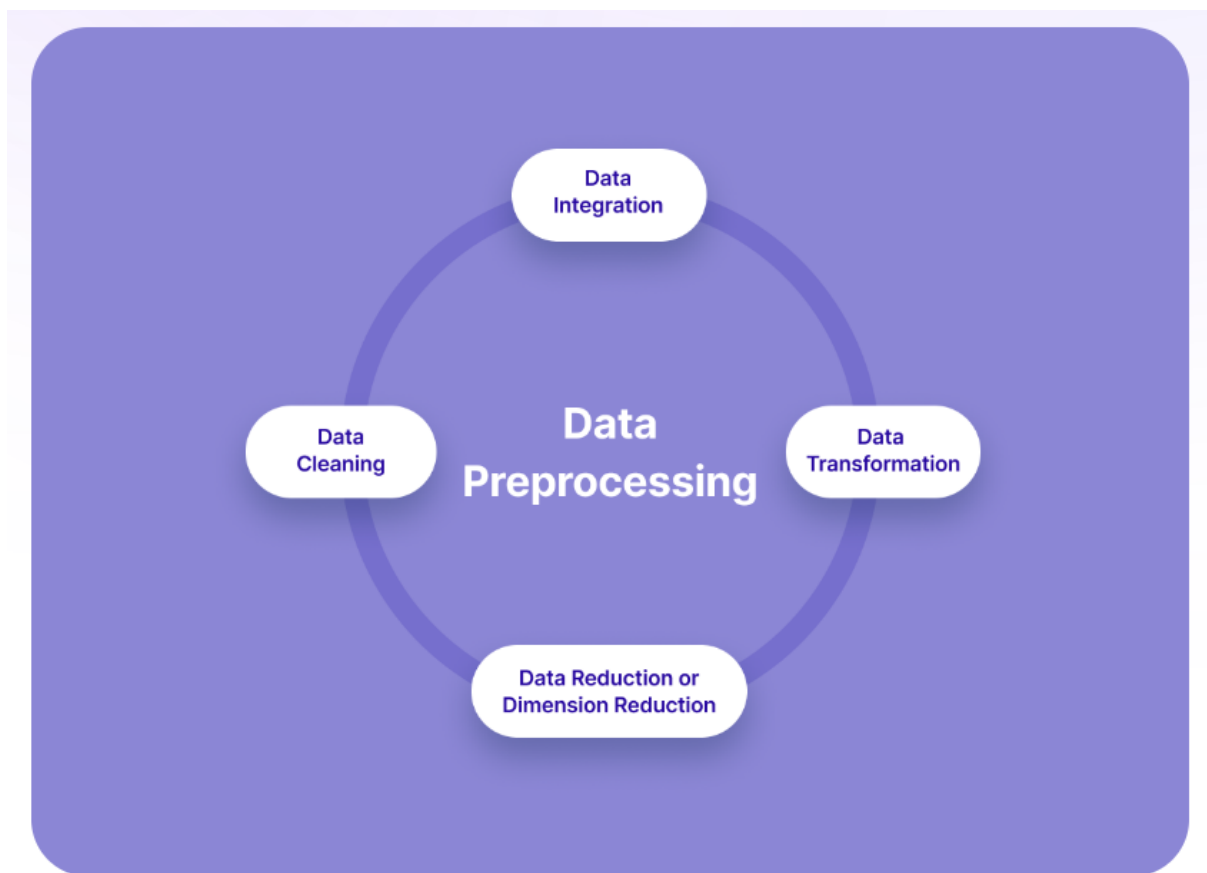


# AI Diabetes Prediction System (PHASE\_3)

## **Dataset preprocessing:**

In the merged dataset, we discovered a few exceptional zero values. For example, skin thickness and Body Mass Index (BMI) cannot be zero. The zero value has been replaced by its corresponding mean value. The training and test dataset has been separated using the holdout validation technique, where 80% is the training data and 20% is the test data.



## **Need of Data Preprocessing:**

For achieving better results from the applied model in Machine Learning projects the format of the data has to be in a proper manner. Some specified Machine Learning model needs information in a specified format, for example, Random Forest algorithm does not support null values, therefore to execute random forest algorithm null values have to be managed from the original raw data set.

## **AI Diabetes Prediction System (PHASE\_3)**

Another aspect is that the data set should be formatted in such a way that more than one Machine Learning and Deep Learning algorithm are executed in one data set, and best out of them is chosen.

### **Data Preprocessing Steps:**

There are four major steps in Data Preprocessing they are:

1. Data quality assessment
2. Data Cleaning
3. Data Transformation
4. Data Reduction

#### **1. Data quality assessment:**

There are a number of data anomalies and inherent problems to look out for in almost any data set, for example:

- Mismatched data types
- Mixed data values
- Data outliers
- Missing data

#### **2. Data Cleaning:**

Data cleaning is the process of adding missing data and correcting, repairing, or removing incorrect or irrelevant data from a data set. Data cleaning is the most important step of preprocessing because it will ensure that your data is ready to go for your downstream needs.

# AI Diabetes Prediction System (PHASE\_3)

## 3. Data Transformation:

With data cleaning, we've already begun to modify our data, but data transformation will begin the process of turning the data into the proper formats.

## 4. Data Reduction:

The more data you're working with, the harder it will be to analyze, even after cleaning and transforming it. Depending on your task at hand, you may actually have more data than you need. Data reduction not only makes the analysis easier and more accurate, but cuts down on data storage.

## Importing Libraries:

### Code:

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

sns.set()

from mlxtend.plotting import plot_decision_regions
import missingno as msno
from pandas.plotting import scatter_matrix
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier

from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import classification_report
import warnings
warnings.filterwarnings('ignore')
%matplotlib inline
```

# AI Diabetes Prediction System (PHASE\_3)

## Reading the dataset:

### Code:

```
diabetes_df = pd.read_csv('diabetes.csv')
diabetes_df.head()
```

### 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

## Columns:

### Code:

```
diabetes_df.columns
```

### Output:

```
Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
       'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
      dtype='object')
```

## Data Cleaning:

## Check the dataset have null values or not:

### Code:

```
diabetes_df.isnull()
```

# AI Diabetes Prediction System (PHASE\_3)

## Output:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False
...	...	...	...	...	...	...	...	...	...
763	False	False	False	False	False	False	False	False	False
764	False	False	False	False	False	False	False	False	False
765	False	False	False	False	False	False	False	False	False
766	False	False	False	False	False	False	False	False	False
767	False	False	False	False	False	False	False	False	False

768 rows × 9 columns

## Replace the 0 value with the NAN:

### Code:

```
diabetes_df_copy = diabetes_df.copy(deep = True)
diabetes_df_copy[['Glucose','BloodPressure','SkinThickness','Insulin',
', 'BMI']] =
diabetes_df_copy[['Glucose','BloodPressure','SkinThickness','Insulin',
', 'BMI']].replace(0,np.NaN)

print(diabetes_df_copy.isnull().sum())
```

## Output:

```
Pregnancies      0
Glucose           5
BloodPressure     35
SkinThickness     227
Insulin           374
BMI               11
DiabetesPedigreeFunction  0
Age              0
Outcome          0
dtype: int64
```

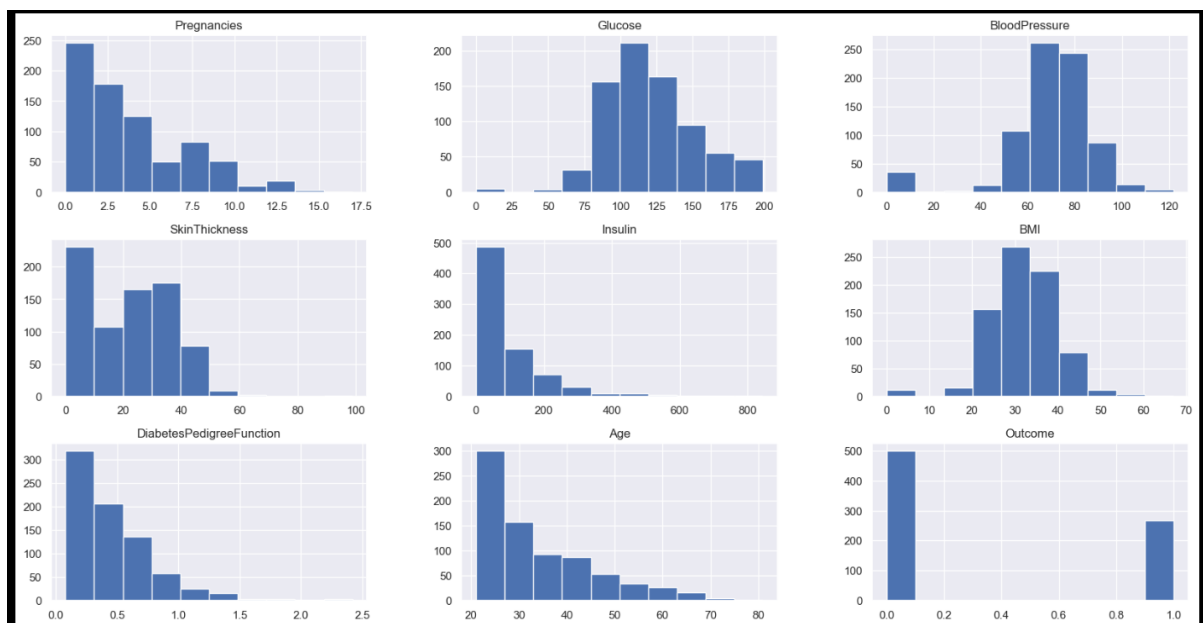
# AI Diabetes Prediction System (PHASE\_3)

## Data Visualization:

### Code:

```
p = diabetes_df.hist(figsize = (10,20))
```

### Output:



## Correlation between all the features:

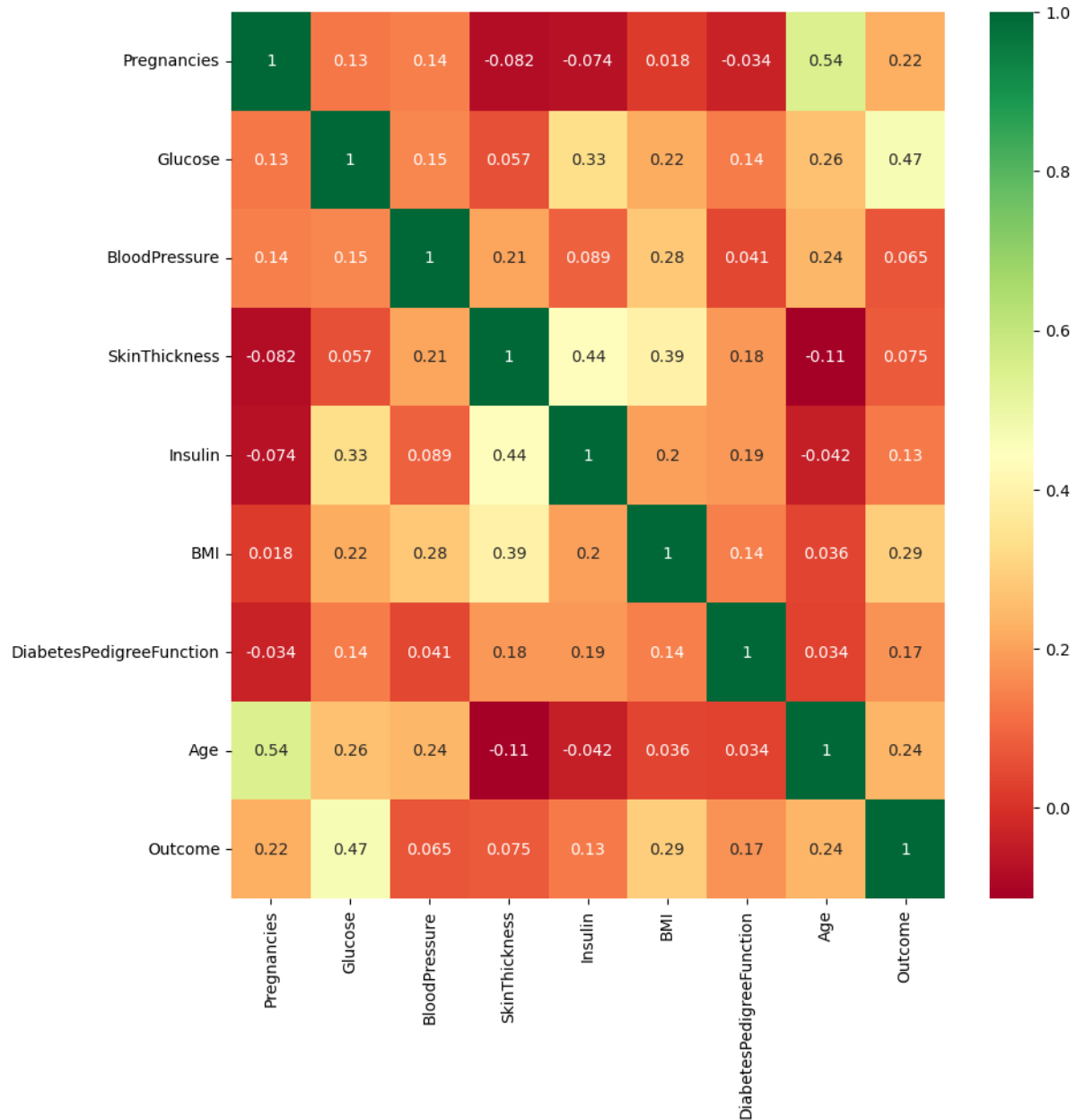
### Correlation between all the features before cleaning:

### Code:

```
plt.figure(figsize=(10,10))  
p = sns.heatmap(diabetes_df.corr(), annot=True,cmap = 'RdYlGn')
```

# AI Diabetes Prediction System (PHASE\_3)

## Output:



## Scaling the Data:

Before scaling down the data:

Code:

```
diabetes_df_copy.head()
```

## AI Diabetes Prediction System (PHASE\_3)

### Output:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148.0	72.0	35.0	NaN	33.6	0.627	50	1
1	1	85.0	66.0	29.0	NaN	26.6	0.351	31	0
2	8	183.0	64.0	NaN	NaN	23.3	0.672	32	1
3	1	89.0	66.0	23.0	94.0	28.1	0.167	21	0
4	0	137.0	40.0	35.0	168.0	43.1	2.288	33	1

### After Standard scaling down the data:

### Code:

```
sc_X = StandardScaler()
X =
pd.DataFrame(sc_X.fit_transform(diabetes_df_copy.drop(["Outcome"],axis = 1)), columns=['Pregnancies',
'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI',
'DiabetesPedigreeFunction', 'Age'])
X.head()
```

### Output:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
0	0.639947	0.862287	-0.032746	0.558557	NaN	0.165097	0.468492	1.425995
1	-0.844885	-1.202229	-0.517645	-0.014657	NaN	-0.846404	-0.365061	-0.190672
2	1.233880	2.009241	-0.679278	NaN	NaN	-1.323254	0.604397	-0.105584
3	-0.844885	-1.071148	-0.517645	-0.587871	-0.518847	-0.629654	-0.920763	-1.041549
4	-1.141852	0.501816	-2.618874	0.558557	0.104968	1.537847	5.484909	-0.020496



## **AI Diabetes Prediction System (PHASE\_3)**

### **Team Members:**

1. M. Keerthivasan
2. P. Thivahar
3. A.K. PraveenKumar
4. S. Muralidharan
5. A. Pachaiyappan