

Problem statement: To predict How Best the data fits

1) Data collection

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn import preprocessing,svm
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
```

```
In [2]: s=pd.read_csv(r"C:\Users\mouni\Downloads\insurance.csv")
s
```

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

2) Data cleaning and Preprocessing

#Exploratory data analysis

```
In [3]: s.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]: s.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]: s.shape
```

Out[5]: (1338, 7)

```
In [6]: s.describe
```

Out[6]: <bound method NDFrame.describe of
region charges

	age	sex	bmi	children	smoker		
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]: s.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]: s.isnull().any()
```

```
Out[8]: age      False
         sex      False
         bmi      False
         children False
         smoker   False
         region   False
         charges  False
dtype: bool
```

```
In [9]: s.isna().sum()
```

```
Out[9]: age      0
         sex      0
         bmi      0
         children 0
         smoker   0
         region   0
         charges  0
dtype: int64
```

```
In [10]: s['region'].value_counts()
```

```
Out[10]: region
         southeast    364
         southwest    325
         northwest    325
         northeast    324
Name: count, dtype: int64
```

```
In [11]: convert={"sex":{"female":1,"male":0}}
s=s.replace(convert)
s
```

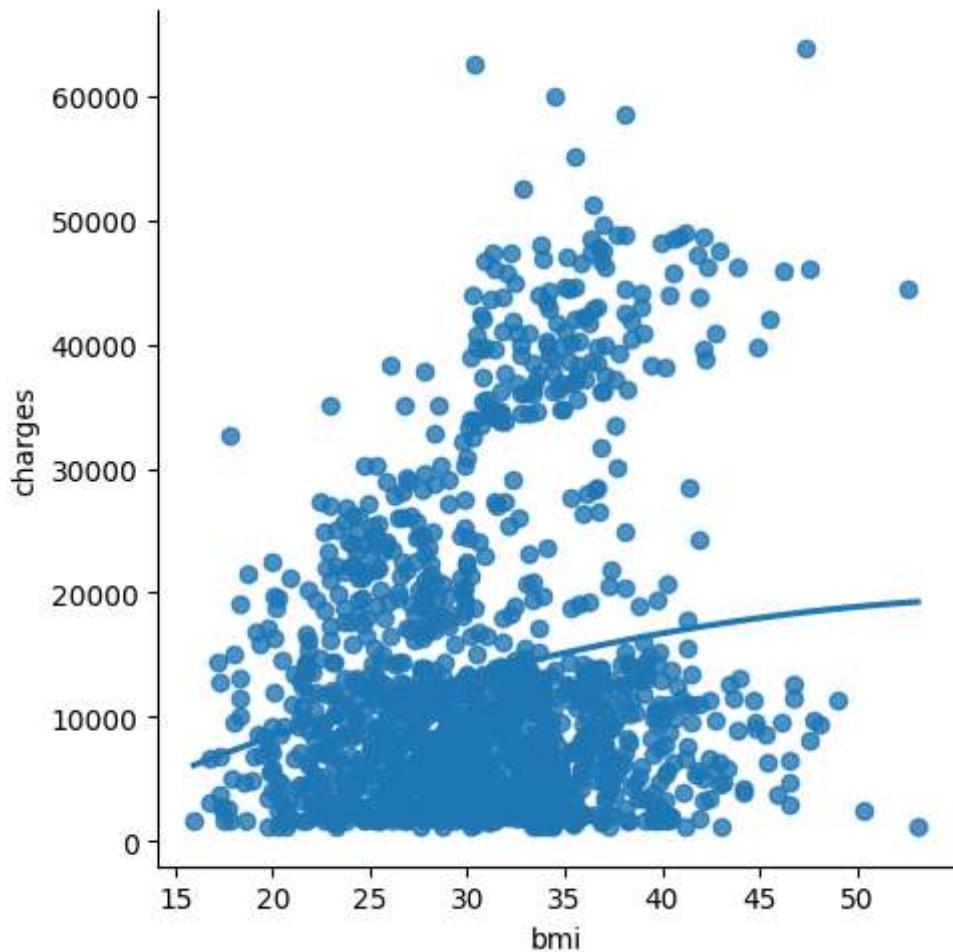
Out[11]:

	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

3)Data Visualization

```
In [12]: sns.lmplot(x='bmi',y='charges',order=2,data=s,ci=None)
plt.show()
```

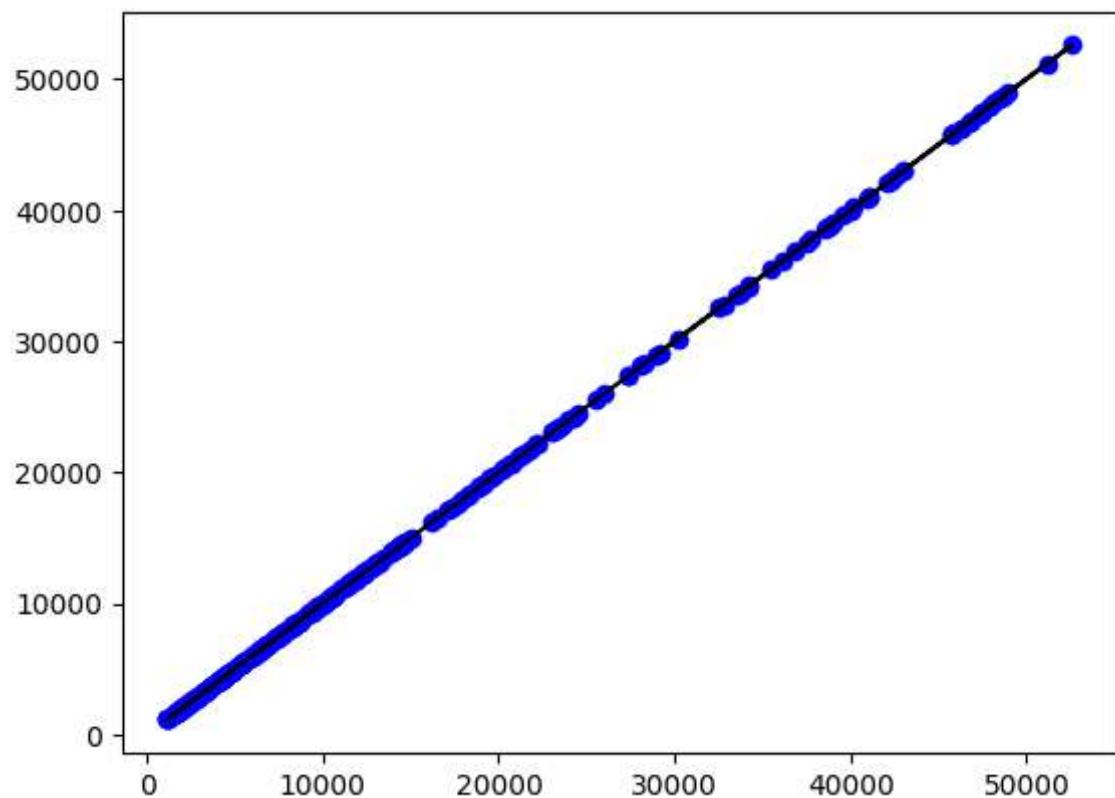


```
In [13]: x=np.array(s['bmi']).reshape(-1,1)
y=x=np.array(s['charges']).reshape(-1,1)
```

```
In [14]: x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25,random_state=42)
lr=LinearRegression()
lr.fit(x_train,y_train)
print(lr.score(x_test,y_test))
```

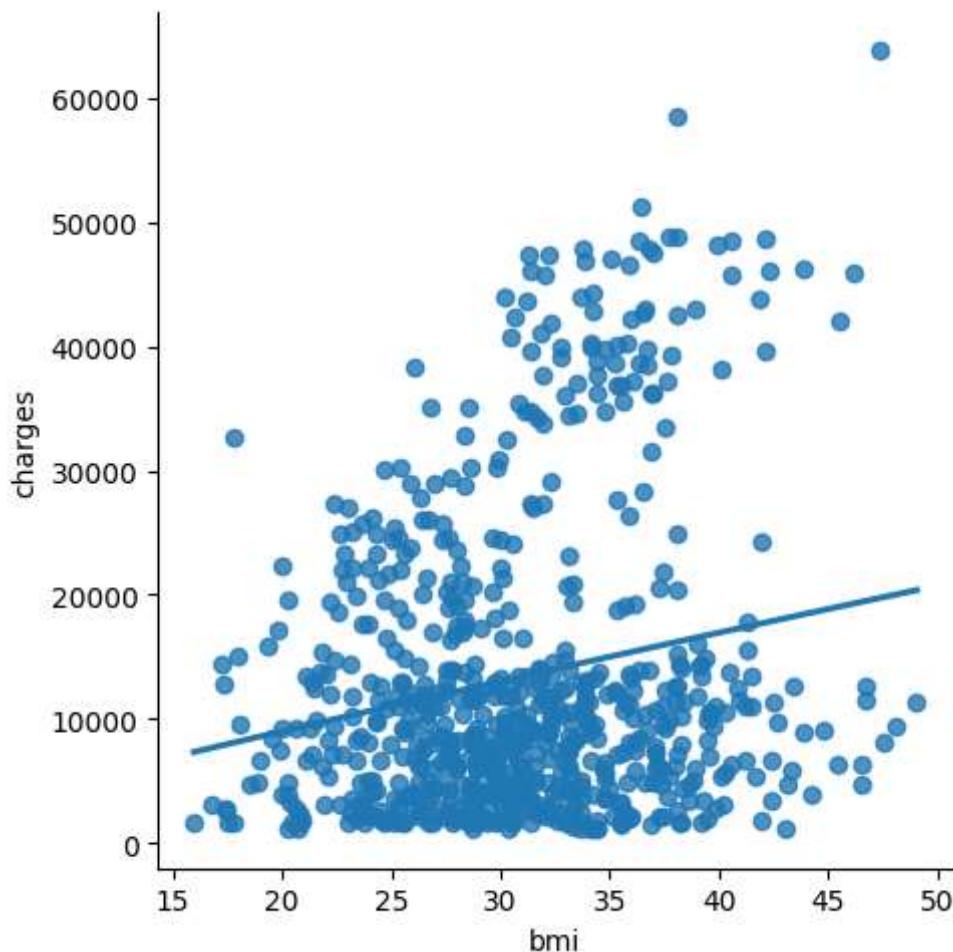
1.0

```
In [15]: y_pred=lr.predict(x_test)
plt.scatter(x_test,y_test,color='b')
plt.plot(x_test,y_pred,color='k')
plt.show()
```



working with subset of data

```
In [67]: s=s[:][:700]
sns.lmplot(x='bmi',y='charges',order=2,ci=None,data=s700)
plt.show()
```



```
In [68]: s.fillna(method='ffill',inplace=True)
```

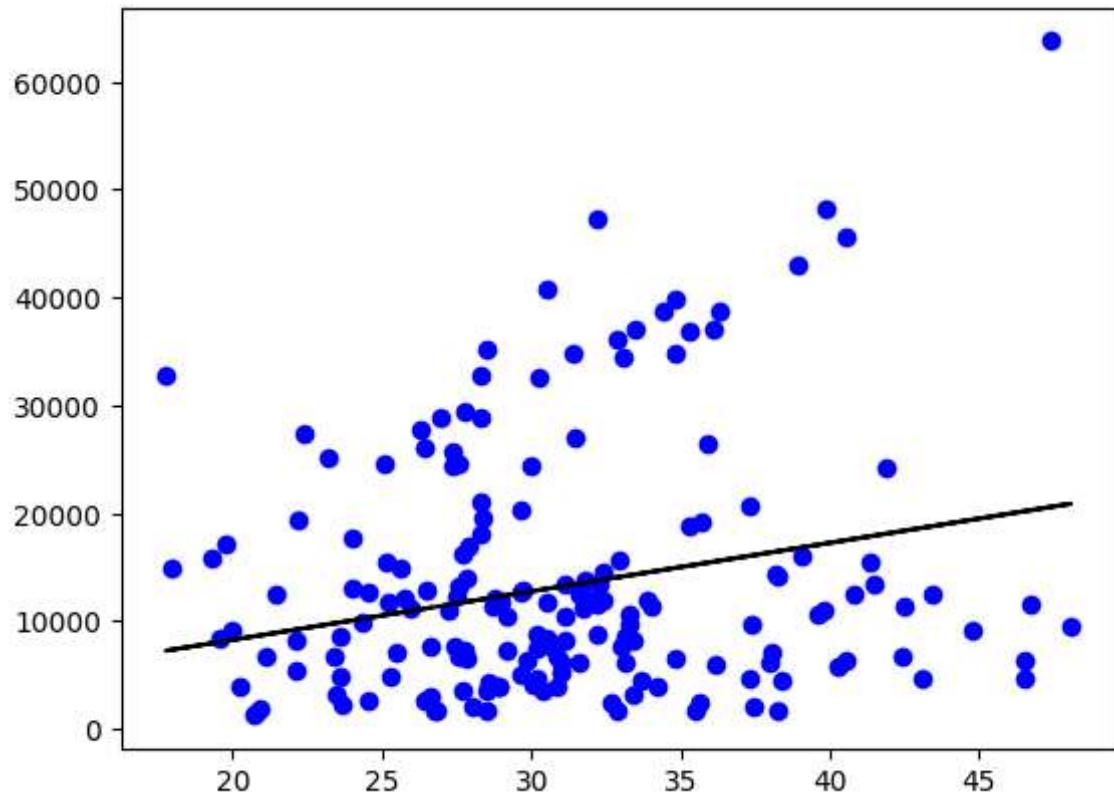
```
In [69]: x=np.array(s["bmi"]).reshape(-1,1)
y=np.array(s['charges']).reshape(-1,1)
```

```
In [70]: s.dropna(inplace=True)
```

```
In [71]: x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25)
lr=LinearRegression()
lr.fit(x_train,y_train)
print(lr.score(x_test,y_test))
```

0.0021400201329184743

```
In [72]: y_pred=lr.predict(x_test)
plt.scatter(x_test,y_test,color='b')
plt.plot(x_test,y_pred,color='k')
plt.show()
```



Evaluation of model

```
In [73]: from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score
```

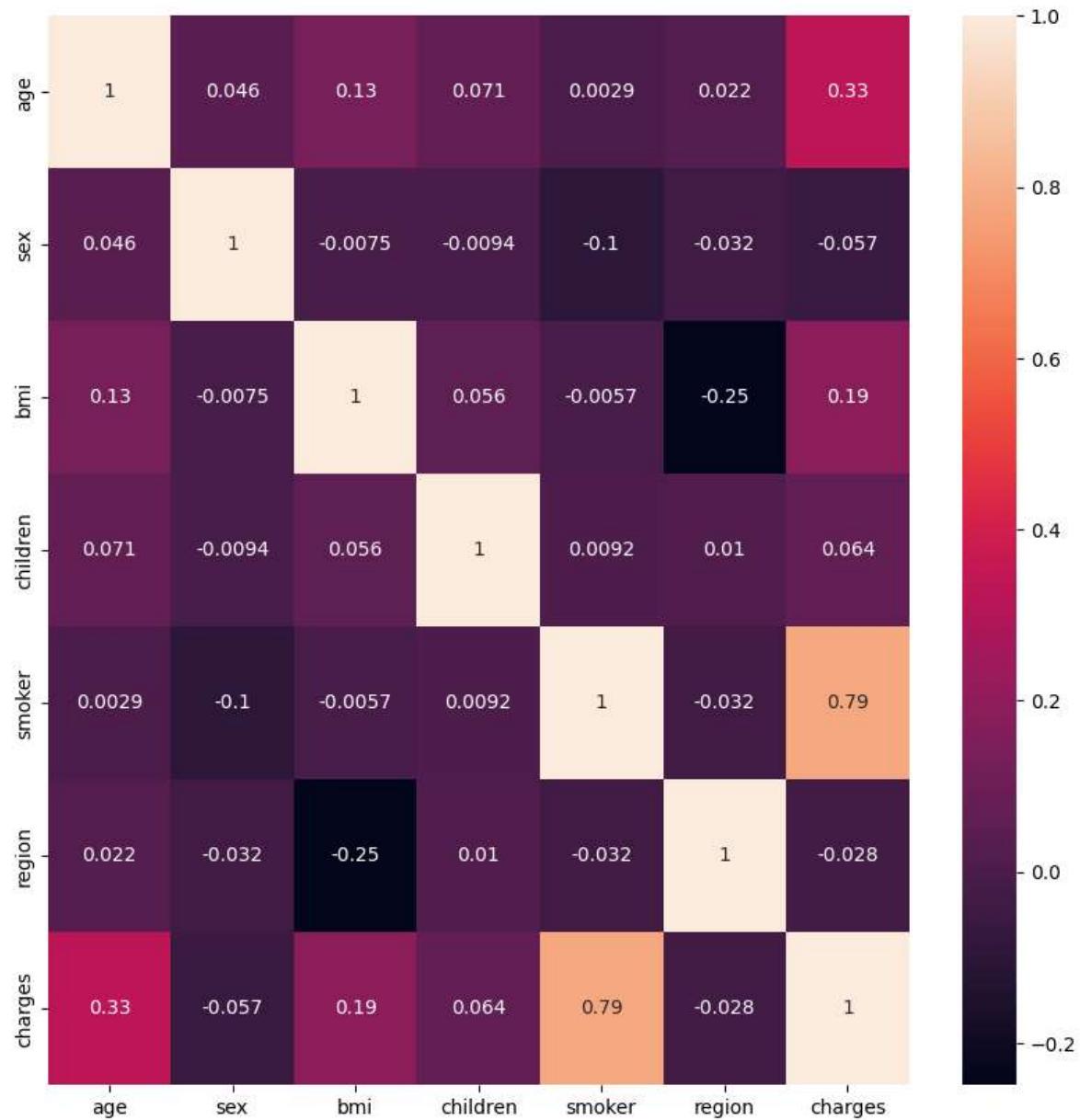
```
In [74]: lr=LinearRegression()
lr.fit(x_train,y_train)
y_pred=lr.predict(x_test)
r2=r2_score(y_test,y_pred)
print(r2)
```

0.0021400201329184743

Ridge Regression

```
In [75]: from sklearn.linear_model import Lasso,Ridge
from sklearn.preprocessing import StandardScaler
```

```
In [76]: plt.figure(figsize=(10, 10))
sns.heatmap(s.corr(), annot=True)
plt.show()
```



```
In [26]: features=s.columns[0:1]
target=s.columns[-1]
```

```
In [27]: x=s[features].values
y=s[target].values
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.30,random_state=42)
print("The dimension of X_train is {}".format(x_train.shape))
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 [28]: lr = LinearRegression()
#Fit model
lr.fit(x_train, y_train)
#predict
actual = y_test
train_score_lr = lr.score(x_train, y_train)
test_score_lr = lr.score(x_test, y_test)
print("\nLinear Regression Model:\n")
print("The train score for lr model is {}".format(train_score_lr))
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 [29]: ridgeReg = Ridge(alpha=10)
ridgeReg.fit(x_train,y_train)
#train and test score for ridge regression
train_score_ridge = ridgeReg.score(x_train, y_train)
test_score_ridge = ridgeReg.score(x_test, y_test)
print("\nRidge Model:\n")
print("The train score for ridge model is {}".format(train_score_ridge))
print("The test score for ridge model is {}".format(test_score_ridge))
```

Ridge Model:

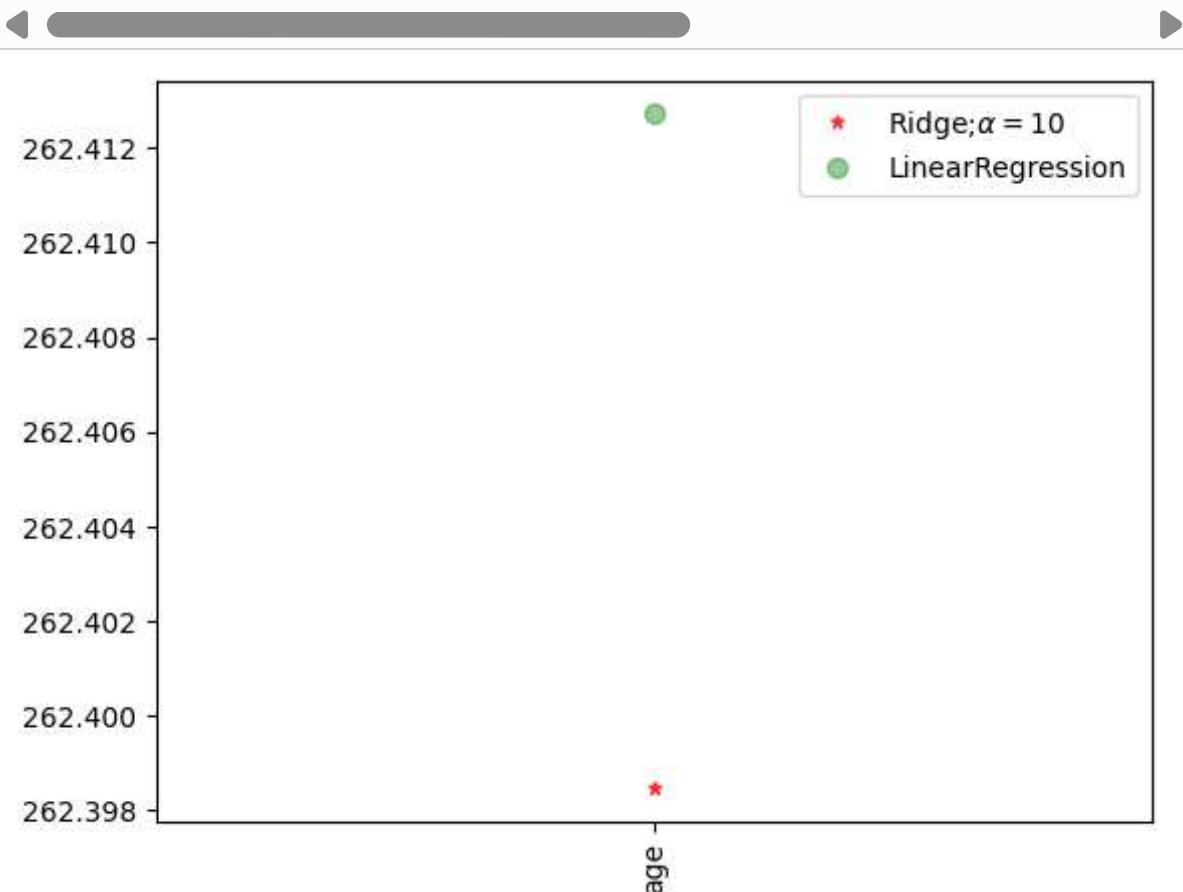
```
The train score for ridge model is 0.09109639711159634
The test score for ridge model is 0.08490538609860176
```

```
In [30]: plt.figure(figsize=(10,10))
```

```
Out[30]: <Figure size 1000x1000 with 0 Axes>
```

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

```
In [31]: plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker="*",markersize=7)
plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker="o",markersize=7,
plt.xticks(rotation=90)
plt.legend()
plt.show()
```



Lasso Regression

```
In [32]: lasso= Lasso(alpha=10)
lasso.fit(x_train,y_train)
#train and test score for ridge regression
train_score_ls = lasso.score(x_train, y_train)
test_score_ls= lasso.score(x_test, y_test)
print("\nRidge Model:")
print("The train score for lasso model is {}".format(train_score_ls))
print("The test score for lasso model is {}".format(test_score_ls))
```

Ridge Model:

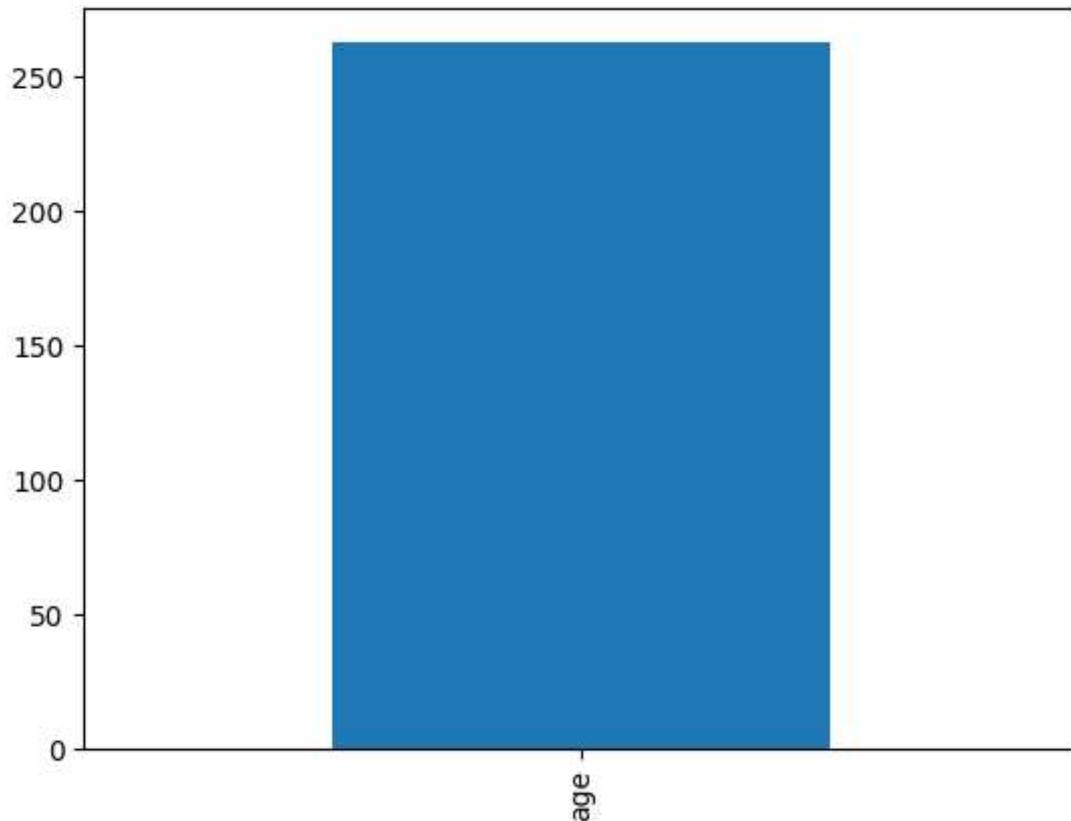
The train score for lasso model is 0.09109639395809055
The test score for lasso model is 0.08490704421828055

```
In [33]: plt.figure(figsize=(10,10))
```

```
Out[33]: <Figure size 1000x1000 with 0 Axes>
```

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

```
In [34]: pd.Series(lasso.coef_, features).sort_values(ascending = True).plot(kind = "bar")  
plt.show()
```



```
In [35]: from sklearn.linear_model import LassoCV
```

```
In [36]: #using the linear cv model  
from sklearn.linear_model import RidgeCV  
#cross validation  
ridge_cv=RidgeCV(alphas =[0.0001,0.001,0.01,0.1,1,10]).fit(x_train,y_train)  
#score  
print(ridge_cv.score(x_train,y_train))  
print(ridge_cv.score(x_test,y_test))
```

```
0.091096397111596
```

```
0.08490538609884224
```

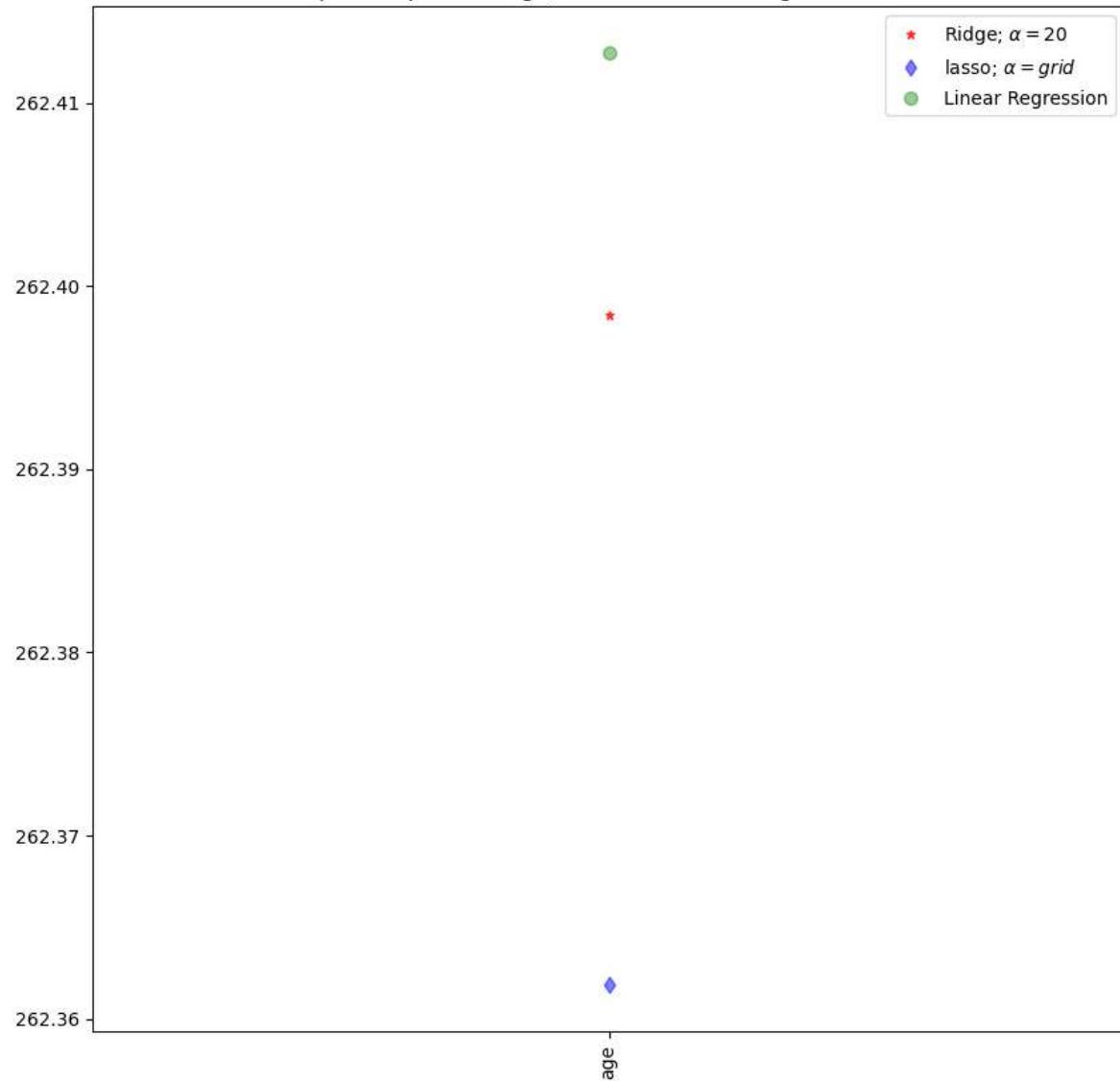
```
In [37]: #using the Linear cv model
from sklearn.linear_model import LassoCV
#cross validation
lasso_cv=LassoCV(alphas =[0.0001,0.001,0.01,0.1,1,10]).fit(x_train,y_train)
#score
print(lasso_cv.score(x_train,y_train))
print(lasso_cv.score(x_test,y_test))
```

```
0.09109639395809055
0.08490704421828055
```

```
In [38]: plt.figure(figsize = (10, 10))
#add plot for ridge regression
plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',markersize=6,color='red')
#add plot for lasso regression
plt.plot(lasso_cv.coef_,alpha=0.5,linestyle='none',marker='d',markersize=6,color='blue')
#add plot for linear model
plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersize=7,color='green')
#rotate axis
plt.xticks(rotation = 90)
plt.legend()
plt.title("Comparison plot of Ridge, Lasso and Linear regression model")
plt.show()
```



Comparison plot of Ridge, Lasso and Linear regression model



ElasticNet Regression

```
In [39]: from sklearn.linear_model import ElasticNet
```

```
In [40]: el=ElasticNet()  
el.fit(x_train,y_train)  
print(el.coef_)  
print(el.intercept_)
```

```
[261.74450967]  
3115.083177426244
```

```
In [41]: y_pred_elastic=el.predict(x_train)
```

```
In [42]: mean_squared_error=np.mean((y_pred_elastic-y_train)**2)  
print(mean_squared_error)
```

```
135077142.70714515
```

```
In [43]: el=ElasticNet()  
el.fit(x_train,y_train)  
print(el.score(x_train,y_train))
```

```
0.09109580670592365
```

Logistic Regression

```
In [44]: import numpy as np  
import pandas as pd  
from sklearn.linear_model import LogisticRegression  
from sklearn.preprocessing import StandardScaler
```

```
In [45]: s=pd.read_csv(r"C:\Users\mouni\Downloads\insurance.csv")
s
```

Out[45]:

	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 [46]: s.shape
```

Out[46]: (1338, 7)

```
In [47]: pd.set_option('display.max_rows',10000000000)
pd.set_option('display.max_columns',10000000000)
pd.set_option('display.width',95)
```

```
In [48]: print('This Dataset has %d rows and %d columns'%(s.shape))
```

This Dataset has 1338 rows and 7 columns

```
In [49]: s.head()
```

Out[49]:

	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 [50]: s.describe
```

```
Out[50]: <bound method NDFrame.describe of
ker      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
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
16      52   female  30.780      1     no northeast  10797.336200
17      ..      ..      ..      ..      ..      ..      ..      ..      ..
```

```
In [51]: s.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 [52]: s.isnull().sum()
```

```
Out[52]: age      0
sex      0
bmi      0
children 0
smoker   0
region   0
charges  0
dtype: int64
```

```
In [53]: convert={"smoker":{"yes":1,"no":0}}
s=s.replace(convert)
s
```

Out[53]:

	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 [54]: convert={"sex":{"female":1,"male":0}}
s=s.replace(convert)
s
```

Out[54]:

	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 [55]: convert={"region":{"southeast":1,"southwest":2,"northeast":3,"northwest":4}}  
s=s.replace(convert)  
s
```

Out[55]:

	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 [56]: features_matrix=s.iloc[:,0:4]
```

```
In [57]: target_vector=s.iloc[:, -3]
```

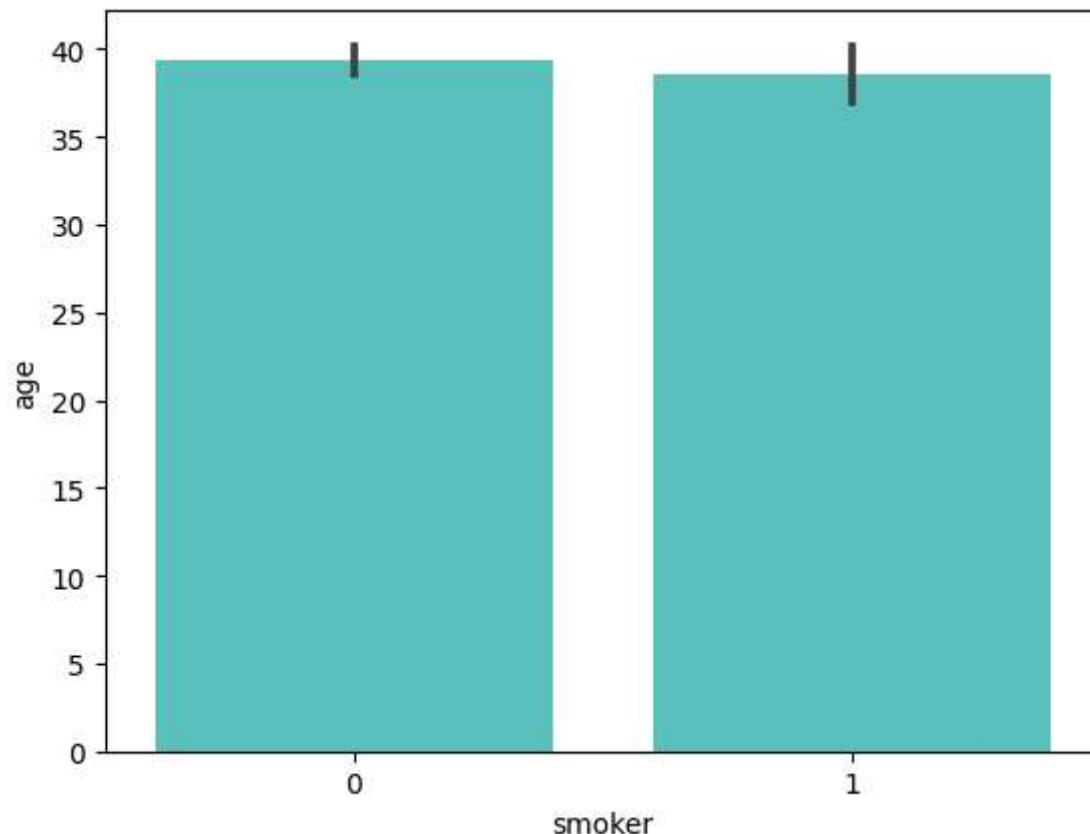
```
In [58]: print('The Feature Matrix has %d Rows and %d columns(s) '%(features_matrix.shape))  
print('The Target Matrix has %d Rows and %d columns(s) '%(np.array(target_vector).shape))
```

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

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

```
In [59]: import matplotlib.pyplot as plt  
import seaborn as sns
```

```
In [60]: sns.barplot(x='smoker', y='age', data=s, color="mediumturquoise")
plt.show()
```



```
In [61]: features_matrix_standardized=StandardScaler().fit_transform(features_matrix)
```

```
In [62]: algorithm=LogisticRegression(max_iter=10000)
```

```
In [63]: Logistic_Regression_Model=algorithm.fit(features_matrix_standardized,target_ve
```

```
In [64]: observation=[[1,0,0.99539,-0.0588]]
```

```
In [65]: predictions=Logistic_Regression_Model.predict(observation)
print('The model predicted the observation to belong to class %s'%(predictions))
```

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

```
In [66]: print('The algorithm was trained to predict one of the two classes:%s'%(algor
```

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

```
In [77]: print(" " "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[0] Is 0.8057075871331396

```
In [80]: print(" " "The Model says the probability of the observation we passed belonging to class['g'] Is 0.19429241286686044
```

The Model says the probability of the observation we passed belonging to class['g'] Is 0.19429241286686044

```
In [82]: x=np.array(s['age']).reshape(-1,1)
y=np.array(s['smoker']).reshape(-1,1)
```

```
In [83]: x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.05)
lo=LogisticRegression()
lo.fit(x_train,y_train)
print(lo.score(x_test,y_test))
```

0.7428571428571429

C:\Users\mouni\AppData\Local\Programs\Python\Python310\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 [84]: import numpy as np
import pandas as pd
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
```

```
In [85]: s=pd.read_csv(r"C:\Users\mouni\Downloads\insurance.csv")  
s
```

Out[85]:

	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 [86]: s.shape
```

Out[86]: (1338, 7)

```
In [87]: s.isnull().any()
```

Out[87]: age False
sex False
bmi False
children False
smoker False
region False
charges False
dtype: bool

```
In [88]: s['region'].value_counts()
```

Out[88]: region
southeast 364
southwest 325
northwest 325
northeast 324
Name: count, dtype: int64

```
In [89]: convert={"sex":{"female":1,"male":0}}
s=s.replace(convert)
s
```

Out[89]:

	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 [90]: convert={"smoker":{"yes":1,"no":0}}
s=s.replace(convert)
s
```

Out[90]:

	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 [91]: x=["bmi","children"]
y=["Yes","No"]
all_inputs=s[x]
all_classes=s["sex"]
```

```
In [92]: (x_train,x_test,y_train,y_test)=train_test_split(all_inputs,all_classes,test_s
```

```
In [93]: clf=DecisionTreeClassifier(random_state=0)
```

```
In [94]: clf.fit(x_train,y_train)
```

```
Out[94]:
```

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

```
In [96]: score=clf.score(x_test,y_test)
print(score)
```

```
0.3902439024390244
```

Random Forest

```
In [97]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt ,seaborn as sns
```

```
In [98]: s=pd.read_csv(r"C:\Users\mouni\Downloads\insurance.csv")
s
```

```
Out[98]:
```

	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 [100]: s.shape
```

```
Out[100]: (1338, 7)
```

```
In [101]: s['region'].value_counts()
```

```
Out[101]: region
southeast    364
southwest    325
northwest    325
northeast    324
Name: count, dtype: int64
```

```
In [102]: s['bmi'].value_counts()
```

```
Out[102]: 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
27.360       7
... ... ...

```

```
In [103]: m={"sex":{"female":1,"male":0}}
          s=s.replace(m)
          print(s)
```

	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	55	0	31.300	2	no	southeast	10660.325000

```
In [105]: n={"smoker":{"yes":1,"no":0}}
s=s.replace(n)
print(s)
```

	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	55	0	32.380	0	0	southeast	16602.325200

```
In [106]: from sklearn.ensemble import RandomForestClassifier
rfc=RandomForestClassifier()
rfc.fit(x_train,y_train)
```

```
Out[106]: RandomForestClassifier()
          |____ RandomForestClassifier()
```

```
In [107]: rf=RandomForestClassifier()
params={'max_depth':[2,3,5,20],
        'min_samples_leaf':[5,10,20,50,100,200],
        'n_estimators':[10,25,30,50,100,200]}
```

```
In [108]: from sklearn.model_selection import GridSearchCV
grid_search=GridSearchCV(estimator=rf, param_grid=params, cv=2, scoring="accuracy")
grid_search.fit(x_train,y_train)
```

```
Out[108]: GridSearchCV()
           |____ estimator: RandomForestClassifier()
                  |____ RandomForestClassifier()
```

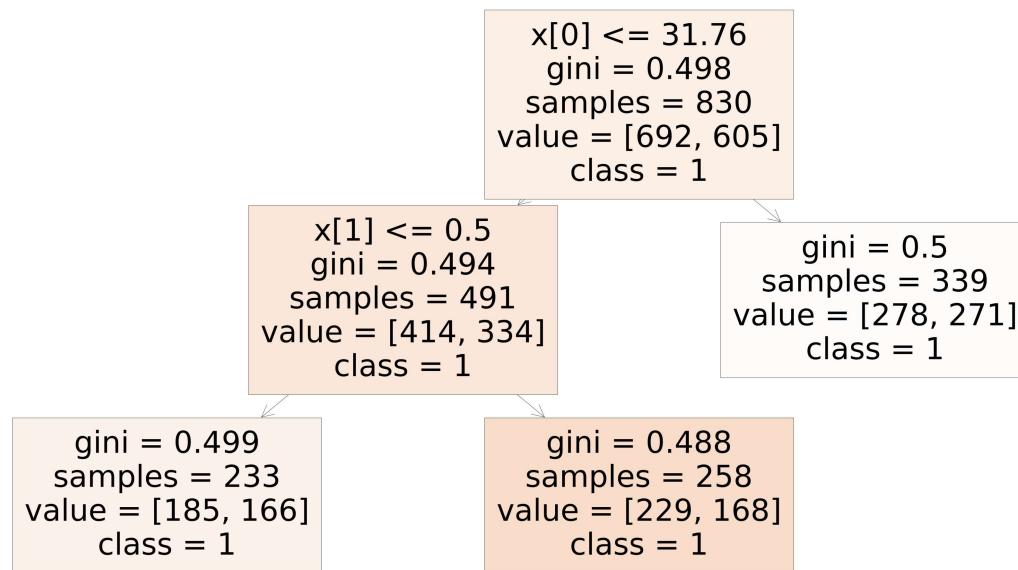
```
In [109]: grid_search.best_score_
```

```
Out[109]: 0.5212066522094772
```

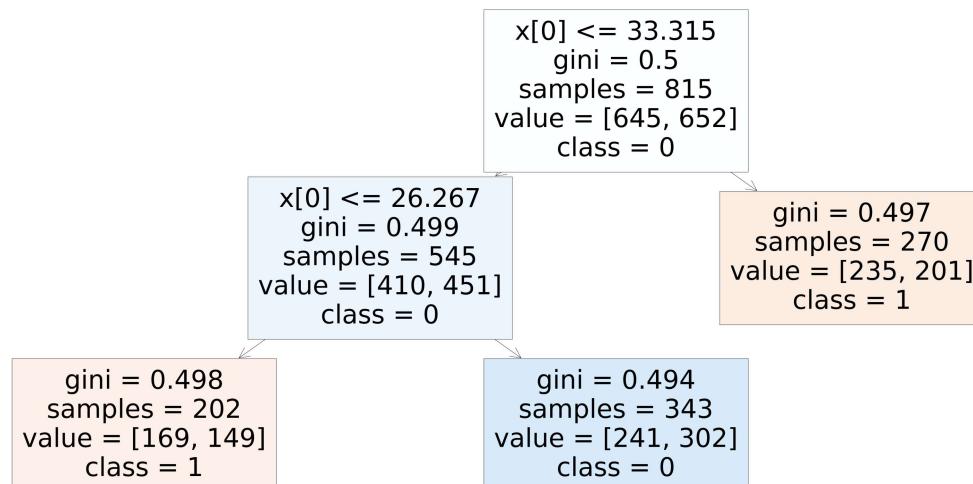
```
In [110]: rf_best=grid_search.best_estimator_
print(rf_best)
```

```
RandomForestClassifier(max_depth=5, min_samples_leaf=200, n_estimators=30)
```

```
In [111]: from sklearn.tree import plot_tree
plt.figure(figsize=(80,40))
plot_tree(rf_best.estimators_[4], class_names=['1','0'], filled=True);
```



```
In [112]: from sklearn.tree import plot_tree
plt.figure(figsize=(70,30))
plot_tree(rf_best.estimators_[6], class_names=["1","0"], filled=True);
```



```
In [113]: rf_best.feature_importances_
```

```
Out[113]: array([0.76473474, 0.23526526])
```

```
In [114]: rf=RandomForestClassifier(random_state=0)
```

Conclusion

```
# For the given insurance data set have performed linear,logistic,random  
forest and decision tree models of regression and classifications.  
#and have conclude that the most accuracy is occured in logistic  
regression,i.e 74percent  
#when compare to other regression models.  
#and concluded that "Logistic Regression" model is fits for the data.
```