# Title :Predictive Modeling with Linear Regression

## 1. Understanding the Dataset:

Before building a predictive model, it's crucial to understand the dataset. This includes examining the structure of the data, identifying the features (independent variables) and the target variable (dependent variable), and understanding the meaning and distribution of each variable.

## 2. Data Preprocessing:

Data preprocessing involves cleaning and transforming the data to make it suitable for modeling. This may include handling missing values, encoding categorical variables (if any), scaling or normalizing numerical features, and dealing with outliers.

## 3. Splitting the Data:

The dataset is typically divided into two subsets: a training set and a test set. The training set is used to train the model, while the test set is used to evaluate its performance. Common split ratios include 70-80% of the data for training and the remaining 20-30% for testing.

## 4. Building the Linear Regression Model: ¶

Linear regression aims to find the best-fitting linear relationship between the input features and the target variable. Mathematically, the model tries to find the coefficients (weights) for each feature that minimize the difference between the actual target values and the predicted values. The model equation is of the form: $y = b_0 + b_1x_1 + b_2x_2 + ... + b_n*x_n$, where y is the predicted target variable, $b_0$ is the intercept, $b_1$ to $b_n$ are the coefficients, and $x_1$ to $x_n$ are the input features.

## 5. Training the Model:

During the training phase, the model is fed with the training data to learn the relationship between the input features and the target variable. The model adjusts its coefficients iteratively using optimization techniques (such as gradient descent) to minimize the prediction error.

## 6. Evaluating the Model:

Once the model is trained, it is evaluated using the test set to assess its performance and generalization ability. Common evaluation metrics for regression models include Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and R-squared (coefficient of determination).

## 7. Model Interpretation:

After evaluation, it's essential to interpret the model coefficients to understand the impact of each feature on the target variable. Positive coefficients indicate a positive relationship with the target variable, while negative coefficients indicate a negative relationship.

## 8. Model Deployment and Monitoring (optional):

Once the model is deemed satisfactory, it can be deployed to make predictions on new, unseen data. Continuous monitoring and periodic retraining may be necessary to maintain the model's performance over time.

**Linear regression is a fundamental and interpretable model widely used for predictive modeling when the relationship between features and the target variable is assumed to be linear. However, it's essential to consider its assumptions and limitations, such as the linearity assumption and sensitivity to outliers, when applying it to real-world datasets.**

**Step 1: Import necessary libraries**

```
In [1]: import pandas as pd
        import numpy as np
        from sklearn.model_selection import train_test_split
        from sklearn.linear_model import LinearRegression
        from sklearn.preprocessing import StandardScaler, OneHotEncoder
        from sklearn.compose import ColumnTransformer
        from sklearn.pipeline import Pipeline
        from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
        import matplotlib.pyplot as plt
```

## Step 2: Load the dataset

```
In [2]: data = pd.read_csv("insurance.csv")
```

## Step 3: Data Preprocessing

```
In [3]: # Handling missing values (if any)
        data.dropna(inplace=True)

        # Encoding categorical variables
        categorical_cols = ['sex', 'smoker', 'region']
        numerical_cols = [col for col in data.columns if col not in categorical_cols + ['charges']]

        # Define preprocessing steps for numerical and categorical features
        numerical_transformer = StandardScaler()
        categorical_transformer = OneHotEncoder(handle_unknown='ignore')

        # Combine preprocessing steps
        preprocessor = ColumnTransformer(
            transformers=[
                ('num', numerical_transformer, numerical_cols),
                ('cat', categorical_transformer, categorical_cols)
            ])
```

## Step 4: Splitting the data into train and test sets
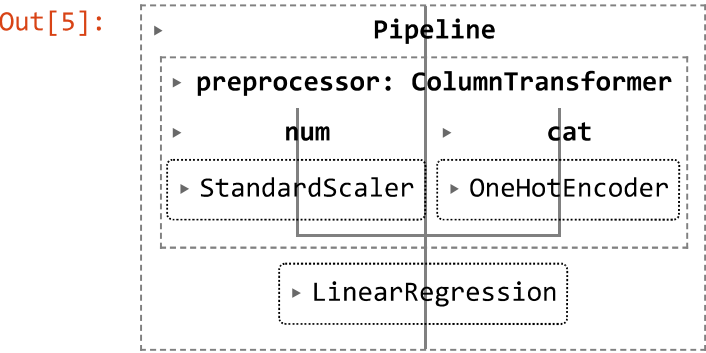
```
In [4]: X = data.drop('charges', axis=1) # Features
        y = data['charges'] # Target variable

        X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

## Step 5: Model building

In [5]:
```python
# Combine preprocessing with model
model = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('regressor', LinearRegression())
])

# Fit the model to the training data
model.fit(X_train, y_train)
```

Out[5]:

```
▸            Pipeline
  ▸ preprocessor: ColumnTransformer
  ▸      num          ▸      cat
  ▸ StandardScaler   ▸ OneHotEncoder

        ▸ LinearRegression
```

## Step 6: Model evaluation

In [6]:
```python
# Model evaluation
y_pred = model.predict(X_test)
mae = mean_absolute_error(y_test, y_pred)
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print("Mean Absolute Error:", mae)
print("Mean Squared Error:", mse)
print("R-squared:", r2)
```
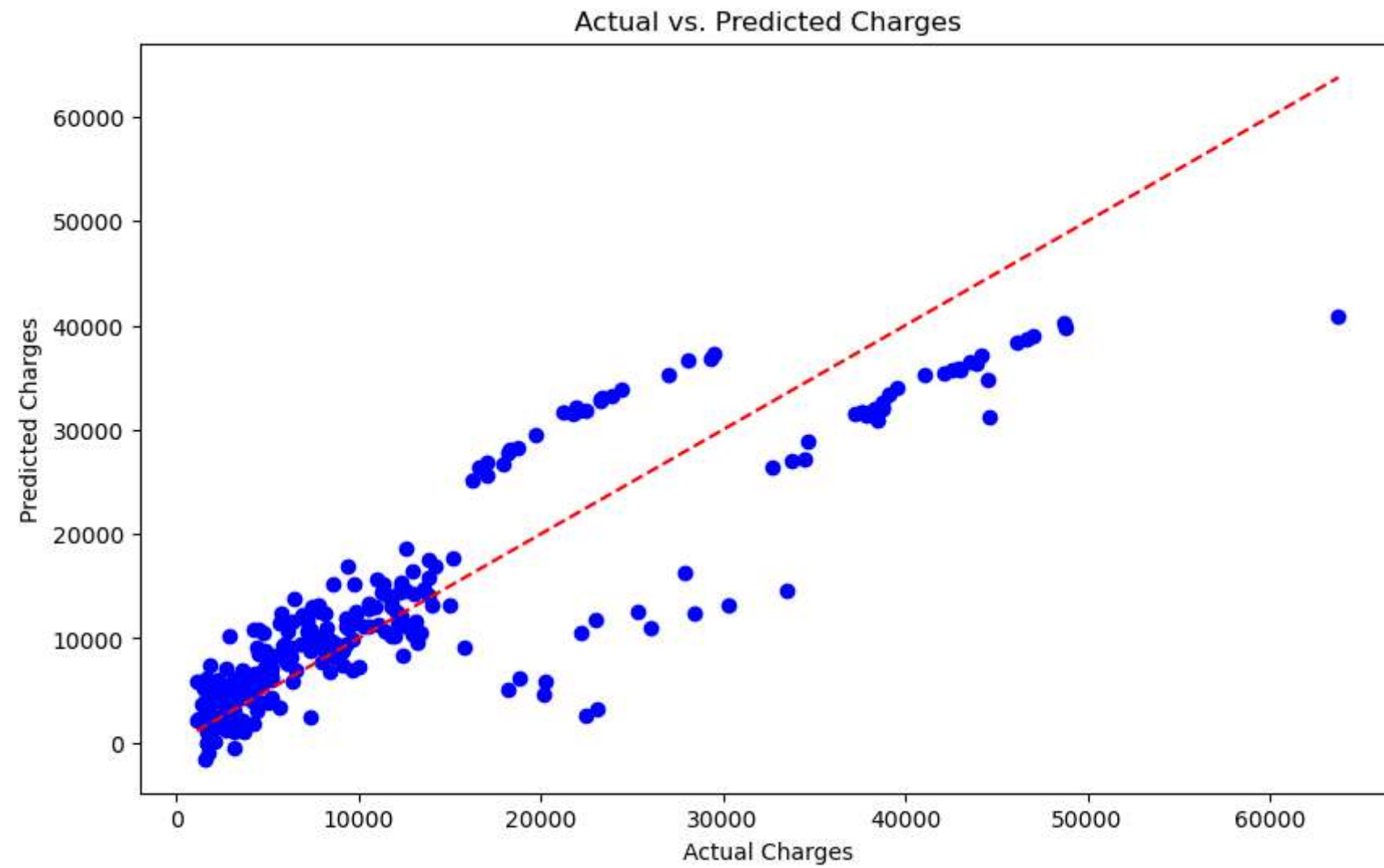
```
Mean Absolute Error: 4181.194473753652
Mean Squared Error: 33596915.85136148
R-squared: 0.7835929767120722
```

**Step 7: Visualization**

In [7]:
```python
# Visualize predictions vs. actual values
plt.figure(figsize=(10, 6))
plt.scatter(y_test, y_pred, color='blue')
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], linestyle='--', color='red')
plt.xlabel('Actual Charges')
plt.ylabel('Predicted Charges')
plt.title('Actual vs. Predicted Charges')
plt.show()
```



In [ ]: