

Problem Statement: To predict The Best Data Fits

Data Collection

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn import preprocessing,svm
6 from sklearn.model_selection import train_test_split
7 from sklearn.linear_model import LinearRegression
```

In [2]:

```
1 df=pd.read_csv(r"C:\Users\HP\Downloads\insurance (1).csv")
2 df
```

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
...
1333	50	male	30.970	3	no	northwest	10600.54830
1334	18	female	31.920	0	no	northeast	2205.98080
1335	18	female	36.850	0	no	southeast	1629.83350
1336	21	female	25.800	0	no	southwest	2007.94500
1337	61	female	29.070	0	yes	northwest	29141.36030

1338 rows × 7 columns

Data Cleaning and Preprocessing

In [3]:

```
1 df.head()
```

Out[3]:

	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 [4]:

```
1 df.tail()
```

Out[4]:

	age	sex	bmi	children	smoker	region	charges
1333	50	male	30.97	3	no	northwest	10600.5483
1334	18	female	31.92	0	no	northeast	2205.9808
1335	18	female	36.85	0	no	southeast	1629.8335
1336	21	female	25.80	0	no	southwest	2007.9450
1337	61	female	29.07	0	yes	northwest	29141.3603

In [5]:

```
1 df.shape
```

Out[5]:

(1338, 7)

In [6]:

```
1 df.describe()
```

Out[6]:

	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
...
1333	50	male	30.970	3	no	northwest	10600.54830
1334	18	female	31.920	0	no	northeast	2205.98080
1335	18	female	36.850	0	no	southeast	1629.83350
1336	21	female	25.800	0	no	southwest	2007.94500
1337	61	female	29.070	0	yes	northwest	29141.36030

[1338 rows x 7 columns]>

In [7]:

```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1338 entries, 0 to 1337
Data columns (total 7 columns):
 #   Column    Non-Null Count  Dtype  
--- 
 0   age        1338 non-null   int64  
 1   sex         1338 non-null   object  
 2   bmi         1338 non-null   float64 
 3   children   1338 non-null   int64  
 4   smoker      1338 non-null   object  
 5   region      1338 non-null   object  
 6   charges     1338 non-null   float64 
dtypes: float64(2), int64(2), object(3)
memory usage: 73.3+ KB
```

In [8]:

```
1 df.isnull().any()
```

Out[8]:

```
age      False
sex      False
bmi      False
children False
smoker   False
region   False
charges  False
dtype: bool
```

In [9]:

```
1 df.isna().sum()
```

Out[9]:

```
age      0
sex      0
bmi      0
children 0
smoker   0
region   0
charges  0
dtype: int64
```

In [10]:

```
1 df['region'].value_counts()
```

Out[10]:

```
region
southeast    364
southwest    325
northwest    325
northeast    324
Name: count, dtype: int64
```

In [11]:

```

1 convert={"sex":{"female":1,"male":0}}
2 df=df.replace(convert)
3 df

```

Out[11]:

	age	sex	bmi	children	smoker	region	charges
0	19	1	27.900	0	yes	southwest	16884.92400
1	18	0	33.770	1	no	southeast	1725.55230
2	28	0	33.000	3	no	southeast	4449.46200
3	33	0	22.705	0	no	northwest	21984.47061
4	32	0	28.880	0	no	northwest	3866.85520
...
1333	50	0	30.970	3	no	northwest	10600.54830
1334	18	1	31.920	0	no	northeast	2205.98080
1335	18	1	36.850	0	no	southeast	1629.83350
1336	21	1	25.800	0	no	southwest	2007.94500
1337	61	1	29.070	0	yes	northwest	29141.36030

1338 rows × 7 columns

In [41]:

```

1 convert={"region":{"southwest":1,"southeast":2,"northwest":3,"northeast":4}}
2 df=df.replace(convert)
3 df

```

Out[41]:

	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
...
1333	50	0	30.970	3	0	3	10600.54830
1334	18	1	31.920	0	0	4	2205.98080
1335	18	1	36.850	0	0	2	1629.83350
1336	21	1	25.800	0	0	1	2007.94500
1337	61	1	29.070	0	1	3	29141.36030

1338 rows × 7 columns

In [42]:

```
1 convert={"smoker":{"yes":1,"no":0}}
```

```
2 df=df.replace(convert)
```

```
3 df
```

Out[42]:

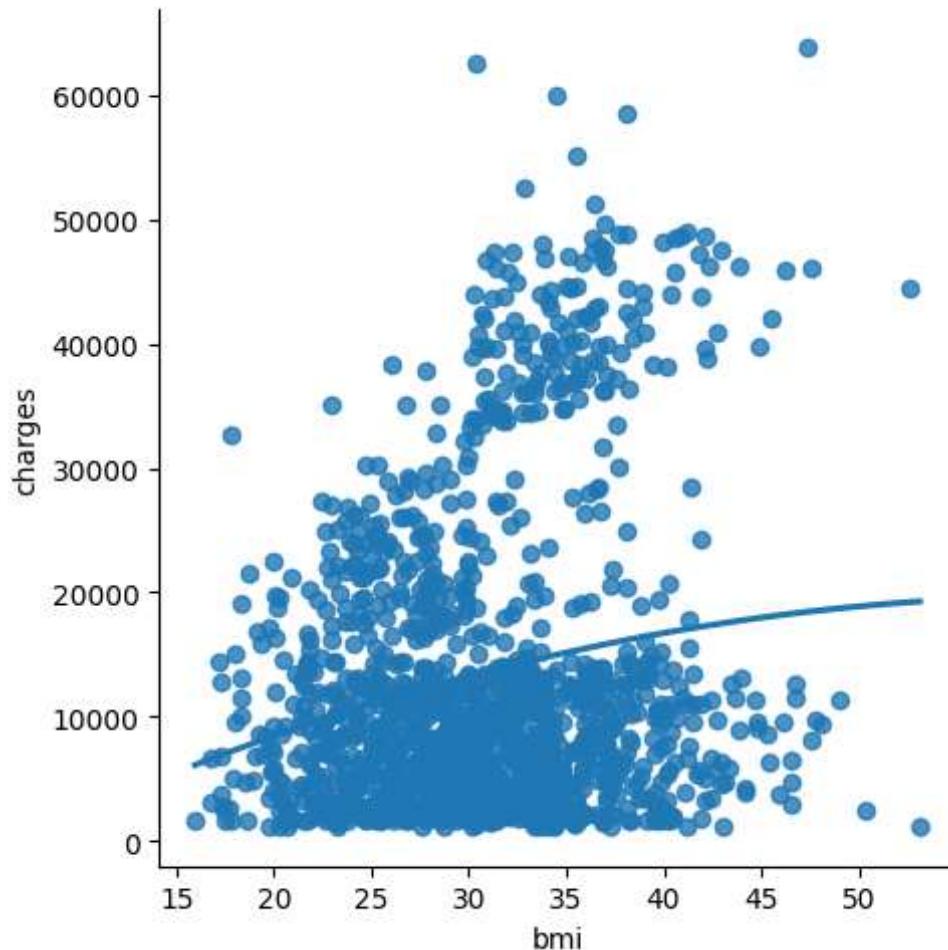
	age	sex	bmi	children	smoker	region	charges
0	19	1	27.900	0	1	1	16884.92400
1	18	0	33.770	1	0	2	1725.55230
2	28	0	33.000	3	0	2	4449.46200
3	33	0	22.705	0	0	3	21984.47061
4	32	0	28.880	0	0	3	3866.85520
...
1333	50	0	30.970	3	0	3	10600.54830
1334	18	1	31.920	0	0	4	2205.98080
1335	18	1	36.850	0	0	2	1629.83350
1336	21	1	25.800	0	0	1	2007.94500
1337	61	1	29.070	0	1	3	29141.36030

1338 rows × 7 columns

Data Visualization

In [43]:

```
1 sns.lmplot(x='bmi',y='charges',order=2,data=df,ci=None)
2 plt.show()
```



In [44]:

```
1 x=np.array(df['bmi']).reshape(-1,1)
2 y=x=np.array(df['charges']).reshape(-1,1)
```

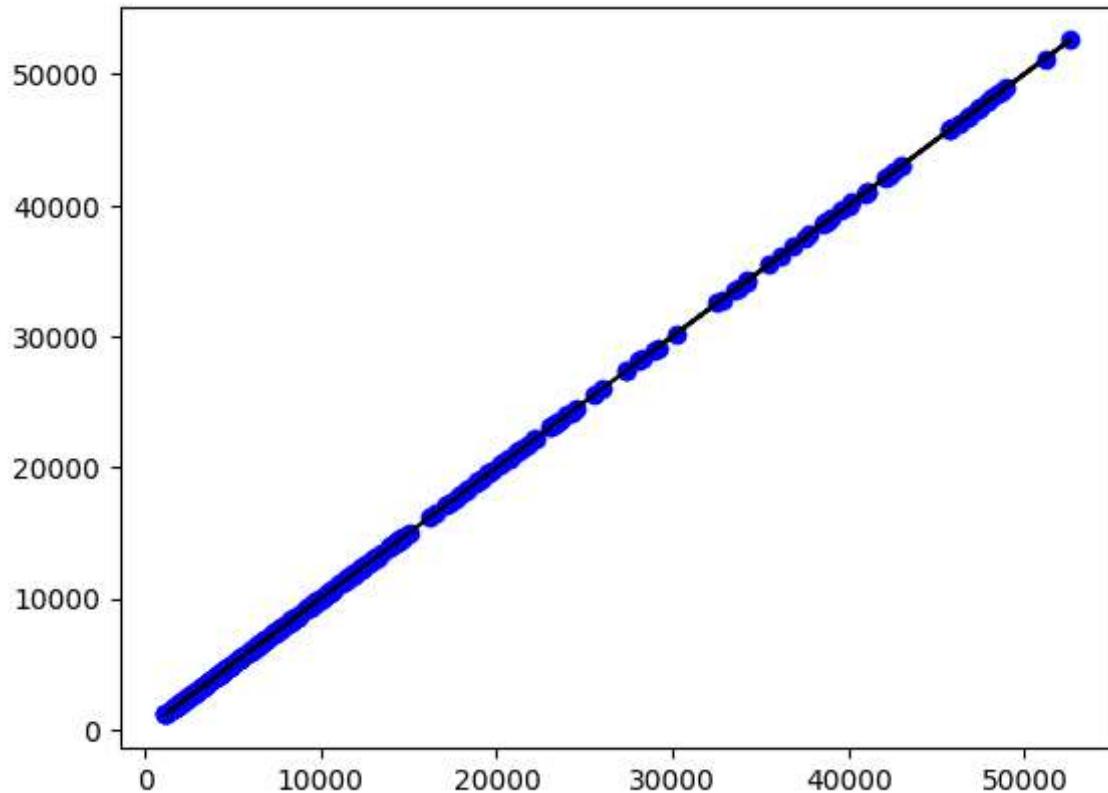
In [45]:

```
1 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25,random_state=0)
2 lr=LinearRegression()
3 lr.fit(x_train,y_train)
4 print(lr.score(x_test,y_test))
```

1.0

In [46]:

```
1 y_pred=lr.predict(x_test)
2 plt.scatter(x_test,y_test,color='b')
3 plt.plot(x_test,y_pred,color='k')
4 plt.show()
```



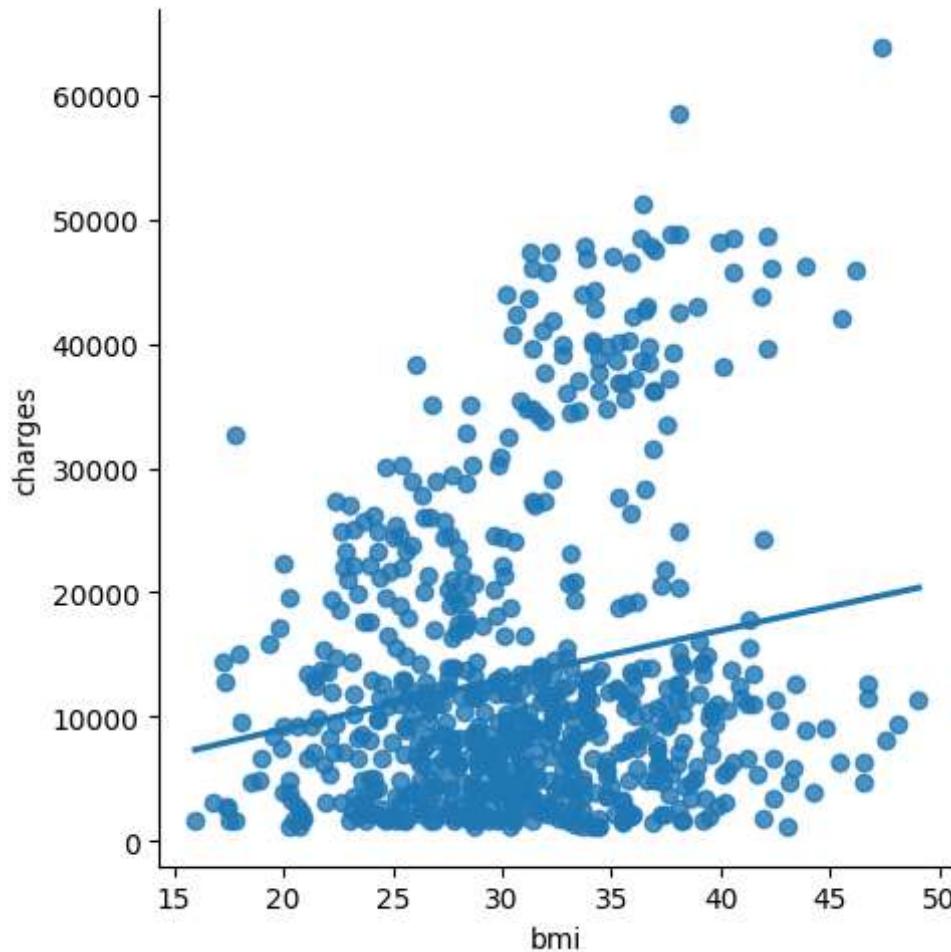
Working of Subset of Data

In [47]:

```

1 df700=df[:][:700]
2 sns.lmplot(x='bmi',y='charges',order=2,ci=None,data=df700)
3 plt.show()

```



In [48]:

```
1 df.fillna(method='ffill',inplace=True)
```

In [49]:

```

1 x=np.array(df700["bmi"]).reshape(-1,1)
2 y=np.array(df700['charges']).reshape(-1,1)

```

In [50]:

```
1 df700.dropna(inplace=True)
```

In [51]:

```

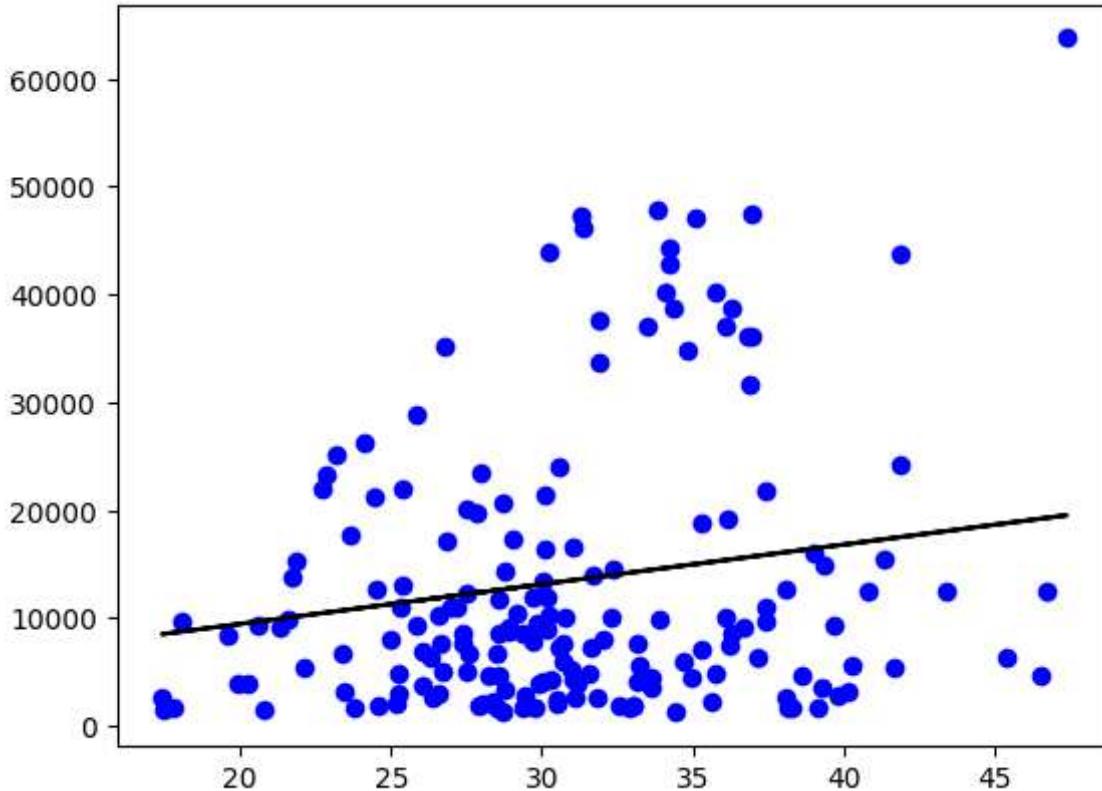
1 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25)
2 lr=LinearRegression()
3 lr.fit(x_train,y_train)
4 print(lr.score(x_test,y_test))

```

0.046811066569145576

In [52]:

```
1 y_pred=lr.predict(x_test)
2 plt.scatter(x_test,y_test,color='b')
3 plt.plot(x_test,y_pred,color='k')
4 plt.show()
```



Evaluation of Model

In [53]:

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.metrics import r2_score
```

In [54]:

```
1 lr=LinearRegression()
2 lr.fit(x_train,y_train)
3 y_pred=lr.predict(x_test)
4 r2=r2_score(y_test,y_pred)
5 print(r2)
```

0.046811066569145576

Ridge Regression

In [55]:

```

1 from sklearn.linear_model import Lasso, Ridge
2 from sklearn.preprocessing import StandardScaler

```

In [56]:

```

1 plt.figure(figsize = (10,10))
2 sns.heatmap(df700.corr(), annot=True)
3 plt.show()
4

```



In [57]:

```

1 features=df.columns[0:1]
2 target=df.columns[-1]

```

In [58]:

```

1 x=df[features].values
2 y=df[target].values
3 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.30,random_state=1)
4 print("The dimension of X_train is {}".format(x_train.shape))
5 print("The dimension of X_test is {}".format(x_test.shape))

```

The dimension of X_train is (936, 1)
The dimension of X_test is (402, 1)

In [59]:

```
1 lr = LinearRegression()
2 #Fit model
3 lr.fit(x_train, y_train)
4 #predict
5 actual = y_test
6 train_score_lr = lr.score(x_train, y_train)
7 test_score_lr = lr.score(x_test, y_test)
8 print("\nLinear Regression Model:\n")
9 print("The train score for lr model is {}".format(train_score_lr))
10 print("The test score for lr model is {}".format(test_score_lr))
```

Linear Regression Model:

The train score for lr model is 0.0910963973805714
The test score for lr model is 0.08490473916580776

In [60]:

```
1 ridgeReg = Ridge(alpha=10)
2 ridgeReg.fit(x_train,y_train)
3 #train and test score for ridge regression
4 train_score_ridge = ridgeReg.score(x_train, y_train)
5 test_score_ridge = ridgeReg.score(x_test, y_test)
6 print("\nRidge Model:\n")
7 print("The train score for ridge model is {}".format(train_score_ridge))
8 print("The test score for ridge model is {}".format(test_score_ridge))
```

Ridge Model:

The train score for ridge model is 0.09109639711159634
The test score for ridge model is 0.08490538609860176

In [61]:

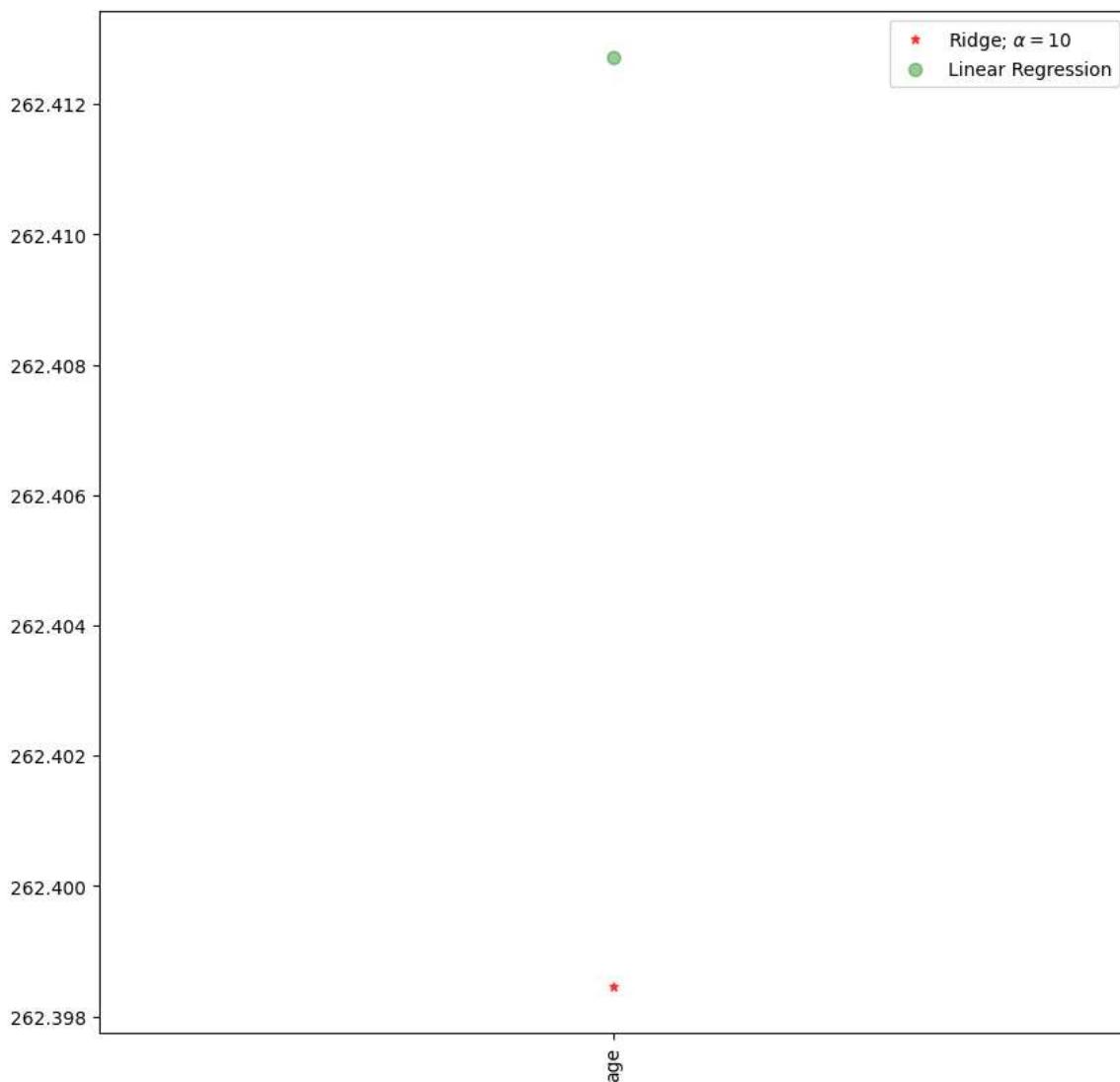
```
1 plt.figure(figsize=(10,10))
```

Out[61]:

```
<Figure size 1000x1000 with 0 Axes>
<Figure size 1000x1000 with 0 Axes>
```

In [62]:

```
1 plt.figure(figsize = (10, 10))
2 plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,color='red')
3 #plt.plot(rr100.coef_,alpha=0.5,linestyle='none',marker='d',markersize=6,color='blue')
4 plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersize=7,color='green')
5 plt.xticks(rotation = 90)
6 plt.legend()
7 plt.show()
```



Lasso Regression

In [63]:

```
1 lasso= Lasso(alpha=10)
2 lasso.fit(x_train,y_train)
3 #train and test score for ridge regression
4 train_score_ls = lasso.score(x_train, y_train)
5 test_score_ls= lasso.score(x_test, y_test)
6 print("\nRidge Model:\n")
7 print("The train score for lasso model is {}".format(train_score_ls))
8 print("The test score for lasso model is {}".format(test_score_ls))
```

Ridge Model:

The train score for lasso model is 0.09109639395809055
The test score for lasso model is 0.08490704421828055

In [64]:

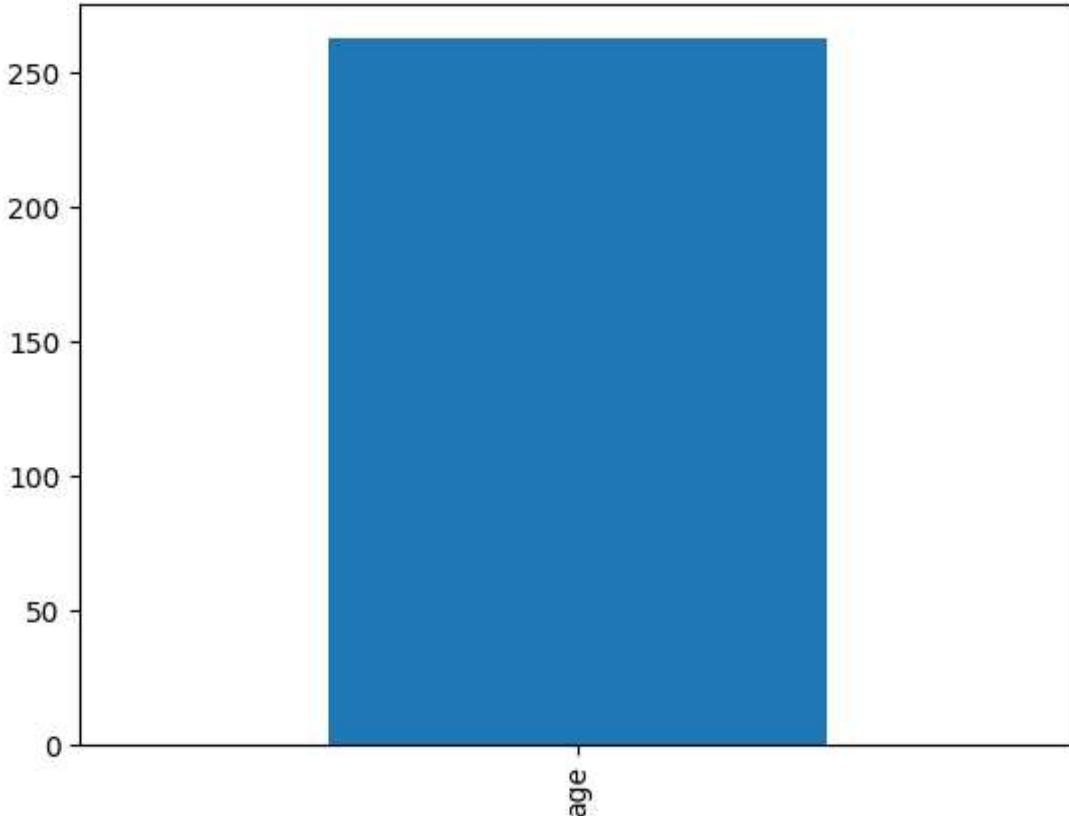
```
1 plt.figure(figsize=(10,10))
```

Out[64]:

<Figure size 1000x1000 with 0 Axes>
<Figure size 1000x1000 with 0 Axes>

In [65]:

```
1 pd.Series(lasso.coef_, features).sort_values(ascending = True).plot(kind = "bar")
2 plt.show()
```



In [66]:

```
1 from sklearn.linear_model import LassoCV
```

In [67]:

```
1 #using the Linear cv model
2 from sklearn.linear_model import RidgeCV
3 #cross validation
4 ridge_cv=RidgeCV(alphas =[0.0001,0.001,0.01,0.1,1,10]).fit(x_train,y_train)
5 #score
6 print(ridge_cv.score(x_train,y_train))
7 print(ridge_cv.score(x_test,y_test))
```

0.09109639711159612

0.08490538609884779

In [68]:

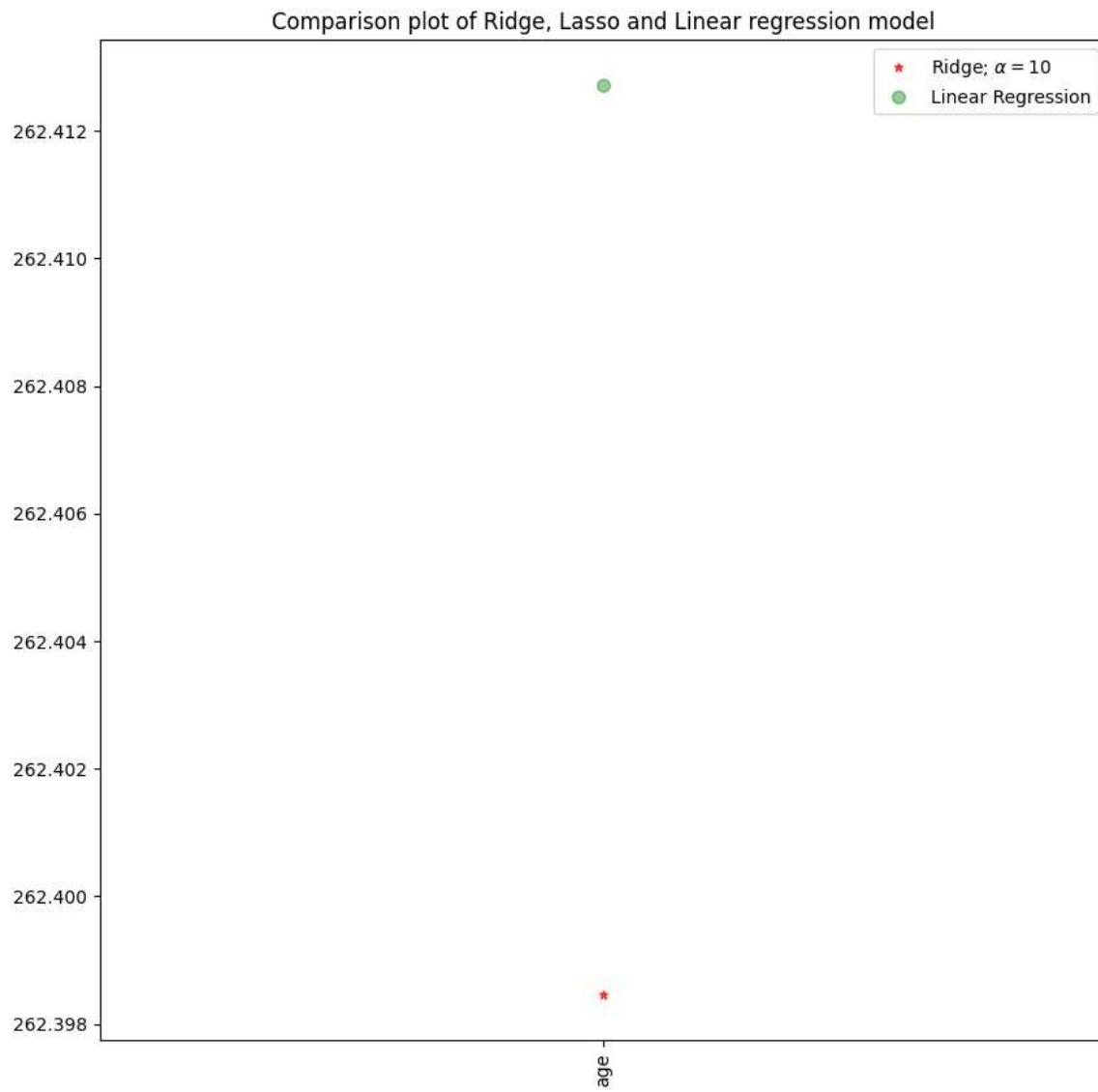
```
1 #using the Linear cv model
2 from sklearn.linear_model import LassoCV
3 #cross validation
4 lasso_cv=LassoCV(alphas =[0.0001,0.001,0.01,0.1,1,10]).fit(x_train,y_train)
5 #score
6 print(lasso_cv.score(x_train,y_train))
7 print(lasso_cv.score(x_test,y_test))
```

0.09109639395809055

0.08490704421828055

In [69]:

```
1 #plot size
2 plt.figure(figsize = (10, 10))
3 #add plot for ridge regression
4 plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,color='red')
5 #add plot for Lasso regression
6 #plt.plot(lasso_cv.coef_,alpha=0.5,linestyle='none',marker='d',markersize=6,color='blue')
7 #add plot for linear model
8 plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersize=7,color='green')
9 #rotate axis
10 plt.xticks(rotation = 90)
11 plt.legend()
12 plt.title("Comparison plot of Ridge, Lasso and Linear regression model")
13 plt.show()
```



Elastic Net Regression

In [70]:

```
1 from sklearn.linear_model import ElasticNet
```

In [71]:

```
1 el=ElasticNet()  
2 el.fit(x_train,y_train)  
3 print(el.coef_)  
4 print(el.intercept_)
```

```
[261.74450967]  
3115.083177426244
```

In [72]:

```
1 y_pred_elastic=el.predict(x_train)
```

In [73]:

```
1 mean_squared_error=np.mean((y_pred_elastic-y_train)**2)  
2 print(mean_squared_error)
```

```
135077142.70714515
```

In [74]:

```
1 el=ElasticNet()  
2 el.fit(x_train,y_train)  
3 print(el.score(x_train,y_train))
```

```
0.09109580670592365
```

Logistic Regression

In [75]:

```
1 import numpy as np  
2 import pandas as pd  
3 from sklearn.linear_model import LogisticRegression  
4 from sklearn.preprocessing import StandardScaler
```

In [76]:

```
1 df=pd.read_csv(r"C:\Users\HP\Downloads\insurance (1).csv")
2 df
```

Out[76]:

	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
...
1333	50	male	30.970	3	no	northwest	10600.54830
1334	18	female	31.920	0	no	northeast	2205.98080
1335	18	female	36.850	0	no	southeast	1629.83350
1336	21	female	25.800	0	no	southwest	2007.94500
1337	61	female	29.070	0	yes	northwest	29141.36030

1338 rows × 7 columns

In [77]:

```
1 df.shape
```

Out[77]:

(1338, 7)

In [78]:

```
1 pd.set_option('display.max_rows',10000000000)
2 pd.set_option('display.max_columns',10000000000)
3 pd.set_option('display.width',95)
```

In [79]:

```
1 print('This Dataset has %d rows and %d columns'%(df.shape))
```

This Dataset has 1338 rows and 7 columns

In [80]:

1 df.head()

Out[80]:

	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 [81]:

1 df.describe

Out[81]:

```
<bound method NDFrame.describe of
  ker      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
5      31   female  25.740      0     no  southeast  3756.621600
6      46   female  33.440      1     no  southeast  8240.589600
7      37   female  27.740      3     no  northwest  7281.505600
8      37     male  29.830      2     no  northeast  6406.410700
9      60   female  25.840      0     no  northwest  28923.136920
10     25     male  26.220      0     no  northeast  2721.320800
11     62   female  26.290      0    yes  southeast  27808.725100
12     23     male  34.400      0     no  southwest  1826.843000
13     56   female  39.820      0     no  southeast  11090.717800
14     27     male  42.130      0    yes  southeast  39611.757700
15     19     male  24.600      1     no  southwest  1837.237000
```

In [82]:

1 df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1338 entries, 0 to 1337
Data columns (total 7 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   age         1338 non-null   int64  
 1   sex         1338 non-null   object  
 2   bmi         1338 non-null   float64 
 3   children    1338 non-null   int64  
 4   smoker      1338 non-null   object  
 5   region      1338 non-null   object  
 6   charges     1338 non-null   float64 
dtypes: float64(2), int64(2), object(3)
memory usage: 73.3+ KB
```

In [83]:

```
1 df.isnull().sum()
```

Out[83]:

```
age      0  
sex      0  
bmi      0  
children  0  
smoker    0  
region    0  
charges   0  
dtype: int64
```

In [84]:

```
1 convert={"smoker":{"yes":1,"no":0}}  
2 df=df.replace(convert)  
3 df
```

Out[84]:

	age	sex	bmi	children	smoker	region	charges
0	19	female	27.900	0	1	southwest	16884.924000
1	18	male	33.770	1	0	southeast	1725.552300
2	28	male	33.000	3	0	southeast	4449.462000
3	33	male	22.705	0	0	northwest	21984.470610
4	32	male	28.880	0	0	northwest	3866.855200
5	31	female	25.740	0	0	southeast	3756.621600
6	46	female	33.440	1	0	southeast	8240.589600
7	37	female	27.740	3	0	northwest	7281.505600
8	37	male	29.830	2	0	northeast	6406.410700
9	60	female	25.840	0	0	northwest	28923.136920
10	25	male	26.220	0	0	northeast	2721.320800

In [85]:

```

1 convert={"sex":{"female":1,"male":0}}
2 df=df.replace(convert)
3 df

```

Out[85]:

	age	sex	bmi	children	smoker	region	charges
0	19	1	27.900	0	1	southwest	16884.924000
1	18	0	33.770	1	0	southeast	1725.552300
2	28	0	33.000	3	0	southeast	4449.462000
3	33	0	22.705	0	0	northwest	21984.470610
4	32	0	28.880	0	0	northwest	3866.855200
5	31	1	25.740	0	0	southeast	3756.621600
6	46	1	33.440	1	0	southeast	8240.589600
7	37	1	27.740	3	0	northwest	7281.505600
8	37	0	29.830	2	0	northeast	6406.410700
9	60	1	25.840	0	0	northwest	28923.136920
10	25	0	26.220	0	0	northeast	2721.320800

In [86]:

```

1 convert={"region":{"southeast":1,"southwest":2,"northeast":3,"northwest":4}}
2 df=df.replace(convert)
3 df

```

Out[86]:

	age	sex	bmi	children	smoker	region	charges
0	19	1	27.900	0	1	2	16884.924000
1	18	0	33.770	1	0	1	1725.552300
2	28	0	33.000	3	0	1	4449.462000
3	33	0	22.705	0	0	4	21984.470610
4	32	0	28.880	0	0	4	3866.855200
5	31	1	25.740	0	0	1	3756.621600
6	46	1	33.440	1	0	1	8240.589600
7	37	1	27.740	3	0	4	7281.505600
8	37	0	29.830	2	0	3	6406.410700
9	60	1	25.840	0	0	4	28923.136920
10	25	0	26.220	0	0	3	2721.320800

In [87]:

```
1 features_matrix=df.iloc[:,0:4]
```

In [88]:

```
1 target_vector=df.iloc[:, -3]
```

In [89]:

```
1 print('The Feature Matrix has %d Rows and %d columns(s)'%(features_matrix.shape))
2 print('The Target Matrix has %d Rows and %d columns(s)'%(np.array(target_vector).resh
```

The Feature Matrix has 1338 Rows and 4 columns(s)

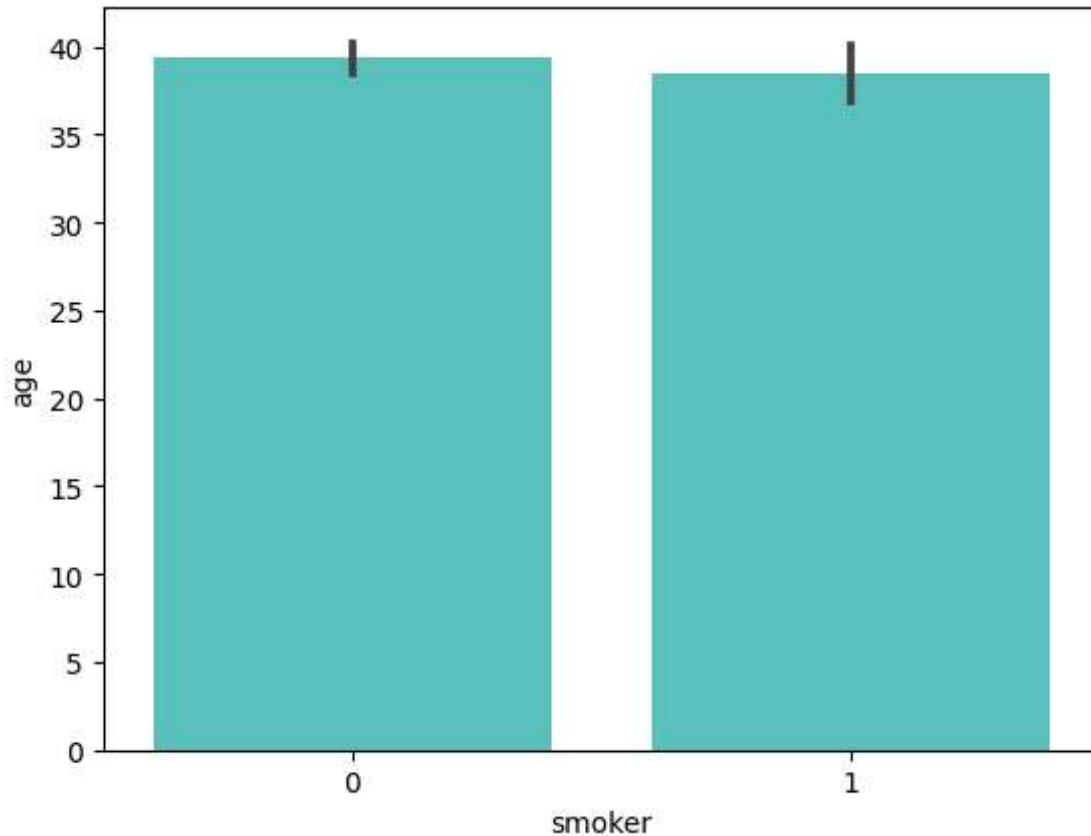
The Target Matrix has 1338 Rows and 1 columns(s)

In [90]:

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
```

In [91]:

```
1 sns.barplot(x='smoker', y='age', data=df, color="mediumturquoise")
2 plt.show()
```



In [92]:

```
1 features_matrix_standardized=StandardScaler().fit_transform(features_matrix)
```

In [93]:

```
1 algorithm=LogisticRegression(max_iter=10000)
```

In [94]:

```
1 Logistic_Regression_Model=algorithm.fit(features_matrix_standardized,target_vector)
```

In [95]:

```
1 observation=[[1,0,0.99539,-0.0588]]
```

In [96]:

```
1 predictions=Logistic_Regression_Model.predict(observation)
2 print('The model predicted the observation to belong to class %s'%(predictions))
```

The model predicted the observation to belong to class [0]

In [97]:

```
1 print('The algorithm was trained to predict one of the two classes:%s'%(algorithm.c))
```

The algorithm was trained to predict one of the two classes:[0 1]

In [98]:

```
1 print("""The Model says The probability of the observation we passed Belonging to class[0] Is 0.8057075871331396
2 print()
3 print("""The Model says The probability of the observation we passed Belonging to class[1] Is 0.19429241286686041
```

The Model says The probability of the observation we passed Belonging to class[0] Is 0.8057075871331396

The Model says The probability of the observation we passed Belonging to class[1] Is 0.19429241286686041

In [99]:

```
1 x=np.array(df['age']).reshape(-1,1)
2 y=np.array(df['smoker']).reshape(-1,1)
```

In [100]:

```
1 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.05)
2 lo=LogisticRegression()
3 lo.fit(x_train,y_train)
4 print(lo.score(x_test,y_test))
```

0.8507462686567164

C:\Users\HP\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\utils\validation.py:1143: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
y = column_or_1d(y, warn=True)

Decision Tree

In [101]:

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 from sklearn.model_selection import train_test_split
5 from sklearn.tree import DecisionTreeClassifier
```

In [102]:

```
1 df=pd.read_csv(r"C:\Users\HP\Downloads\insurance (1).csv")
2 df
```

Out[102]:

	age	sex	bmi	children	smoker	region	charges
0	19	female	27.900	0	yes	southwest	16884.924000
1	18	male	33.770	1	no	southeast	1725.552300
2	28	male	33.000	3	no	southeast	4449.462000
3	33	male	22.705	0	no	northwest	21984.470610
4	32	male	28.880	0	no	northwest	3866.855200
5	31	female	25.740	0	no	southeast	3756.621600
6	46	female	33.440	1	no	southeast	8240.589600
7	37	female	27.740	3	no	northwest	7281.505600
8	37	male	29.830	2	no	northeast	6406.410700
9	60	female	25.840	0	no	northwest	28923.136920
10	25	male	26.220	0	no	northeast	2721.320800

In [103]:

```
1 df.shape
```

Out[103]:

(1338, 7)

In [104]:

```
1 df.isnull().any()
```

Out[104]:

```
age      False
sex      False
bmi      False
children  False
smoker   False
region   False
charges  False
dtype: bool
```

In [105]:

```
1 df['region'].value_counts()
```

Out[105]:

```
region
southeast    364
southwest    325
northwest    325
northeast    324
Name: count, dtype: int64
```

In [106]:

```
1 convert={"sex":{"female":1,"male":0}}
2 df=df.replace(convert)
3 df
```

Out[106]:

	age	sex	bmi	children	smoker	region	charges
0	19	1	27.900	0	yes	southwest	16884.924000
1	18	0	33.770	1	no	southeast	1725.552300
2	28	0	33.000	3	no	southeast	4449.462000
3	33	0	22.705	0	no	northwest	21984.470610
4	32	0	28.880	0	no	northwest	3866.855200
5	31	1	25.740	0	no	southeast	3756.621600
6	46	1	33.440	1	no	southeast	8240.589600
7	37	1	27.740	3	no	northwest	7281.505600
8	37	0	29.830	2	no	northeast	6406.410700
9	60	1	25.840	0	no	northwest	28923.136920
10	25	0	26.220	0	no	northeast	2721.320800

In [107]:

```

1 convert={"smoker":{"yes":1,"no":0}}
2 df=df.replace(convert)
3 df

```

Out[107]:

	age	sex	bmi	children	smoker	region	charges
0	19	1	27.900	0	1	southwest	16884.924000
1	18	0	33.770	1	0	southeast	1725.552300
2	28	0	33.000	3	0	southeast	4449.462000
3	33	0	22.705	0	0	northwest	21984.470610
4	32	0	28.880	0	0	northwest	3866.855200
5	31	1	25.740	0	0	southeast	3756.621600
6	46	1	33.440	1	0	southeast	8240.589600
7	37	1	27.740	3	0	northwest	7281.505600
8	37	0	29.830	2	0	northeast	6406.410700
9	60	1	25.840	0	0	northwest	28923.136920
10	25	0	26.220	0	0	northeast	2721.320800

In [108]:

```

1 x=["bmi","children"]
2 y=["Yes","No"]
3 all_inputs=df[x]
4 all_classes=df["sex"]

```

In [109]:

```

1 (x_train,x_test,y_train,y_test)=train_test_split(all_inputs,all_classes,test_size=0.6)

```

In [110]:

```

1 clf=DecisionTreeClassifier(random_state=0)

```

In [111]:

```

1 clf.fit(x_train,y_train)

```

Out[111]:

```
DecisionTreeClassifier(random_state=0)
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [112]:

```
1 score=clf.score(x_test,y_test)
2 print(score)
```

0.4878048780487805

Random Forest

In [113]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt ,seaborn as sns
```

In [114]:

```
1 df=pd.read_csv(r"C:\Users\HP\Downloads\insurance (1).csv")
2 df
```

Out[114]:

	age	sex	bmi	children	smoker	region	charges
0	19	female	27.900	0	yes	southwest	16884.924000
1	18	male	33.770	1	no	southeast	1725.552300
2	28	male	33.000	3	no	southeast	4449.462000
3	33	male	22.705	0	no	northwest	21984.470610
4	32	male	28.880	0	no	northwest	3866.855200
5	31	female	25.740	0	no	southeast	3756.621600
6	46	female	33.440	1	no	southeast	8240.589600
7	37	female	27.740	3	no	northwest	7281.505600
8	37	male	29.830	2	no	northeast	6406.410700
9	60	female	25.840	0	no	northwest	28923.136920
10	25	male	26.220	0	no	northeast	2721.320800

In [115]:

```
1 df.shape
```

Out[115]:

(1338, 7)

In [116]:

```
1 df['region'].value_counts()
```

Out[116]:

```
region
southeast    364
southwest    325
northwest    325
northeast    324
Name: count, dtype: int64
```

In [117]:

```
1 df['bmi'].value_counts()
```

Out[117]:

```
bmi
32.300      13
28.310       9
30.495       8
30.875       8
31.350       8
30.800       8
34.100       8
28.880       8
33.330       7
35.200       7
25.800       7
32.775       7
27.645       7
32.110       7
38.060       7
25.460       7
30.590       7
```

In [118]:

```

1 m={"sex":{"female":1,"male":0}}
2 df=df.replace(m)
3 print(df)

```

	age	sex	bmi	children	smoker	region	charges
0	19	1	27.900	0	yes	southwest	16884.924000
1	18	0	33.770	1	no	southeast	1725.552300
2	28	0	33.000	3	no	southeast	4449.462000
3	33	0	22.705	0	no	northwest	21984.470610
4	32	0	28.880	0	no	northwest	3866.855200
5	31	1	25.740	0	no	southeast	3756.621600
6	46	1	33.440	1	no	southeast	8240.589600
7	37	1	27.740	3	no	northwest	7281.505600
8	37	0	29.830	2	no	northeast	6406.410700
9	60	1	25.840	0	no	northwest	28923.136920
10	25	0	26.220	0	no	northeast	2721.320800
11	62	1	26.290	0	yes	southeast	27808.725100
12	23	0	34.400	0	no	southwest	1826.843000
13	56	1	39.820	0	no	southeast	11090.717800
14	27	0	42.130	0	yes	southeast	39611.757700
15	19	0	24.600	1	no	southwest	1837.237000
16	52	1	30.780	1	no	northeast	10797.336200
17	23	0	23.845	0	no	northeast	2395.171550
18	56	0	40.200	0	no	southwest	10600.385000

In [119]:

```

1 n={"smoker":{"yes":1,"no":0}}
2 df=df.replace(n)
3 print(df)

```

	age	sex	bmi	children	smoker	region	charges
0	19	1	27.900	0	1	southwest	16884.924000
1	18	0	33.770	1	0	southeast	1725.552300
2	28	0	33.000	3	0	southeast	4449.462000
3	33	0	22.705	0	0	northwest	21984.470610
4	32	0	28.880	0	0	northwest	3866.855200
5	31	1	25.740	0	0	southeast	3756.621600
6	46	1	33.440	1	0	southeast	8240.589600
7	37	1	27.740	3	0	northwest	7281.505600
8	37	0	29.830	2	0	northeast	6406.410700
9	60	1	25.840	0	0	northwest	28923.136920
10	25	0	26.220	0	0	northeast	2721.320800
11	62	1	26.290	0	1	southeast	27808.725100
12	23	0	34.400	0	0	southwest	1826.843000
13	56	1	39.820	0	0	southeast	11090.717800
14	27	0	42.130	0	1	southeast	39611.757700
15	19	0	24.600	1	0	southwest	1837.237000
16	52	1	30.780	1	0	northeast	10797.336200
17	23	0	23.845	0	0	northeast	2395.171550
18	56	0	40.200	0	0	southwest	10600.385000

In [120]:

```
1 from sklearn.ensemble import RandomForestClassifier
2 rfc=RandomForestClassifier()
3 rfc.fit(x_train,y_train)
```

Out[120]:

RandomForestClassifier()

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [121]:

```
1 rf=RandomForestClassifier()
2 params={'max_depth':[2,3,5,20],
3         'min_samples_leaf':[5,10,20,50,100,200],
4         'n_estimators':[10,25,30,50,100,200]}
```

In [122]:

```
1 from sklearn.model_selection import GridSearchCV
2 grid_search=GridSearchCV(estimator=rf,param_grid=params,cv=2,scoring="accuracy")
3 grid_search.fit(x_train,y_train)
```

Out[122]:

```
GridSearchCV(cv=2, estimator=RandomForestClassifier(),
             param_grid={'max_depth': [2, 3, 5, 20],
                         'min_samples_leaf': [5, 10, 20, 50, 100, 200],
                         'n_estimators': [10, 25, 30, 50, 100, 200]},
             scoring='accuracy')
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [123]:

```
1 grid_search.best_score_
```

Out[123]:

0.5235250337651468

In [124]:

```
1 rf_best=grid_search.best_estimator_
2 print(rf_best)
```

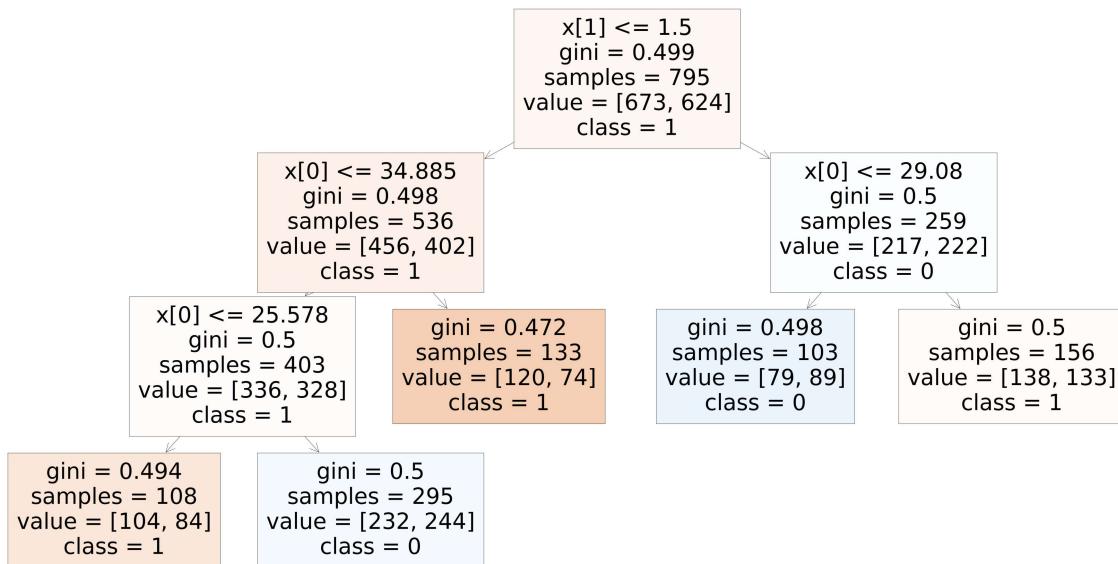
RandomForestClassifier(max_depth=3, min_samples_leaf=100, n_estimators=30)

In [125]:

```

1 from sklearn.tree import plot_tree
2 plt.figure(figsize=(80,40))
3 plot_tree(rf_best.estimators_[4], class_names=['1','0'], filled=True);

```

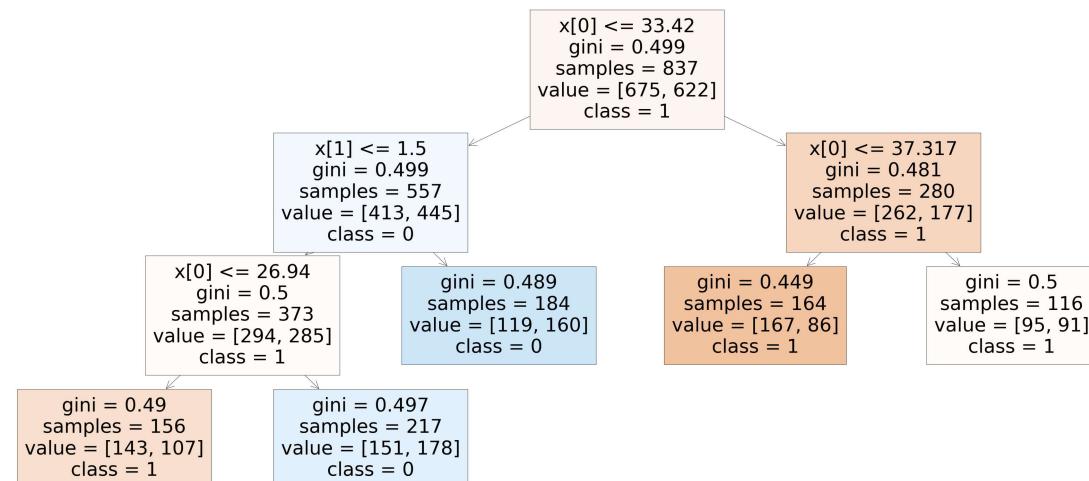


In [126]:

```

1 from sklearn.tree import plot_tree
2 plt.figure(figsize=(70,30))
3 plot_tree(rf_best.estimators_[6], class_names=["1","0"], filled=True);

```



In [127]:

```
1 rf_best.feature_importances_
```

Out[127]:

```
array([0.84852841, 0.15147159])
```

In [128]:

```
1 rf=RandomForestClassifier(random_state=0)
```

In [129]:

```
1 rf.fit(x_train,y_train)
```

Out[129]:

```
RandomForestClassifier(random_state=0)
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [130]:

```
1 score=rf.score(x_test,y_test)
2 print(score)
```

```
0.5365853658536586
```

CONCLUSION:

For the given insurance data set have performed linear, logistic, random forest and decision tree models of regression and classifications. #conclude that the most accuracy is occurred in logistic regression, i.e 85percent #conclude that the Logistic Regression model is best fit for given data

In []:

```
1
```

In []:

```
1
```

In []:

```
1
```

In []:

```
1
```