

1.Data Collection

```
In [65]: 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 [66]: df=pd.read_csv(r"C:\Users\Ajay Reddy\Downloads\insurance.csv")
df
```

Out[66]:

	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 Preprocesing

```
In [67]: df.head()
```

Out[67]:

	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 [68]: df.tail()
```

Out[68]:

	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 [69]: df.shape
```

Out[69]: (1338, 7)

In [70]: df.describe

```
Out[70]: <bound method NDFrame.describe of
   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 [71]: 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 [72]: df.isnull().any()

```
Out[72]: age      False
          sex      False
          bmi      False
          children False
          smoker   False
          region   False
          charges  False
dtype: bool
```

In [73]: df.isna().sum()

```
Out[73]: age      0
          sex      0
          bmi      0
          children 0
          smoker   0
          region   0
          charges  0
dtype: int64
```

In [74]: df['region'].value_counts()

```
Out[74]: region
          southeast    364
          southwest    325
          northwest    325
          northeast    324
Name: count, dtype: int64
```

```
In [75]: convert={"sex":{"female":1,"male":0}}
df=df.replace(convert)
df
```

Out[75]:

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

Out[76]:

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

1338 rows × 7 columns

```
In [77]: convert={"region":{"southwest":1, "southeast":2, "northwest":3, "northeast":4}}
df=df.replace(convert)
df
```

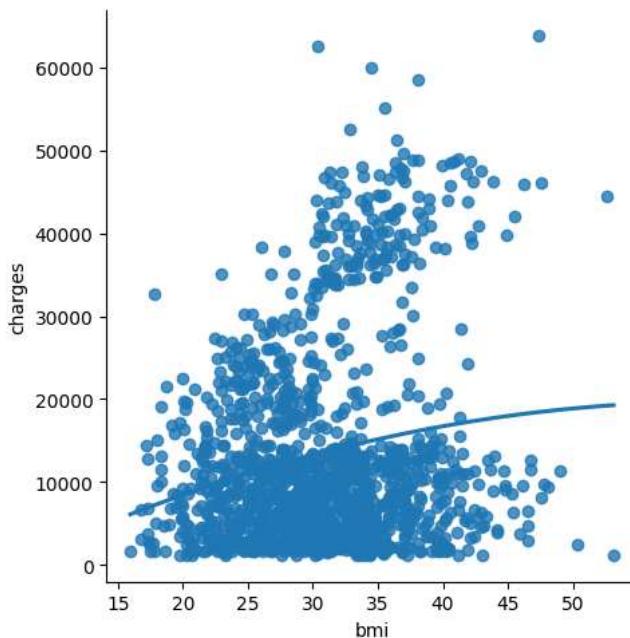
Out[77]:

	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

3.Data Visualization

```
In [78]: sns.lmplot(x='bmi',y='charges',order=2,data=df,ci=None)
plt.show()
```

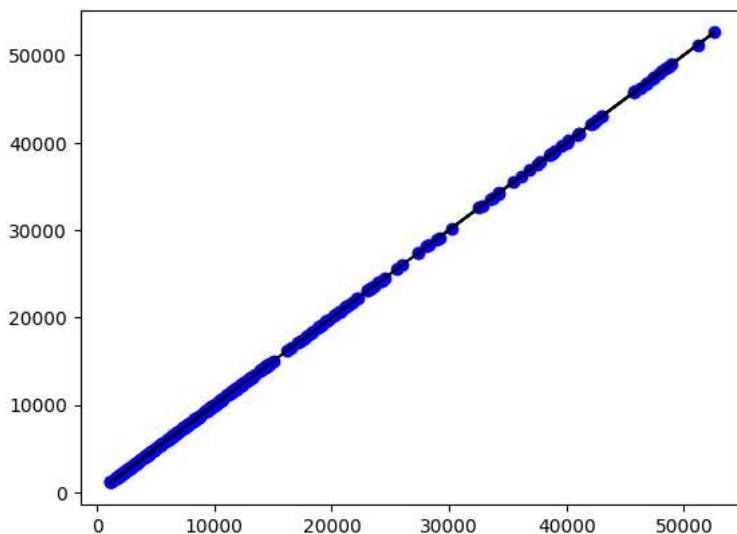


```
In [79]: x=np.array(df['bmi']).reshape(-1,1)
y=np.array(df['charges']).reshape(-1,1)
```

```
In [80]: x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25,random_state=0)
lr=LinearRegression()
lr.fit(x_train,y_train)
print(lr.score(x_test,y_test))
```

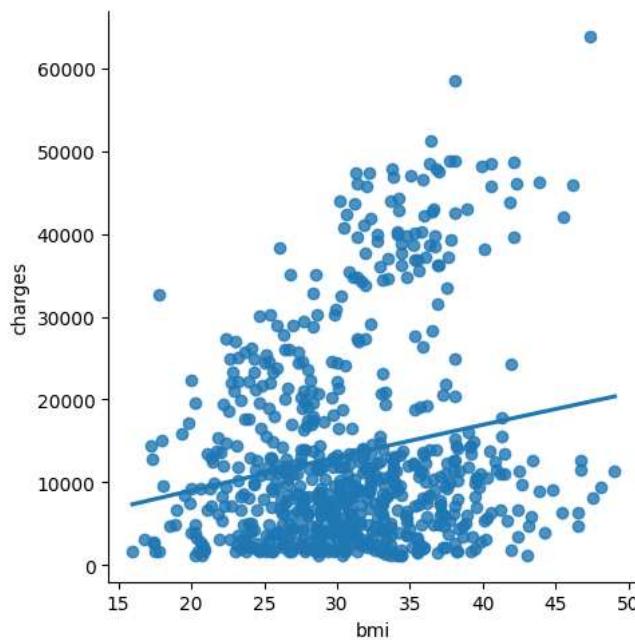
```
1.0
```

```
In [81]: 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 [82]: df700=df[:][:700]
sns.lmplot(x='bmi',y='charges',order=2,ci=None,data=df700)
plt.show()
```



```
In [83]: df700.fillna(method='ffill',inplace=True)
```

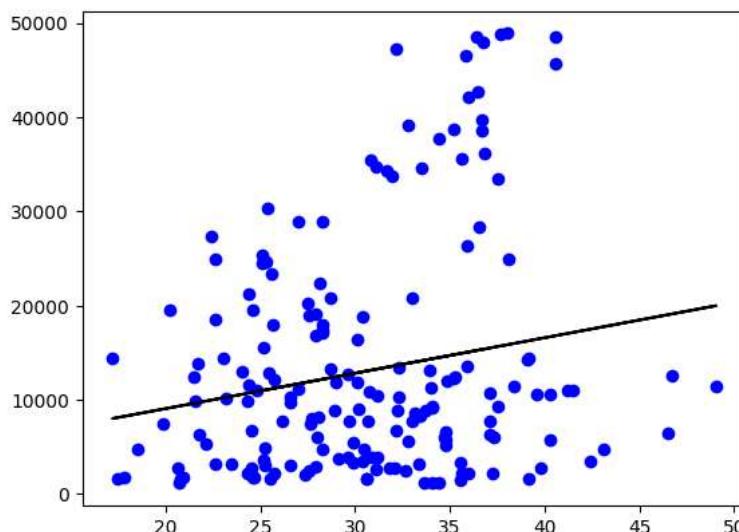
```
In [84]: x=np.array(df700["bmi"]).reshape(-1,1)
y=np.array(df700['charges']).reshape(-1,1)
```

```
In [85]: df700.dropna(inplace=True)
```

```
In [86]: 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))
```

```
0.04307515742285595
```

```
In [87]: 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 [88]: └─▶ from sklearn.linear_model import LinearRegression
      from sklearn.metrics import r2_score
      lr=LinearRegression()
      lr.fit(x_train,y_train)
      y_pred=lr.predict(x_test)
      r2=r2_score(y_test,y_pred)
      print(r2)
```

0.04307515742285595

Ridge Regression

```
In [89]: └─▶ from sklearn.linear_model import Lasso,Ridge
      from sklearn.preprocessing import StandardScaler
```

```
In [90]: └─▶ plt.figure(figsize= (10,10))
      sns.heatmap(df700.corr(),annot=True)
      plt.show()
```



```
In [91]: └─▶ features=df.columns[0:1]
      target=df.columns[-1]
```

```
In [92]: x=df[features].values
y=df[target].values
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.30,random_state=1)
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 [93]: 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 [94]: ridgeReg = Ridge(alpha=10)
ridgeReg.fit(x_train,y_train)
#train and test scorefor 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))
```

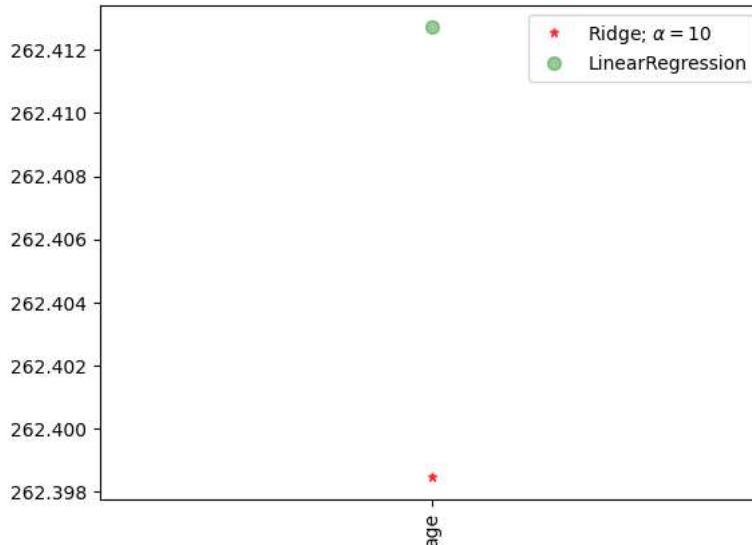
Ridge Model:

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

```
In [95]: plt.figure(figsize=(10,10))
```

```
Out[95]: <Figure size 1000x1000 with 0 Axes>
<Figure size 1000x1000 with 0 Axes>
```

```
In [97]: plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,color='red',label=r'Ridge; $\alpha=10$',zorder=10)
plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersize=7,color='green',label='LinearRegression')
plt.xticks(rotation=90)
plt.legend()
plt.show()
```



Lasso Regression

```
In [98]: 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:\n")
print("The train score for lasso model is {}".format(train_score_ls))
print("The test score for lasso model is {}".format(test_score_ls))
```

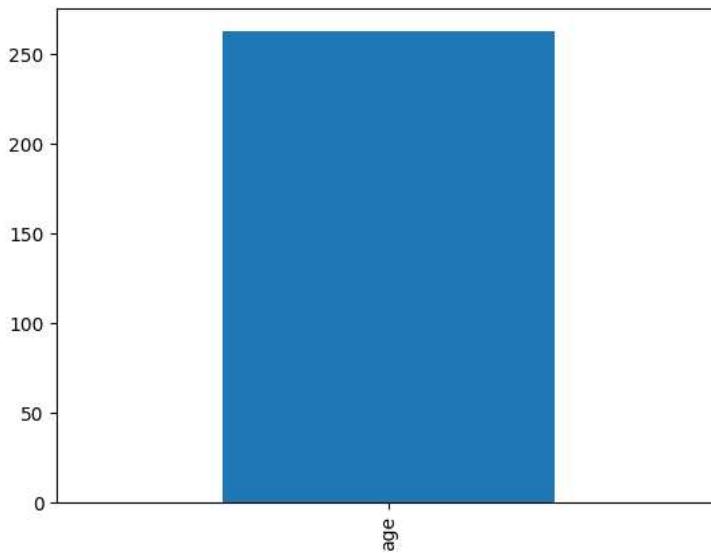
Ridge Model:

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

```
In [99]: plt.figure(figsize=(10,10))
```

```
Out[99]: <Figure size 1000x1000 with 0 Axes>
<Figure size 1000x1000 with 0 Axes>
```

```
In [100]: pd.Series(lasso.coef_, features).sort_values(ascending = True).plot(kind = "bar")
plt.show()
```



```
In [105]: from sklearn.linear_model import LassoCV
```

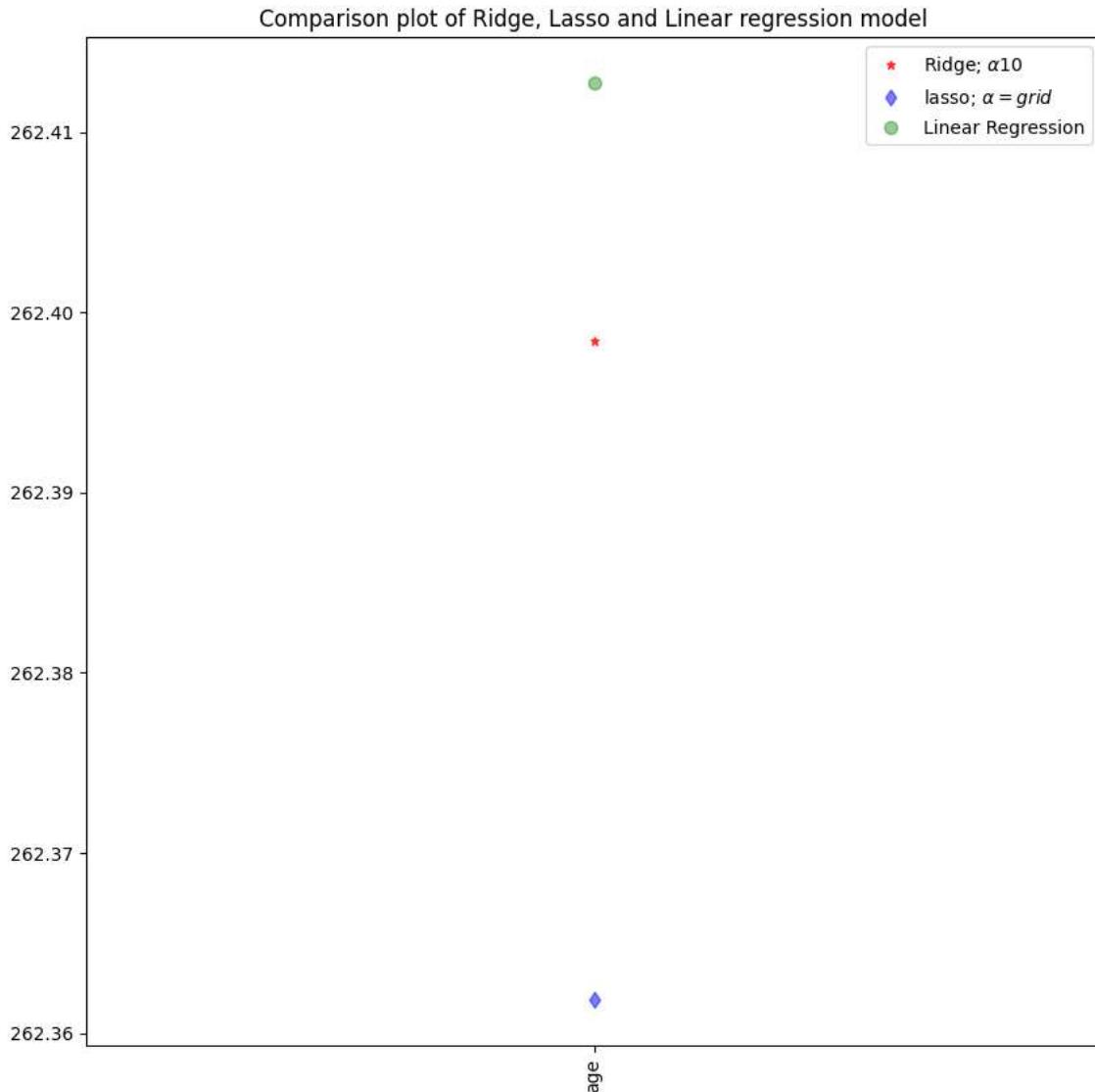
```
In [106]: #using the Linear cv model
from sklearn.linear_model import RidgeCV
#cross validation
ridge_cv=RidgeCV(alphas =[0.0001,0.001,0.01,0.1,1,10]).fit(x_train,y_train)
#score
print(ridge_cv.score(x_train,y_train))
print(ridge_cv.score(x_test,y_test))
```

0.09109639711159612
 0.08490538609884779

```
In [109]: #using the Linear CV model
from sklearn.linear_model import LassoCV
#Ridge Cross validation
lasso_cv=LassoCV(alphas =[0.0001,0.001,0.01,0.1,1,10]).fit(x_train,y_train)
#score
print("The train score for ls model is {}".format(ridge_cv.score(x_train,y_train)))
print("The test score for ls model is {}".format(ridge_cv.score(x_test,y_test)))
```

The train score for ls model is 0.09109639711159612
 The test score for ls model is 0.08490538609884779

```
In [110]: # plt.figure(figsize = (10, 10))
#add plot for ridge regression
plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,color='red',label=r'Ridge; $\alpha$10$',zorder=1)
#add plot for lasso regression
plt.plot(lasso_cv.coef_,alpha=0.5,linestyle='none',marker='d',markersize=6,color='blue',label=r'lasso; $\alpha = \text{grid}$')
#add plot for Linear model
plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersize=7,color='green',label='Linear Regression')
#rotate axis
plt.xticks(rotation = 90)
plt.legend()
plt.title("Comparison plot of Ridge, Lasso and Linear regression model")
plt.show()
```



Elastic Net Regression

```
In [111]: from sklearn.linear_model import ElasticNet
```

```
In [112]: el=ElasticNet()
el.fit(x_train,y_train)
print(el.coef_)
print(el.intercept_)
```

```
[261.74450967]
3115.083177426244
```

```
In [113]: y_pred_elastic=el.predict(x_train)
```

```
In [114]: mean_squared_error=np.mean((y_pred_elastic-y_train)**2)
print(mean_squared_error)
```

135077142.70714515

```
In [115]: el=ElasticNet()
el.fit(x_train,y_train)
print(el.score(x_train,y_train))
```

0.09109580670592365

Logistic Regression

```
In [116]: import numpy as np
import pandas as pd
from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import StandardScaler
```

```
In [117]: df=pd.read_csv(r"C:\Users\Ajay Reddy\Downloads\insurance.csv")
df
```

Out[117]:

	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 [118]: df.shape
```

Out[118]: (1338, 7)

```
In [119]: pd.set_option('display.max_rows',10000000000)
pd.set_option('display.max_columns',10000000000)
pd.set_option('display.width',95)
```

```
In [120]: print('This Dataset has %d rows and %d columns'%(df.shape))
```

This Dataset has 1338 rows and 7 columns

```
In [121]: df.head()
```

Out[121]:

	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 [122]: df.describe

52	48	male	26.000	1	yes	southwest	25508.272000
53	36	male	34.430	0	yes	southeast	37742.575700
54	40	female	28.690	3	no	northwest	8059.679100
55	58	male	36.955	2	yes	northwest	47496.494450
56	58	female	31.825	2	no	northeast	13607.368750
57	18	male	31.680	2	yes	southeast	34303.167200
58	53	female	22.880	1	yes	southeast	23244.790200
59	34	female	37.335	2	no	northwest	5989.523650
60	43	male	27.360	3	no	northeast	8606.217400
61	25	male	33.660	4	no	southeast	4504.662400
62	64	male	24.700	1	no	northwest	30166.618170
63	28	female	25.935	1	no	northwest	4133.641650
64	20	female	22.420	0	yes	northwest	14711.743800
65	19	female	28.900	0	no	southwest	1743.214000
66	61	female	39.100	2	no	southwest	14235.872000
67	40	male	26.315	1	no	northwest	6389.377850
68	40	female	36.190	0	no	southeast	5920.104100
69	28	male	23.980	3	yes	southeast	17663.144200
70	27	female	24.750	0	yes	southeast	16577.779500
71	31	male	28.500	5	no	northeast	6799.458000

In [123]: 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 [124]: df.isnull().sum()

```
Out[124]: age      0
          sex      0
          bmi      0
          children  0
          smoker   0
          region   0
          charges   0
          dtype: int64
```

In [125]: convert={"smoker":{"yes":1,"no":0}}
df=df.replace(convert)
df

78	22	female	39.805	0	0	northeast	2755.020950
79	41	female	32.965	0	0	northwest	6571.024350
80	31	male	26.885	1	0	northeast	4441.213150
81	45	female	38.285	0	0	northeast	7935.291150
82	22	male	37.620	1	1	southeast	37165.163800
83	48	female	41.230	4	0	northwest	11033.661700
84	37	female	34.800	2	1	southwest	39836.519000
85	45	male	22.895	2	1	northwest	21098.554050
86	57	female	31.160	0	1	northwest	43578.939400
87	56	female	27.200	0	0	southwest	11073.176000
88	46	female	27.740	0	0	northwest	8026.666600
89	55	female	26.980	0	0	northwest	11082.577200
90	21	female	39.490	0	0	southeast	2026.974100

```
In [126]: convert={"sex":{"female":1,"male":0}}
df=df.replace(convert)
df
```

Out[126]:

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

Out[127]:

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

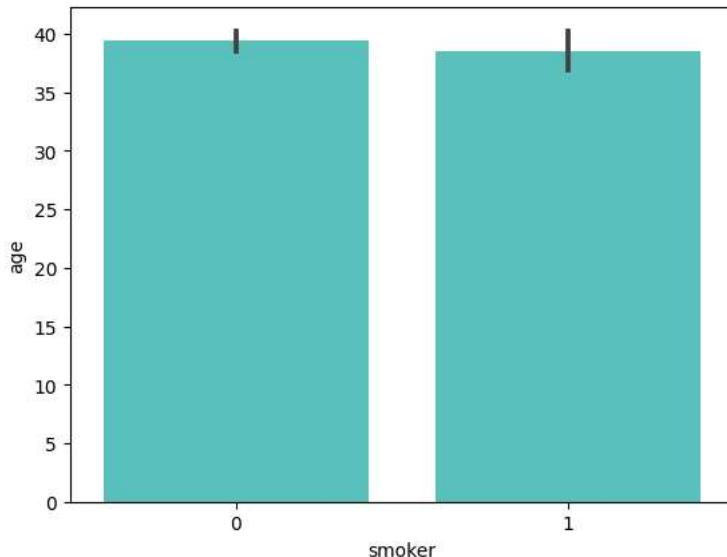
```
In [129]: target_vector=df.iloc[:,+3]
```

```
In [130]: 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).reshape(-1,1).shape))
```

The Feature Matrix has 1338 Rows and 4 columns(s)
The Target Matrix has 1338 Rows and 1 columns(s)

```
In [131]: import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [132]: sns.barplot(x='smoker', y='age', data=df, color="mediumturquoise")
plt.show()
```



```
In [133]: features_matrix_standardized=StandardScaler().fit_transform(features_matrix)
```

```
In [134]: algorithm=LogisticRegression(max_iter=10000)
```

```
In [135]: Logistic_Regression_Model=algorithm.fit(features_matrix_standardized,target_vector)
```

```
In [136]: observation=[[1,0,0.99539,-0.0588]]
```

```
In [137]: 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 [138]: print('The algorithm was trained to predict one of the two classes:%s'%(algorithm.classes_))
```

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

```
In [140]: print(" " "The Model says the probability of the observation we passed belonging to class[0] %s" " "%(algorithm.predict_proba(observation)[0][0])
```

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

```
In [141]: print(" " "The probability of the observation we passed belonging to class['g'] Is %s" " "%(algorithm.predict_proba(observation)[0][1]))
```

The Model says the probability of the observation we passed belonging to class['g'] Is 0.19429241286686041

```
In [142]: x=np.array(df['age']).reshape(-1,1)
y=np.array(df['smoker']).reshape(-1,1)
```

```
In [143]: 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))
```

0.7761194029850746

```
C:\Users\Ajay Reddy\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 [144]: 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 [145]: df=pd.read_csv(r"C:\Users\Ajay Reddy\Downloads\insurance.csv")
df
```

	age	sex	bmi	children	smoker	region	charges
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	23	male	23.845	0	no	northeast	2395.171550
18	56	male	40.300	0	no	southwest	10602.385000
19	30	male	35.300	0	yes	southwest	36837.467000
20	60	female	36.005	0	no	northeast	13228.846950
21	30	female	32.400	1	no	southwest	4149.736000
22	18	male	34.100	0	no	southeast	1137.011000
23	34	female	31.920	1	yes	northeast	37701.876800
24	37	male	28.025	2	no	northwest	6203.901750

```
In [146]: df.shape
```

```
Out[146]: (1338, 7)
```

```
In [147]: df.isnull().any()
```

```
Out[147]: age      False
          sex      False
          bmi      False
          children  False
          smoker   False
          region   False
          charges   False
          dtype: bool
```

```
In [148]: df['region'].value_counts()
```

```
Out[148]: region
          southeast    364
          southwest    325
          northwest    325
          northeast    324
          Name: count, dtype: int64
```

```
In [149]: convert={"sex":{"female":1,"male":0}}
df=df.replace(convert)
df
```

```
Out[149]:
```

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

Out[150]:

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

```
In [152]: (x_train,x_test,y_train,y_test)=train_test_split(all_inputs,all_classes,test_size=0.03)
```

```
In [153]: clf=DecisionTreeClassifier(random_state=0)
```

```
In [154]: clf.fit(x_train,y_train)
```

```
Out[154]: DecisionTreeClassifier
DecisionTreeClassifier(random_state=0)
```

```
In [155]: score=clf.score(x_test,y_test)
print(score)
```

0.43902439024390244

Random Forest

```
In [156]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt ,seaborn as sns
```

```
In [158]: df=pd.read_csv(r"C:\Users\Ajay Reddy\Downloads\insurance.csv")
df
```

Out[158]:

	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 [159]: df.shape
```

```
Out[159]: (1338, 7)
```

```
In [160]: df['region'].value_counts()
```

```
Out[160]: region
southeast    364
southwest   325
northwest   325
northeast   324
Name: count, dtype: int64
```

```
In [161]: df['bmi'].value_counts()
```

```
Out[161]: 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 [162]: m={"sex":{"female":1,"male":0}}
df=df.replace(m)
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
..

```
In [163]: n={"smoker":{"yes":1,"no":0}}
df=df.replace(n)
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
..

```
In [164]: from sklearn.ensemble import RandomForestClassifier
rfc=RandomForestClassifier()
rfc.fit(x_train,y_train)
```

Out[164]:

```
RandomForestClassifier()
RandomForestClassifier()
```

```
In [165]: 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 [171]: 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[171]:

```
GridSearchCV
  estimator: RandomForestClassifier
    RandomForestClassifier()
```

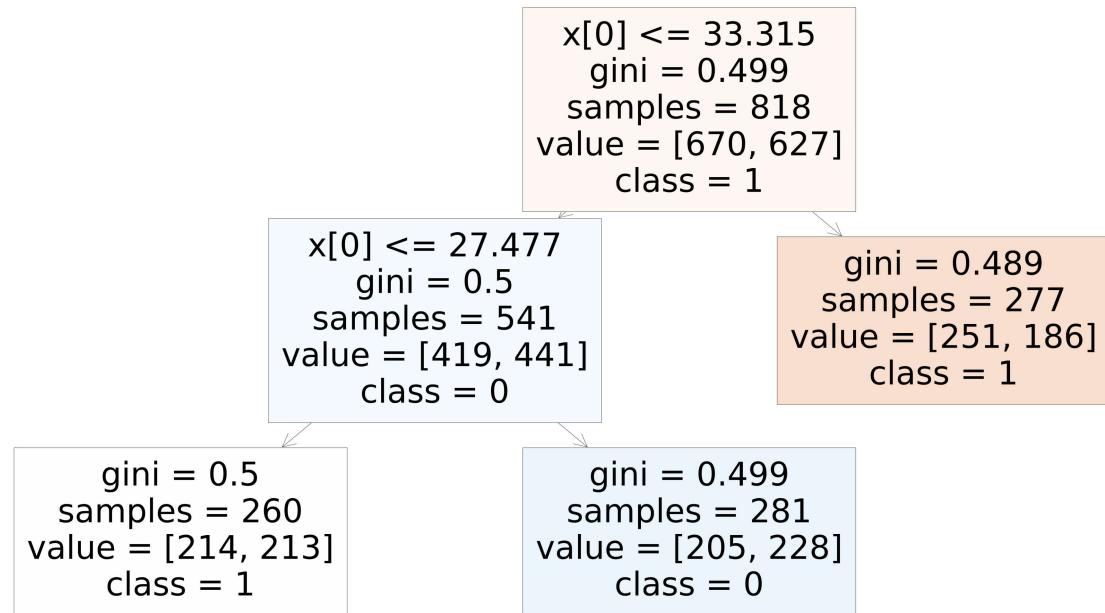
```
In [169]: grid_search.best_score_
```

Out[169]: 0.5235190891970554

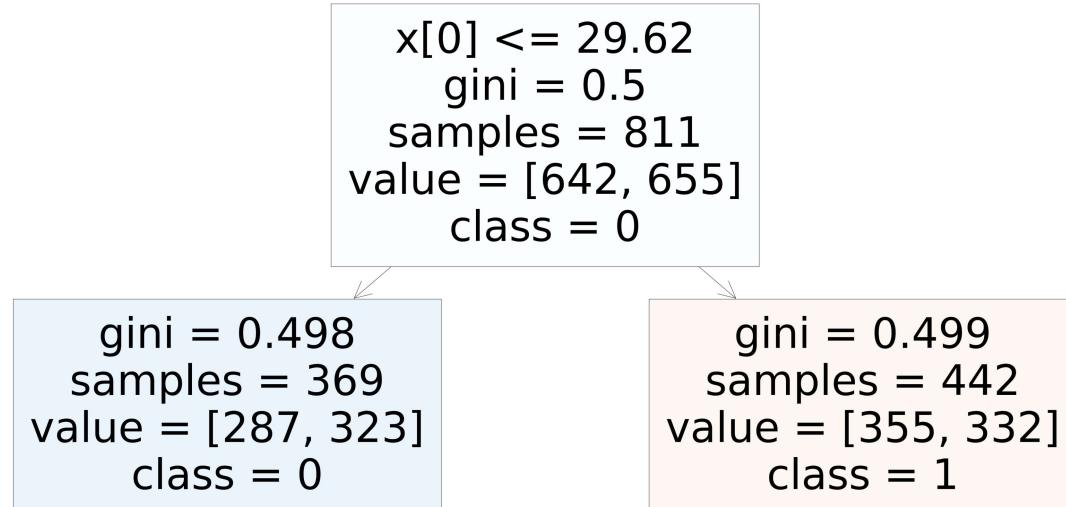
```
In [173]: rf_best=grid_search.best_estimator_
print(rf_best)
```

RandomForestClassifier(max_depth=3, min_samples_leaf=200, n_estimators=10)

```
In [174]: ┌─┐ from sklearn.tree import plot_tree
      plt.figure(figsize=(80,40))
      plot_tree(rf_best.estimators_[4], class_names=['1', '0'], filled=True);
```



```
In [175]: ┌─┐ from sklearn.tree import plot_tree
      plt.figure(figsize=(70,30))
      plot_tree(rf_best.estimators_[6], class_names=["1", "0"], filled=True);
```



```
In [176]: ┌─┐ rf_best.feature_importances_
```

```
Out[176]: array([0.77540327, 0.22459673])
```

```
In [177]: ┌─┐ rf=RandomForestClassifier(random_state=0)
```

```
In [178]: ┌─┐ rf.fit(x_train,y_train)
```

```
Out[178]: RandomForestClassifier
          RandomForestClassifier(random_state=0)
```

```
In [179]: ┌─┐ score=rf.score(x_test,y_test)
      print(score)
```

```
0.3902439024390244
```

In []: