```
In [ ]:  1  # ===============================================================
         2  # This code is part of Assignment 3 of ML lab (Executive M-Tech -ML Ass
         3  # Submitted by:
         4  #        IITP001300: Sukhvinder Singh  (email id: sukhvinder.malik13@gm
         5  #        IITP001316: Manjit Singh Duhan (email id: duhan.manjit@gmail.c
         6  #        IITP001508: Atul Singh (email id: atulsingh.xcvi@gmail.com)
         7  #===============================================================
```

```
In [1]:  1  import pandas as pd
         2  import numpy as np
         3  import matplotlib.pyplot as plt
         4  from sklearn.preprocessing import LabelEncoder
         5  from math import sqrt
```

```
In [2]:  1  # Load the dataset from CSV
         2  df = pd.read_csv('./insurance.csv')
         3
         4
         5  # Remove the duplicate entries and Do the re-indexing of all the data
         6  df.drop_duplicates(inplace=True)
         7  df.reset_index(drop=True, inplace=True)
         8
         9  df.head(5)
```

Out[2]:

|   | age | sex | bmi | children | smoker | region | charges |
|---|-----|-----|-----|----------|--------|--------|---------|
| 0 | 19 | female | 27.900 | 0 | yes | southwest | 16884.92400 |
| 1 | 18 | male | 33.770 | 1 | no | southeast | 1725.55230 |
| 2 | 28 | male | 33.000 | 3 | no | southeast | 4449.46200 |
| 3 | 33 | male | 22.705 | 0 | no | northwest | 21984.47061 |
| 4 | 32 | male | 28.880 | 0 | no | northwest | 3866.85520 |

```
In [3]:  1  # Convert strings to digits we need to know the unique values
         2  print("sex:    ", df['sex'].unique())
         3  print("smoker: ", df['smoker'].unique())
         4  print("region: ", df['region'].unique())
         5
```
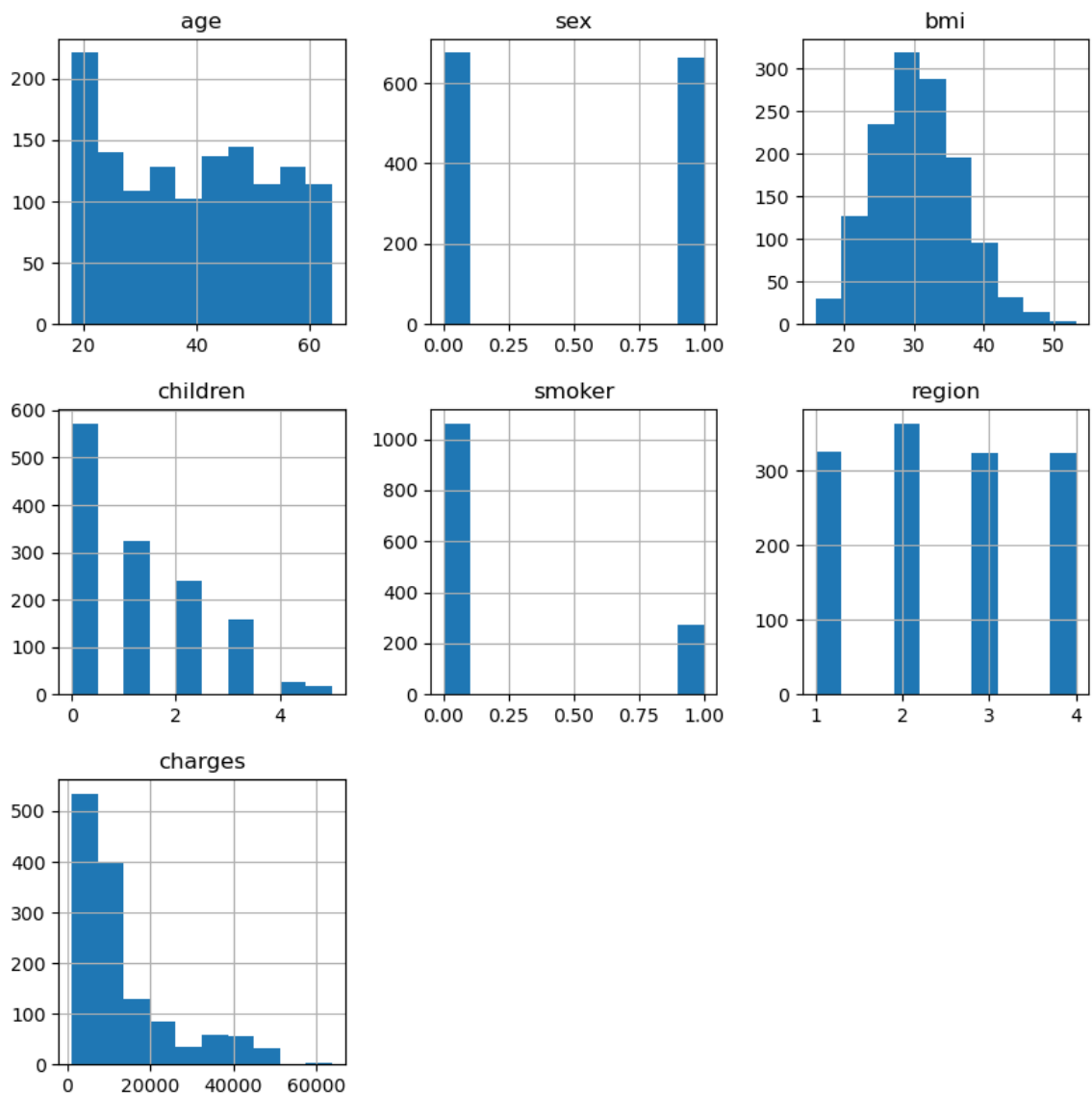
```
sex:    ['female' 'male']
smoker: ['yes' 'no']
region: ['southwest' 'southeast' 'northwest' 'northeast']
```

```
1  #update the strings with digit values
2  df['sex'] = df['sex'].apply({'male':0, 'female':1}.get)
3  df['smoker'] = df['smoker'].apply({'yes':1, 'no':0}.get)
4  df['region'] = df['region'].apply({'southwest':1, 'southeast':2, 'north
5
6  df.head(5)
```
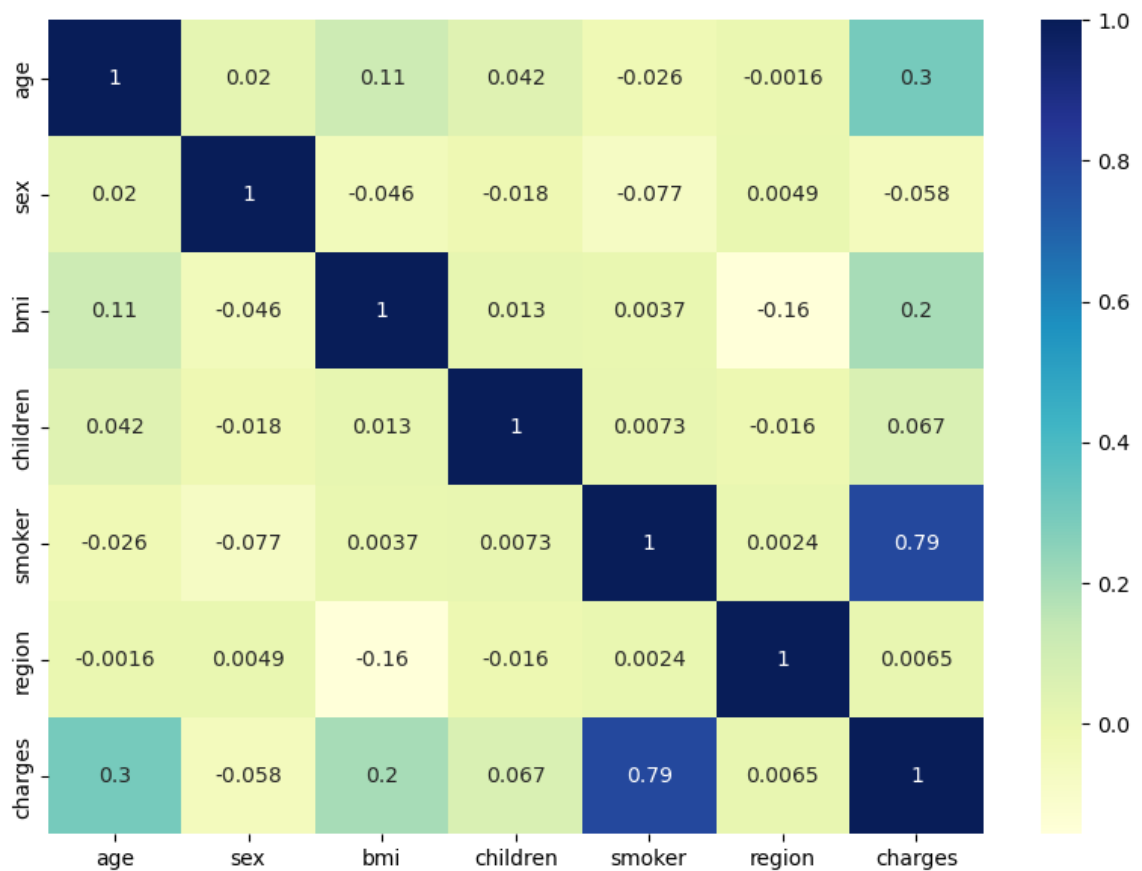
| | age | sex | bmi | children | smoker | region | charges |
|---|---|---|---|---|---|---|---|
| 0 | 19 | 1 | 27.900 | 0 | 1 | 1 | 16884.92400 |
| 1 | 18 | 0 | 33.770 | 1 | 0 | 2 | 1725.55230 |
| 2 | 28 | 0 | 33.000 | 3 | 0 | 2 | 4449.46200 |
| 3 | 33 | 0 | 22.705 | 0 | 0 | 3 | 21984.47061 |
| 4 | 32 | 0 | 28.880 | 0 | 0 | 3 | 3866.85520 |

```
1  #Let us see the data and uniformity of the data
2  df.hist(bins=10, figsize=(10, 10))
3  plt.show()
```

```
1  #Let us check correlation of the data
2  import seaborn as sns
3
4  plt.figure(figsize=(10,7))
5  sns.heatmap(df.corr(), annot = True, cmap="YlGnBu")
6  plt.show()
```



```
1  #We can see that the "smoker" feature is the most affecting feature to
   the charges.
2  #Beside Smoker, BMI and age is also hold very good corr. in deciding
   the charges.
3  #Lets take only BMI feature and refine it further.
```

```
1  print('bmi correlation : ', df['charges'].corr(df['bmi']))
2  print('smoker correlation : ', df['charges'].corr(df['smoker']))
3  print('age correlation : ', df['charges'].corr(df['age']))
```

```
bmi correlation :  0.1984008312262494
smoker correlation :  0.787234367280032
age correlation :  0.2983082125097864
```

Besed on these we can say that BMI is more accurate to move ahead so, we will take only BMI & charges

```python
In [9]:  1  from sklearn.model_selection import train_test_split
         2
         3  x_data = df["bmi"].values
         4  y_data = df["charges"].values
         5
         6  #Let us split the dataset to training and testing data
         7  x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, tes
```

```python
In [10]:  1  #Let us see the number of rows in each i.e. training & test data
          2  print("X_train shape: ", x_train.shape)
          3  print("X_test shape: ", x_test.shape)
          4  print("y_train shape: ", y_train.shape)
          5  print("y_test shape: ", y_test.shape)
```

```
X_train shape:  (935,)
X_test shape:  (402,)
y_train shape:  (935,)
y_test shape:  (402,)
```

```python
In [11]:   1  class myLinearRegression() :
           2
           3      def __init__( self ):
           4          self.b0 = 0
           5          self.b1 = 0
           6          self.predictions = list()
           7
           8      def fit(self, x, y):
           9          x_mean = sum(x) / float(len(x))
          10          y_mean = sum(y) / float(len(y))
          11          n = len(x)
          12          numerator = 0
          13          denominator = 0
          14
          15          for i in range(n):
          16              numerator += (x[i] - x_mean) * (y[i] - y_mean)
          17              denominator += (x[i] - x_mean) ** 2
          18
          19          b1 = numerator / denominator
          20          b0 = y_mean - (b1 * x_mean)
          21
          22          self.b0 = b0
          23          self.b1 = b1
          24
          25      def predict(self, x_test):
          26          self.predictions.clear()
          27          for row in x_test:
          28              yhat = self.b0 + self.b1 * x_test
          29              self.predictions.append(yhat)
          30          return self.predictions
          31
          32      def plot(self, X, Y):
          33          #plotting values
          34          x_max = np.max(X) + 5
          35          x_min = np.min(X) - 5
          36
          37          #calculating line values of x and y
          38          x = np.linspace(x_min, x_max, 1000)
          39          y = self.b0 + self.b1 * x
          40
          41          #plotting line
          42          plt.plot(x, y, color='red', label='Linear Regression')
          43
          44          #plot the data point
          45          plt.scatter(X, Y, color='blue', label='Data Point')
          46
          47          plt.xlabel('BMI')
          48          plt.ylabel('Charges')
          49
          50          plt.legend()
          51          plt.show()
          52
          53
          54      def r_square(self, x_test, y_test):
          55          y_mean = np.mean(y_test)
          56
          57          sumofsquares = 0
          58          sumofresiduals = 0
          59          n = len(x_test)
          60
          61          for i in range(n) :
```

```
62              y_pred = self.b0 + self.b1 * x_test[i]
63              sumofsquares += (y_test[i] - y_mean) ** 2
64              sumofresiduals += (y_test[i] - y_pred) **2
65
66          score  = 1 - (sumofresiduals/sumofsquares)
67
68          return score
```

In [12]:
```
1  my_lr = myLinearRegression()
2
3  my_lr.fit(x_train, y_train)
4  test_prediction = my_lr.predict(x_test)
5
6  score = my_lr.r_square(x_test, y_test)
7  print(score)
```
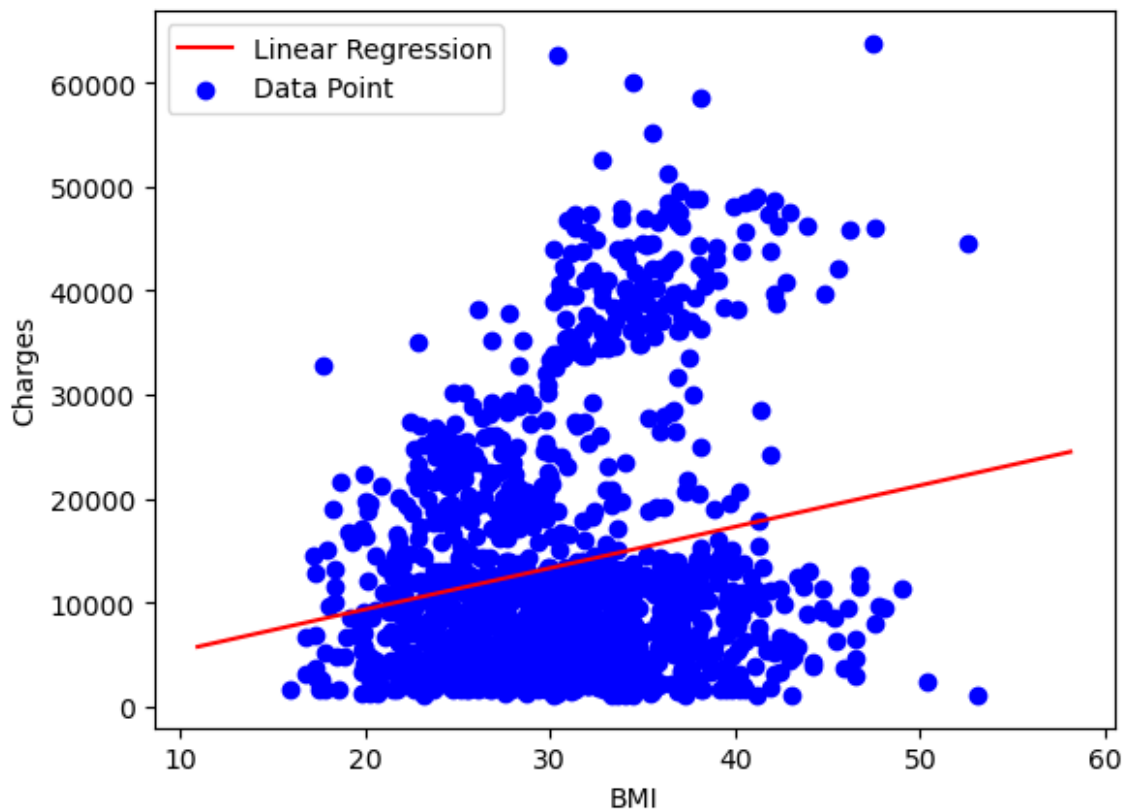
0.7275704833996746

In [13]:
```
1  my_lr.plot(x_data, y_data)
```



In [ ]:
```
1
```