

## 1 Understanding data

```
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns
```

Create dataframe

```
df = pd.read_csv("diabetes.csv")  
df.head()
```

`df.columns` → Return column name / features

`df.dtypes` → Data types of each column

Missing values in given dataset.

```
df.isnull().sum()
```

```
df.head(20)
```

⇒ There are many zeros like in blood pressure, insulin, skin thickness, BMI columns, which are not possible (problem).

\* Data Imputation of 0's in every feature

\* Size of the data

`df.shape`

\* Target column: Outcome [0,1] (Binary classification Task)

\* As the target column is available in the dataset, supervised machine learning algorithm

\* Records : 768

⇒ Data does not contain any null values but having some weird zeros.

## 2. EDA : Feature Selection and Data Imputation

Correlation coeff. between the features

`df.corr()`

Heat map

`plt.figure(figsize=(15,15))`

`ax = sns.heatmap(df.corr(), annot=True)`

`plt.savefig('correlation-coeff.jpg')`

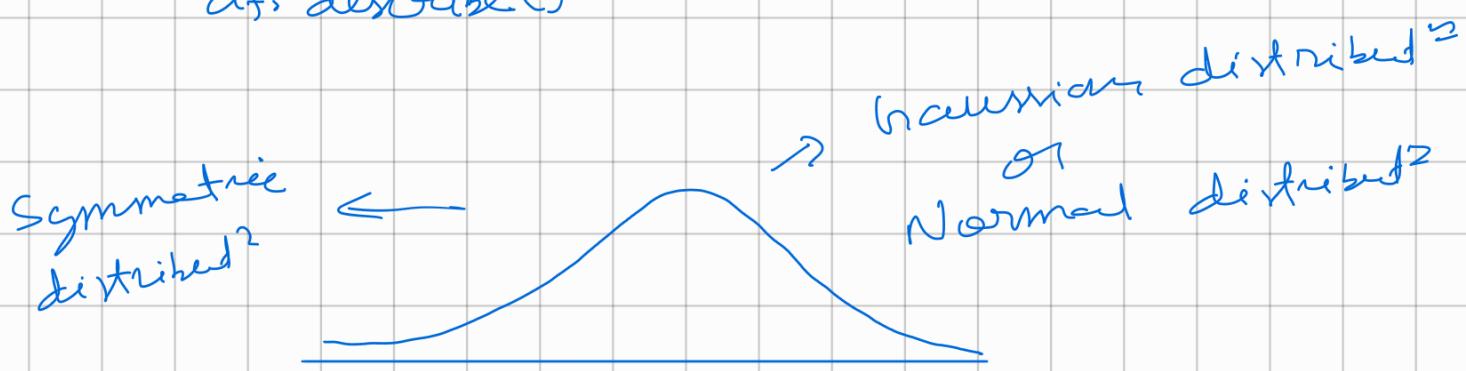
`plt.show()`

To add value  
to heat map

$\Rightarrow$  No features are highly correlated  
so we cannot remove any feature.

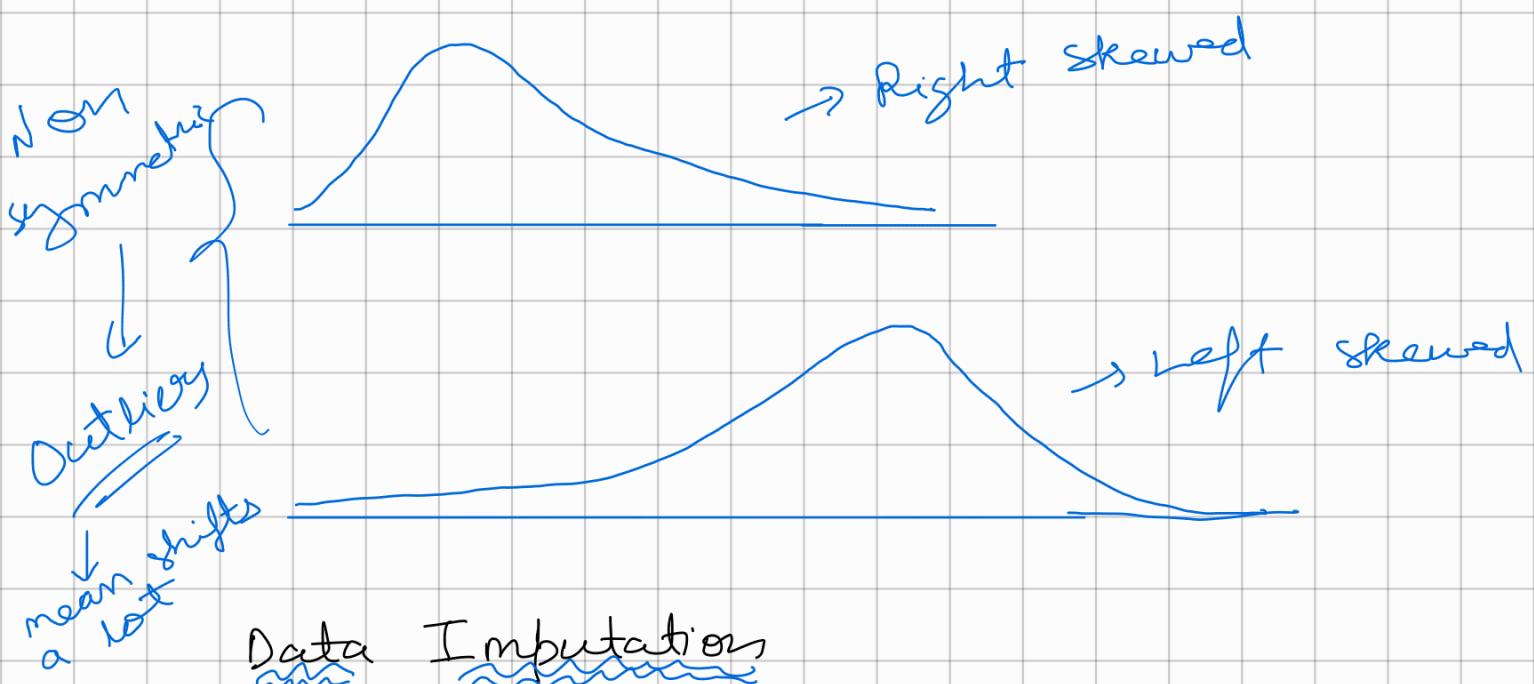
Descriptive Statistics of the given data

`df.describe()`

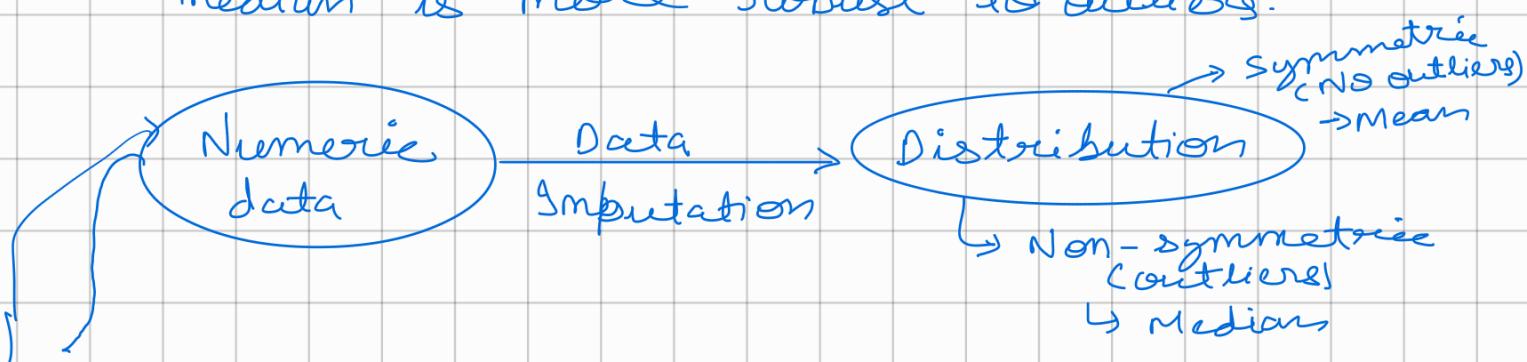


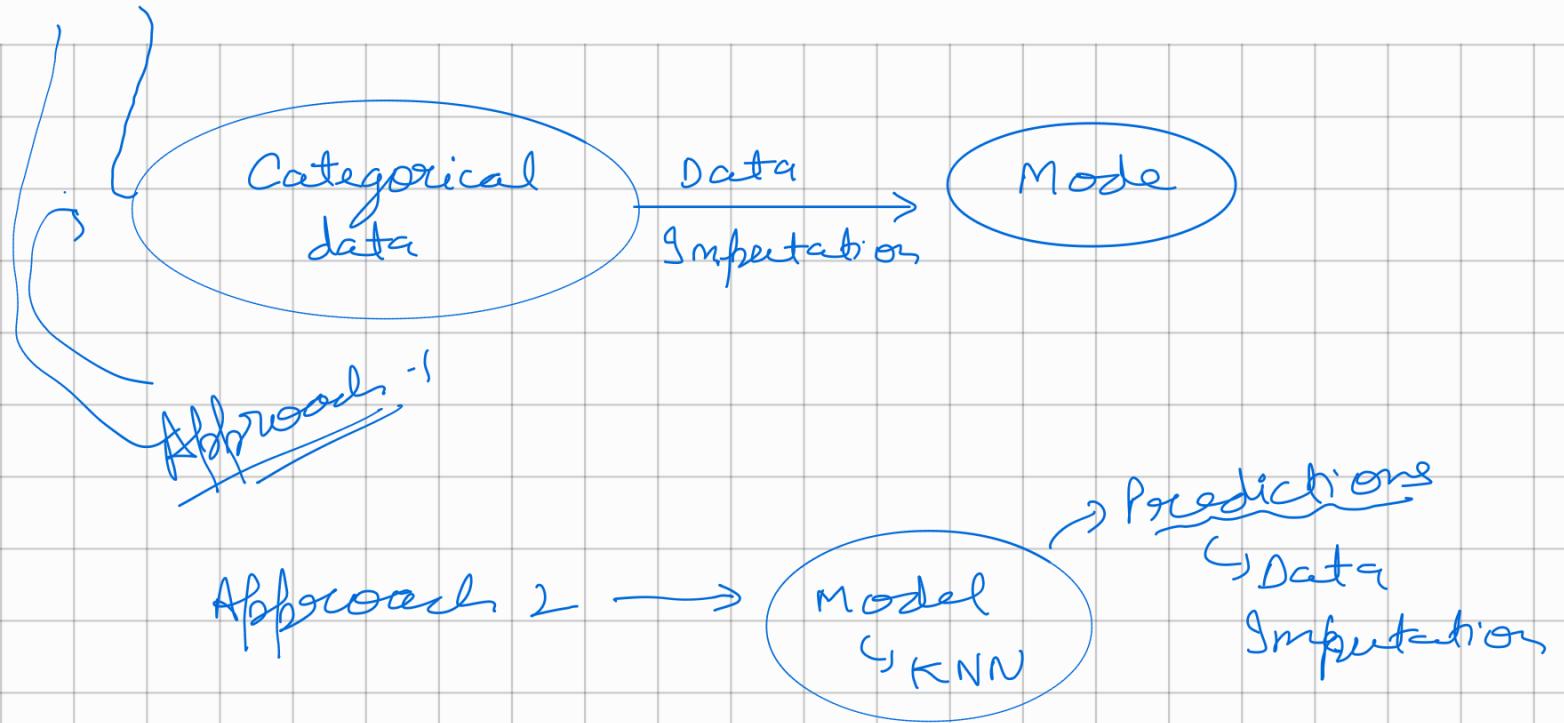
Standard Normal Distribution

$$\text{↳ } \mu=0, \sigma=1$$



median is more robust to outliers.





Distribution

`sns.distplot(df.Pregnancies)` # median

`sns.distplot(df.BloodPressure)` # mean

`sns.distplot(df.Insulin)` # median

# Data imputation of insulin → median

# Insulin → Right skewed distribution

`df['Insulin'] = df['Insulin'].replace(0, df['Insulin'].median())`

# Blood Pressure → mean

# Pregnancies → median

# skin thickness →

Data Imputation :- Data imputation is the process of replacing missing data values in a dataset with substituted values, a crucial step in data preprocessing to ensure data integrity and enable effective analysis.

Why Impute Data?

i) Handling Missing Data :- Real-world datasets often contain missing values due to various reasons, such as data entry errors, incomplete data recording, or sensor failures.

ii) Maintaining Data Integrity :- Ignoring or deleting rows with missing values can

lead to a significant reduction in dataset size and potentially introduce bias, impacting the reliability of analysis.

iii) Enabling Analysis :- Many statistical and machine learning algorithms cannot handle missing data, making imputation necessary for data analysis and model training.

### Common Imputation Techniques:

i) Mean / Median / Mode Imputation :- Replacing missing values with the mean, median or mode of the non-missing values in the same column.

ii) Constant Value Imputation :- Replacing missing values with a constant value, such as 0 or a specific placeholder.

iii) Regression Imputation :- Using a regression model to predict missing values based on other variables in the dataset.

iv) K - Nearest Neighbors (KNN) Imputation:- Using the values of similar data points (neighbors) to estimate missing values.

v) Multiple Imputation :- Generating multiple imputed datasets and analyzing them separately, allowing for uncertainty in the imputed values to be reflected in the analysis.

### Considerations When Choosing an Imputation Method:

i) Type of Missing Data :- Understanding the reason for missing

values is important, as different imputation methods are suitable for different types of missing data.

e.g.: - Missing Completely at Random (MCAR),  
Missing at Random (MAR), Missing Not at Random (MNAR).

ii) Data Characteristics :-

The nature of the data (e.g., numerical, categorical) and its distribution can influence the choice of imputation method.

iii) Potential Bias :- Imputation methods can introduce bias, so it's important to choose a method that minimizes bias and preserves the integrity of the data.

iv) Computational cost :- Some imputation methods, like multiple imputation, can be computationally expensive.

⇒ Data integrity is the accuracy, consistency, and reliability of data over its entire lifecycle. It can also refer to the process of ensuring that data is valid and accurate.



# EDA : Outlier Detection and Normalization

#  $x \rightarrow$  input features

#  $y \rightarrow$  target value

$x = df.drop(columns = "Outcome", axis=1)$

↳ Deleted the outcome column for  $x$ .

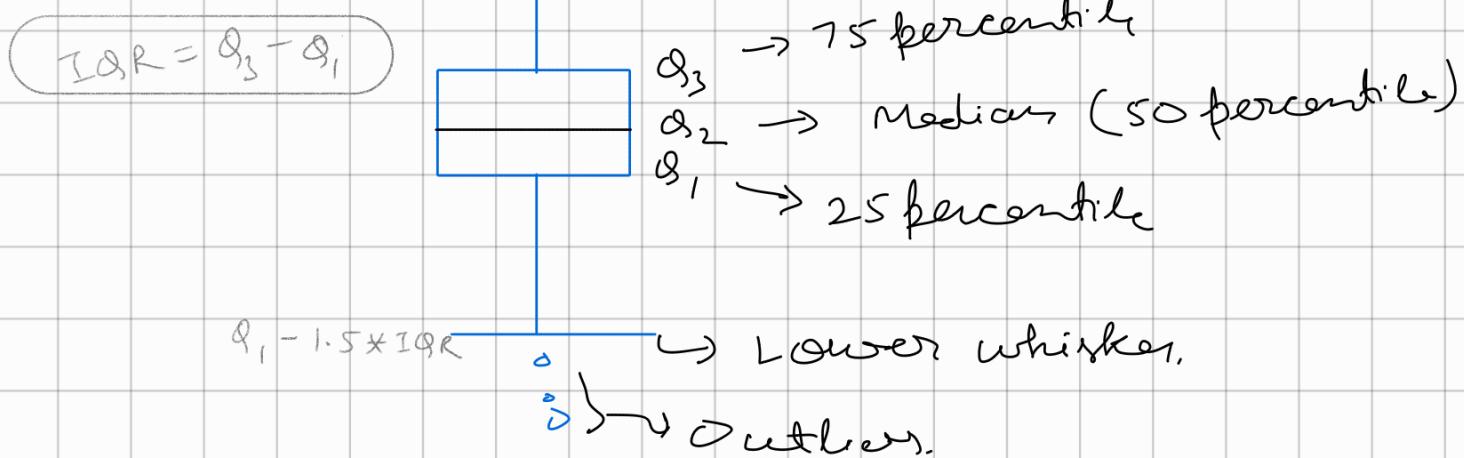
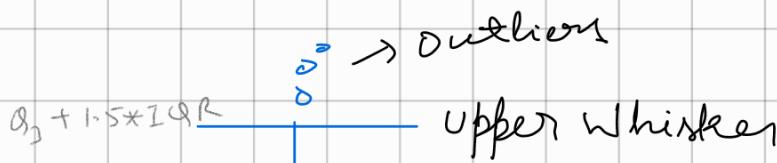
$y = df["Outcome"]$

## 1: Outlier Detection

`fig, ax = plt.subplots(figsize=(15, 15))`

`sns.boxplot(data=x, ax=ax)`

`plt.savefig("BoxPlot.jpg")`



`cols = ['Pregnancies', 'Glucose', 'BloodPressure',  
 'SkinThickness', 'Insulin', 'BMI',  
 'DiabetesPedigreeFunction', 'Age']`

for col in cols:

$Q1 = x[col].quantile(0.25)$

$Q3 = x[col].quantile(0.75)$

$IQR = Q3 - Q1$

$lower\_bound = Q1 - 1.5 * IQR$

$upper\_bound = Q3 + 1.5 * IQR$

$mask = (x[col] \geq lower\_bound) \& (x[col] \leq upper\_bound)$

$x_{outlier\_detection} = x[mask]$

$y_{outlier\_detection} = y[mask]$

} After removing outlier

## Standardization / Normalization

↳ To reduce biasness ( $\mu=0$ ,  $\sigma=1$ )

from sklearn.preprocessing import  
StandardScaler

scaler = StandardScaler()

x\_scaled = scaler.fit\_transform(x\_outlier\_detection)

⇒ After this outliers are found

x\_scaled = pd.DataFrame(x\_scaled)

x\_scaled.describe()

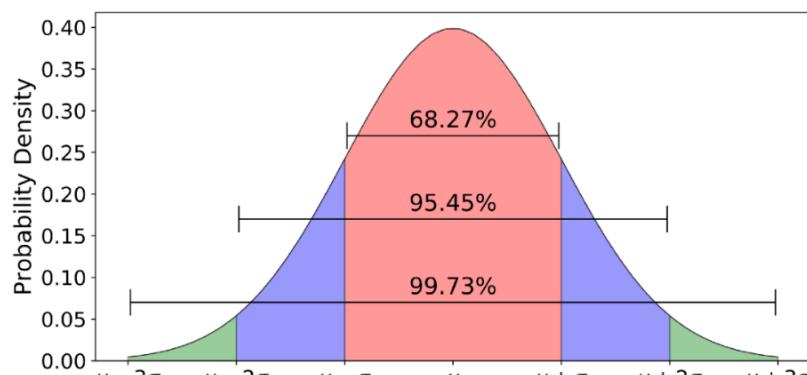
y\_outlier\_detection.value\_counts()

↳ Outcome

0 : 493      } → Biased  
1 : 263

## EDA :- Outlier Handling Techniques

### Normal / Gaussian Distribution



cols = ['Pregnancies', 'Glucose', 'BloodPressure',  
'SkinThickness', 'Insulin', 'BMI',  
'DiabetesPedigreeFunction', 'Age']

x\_scaled = pd.DataFrame(x\_scaled, columns=cols)  
x\_scaled.describe()

fig, ax = plt.subplots(figsize=(15, 15))

sns.boxplot(data=x\_scaled, ax=ax)

plt.savefig('boxplot-1.jpg')

## 95 percentile

$q = \text{x-scaled}["\text{Insulin}"].\text{quantile}(0.95)$   
 $\text{datanew} = \text{x-scaled}[\text{x-scaled}["\text{Insulin}"] < q]$

```
fig, ax = plt.subplots(figsize=(15, 15))
sns.boxplot(data = datanew, ax=ax)
plt.savefig('boxplot2.jpg')
```

## 6 Model Training : Splitting data and Handling of Imbalanced Data

Sklearn test train split

### Model Training

Splitting of data into training and testing

```
from sklearn.model_selection import
train_test_split
```

$X_{\text{train}}, X_{\text{test}}, Y_{\text{train}}, Y_{\text{test}} =$   
 $\text{train\_test\_split}(X_{\text{scaled}}, Y_{\text{cleaned}},$   
 $\text{test\_size} = 0.33, \text{random\_state} = 42)$

↳ 33%

repeat with  
less random  
everytime

$Y_{\text{train}}.\text{value_counts}()$

⇒ Outcome

0	293	}	Data Imbalance
1	135		

### Data Imbalancing

→ oversampling :- Minority class  
and increase that number to the  
majority class.

→ undersampling :- Majority class  
and decrease that number to  
the minority class

→ SMOTE :- Synthetic data and

increase the number of samples  
to the majority class.

Imblearn documentation

### SMOTE Technique

```
from imblearn.over_sampling import SMOTE  
smote = SMOTE(random_state=42)  
X_train_resampled, Y_train_resampled =  
smote.fit_resample(X_train, Y_train)
```

```
# Check resampled class distribution  
print("\nResampled class distribution :")  
print(pd.Series(Y_train_resampled).value_  
counts())
```

## 7. Model Training: Implementation of Logistic Regression

sklearn logistic regression

```
from sklearn.linear_model import LogisticRegression  
classification = LogisticRegression()  
classification.fit(X_train_resampled, Y_train_  
resampled)
```

## 8. Model Prediction and Evaluation

### Model Prediction

```
Y_predictions = classification.predict(X_test)  
print(Y_predictions)
```

### Model Evaluation

sklearn accuracy score

```
from sklearn.metrics import accuracy_score  
accuracy_score(Y_test, Y_predictions)
```

Prediction is 75.943%.

### Classification report

sklearn classification report

Health care :- Recall is very important metric

```
from sklearn.metrics import classification_report  
target_names = ['Non-Diabetic', 'Diabetic']  
print(classification_report(Y_test, Y_predictions,  
                            target_names=target_names))
```

## 9. Saving the Model

```
import pickle  
pickle.dump(classification,  
            open("classification_model.pkl", "wb"))
```

Use saved file

```
classification_model = pickle.load(open(  
                                    "classification_model.pkl", "rb"))
```

```
classification_model.predict(X_test)
```

## 10 Project Outro