting Glucose - 8 hrs of fasting
3. Age - In years -> 21 and up
4. BMI - weight in kg / Height in meters^2
5. Family History - 0 for no Family history in diabetes, 1 if there is family history
6. Waist to hip ratio = Waist / Hips in cm
7. Bloodpressure - Diastolic blood pressure (mm Hg)
8. Outcome - 0 for not diabetic and 1 if diabetic (Target Variable)

BloodPressure	1.000000	LIME Feature Importance Matrix for B	
Age	0.281644		
BMI	0.227042		
WaistToHipRatio	0.205303		
Pregnancies	0.186144		
Glucose	0.111212		
FamilyHistory	0.071707		
Outcome	0.043079		
Name: BloodPressure, dtype: float64			

	Feature	LIME Importance
5	Age	2.464679
2	BMI	1.783704
4	FamilyHistory	1.215964
3	WaistToHipRatio	0.608275
1	Glucose	0.287391
0	Pregnancies	0.191928



Feature	Correlation Rank	LIME Importance Rank
Age	1	1
BMI	2	2
WaistToHipRatio	3	4
FamilyHistory	6	3
Glucose	5	5
Pregnancies	4	6

```
feature_sets = {  
    "Set 1 (Age,BMI,WaistToHipRatio)": ['Age', 'BMI', 'WaistToHipRatio'],  
    "Set 2 (Age Only)": ['Age'],  
    "Set 3 (All Features Except Outcome)": ['Age', 'BMI', 'WaistToHipRatio', 'Pregnancies', 'Glucose', 'FamilyHistory']  
}
```

```
# Step 1: Perform 80/20 split  
BP_X = data[features]  
BP_y = data['BloodPressure']  
BP_X_train_80, BP_X_test_20, BP_y_train_80, BP_y_test_20 = train_test_split(BP_X, BP_y, test_size=0.2, random_state=42)  
  
# Step 2: Perform 60/20 split on 80% training data  
BP_X_train_60, BP_X_val_20, BP_y_train_60, BP_y_val_20 = train_test_split(BP_X_train_80, BP_y_train_80, test_size=0.25, random_state=42)  
  
# Step 3: Scale data  
scaler = StandardScaler()  
BP_X_train_60_scaled = scaler.fit_transform(BP_X_train_60)  
BP_X_val_20_scaled = scaler.transform(BP_X_val_20)  
BP_X_test_20_scaled = scaler.transform(BP_X_test_20)  
  
# Step 4: Train model on 60% training data (only rows with non-missing y)  
BP_train_complete = BP_X_train_60_scaled[BP_y_train_60.notna()]  
BP_y_train_complete = BP_y_train_60[BP_y_train_60.notna()]
```

Missing values successfully filled!

Chosen Feature Set for Imputation of BloodPressure:
['Age', 'BMI', 'WaistToHipRatio']

Feature Set	R-squared (Validation)	MSE (Validation)	R-squared (Test)	MSE (Test)
Set 1 (Age,BMI,WaistToHipRatio)	0.094433	100.519790	0.086290	102.787343
Set 2 (Age Only)	0.073139	102.883468	0.066154	105.052557
Set 3 (All Features Except Outcome)	0.094287	100.535924	0.088351	102.555416

MODEL INTERPRETABILITY AND EMPERICAL TUNING

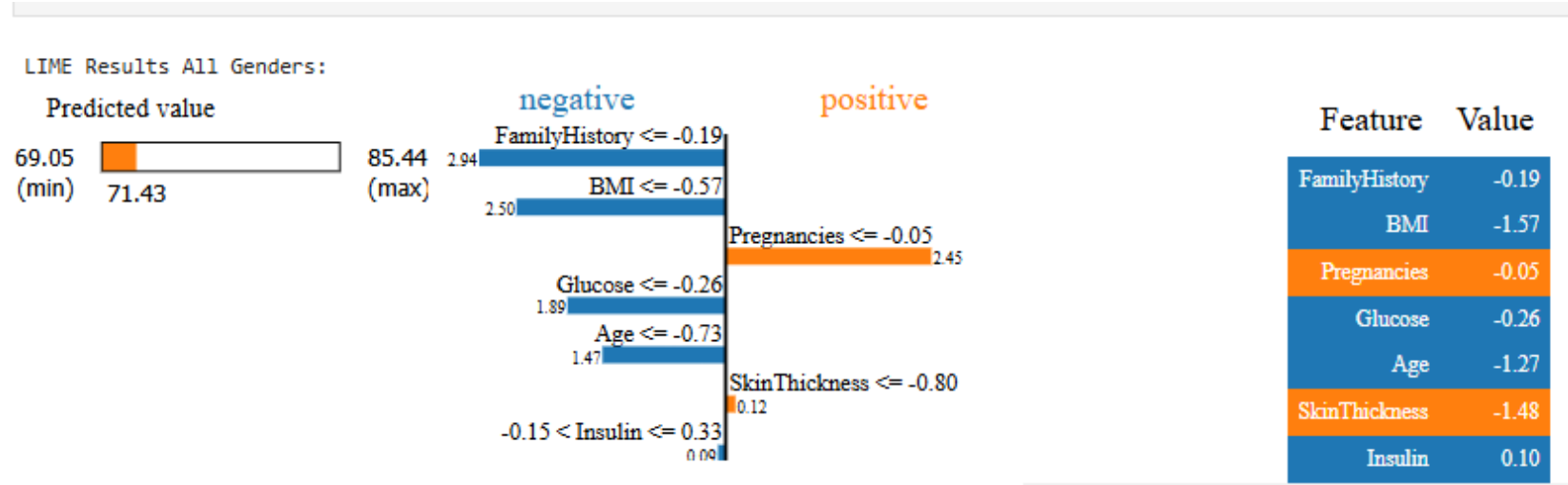
BLOOD PRESSURE IMPUTATION - ALL GENDER

VISUALIZE RELATIONSHIPS

Features and Their Descriptions

- **Pregnancies** - Number of times the person conceived.
- **Glucose** - 2 hours of fasting.
- **BloodPressure** - Diastolic blood pressure (mm Hg).
- **SkinThickness** - Measurement of Triceps Skinfold Thickness (mm):
- **Insulin** - Fasting insulin levels (μU/mL). Normal range: 2–25 μU/mL.
- **BMI** - Body Mass Index (weight in kg / height in meters²):
- **Family History** - 0 for no family history in diabetes, 1 if there is family history.
- **Age** - In years, ranging from 1 to 85.
- **Outcome** - 0 for not diabetic and 1 if diabetic (Target Variable).

BloodPressure	1.000000	LIME Results All Genders:	
Outcome	0.335337	Feature	LIME Importance
BMI	0.283767	5 FamilyHistory	2.865276
Glucose	0.223614	6 Pregnancies	2.272630
Insulin	0.204414	1 Glucose	1.854726
Age	0.201559	0 BMI	1.396009
SkinThickness	0.100627	3 Age	0.696293
FamilyHistory	0.097086	2 Insulin	0.460933
Pregnancies	-0.022024	4 SkinThickness	0.063152
Name: BloodPressure, dtype: float64			

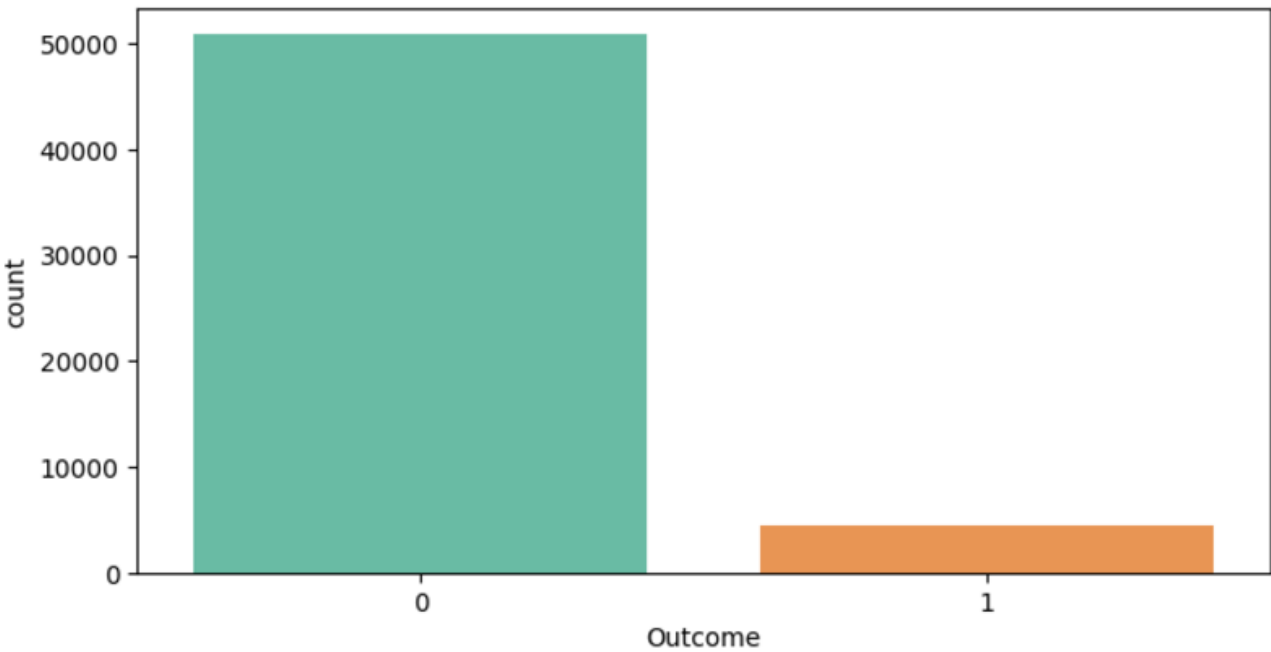


ORIGINAL FEATURES - FEMALE ONLY

Features and their description

- 1. Pregnancies - Number of times the person conceived
- 2. Fasting Glucose - 8 hrs of fasting
- 3. Age - In years -> 21 and up
- 4. BMI - weight in kg / Height in meters^2
- 5. Family History - 0 for no Family history in diabetes, 1 if there is family history
- 6. Waist to hip ratio = Waist / Hips in cm
- 7. Bloodpressure - Diastolic blood pressure (mm Hg)
- 8. Outcome - 0 for not diabetic and 1 if diabetic (Target Variable)

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 55299 entries, 0 to 55298  
Data columns (total 8 columns):  
#   Column          Non-Null Count  Dtype  
---  -  
0   Pregnancies     55299 non-null  int64  
1   Glucose         55299 non-null  int64  
2   BloodPressure   55299 non-null  float64  
3   BMI             55299 non-null  float64  
4   WaistToHipRatio 55299 non-null  float64  
5   FamilyHistory   55299 non-null  int64  
6   Age             55299 non-null  int64  
7   Outcome         55299 non-null  int64  
dtypes: float64(3), int64(5)  
memory usage: 3.4 MB
```



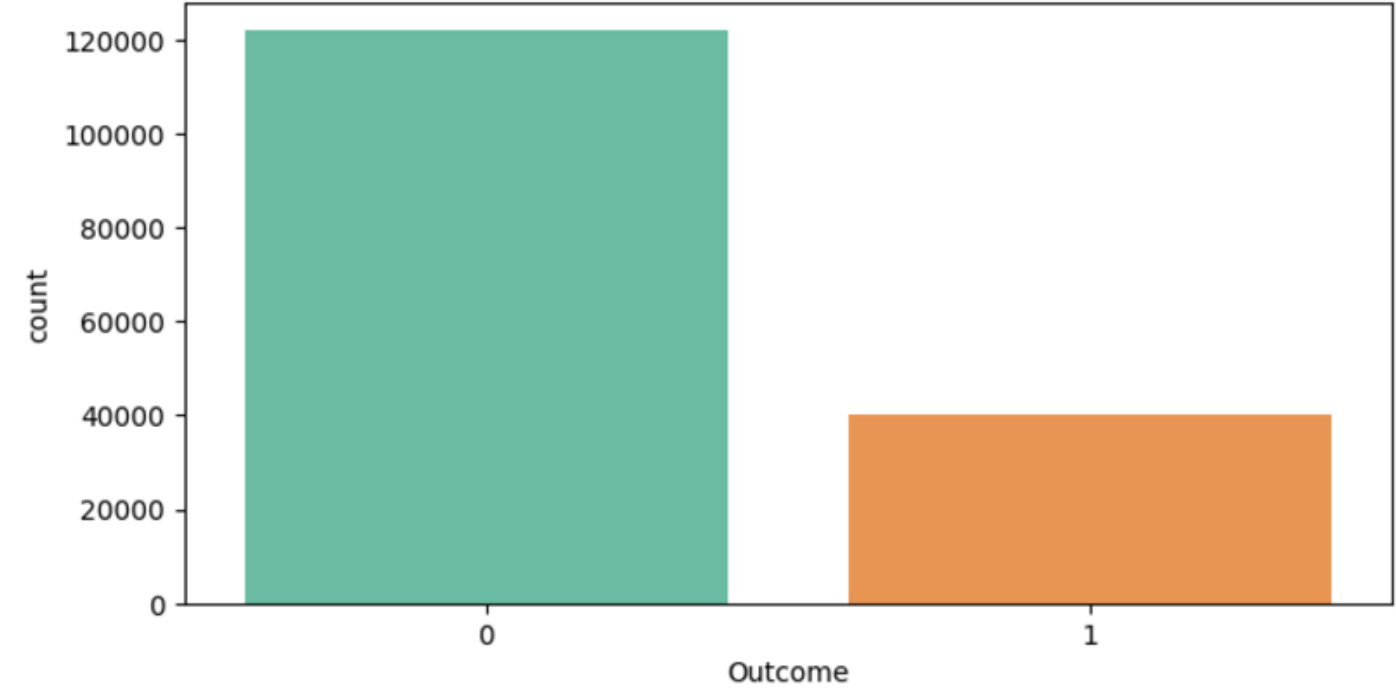
```
Class distribution after undersampling: 0    4463  
1    4463  
Name: Outcome, dtype: int64  
  
Training Accuracy of Logistic Regression: 95.30812324929971  
Accuracy (Test) score of Logistic Regression: 95.07278835386337  
Accuracy score of Logistic Regression: 95.07278835386337  
AUC: 0.9849910512569402  
  
Training Accuracy of SVM: 95.92436974789916  
Accuracy (Test) score of SVM: 95.18477043673013  
Accuracy score of SVM: 95.18477043673013  
AUC: 0.9837198501721406  
  
Training Accuracy of Random Forest: 99.9859943977591  
Accuracy (Test) score of Random Forest: 95.91265397536394  
Accuracy score of Random Forest: 95.91265397536394  
AUC: 0.9899184419828658
```

ALL GENDER RESULTS

Features and Their Descriptions

- **Pregnancies** - Number of times the person conceived.
- **Glucose** - 2 hours of fasting.
- **BloodPressure** - Diastolic blood pressure (mm Hg).
- **SkinThickness** - Measurement of Triceps Skinfold Thickness (mm):
- **Insulin** - Fasting insulin levels (μU/mL). Normal range: 2–25 μU/mL.
- **BMI** - Body Mass Index (weight in kg / height in meters²):
- **Family History** - 0 for no family history in diabetes, 1 if there is family history.
- **Age** - In years, ranging from 1 to 85.
- **Outcome** - 0 for not diabetic and 1 if diabetic (Target Variable).

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 162248 entries, 0 to 162247
Data columns (total 9 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Age             162248 non-null int64
1   Pregnancies     162248 non-null int64
2   Glucose         162248 non-null float64
3   BloodPressure   136239 non-null float64
4   SkinThickness   162248 non-null float64
5   Insulin         162248 non-null float64
6   BMI             162248 non-null float64
7   FamilyHistory   162248 non-null int64
8   Outcome         162248 non-null int64
dtypes: float64(5), int64(4)
memory usage: 11.1 MB
```



Class_weight='balanced'

Model: Logistic Regression
Training Accuracy: 94.64%
Test Accuracy: 94.51%
Accuracy Score: 94.51%
AUC: 0.98

Model: Support Vector Machine (SVM)
Training Accuracy: 95.14%
Test Accuracy: 95.00%
Accuracy Score: 95.00%
AUC: 0.98

Model: Random Forest
Training Accuracy: 100.00%
Test Accuracy: 97.61%
Accuracy Score: 97.61%
AUC: 0.99

SUBSEQUENT MODEL BUILDING

RISK LEVEL CLASSIFICATION:

High Risk:
Glucose > 180 or BloodPressure > 140

Moderate Risk:
BMI > 30

Low Risk:
All other cases

Classification applied to Outcome == 1

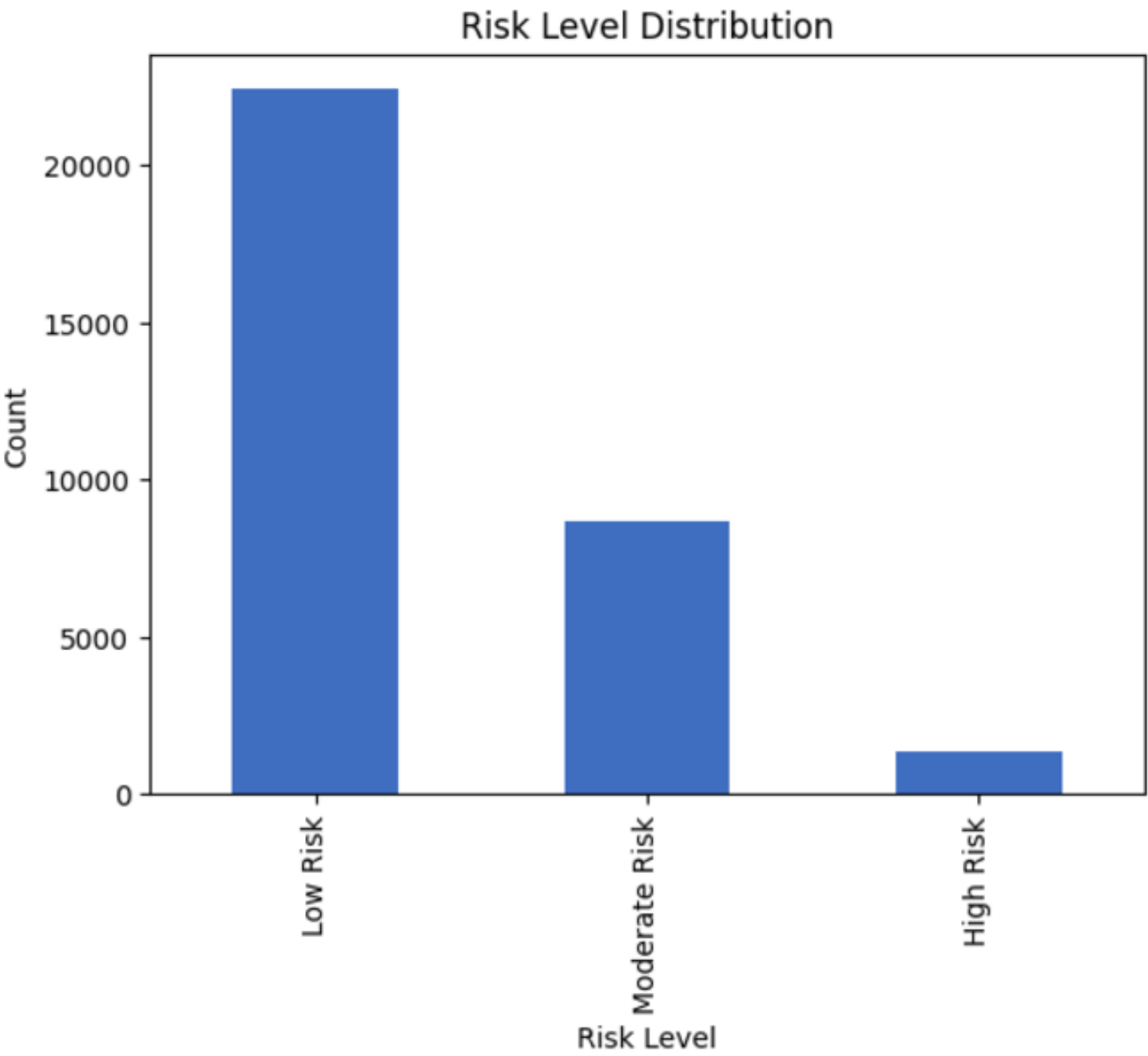
```
Model: Random Forest
Training Accuracy: 100.00%
Test Accuracy: 97.61%
Accuracy Score: 97.61%
AUC: 0.99

Confusion Matrix:
[[24143  121]
 [  656 7530]]

Definitions:
True Positives (TP): 7530 - Correctly predicted positive cases.
False Negatives (FN): 656 - Positive cases incorrectly predicted as negative.
False Positives (FP): 121 - Negative cases incorrectly predicted as positive.
True Negatives (TN): 24143 - Correctly predicted negative cases.
```

Classification Report:

	precision	recall	f1-score	support
0	0.97	1.00	0.98	24264
1	0.98	0.92	0.95	8186
accuracy			0.98	32450
macro avg	0.98	0.96	0.97	32450
weighted avg	0.98	0.98	0.98	32450



GIT REPOSITORY

CodeIssuesPull requestsActionsProjectsWikiSecurityInsights

diabetes-predictionPublic

Unwatch2Fork0Star0

main3 Branches0 Tags

Go to file

Add fileCode

sunny-aiden flask855c5cf · 4 days ago58 Commits

.github/workflows	cicd	4 days ago
app	flask	4 days ago
config	connection	4 days ago
reaserch	connection	4 days ago
.env	setup	4 days ago
.gitignore	setup	4 days ago
Dockerfile	docker	4 days ago
README.md	setup	4 days ago
docker-compose.yml	docker	4 days ago
init_setup.sh	start files	4 days ago
requirement.txt	setup	4 days ago
run.py	start files	4 days ago
test_app.py	start files	4 days ago

README

About

No description, website, or topics provided.

ReadmeActivity0 stars2 watching0 forksReport repository

Releases

No releases publishedCreate a new release

Packages

No packages publishedPublish your first package

Contributors8

Languages

Jupyter Notebook49.7%Python39.3%

HTML7.0%CSS2.0%

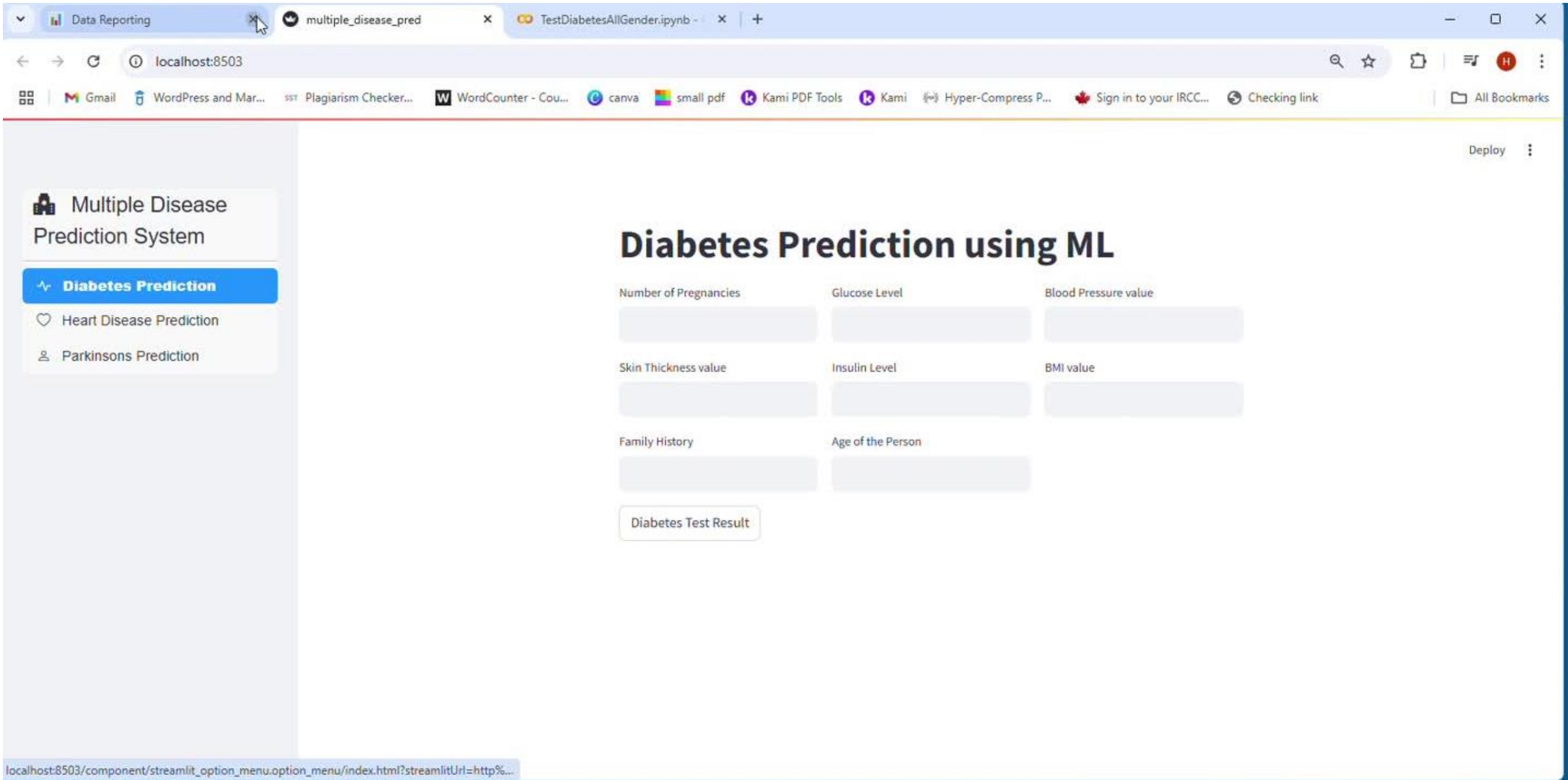
https://github.com/krunalpatel355/diabetes-prediction/graphs/contributors

A photograph showing a female doctor in a white lab coat and stethoscope examining the arm of a young girl. A male doctor is also present, looking on. They are in a clinical or hospital setting.

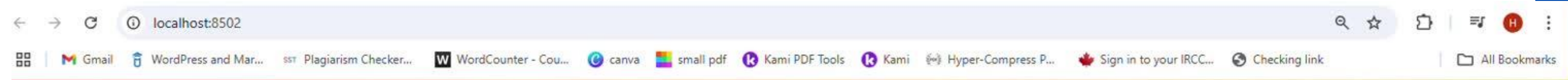
PITCH



UI APPLICATION



UI APPLICATION - DATA REPORTING



Data Reporting

Select a file

diabetes.csv

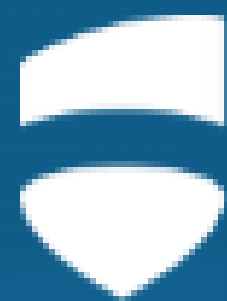
heart.csv

parkinsons.csv



THANK YOU

FROM: GROUP 6



Lambton
College