



Motivations in Predicting Diabetes

Early Diagnosis & Intervention: Catching diabetes early for better outcomes.

Personalized Care: Tailored plans based on individual risk factors. Resource

Allocation: Targeting high-risk populations for efficient care.

Public Health Planning: Predicting future trends to guide policy and programs.

Research & Drug Development: Identifying targets for new treatments.

Empowering Individuals: Taking control of health with risk assessments.

Economic Impact: Reducing healthcare costs through better management.







About Dataset

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 or not a patient has diabetes, based on certain diagnostic measurements included in the dataset. Several constraints were placed on the selection of these instances from a larger database. In particular, all patients here are females at least 21 years old of Pima Indian heritage.











df.head()

✓ 0.0:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	ВМІ	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

df.tail()

✓ 0.0

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	вмі	DiabetesPedigreeFunction	Age	Outcome
763	10	101	76	48	180	32.9	0.171	63	0
764	2	122	70	27	0	36.8	0.340	27	0
765	5	121	72	23	112	26.2	0.245	30	0
766	1	126	60	0	0	30.1	0.349	47	1
767	1	93	70	31	0	30.4	0.315	23	0



Import Dataset



Dataset Content



```
df.info()
✓ 0.0s
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
# Column
                              Non-Null Count Dtype
    Pregnancies
                              768 non-null
                                              int64
    Glucose
                              768 non-null
                                              int64
    BloodPressure
                              768 non-null
                                              int64
    SkinThickness
                              768 non-null
                                              int64
    Insulin
                              768 non-null
                                              int64
    BMI
                              768 non-null
                                              float64
    DiabetesPedigreeFunction 768 non-null
                                              float64
                              768 non-null
                                              int64
8 Outcome
                              768 non-null
                                              int64
dtypes: float64(2), int64(7)
memory usage: 54.1 KB
   df.columns
 ✓ 0.0s
Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
       'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
     dtype='object')
```



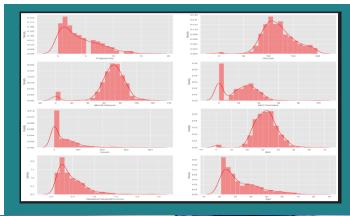


- **Pregnancies**: Number of times pregnant.
- Glucose: Plasma glucose concentration a 2 hours in an oral glucose tolerance test.
- **Blood Pressure**: Diastolic blood pressure (mm Hg).
- **Skin Thickness**: Triceps skin fold thickness (mm).
- **Insulin**: 2-Hour serum insulin (mu U/ml).
- **BMI**: Body mass index (weight in kg/(height in m)^2).
- **DiabetesPedigreeFunction**: Diabetes pedigree function
- Age: Age (years).
- Outcome: '1' denotes patient having diabetes and '0' denotes patient not having diabetes.





Exploratory Data Analysis

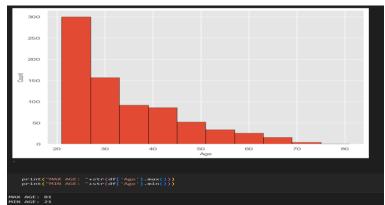


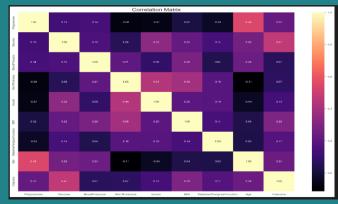






df.gro	upby("Out	come	").agg({'Pregnancies':'max'})
	Pregnanc	ies	
Outcome			
0		13	
1		17	
df.gro	upby("Out	come	").agg({'Glucose':'max'})
	Glucose		
Outcome			
0	197		
1	199		

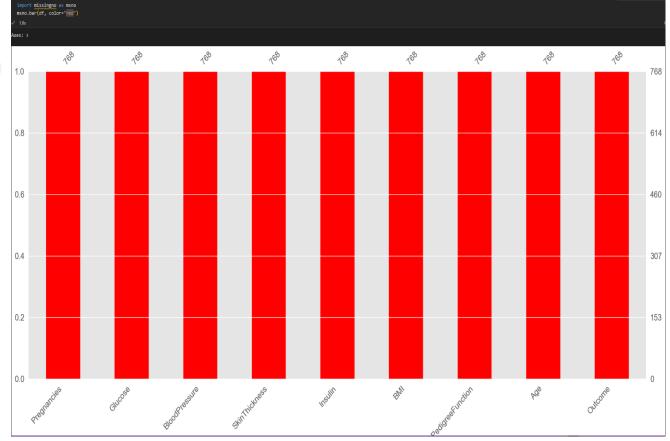






1. Missing Observation Analysis





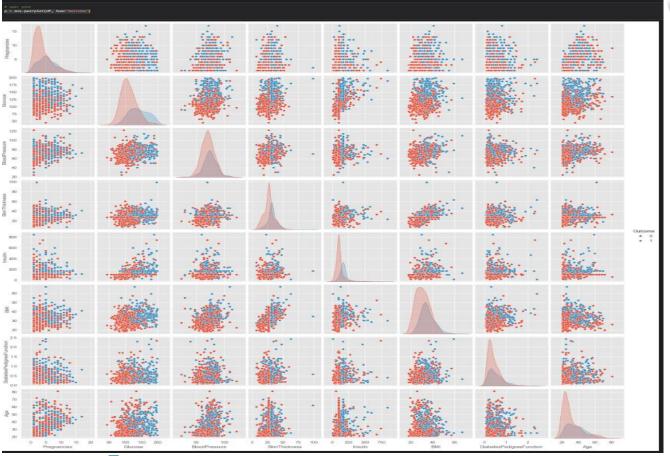


Data Preprocessing





2. Outlier Observation Analysis





Data Preprocessing



3. Local Outlier Factor (LOF)



```
from sklearn.neighbors import LocalOutlierFactor
   lof = LocalOutlierFactor(n_neighbors=10)
   lof.fit_predict(df)
Output is truncated. View as a <u>scrollable element</u> or open in a <u>text editor</u>. Adjust cell output <u>settings</u>...
```





Data Preprocessing





BMI



Feature Engineering

Creating new variables is important for models. But you need to create a logical new variable. For this data set, some new variables were created according to BMI, Insulin and glucose variables.

```
NewBMI
     Underweight
          Normal
      Overweight
       Obesity 1
       Obesity 2
       Obesity 3
dtype: category
Categories (6, object): ['Normal', 'Obesity 1', 'Obesity 2', 'Obesity 3', 'Overweight', 'Underweight']
   df['NewBMI'] = NewBMI
   df.loc[df["BMI"]<18.5, "NewBMI"] = NewBMI[0]</pre>
   df.loc[(df["BMI"]>18.5) & df["BMI"]<=24.9, "NewBMI"] = NewBMI[1]</pre>
   df.loc[(df["BMI"]>24.9) & df["BMI"]<=29.9, "NewBMI"] = NewBMI[2]</pre>
   df.loc[(df["BMI"]>29.9) & df["BMI"]<=34.9, "NewBMI"] = NewBMI[3]</pre>
   df.loc[(df["BMI"]>34.9) & df["BMI"]<=39.9, "NewBMI"] = NewBMI[4]</pre>
   df.loc[df["BMI"]>39.9, "NewBMI"] = NewBMI[5]
```



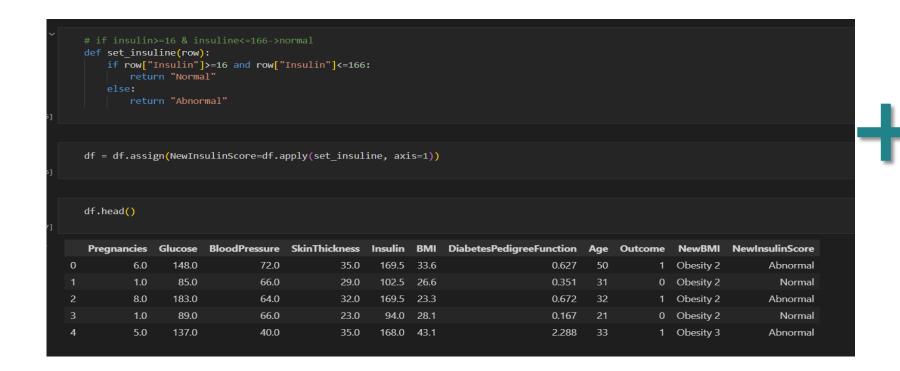
Feature Engineering

Creating new variables is important for models. But you need to create a logical new variable. For this data set, some new variables were created according to BMI, Insulin and glucose variables.



Insulin







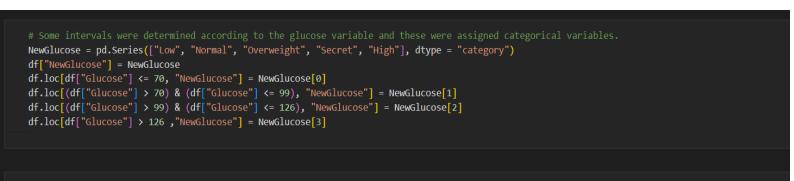
Feature Engineering

Creating new variables is important for models. But you need to create a logical new variable. For this data set, some new variables were created according to BMI, Insulin and glucose variables.



Glucose





df.head()

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	ВМІ	DiabetesPedigreeFunction	Age	Outcome	NewBMI	NewInsulinScore	NewGlucose
0	6.0	148.0	72.0	35.0	169.5	33.6	0.627	50	1	Obesity 2	Abnormal	Secret
1	1.0	85.0	66.0	29.0	102.5	26.6	0.351	31	0	Obesity 2	Normal	Normal
2	8.0	183.0	64.0	32.0	169.5	23.3	0.672	32	1	Obesity 2	Abnormal	Secret
3	1.0	89.0	66.0	23.0	94.0	28.1	0.167	21	0	Obesity 2	Normal	Normal
4	5.0	137.0	40.0	35.0	168.0	43.1	2.288	33	1	Obesity 3	Abnormal	Secret





One Hot Encoding



Categorical variables in the data set should be converted into numerical values. For this reason, these transformation processes are performed with Label Encoding and One Hot Encoding method.

by making One Hot Encoding transformation, categorical variables were converted into numerical values. It is also protected from the Dummy variable trap.

	hot enco		f, columns = ["N	ewBMI", "Nev	vInsulin	Score",	"NewGlucose"], drop	_first	=True)											
df.he		5 1																		
Preg	nancies 6	Glucose	BloodPressure S	kinThickness	Insulin	BMI	DiabetesPedigreeFunct	tion A	ge Outcon	ne NewBMI_Obesity 1	NewBMI_Obesity 2	2 NewBMI_Obesity :	3 NewBMI_Overwei	ight NewBMI_Under	weight New	InsulinScore_Normal	NewGlucose_Low	NewGlucose_Norma	I NewGlucose_Overweight	NewGlucose_Secret
0	6.0	148.0	72.0	35.0	169.5	33.6	0.	.627	50	1 False	: True	e False	e F	False	False	False	False	False	e False	True
	1.0	85.0	66.0	29.0	102.5	26.6	0.	351	31	0 False	: True	e Falsi	ie F	False	False	True	False	True	e False	False
2	8.0	183.0	64.0	32.0	169.5	23.3	0.	.672	32	1 False	: True	e False	e F	False	False	False	False	False	e False	True
3	1.0	89.0	66.0	23.0	94.0	28.1	0.	.167	21	0 False	: True	e Falsı	ie F	False	False	True	False	True	e False	False
4	5.0	137.0	40.0	35.0	168.0	43.1	2.	288	33	1 False	e Falsı	e Tru	e F	False	False	False	False	False	e False	True





One Hot Encoding



Categorical variables in the data set should be converted into numerical values. For this reason, these transformation processes are performed with Label Encoding and One Hot Encoding method.

Make new data frame in which all newly created variable are present

	'NewBMI_Und	sity 2', 'NewBMI_O	besity 3', 'NewBMI ulinScore_Normal',	_Overweight', 'NewGlucose_Low', NewGlucose_Secret']]						
	categorical_df.hea	d()								
	NewBMI_Obesity 1	NewBMI_Obesity 2	NewBMI_Obesity 3	NewBMI_Overweight	NewBMI_Underweight	NewInsulinScore_Normal	NewGlucose_Low	NewGlucose_Normal	NewGlucose_Overweight	NewGlucose_Secret
0	False	True	False	False	False	False	False	False	False	True
1	False	True	False	False	False	True	False	True	False	False
2	False	True	False	False	False	False	False	False	False	True
	False	True	False	False	False	True	False	True	False	False
3										





One Hot Encoding



Categorical variables in the data set should be converted into numerical values. For this reason, these transformation processes are performed with Label Encoding and One Hot Encoding method.

The variables in the data set are an effective factor in increasing the performance of the models by standardization. There are multiple standardization methods. These are methods such as "Normalize"," MinMax"," Robust" and "Scale".

tra X=t	nsformer ransforme	= RobustS r.transfo	essing import Ro caler().fit(X) erm(X) ummns = cols, in															
X.h	ead()																	
Pro	gnancies	Glucose	BloodPressure	SkinThickness	Insulin	ВМІ	DiabetesPedigreeFunction	Age										
0	0.75	0.775	0.000	1.000000	1.000000	0.177778	0.669707	1.235294										
	-0.50	-0.800	-0.375	0.142857	0.000000	-0.600000	-0.049511	0.117647										
	1.25	1.650	-0.500	0.571429	1.000000	-0.966667	0.786971	0.176471										
	-0.50	-0.700	-0.375	-0.714286	-0.126866	-0.433333	-0.528990	-0.470588										
	0.50	0.500	-2.000	1.000000	0.977612	1.233333	4.998046	0.235294										
X =	pd.conca	t([X, cat	egorical_df], a	exis=1)														
X.h	ead()																	
Pre	gnancies	Glucose	BloodPressure	SkinThickness	Insulin	ВМІ	DiabetesPedigreeFunction	Age	NewBMI_Obesity 1	NewBMI_Obesity 2	NewBMI_Obesity 3	NewBMI_Overweight	NewBMI_Underweight	NewInsulinScore_Normal	NewGlucose_Low	NewGlucose_Normal	NewGlucose_Overweight	NewGlucose_Secret
	0.75	0.775	0.000	1.000000	1.000000	0.177778	0.669707	1.235294	False	True	False	False	False	False	False	False	False	True
	-0.50	-0.800	-0.375	0.142857	0.000000	-0.600000	-0.049511	0.117647	False	True	False	False	False	True	False	True	False	False
	1.25	1.650	-0.500	0.571429	1.000000	-0.966667	0.786971	0.176471	False	True	False	False	False	False	False	False	False	True
	-0.50	-0.700	-0.375	-0.714286	-0.126866	-0.433333	-0.528990	-0.470588	False	True	False	False	False	True	False	True	False	False
4	0.50	0.500	-2.000	1.000000	0.977612	1.233333	4.998046	0.235294	False	False	True	False	False	False	False	False	False	True



Splitting and Scaling Data for Modelling



Splitting Data (train_test_split):



- This code splits our data (features X and target variable y) into training and testing sets.
- test_size=0.2 allocates 20% of the data for testing and the remaining 80% for training.
- random_state=0 sets a seed for reproducibility (ensuring the same split each time you run the code).

Scaling Data (StandardScaler):

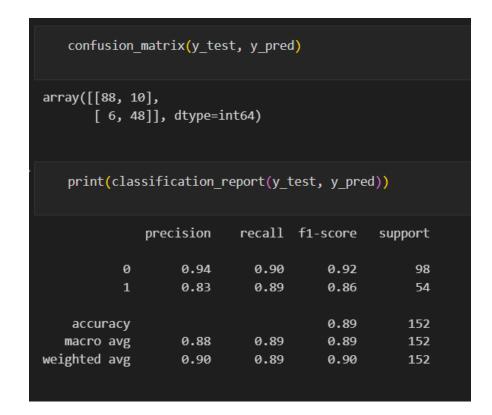
- A StandardScaler object is created (scaler).
- scaler.fit transform(X train) standardizes the training data (X train).
 - This removes the mean and scales the features to unit variance.
- The same scaler (scaler) is then used to transform the testing data (X test).
 - It applies the same mean and scaling factors learned from the training data.

```
X_train, X_test, y_train , y_test = train_test_split(X,y, test_size=0.2, random_state=0)

scaler =StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```









Modelling













SVM

[[90 8] [6 48]]	precision	recall	f1-score	support	
0 1	0.94 0.86	0.92 0.89	0.93 0.87	98 54	
accuracy macro avg weighted avg	0.90 0.91	0.90 0.91	0.91 0.90 0.91	152 152 152	

















[[87 11]					
[9 45]]	precision	recall	f1-score	support	
(0.91	0.89	0.90	98	
1	1 0.80	0.83	0.82	54	
accuracy	y		0.87	152	
macro av	g 0.85	0.86	0.86	152	
weighted av	g 0.87	0.87	0.87	152	









Model Comparison









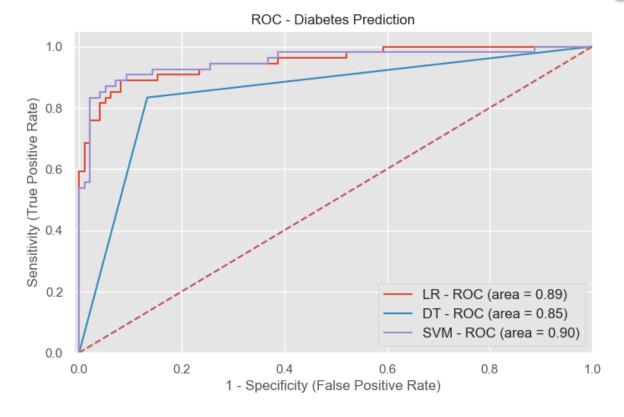
2 Decision Tree Classifier

86.84

ROC
Receiver Operating Characteristic











Performance Evaluation





