

AI Based Diabetes Prediction System

TEAM MEMBER

311521106081- SAI HARSHAVARADHAN

PHASE 1 DOCUMENT SUBMISSION



OBJECTIVE:

The objective of this project is to develop a machine learning model that analyze medical data and predict the likelihood of an individual developing diabetes and to provide early risk assessment and personalized preventive measures, allowing individuals to take proactive actions to manage their health.

1.DATA SOURCE:

A good data source containing medical features such as glucose levels, blood pressure, BMI, etc., along with information about whether the individual has diabetes or not.

Dataset Link: <https://www.kaggle.com/datasets/mathchi/diabetes-data-set>

Pregnanci	Glucose	BloodPres	SkinThickn	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
4	110	92	0	0	37.6	0.191	30	0
10	168	74	0	0	38	0.537	34	1
10	139	80	0	0	27.1	1.441	57	0
1	189	60	23	846	30.1	0.398	59	1
5	166	72	19	175	25.8	0.587	51	1
7	100	0	0	0	30	0.484	32	1
0	118	84	47	230	45.8	0.551	31	1
7	107	74	0	0	29.6	0.254	31	1
1	103	30	38	83	43.3	0.183	33	0
1	115	70	30	96	34.6	0.529	32	1
3	126	88	41	235	39.3	0.704	27	0
8	99	84	0	0	35.4	0.388	50	0
7	196	90	0	0	39.8	0.451	41	1
9	119	80	35	0	29	0.263	29	1
11	143	94	33	146	36.6	0.254	51	1

2.DATA PREPROCESSING:

Preprocessing steps involve data cleaning,feature extraction and handling missing values to ensure data quality and consistency.

a) DUPLICATE REMOVAL:

We will identify and remove duplicate's typically by sorting the dataset based on a unique identifier and then eliminating consecutive rows with the same identifier.

b) HANDLING MISSING VALUES:

Missing data is common and needs to be addressed.We will implement suitable methods such as:

Mean imputation: Replace missing values with the mean of the feature for the remaining rows. This is appropriate for numerical features.

Median imputation: If data contains outliers, median imputation can be more robust as it is less sensitive to extreme values.

3. MODEL SELECTION:

LOGISTIC REGRESSION ALGORITHM

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, confusion_matrix
import matplotlib.pyplot as plt
import seaborn as sns

#read the data file
data = pd.read_csv("/kaggle/input/diabetes-data-set/diabetes.csv")
data.head()

data.describe()

data.isnull().sum()
```

OUTPUT:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148	72	35	0	33.6	0.627	50	1
1	1	85	66	29	0	26.6	0.351	31	0
2	8	183	64	0	0	23.3	0.672	32	1
3	1	89	66	23	94	28.1	0.167	21	0
4	0	137	40	35	168	43.1	2.288	33	1

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000
mean	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	0.471876	33.240885	0.348958
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	0.331329	11.760232	0.476951
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.078000	21.000000	0.000000
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	0.243750	24.000000	0.000000
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	0.372500	29.000000	0.000000
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.626250	41.000000	1.000000
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.420000	81.000000	1.000000

In conclusion, AI-powered diabetes prediction systems using machine learning have the potential to revolutionize the way that diabetes is diagnosed and managed. By predicting which individuals are at risk of developing diabetes, these systems can help to prevent the onset of the disease and its associated complications.