### **Distributed and Scalable Data Engineering Final project**

**Group4-Phoenix Attackers** 

# **TECHNICAL REPORT**



### **CONTENTS**

Project Title	2
Executive Summary	
Highlights of Project	4
Submitted on:	4
Abstract	5
Methodology	
Problems or Challenges	
Tools Used	
Exploratory Data Analysis	9
Analysis of Relationaship Between Variables	
Data Deployment and Data Modeling	
Model Evaluation	
Results Section	
Building a Flask Web App	
Conclusion	
Contributions/Poforonces	



# **The Early Prediction of Diabetes**

ARDIANA SULA
PROFESSOR, ACADEMICADVISOR

# **GROUP4-PHOENIX ATTACKERS**

NAME	E-MAIL ID
RAHUL MUVVALA	rmuvv1@unh.newhaven.edu
KALPANA CHAMALA	kcham5@unh.newhaven.edu
RISHITHA VENKANNAGARI	rvenk3@unh.newhaven.edu
MASTANVALI SHAIK	mshai32@unh.newhaven.edu
FARHAN SHAIK	fshai2@unh.newhaven.edu



### **Executive Summary**

Diabetes is one of the deadliest diseases in the world which is rapidly increasing and the cause of which vary depending on the genetic makeup family, history, ethnicity, health etc. It is an illness caused because of high glucose level in the body. Diabetes, untreated, may cause some major issues in like heart related problems, kidney problem, blood pressure, eye damage, it can also affect other organs of the body. Good thing is diabetes can be controlled if it is predicted earlier. So, the aim of this project is to design and implement diabetes prediction using machine learning methods and performance analysis of those methods. This helps a lot in saving the valuable time and money and is also a contribution for developing the technology in medical history.



Title of Project
------------------

The Early Prediction of Diabetes

## **Highlights of Project:**

~	various features such as				
	٥	Pregnancies	☐ Age		
		Insulin Level	□ вмі		
>	We have also predicted the predict the Diabetes.	e best model among multi	ple Machine Learning Algorithms, to		

### **Submitted on:**

12 Dec 2022.

#### **Abstract**

The aim of our project is to design and implement diabetes prediction using machine learning methods and performance analysis of those methods. The techniques we are using are support vector machine-nearest Neighbor, decision Tree, logistic regression, random forest. Applying machine learning techniques to analyze the performance of these methods will help to find accuracy of the model. It will help to figure out the responsible/important feature which play a major role in prediction. Our result shows that random forest achieved higher accuracy compared to other machine learning techniques. Diabetes can be controlled if it is predicted earlier. Machine learning techniques will provide better result for prediction by constructing models from datasets collected from patients.

The main goal of the project is to decrease the diabetic related issues faced by the public. As diabetics are very common these days that most of the people are suffering from, the early prediction can be very helpful.

**GitHub Link:** https://github.com/You-Can-Do/DSCI-6007-02-Phonix-Attackers

### **Methodology**

We imported the data sets into Jupyter notebook and performed the analysis. In this project, the objective is to predict whether the person has Diabetes or not based on various features such as:

- 1. Pregnancies
- 2. Insulin Level
- 3. Age
- 4. BMI

In this project we have used CRISP-DM methodology.

### **Business Understanding**

When it comes to health issues, there are many deadly diseases which effects a lot of people due to their inappropriate health habits and maintenance. Many people across the world are suffering a lot even some are losing their lives which is sad to know. There are many reasons for this like the feeling of getting their problems exposed, no proper guidance and knowledge, not enough money to consult a doctor and many more. So, in this project we developed a medical based application which helps at least a few in their medical life.

# The Problem challenges



Problem
Cause of
diabetes vary
depending on
the genetic
makeup family,
history, ethnicity,
health etc.,



Problem
Diabetes is one of the deadliest diseases in the world which is rapidly increasing.



Problem
It is an illness
caused because
of high glucose;
level in a human
body.



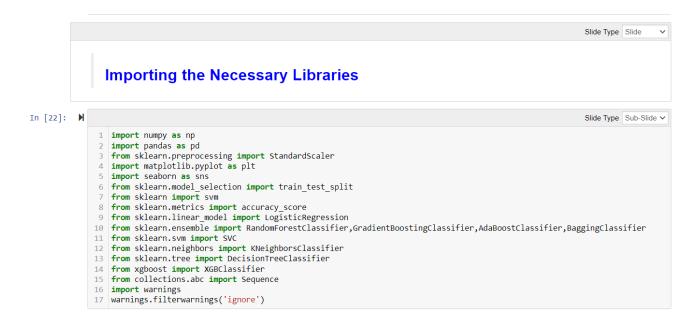




## **TOOLS USED**



#### **IMPORTING LIBRARIES**

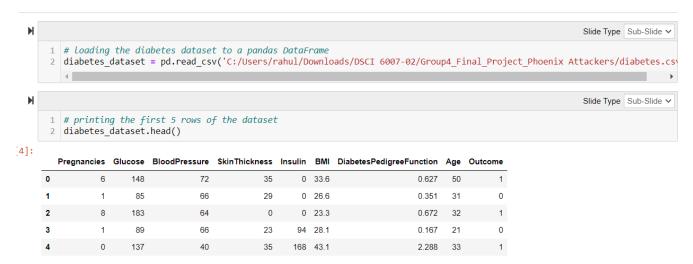


#### **EXPLORATORY DATA ANALYSIS**

1 UnderStanding the Da 2 Clean The Da		
2 Clean The Da	No	Steps
	1	UnderStanding the Data
3 Analysis of Relationship between variable	2	Clean The Data
	3	Analysis of Relationship between variables

#### 1. UNDERSTANDING THE DATA

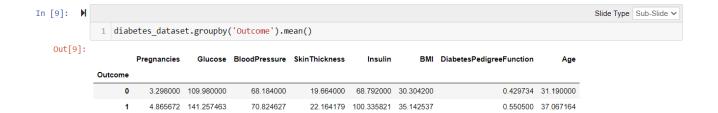
The data set that has used in this project has taken from the Kaggle. "This dataset is originally from the National Institute of Diabetes and Digestive and Kidney Diseases. The objective of the dataset is to diagnostically predict whether a patient has diabetes or not, based on certain diagnostic measurements included in the dataset. Several constraints were placed on the selection of these instances from a larger database. All patients here are females at least 21 years old of Pima Indian heritage." and used a simple random forest classifier.



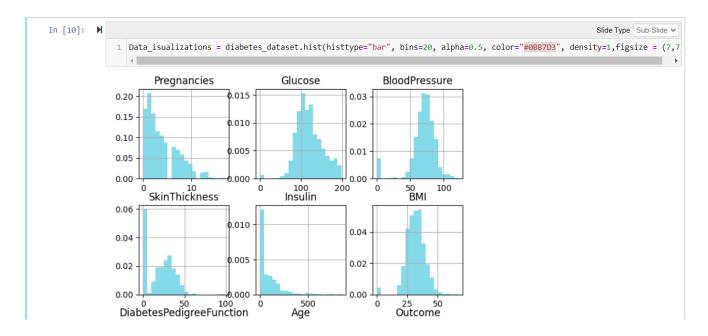
		Slide Type Sub-Slide V
	# number of rows and Columns in this dataset diabetes_dataset.shape	
(76	68, 9)	

#### 2. CLEANING THE DATA

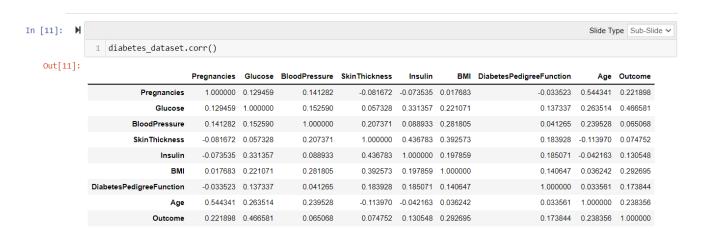
```
Slide Type Sub-Slide ✓
      1 diabetes_dataset.info()
     <class 'pandas.core.frame.DataFrame'>
     RangeIndex: 768 entries, 0 to 767
     Data columns (total 9 columns):
      # Column
                                        Non-Null Count Dtype
                                        768 non-null
                                                          int64
          Pregnancies
                                        768 non-null
                                                          int64
          Glucose
      1
      2
          BloodPressure
                                        768 non-null
                                                          int64
      3
          SkinThickness
                                        768 non-null
                                                          int64
      4
          Insulin
                                        768 non-null
                                                           int64
                                                           float64
          BMI
                                        768 non-null
      6
          DiabetesPedigreeFunction
                                        768 non-null
                                                           float64
                                        768 non-null
                                                          int64
          Age
      8
          Outcome
                                        768 non-null
                                                          int64
     dtypes: float64(2), int64(7)
     memory usage: 54.1 KB
In [7]: ▶
                                                                                                                              Slide Type Sub-Slide V
              diabetes_dataset.isnull().sum(axis=0)
   Out[7]: Pregnancies
            Glucose
             BloodPressure
                                          0
             SkinThickness
             Insulin
                                          0
            BMI
            DiabetesPedigreeFunction
            Age
            Outcome
            dtype: int64
                                                                                                                              Slide Type Sub-Slide ✓
         There is no empty data, so nothing is there to clean
In [8]: ▶
                                                                                                                                Slide Type Sub-Slide ✓
               1 # getting the statistical measures of the data
               2 diabetes_dataset.describe()
    Out[8]:
                    Pregnancies
                                  Glucose BloodPressure SkinThickness
                                                                         Insulin
                                                                                      BMI DiabetesPedigreeFunction
                                                                                                                              Outcome
                                                                                                                            768.000000
                     768.000000 768.000000
                                             768.000000
                                                           768.000000 768.000000 768.000000
                                                                                                       768.000000 768.000000
              count
              mean
                       3.845052 120.894531
                                              69.105469
                                                            20.536458
                                                                      79.799479
                                                                                 31.992578
                                                                                                         0.471876
                                                                                                                  33.240885
                                                                                                                              0.348958
                       3.369578 31.972618
                                              19.355807
                                                                                                                              0.476951
                                                            15.952218 115.244002
                                                                                  7 884160
                                                                                                         0.331329
                                                                                                                  11.760232
                std
                       0.000000
                                 0.000000
                                               0.000000
                                                             0.000000
                                                                       0.000000
                                                                                  0.000000
                                                                                                         0.078000
                                                                                                                  21.000000
                                                                                                                              0.000000
               min
               25%
                       1.000000
                                99.000000
                                              62.000000
                                                             0.000000
                                                                       0.000000
                                                                                27.300000
                                                                                                         0.243750
                                                                                                                  24.000000
                                                                                                                              0.000000
                       3.000000 117.000000
                                              72.000000
                                                                      30.500000
                                                                                                                  29.000000
                                                                                                                              0.000000
                                                            23.000000
                                                                                                         0.372500
                                                                                36.600000
               75%
                       6.000000 140.250000
                                              80.000000
                                                            32.000000 127.250000
                                                                                                         0.626250 41.000000
                                                                                                                              1.000000
                      17.000000 199.000000
                                              122.000000
                                                            99.000000 846.000000
                                                                                                         2.420000
                                                                                                                  81.000000
                                                                                                                              1.000000
```

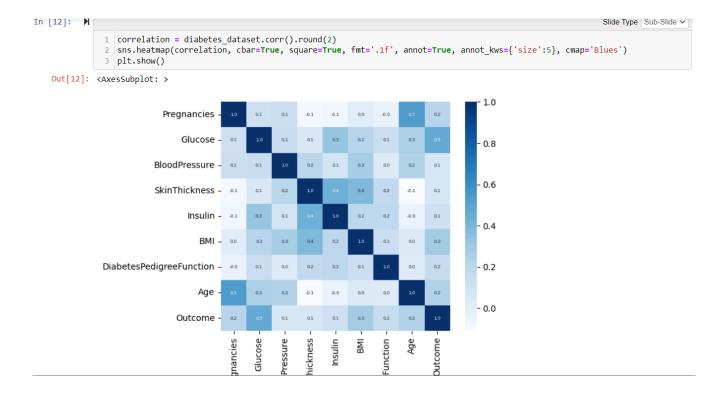


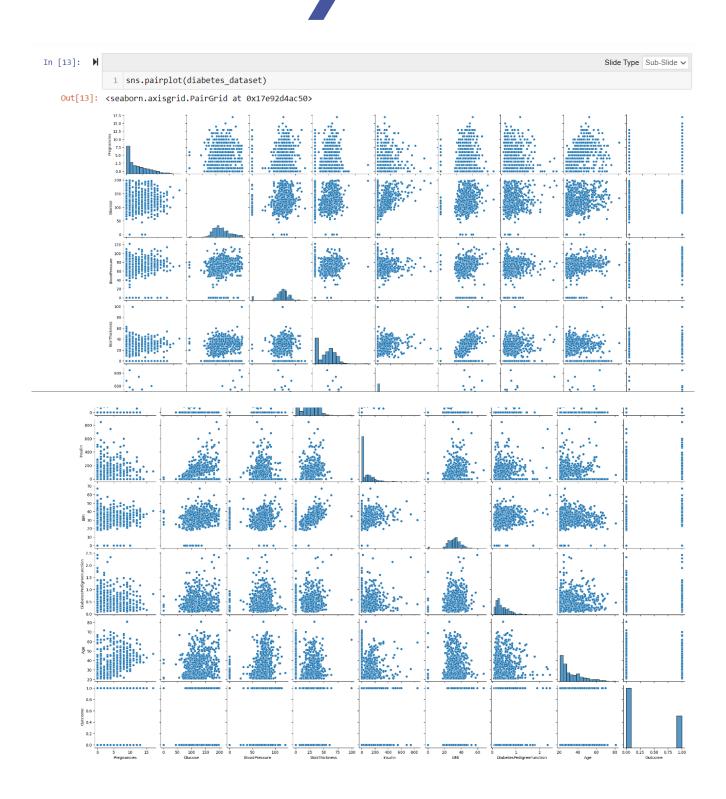
#### 3. DATA VISUALIZATION

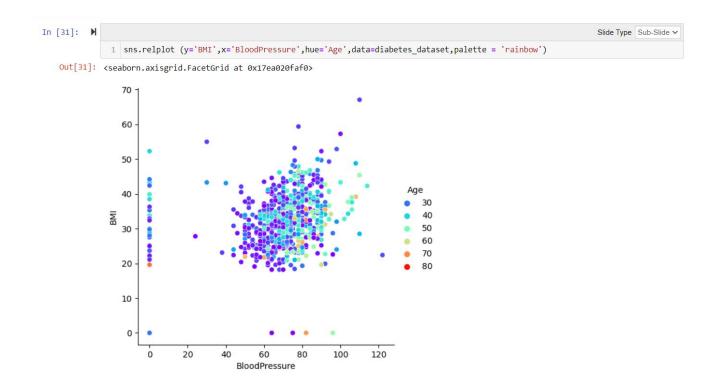


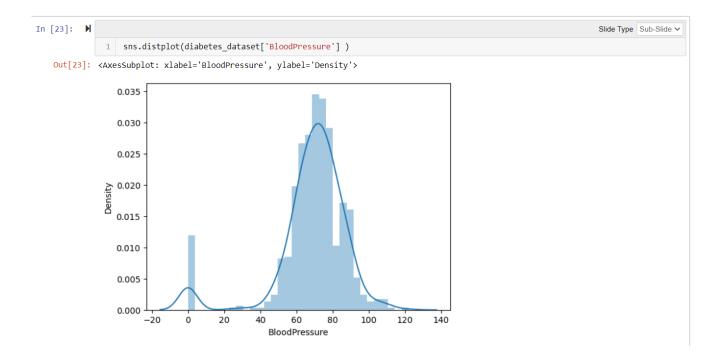
#### 4. ANALYSIS OF RELATIONALSHIP BETWEEN VARIABLES











#### 5. DATA DEPLOYMENT

#### 6. DATA MODELING

```
In [25]: ▶
                                                                                                                                                Slide Type Sub-Slide V
                 1 # RandomForestClassifier
                 2 rndf = RandomForestClassifier()
                 3 rndf.fit(X_train, Y_train)
                 4 # LogisticRegression
                 5 log=LogisticRegression()
                 6 log = LogisticRegression(solver='lbfgs', max_iter=1000)
7 log.fit(X_train,Y_train)
                 8 # Support Vector Machine
9 s=SVC()
                s.fit(X_train,Y_train)
# XgBoost
                12     xg=XGBClassifier()
13     xg.fit(X_train,Y_train)
14     # k-Nearest-Neighbor
                15 k=KNeighborsClassifier()
                16 k.fit(X_train,Y_train)
    Out[25]: RandomForestClassifier()
    Out[25]: LogisticRegression(max_iter=1000)
    Out[25]: SVC()
```

#### 7. MODEL EVALUATION

```
In [19]: ▶
                                                                                                                                                           Slide Type Sub-Slide V
                   1 # RandomForestClassifier
                   pred=rndf.predict(X_test)
                   3 R = round(accuracy_score(Y_test, pred),2)
                   4 # LogisticRegression
                  5 pred1=log.predict(X_test)
                  R1 = round(accuracy_score(Y_test, pred1),2)

# Support Vector Machine

pred2=s.predict(X_test)

R2 = round(accuracy_score(Y_test, pred2),2)
                  10 # XgBoost
                  11 pred3=xg.predict(X_test)
12 R3 = round(accuracy_score(Y_test, pred3),2)
13 # k-Nearest-Neighbor
                  14 pred4=k.predict(X_test)
                  15 R4 = round(accuracy_score(Y_test, pred4),2)
In [20]: ▶
                                                                                                                                                           Slide Type Sub-Slide V
                   1 scores=[]
                      scores.append(R)
                       scores.append(R1)
                   4 scores.append(R2)
                   5 scores.append(R3)
                  6 scores.append(R4)
7 Accuracy=pd.DataFrame(scores)
```

#### **Results Section**

```
In [21]: ▶
                                                                                                                                            Slide Type Sub-Slide V
                 1 Accuracy.index= ["RandomForest","LogisticRegression ","SVM","XgBoost","KNN"]
                 2 Accuracy.columns=["Accuracy"]
3 Accuracy.index.name='Methods'
                 4 Accuracy
   Out[21]:
                                   Accuracy
                         Methods
                    RandomForest
                                       0.76
                LogisticRegression
                                       0.76
                             SVM
                                       0.79
                                       0.75
                         XgBoost
                             KNN
                                       0.73
```

#### **Building Flask Web App**

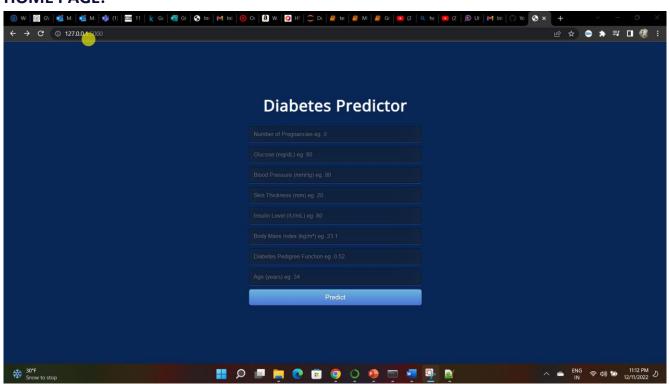
#### For this we developed

- 1. App.py file
- 2. Pickle file
- 3. Static files
- 4. Templates
  - a. Index.html
  - b. Result.html

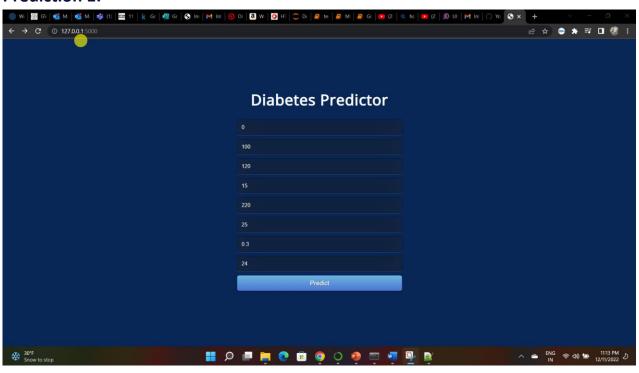
#### **Anaconda Prompt:**

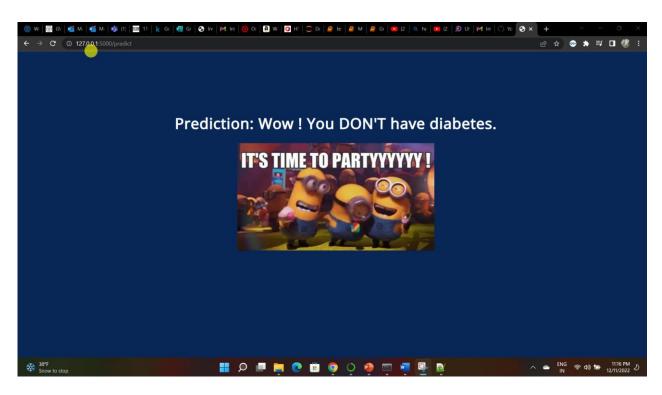
```
(rahul) C:\Users\rahul\Downloads\DSCI 6007-02\Group4_Final_Project_Phoenix Attackers>python app.py
* Serving Flask app 'app' (lazy loading)
* Environment: production
WARNING: This is a development server. Do not use it in a production deployment.
Use a production WSGI server instead.
* Debug mode: on
* Restarting with stat
* Debugger is active!
* Debugger PIN: 733-395-232
* Running on http://127.0.0.1:5000/ (Press CTRL+C to quit)
127.0.0.1 - - [11/Dec/2022 23:08:32] "GET / HTTP/1.1" 200 -
127.0.0.1 - - [11/Dec/2022 23:08:32] "GET / static/css/style.css HTTP/1.1" 404 -
```

#### **HOME PAGE:**

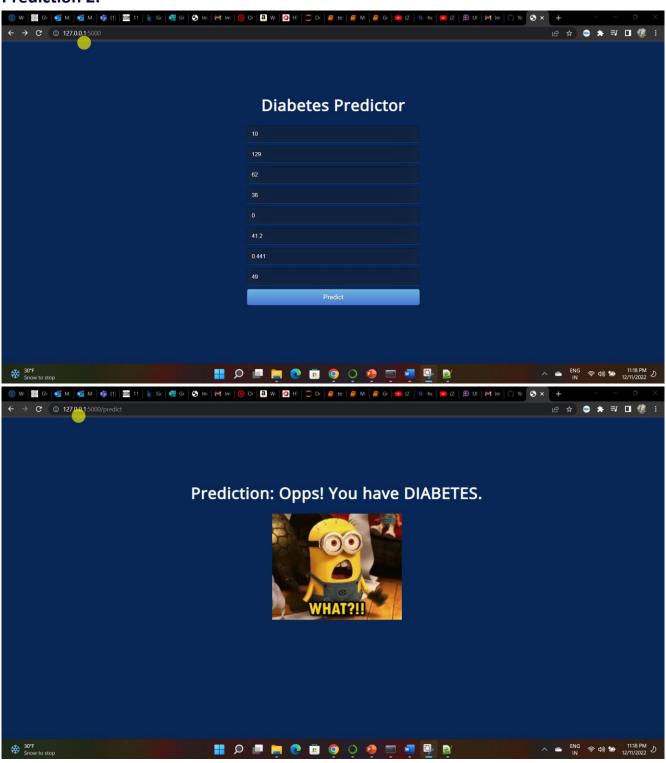


#### **Prediction 1:**





#### **Prediction 2:**



#### **Conclusion**

Based on the above analysis, we can conclude that "Support Vector Machine (SVM)" is the best model to predict the diabetes with most accurate value compared to other machine learning algorithms.

#### **Contributions/References:**

- 1. Data sets We have taken from the Kaggle. This dataset is originally from the National Institute of Diabetes and Digestive and Kidney Diseases.
- 2. Some Information is taken from Internet.