



Scikit Learn

Professor: Gopinath Panda

Group No: 19

Kathan Vora - 202411073
Aniruddha Shinde - 202411065
Monson Reji Verghese - 202411039

Date: October 24, 2024

Contents

1	Introduction	2
2	1) Importing Libraries & Data	2
2.1	Purpose	2
2.2	How it Works	2
2.3	Key Parameters	2
3	2) Outlier Detection and Treatment Using IQR	2
3.1	Purpose	2
3.2	How it Works	2
3.3	Key Parameters	3
3.4	Result	3
4	3) Numerical Feature Scaling Using Scikit-learn	3
4.1	Purpose	3
4.2	How it Works	3
4.3	Key Parameters	3
4.4	Code	3
5	4) Feature Selection Using SelectKBest and f_regression	4
5.1	Purpose	4
5.2	How it Works	4
5.3	Key Parameters	4
5.4	Code	4
5.5	Result	4
6	5) Principal Component Analysis (PCA)	4
6.1	Purpose	4
6.2	How it Works	4
6.3	Code	5
6.4	Result	5
7	6) Recursive Feature Elimination (RFE)	5
7.1	Purpose	5
7.2	How it Works	5
7.3	Code	5
7.4	Result	6
8	7) Feature Importance	6
8.1	Purpose	6
8.2	How it Works	6
8.3	Code	6
8.4	Result	6

1 Introduction

This report provides a step-by-step analysis of an insurance dataset using scikit-learn. The aim is to predict insurance charges based on key features such as age, BMI, and children. The workflow includes data loading, outlier detection and treatment, numerical feature scaling, feature selection, and model training. Scikit-learn is used throughout the process for seamless machine learning tasks.

2 1) Importing Libraries & Data

2.1 Purpose

The required Python libraries are imported, and the dataset is loaded for analysis.

2.2 How it Works

Libraries like pandas and numpy are used for data handling, while scikit-learn modules support data splitting, scaling, and modeling. The dataset is read from a URL into a DataFrame for further operations.

2.3 Key Parameters

- `pd.read_csv(url)`: Reads the dataset from a specified URL.
- `df.head()`: Displays the first five rows of the dataset.
- `df.describe()`: Provides summary statistics of the data.

3 2) Outlier Detection and Treatment Using IQR

3.1 Purpose

Outliers are detected using the Interquartile Range (IQR) method to improve model performance.

3.2 How it Works

IQR identifies outliers as data points that lie beyond 1.5 times the IQR from Q1 or Q3. After detection, rows with outliers are removed from the dataset. The following function is used for outlier removal:

```
def remove_outliers_iqr(df):
    numerical_columns = ['age', 'bmi', 'children', 'charges']
    for col in numerical_columns:
        q25, q75 = np.quantile(df[col], 0.25), np.quantile(df[col], 0.75)
        iqr = q75 - q25
        cut_off = iqr * 1.5
        lower, upper = q25 - cut_off, q75 + cut_off

        # Remove rows with outliers
        df = df[(df[col] >= lower) & (df[col] <= upper)]

    return df

# Apply the function to remove outliers
df_cleaned = remove_outliers_iqr(df)
```

3.3 Key Parameters

- Q1, Q3: Boundaries for the 25th and 75th percentiles.
- IQR: Difference between Q3 and Q1.
- Cut-off: 1.5 times the IQR for identifying outliers.

3.4 Result

- Total outliers detected: 148
- Outlier distribution:
 - Age: 0
 - BMI: 9
 - Children: 0
 - Charges: 139

After applying the outlier removal function, the dataset shape is (1191, 7), and the updated outlier distribution shows:

- Total outliers across all columns: 0
- Age: 0, BMI: 0, Children: 0, Charges: 0

4 3) Numerical Feature Scaling Using Scikit-learn

4.1 Purpose

Feature scaling ensures all numerical features are on a similar scale to prevent bias during model training.

4.2 How it Works

The StandardScaler from scikit-learn standardizes the numerical features by subtracting the mean and dividing by the standard deviation.

4.3 Key Parameters

- `fit()`: Computes the mean and standard deviation.
- `transform()`: Scales the data using the computed statistics.

4.4 Code

```
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
df_cleaned[['age', 'bmi']] = scaler.fit_transform(df_cleaned[['age', 'bmi']])
```

5 4) Feature Selection Using SelectKBest and f_regression

5.1 Purpose

SelectKBest identifies the top features correlated with the target variable using statistical tests.

5.2 How it Works

f_regression evaluates the relationship between each feature and the target, retaining only the top k features with the highest scores.

5.3 Key Parameters

- k: Number of features to select.
- Score function: f-statistics for feature relevance.

5.4 Code

```
from sklearn.feature_selection import SelectKBest, f_regression

X = df_cleaned[['age', 'bmi', 'children']]
y = df_cleaned['charges']

select_k_best = SelectKBest(score_func=f_regression, k=3)
X_selected = select_k_best.fit_transform(X, y)

# Get selected feature names
selected_features = X.columns[select_k_best.get_support()]
```

5.5 Result

- Selected features: ['age', 'bmi', 'children']
- Mean Squared Error: 20846598.562740713

6 5) Principal Component Analysis (PCA)

6.1 Purpose

PCA reduces the dimensionality of the dataset while preserving as much variance as possible.

6.2 How it Works

PCA transforms the original features into a new set of uncorrelated variables (principal components) ordered by the amount of variance they capture.

6.3 Code

```
from sklearn.decomposition import PCA
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error

# Define features (X) and target (y)
X = df_cleaned[['age', 'bmi', 'children']]
y = df_cleaned['charges']

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Initialize PCA for feature reduction
n_components = 3
pca = PCA(n_components=n_components)

# Fit PCA on the training data and transform both train and test sets
X_train_pca = pca.fit_transform(X_train)
X_test_pca = pca.transform(X_test)

# Print the explained variance ratio for each component
explained_variance = pca.explained_variance_ratio_
```

6.4 Result

- Explained variance ratio by each component: [0.37999397, 0.33316148, 0.28684455]
- Mean Squared Error: 20846598.562740713

7 6) Recursive Feature Elimination (RFE)

7.1 Purpose

RFE selects the most important features by recursively removing the least important ones.

7.2 How it Works

A model is trained on the dataset, and features are ranked based on their importance. The least important features are removed iteratively.

7.3 Code

```
from sklearn.feature_selection import RFE
from sklearn.linear_model import LinearRegression

# Define a model
model = LinearRegression()

# Perform RFE
rfe = RFE(model, n_features_to_select=3)
X_rfe = rfe.fit_transform(X, y)
```

```
# Get ranking of features
ranking = rfe.ranking_
```

7.4 Result

- Selected features: ['age', 'bmi', 'children']
- Mean Squared Error: 20846598.562740713

8 7) Feature Importance

8.1 Purpose

Feature importance indicates the contribution of each feature in predicting the target variable.

8.2 How it Works

Various models (like decision trees) provide built-in feature importance metrics, indicating the significance of each feature in the decision-making process.

8.3 Code

```
from sklearn.ensemble import RandomForestRegressor

# Fit a Random Forest model
rf_model = RandomForestRegressor()
rf_model.fit(X_train, y_train)

# Get feature importance
importances = rf_model.feature_importances_
```

8.4 Result

- Selected features: ['age', 'bmi']
- Mean Squared Error: 32634129.98330079

You can access our work here :[GitHub Repository](#)
[Google Colab Link](#) (please access using college ID): [Colab Link](#)