



In [1]:

```
1 #importing packages
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.preprocessing import MinMaxScaler
5 from sklearn.decomposition import PCA
6 from sklearn.pipeline import Pipeline
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.tree import DecisionTreeClassifier
9 from sklearn.ensemble import RandomForestClassifier
10
11 #Filtering Warnings
12 import warnings
13 warnings.filterwarnings('ignore')
```

In [2]:

```
1 #Importing Diabetes Data
2 DiabetesData = pd.read_csv("pima-indians-diabetes.csv")
3 DiabetesData.head()
```

Out[2]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction
0	6	148	72	35	0	33.6	0.62
1	1	85	66	29	0	26.6	0.35
2	8	183	64	0	0	23.3	0.67
3	1	89	66	23	94	28.1	0.16
4	0	137	40	35	168	43.1	2.28

In [3]:

```
1 X=DiabetesData.iloc[:,0:8]
2 y=DiabetesData.iloc[:,8]
```

In [4]:

```
1 #Dividing Data in test and train
2 X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=1)
```

## Creating pipelines

Creating pipelines for Logistic regression, Decision Tree and Random Forest models

## Pipeline steps will include

### 1. Data Preprocessing using MinMax Scaler



## 2. Reducing Dimensionality using PCA

## 3. Training respective models

In [5]:

```
1 #Logistic Regression Pipeline
2 LogisticRegressionPipeline=Pipeline([('myscaler',MinMaxScaler()),
3                                     ('mypca',PCA(n_components=3)),
4                                     ('logistic_classifier',LogisticRegression())])
```

In [6]:

```
1 #Decision tree Pipeline
2 DecisionTreePipeline=Pipeline([('myscaler',MinMaxScaler()),
3                                ('mypca',PCA(n_components=3)),
4                                ('decisiontree_classifier',DecisionTreeClassifier())])
```

In [7]:

```
1 #Random Forest Pipeline
2 RandomForestPipeline=Pipeline([('myscaler',MinMaxScaler()),
3                                ('mypca',PCA(n_components=3)),
4                                ('randomforest_classifier',RandomForestClassifier())])
```

# Modelling and Evaluation

In [8]:

```
1 ## Defining the pipelines in a list
2 mypipeline = [LogisticRegressionPipeline, DecisionTreePipeline, RandomForestPipeline]
```

In [9]:

```
1 # Fit the pipelines
2 for mypipe in mypipeline:
3     mypipe.fit(X_train, y_train)
```

In [10]:

```
1 #Defining variables for choosing best model
2 accuracy=0.0
3 classifier=0
4 pipeline=""
```

In [11]:

```
1 # Creating dictionary of pipelines and training models
2 PipelineDict = {0: 'Logistic Regression', 1: 'Decision Tree', 2: 'Random Forest'}
```



In [12]:

```
1 #getting test accuracy for all classifiers
2 for i,model in enumerate(mypipeline):
3     print("{} Test Accuracy: {}".format(PipelineDict[i],model.score(X_test,y_test)))
```

Logistic Regression Test Accuracy: 0.7597402597402597

Decision Tree Test Accuracy: 0.6883116883116883

Random Forest Test Accuracy: 0.7727272727272727

In [13]:

```
1 #Choosing best model for the given data
2 for i,model in enumerate(mypipeline):
3     if model.score(X_test,y_test)>accuracy:
4         accuracy=model.score(X_test,y_test)
5         pipeline=model
6         classifier=i
7 print('Classifier with best accuracy:{}'.format(PipelineDict[classifier]))
```

Classifier with best accuracy:Random Forest

"Data Science & AI"  
by  
Siva Rama Krishna