# Practical 01

**Aim:**

To build a logistic regression model to predict whether a person has diabetes using the Pima Indians Diabetes dataset and to compare the model performance before and after selecting the top 5 most important features using feature selection techniques.

**Theory:**

**Feature Selection and Filtering**

Feature selection is the process of identifying and selecting the most relevant features (variables) from a dataset that contribute the most to predicting the target variable. This reduces model complexity, speeds up training, and can improve model performance. Common methods include:

- **Filter Methods** (e.g., SelectKBest)

- **Wrapper Methods** (e.g., Recursive Feature Elimination)

- **Embedded Methods** (e.g., feature importance from tree models)

In this practical, we use **SelectKBest** with the **f_classif** scoring function, which uses ANOVA F-statistics to select the features most associated with the target outcome.

**Logistic Regression Model**

Logistic Regression is a supervised machine learning algorithm used for classification problems. It predicts the probability of a categorical dependent variable (binary outcome).
Key points:

- It uses the logistic (sigmoid) function to map predicted values between 0 and 1.

- Suitable for binary classification (e.g., diabetes: yes/no).

- The decision boundary is based on probability thresholds (typically 0.5).

## Code:

Step 1: Import necessary libraries

```python
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.metrics import accuracy_score
5 from sklearn.feature_selection import SelectKBest, f_classif
```

Step 2: Load dataset

```python
1 url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabetes.data.csv"
2 columns = ['Pregnancies','Glucose','BloodPressure','SkinThickness','Insulin',
3          'BMI','DiabetesPedigreeFunction','Age','Outcome']
4 df = pd.read_csv(url, names=columns)
```

Step 3: Prepare input (X) and output (y)

```python
1 X = df.drop('Outcome', axis=1)
2 y = df['Outcome']
```

Step 4: Train-test split

```python
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=1)
```

Step 5: Train Logistic Regression on all features

```python
1 model = LogisticRegression(max_iter=200)
2 model.fit(X_train, y_train)
```

```
     LogisticRegression        ⓘ ❓

LogisticRegression(max_iter=200)
```

Step 6: Predict and evaluate

```python
1 y_pred = model.predict(X_test)
2
3 print("\n" + "="*50)
4 print("📌 STEP 1: MODEL WITH ALL FEATURES")
5 print("="*50)
6 print("Accuracy with all features:", accuracy_score(y_test, y_pred))
```

```
==================================================
📌 STEP 1: MODEL WITH ALL FEATURES
==================================================
Accuracy with all features: 0.7792207792207793
```

Step 7: Feature Selection (top 5 features)

```python
1 selector = SelectKBest(score_func=f_classif, k=5)
2 X_new = selector.fit_transform(X, y)
3
4 selected_features = selector.get_support(indices=True)
5 selected_feature_names = [columns[i] for i in selected_features]
6
7 print("\n" + "="*50)
8 print("📌 STEP 2: FEATURE SELECTION RESULTS")
9 print("="*50)
10 print("Selected Feature Indices:", selected_features)
11 print("Selected Feature Names:", selected_feature_names)
```

```
==================================================
📌 STEP 2: FEATURE SELECTION RESULTS
==================================================
Selected Feature Indices: [0 1 5 6 7]
Selected Feature Names: ['Pregnancies', 'Glucose', 'BMI', 'DiabetesPedigreeFunction', 'Age']
```

Step 8: Train model with selected features

```
1 X_train_new, X_test_new, y_train_new, y_test_new = train_test_split(
2     X_new, y, test_size=0.2, random_state=1)
3
4 model2 = LogisticRegression(max_iter=200)
5 model2.fit(X_train_new, y_train_new)
```

```
▼    LogisticRegression    ⓘ ⓘ
LogisticRegression(max_iter=200)
```

Step 9: Evaluate new model

```
1 y_pred_new = model2.predict(X_test_new)
2
3 print("\n" + "="*50)
4 print("📌 STEP 3: MODEL WITH SELECTED TOP 5 FEATURES")
5 print("="*50)
6 print("Accuracy with selected features:", accuracy_score(y_test_new, y_pred_new))
7 print("="*50)
```

```
==================================================
📌 STEP 3: MODEL WITH SELECTED TOP 5 FEATURES
==================================================
Accuracy with selected features: 0.7662337662337663
==================================================
```

## Output Screenshots:

- **First accuracy (before feature selection)**

```
==================================================
📌 STEP 1: MODEL WITH ALL FEATURES
==================================================
Accuracy with all features: 0.7792207792207793
```

- **List of selected features**

```
==================================================
📌 STEP 2: FEATURE SELECTION RESULTS
==================================================
Selected Feature Indices: [0 1 5 6 7]
Selected Feature Names: ['Pregnancies', 'Glucose', 'BMI', 'DiabetesPedigreeFunction', 'Age']
```

- **Second accuracy (after feature selection)**

```
==================================================
📌 STEP 3: MODEL WITH SELECTED TOP 5 FEATURES
==================================================
Accuracy with selected features: 0.7662337662337663
==================================================
```

## Conclusion:

From the above practical, we successfully built a logistic regression model to predict diabetes using all features and then after applying feature selection using SelectKBest.