

# To predict how best the data fits

## 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]: df=pd.read_csv(r"C:\Users\sudheer\Downloads\insurance.csv")
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]: 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]: 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]: df.shape

Out[5]: (1338, 7)

In [6]: 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]: 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]: df.isnull().any()
```

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

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

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

```
In [10]: df['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}}  
df=df.replace(convert)  
df
```

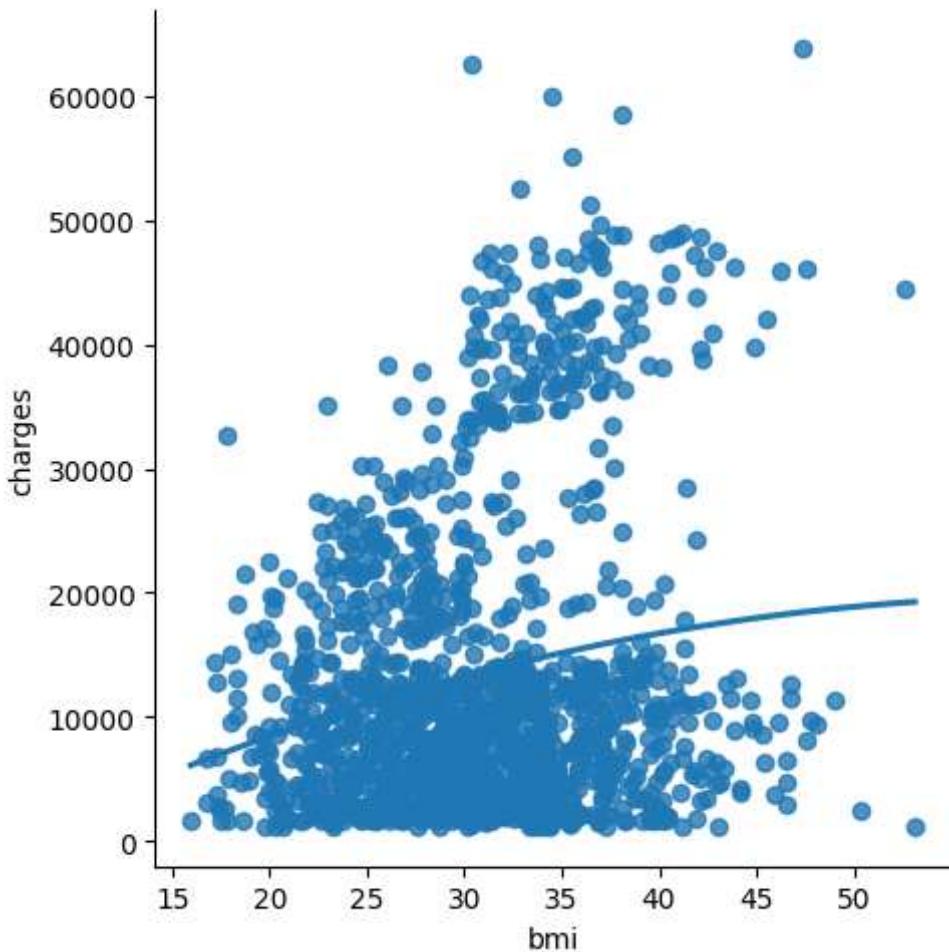
```
Out[11]:
```

	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

## Data Visualisation

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

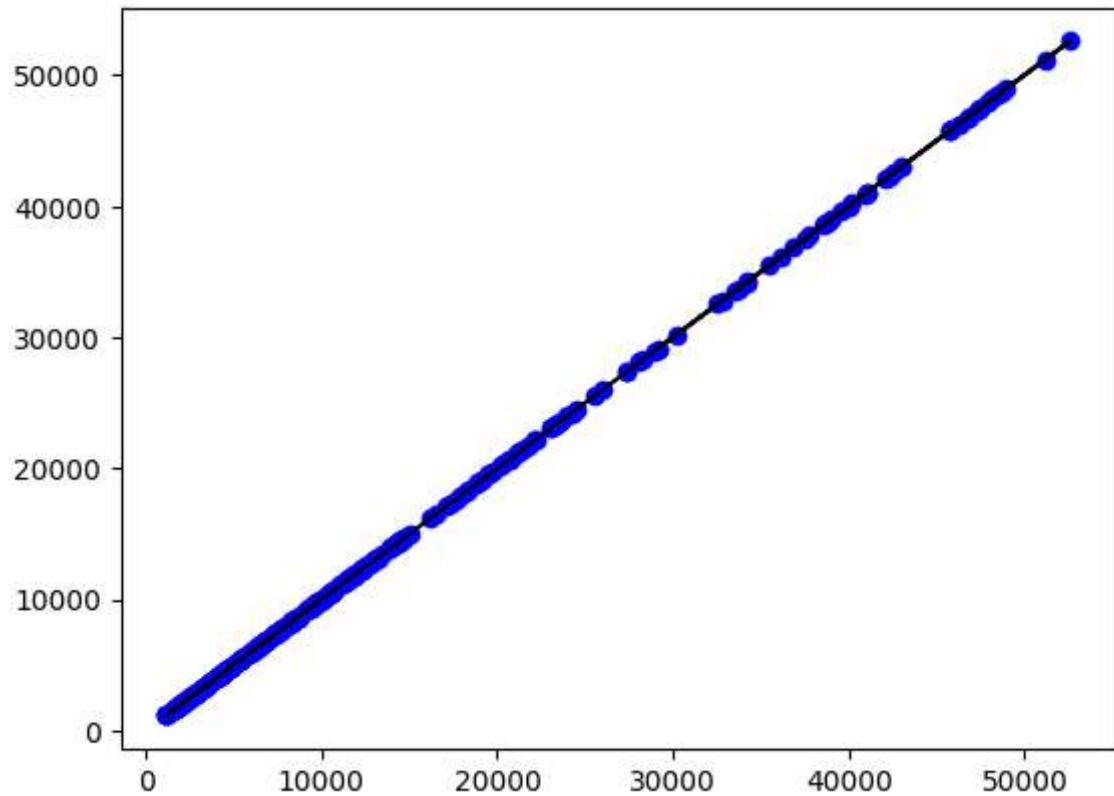


```
In [13]: x=np.array(df['bmi']).reshape(-1,1)
y=x=np.array(df['charges']).reshape(-1,1)
```

```
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))
```

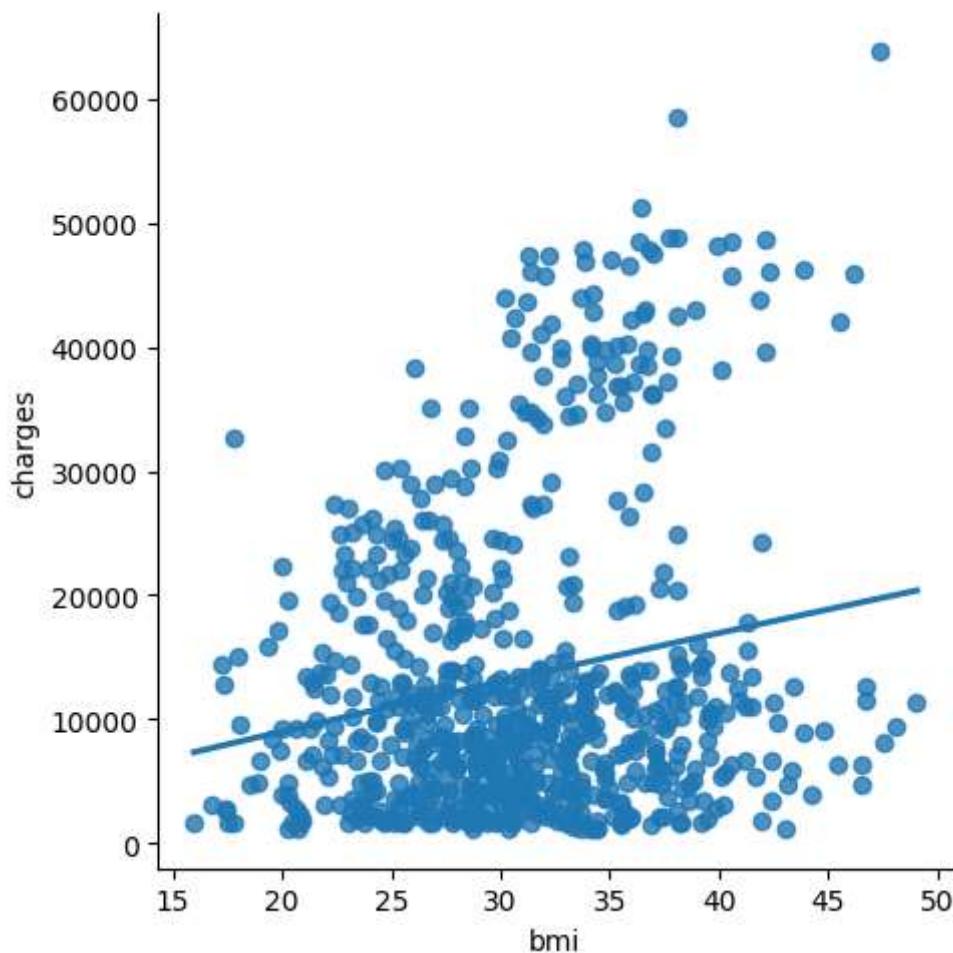
1.0

```
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 data subset

```
In [16]: df700=df[:][:700]
sns.lmplot(x='bmi',y='charges',order=2,ci=None,data=df700)
plt.show()
```



```
In [17]: df700.fillna(method='ffill',inplace=True)
```

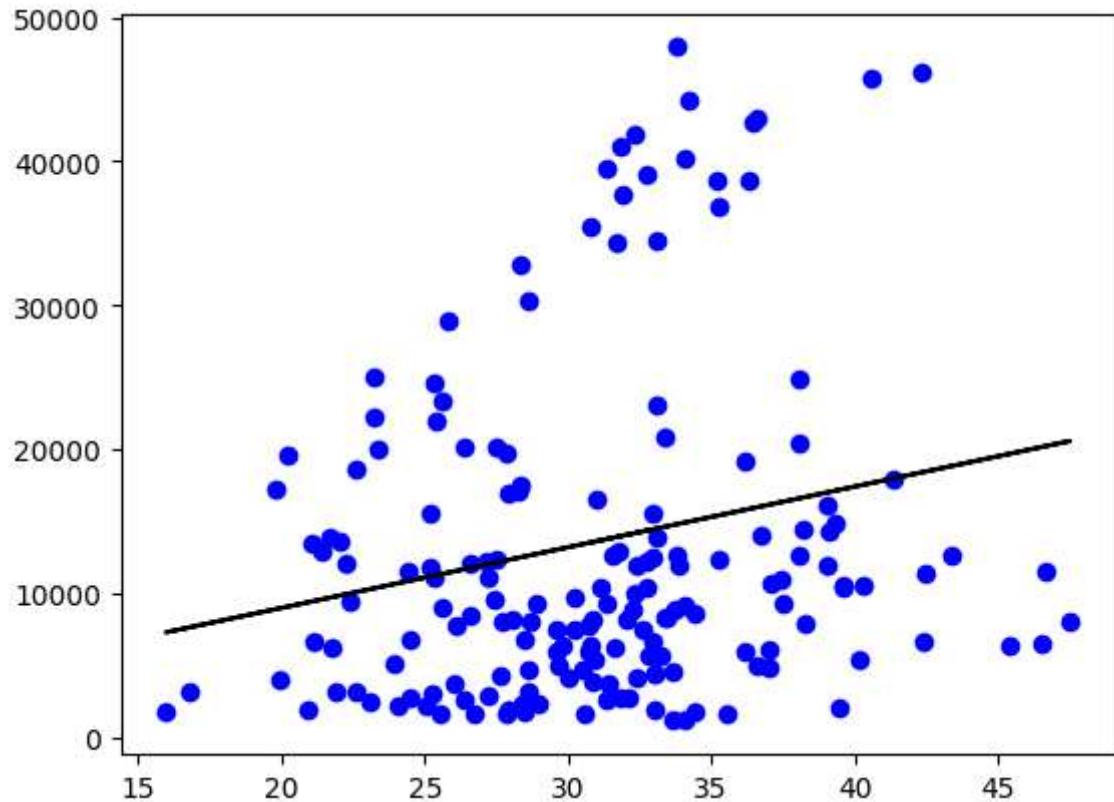
```
In [18]: x=np.array(df700["bmi"]).reshape(-1,1)
y=np.array(df700['charges']).reshape(-1,1)
```

```
In [19]: df700.dropna(inplace=True)
```

```
In [20]: 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))
```

0.02101977957759371

```
In [21]: 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 [22]: from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score
```

```
In [23]: lr=LinearRegression()
lr.fit(x_train,y_train)
y_pred=lr.predict(x_test)
r2=r2_score(y_test,y_pred)
print(r2)
```

0.02101977957759371

## Ridge Regression

```
In [24]: from sklearn.linear_model import Lasso,Ridge
from sklearn.preprocessing import StandardScaler
```

```
In [28]: features=df.columns[0:1]
target=df.columns[-1]
```

```
In [29]: 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=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 [30]: 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 [32]: 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 [33]: plt.figure(figsize=(10,10))
```

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

## Lasso Regression

```
In [34]: 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))
```

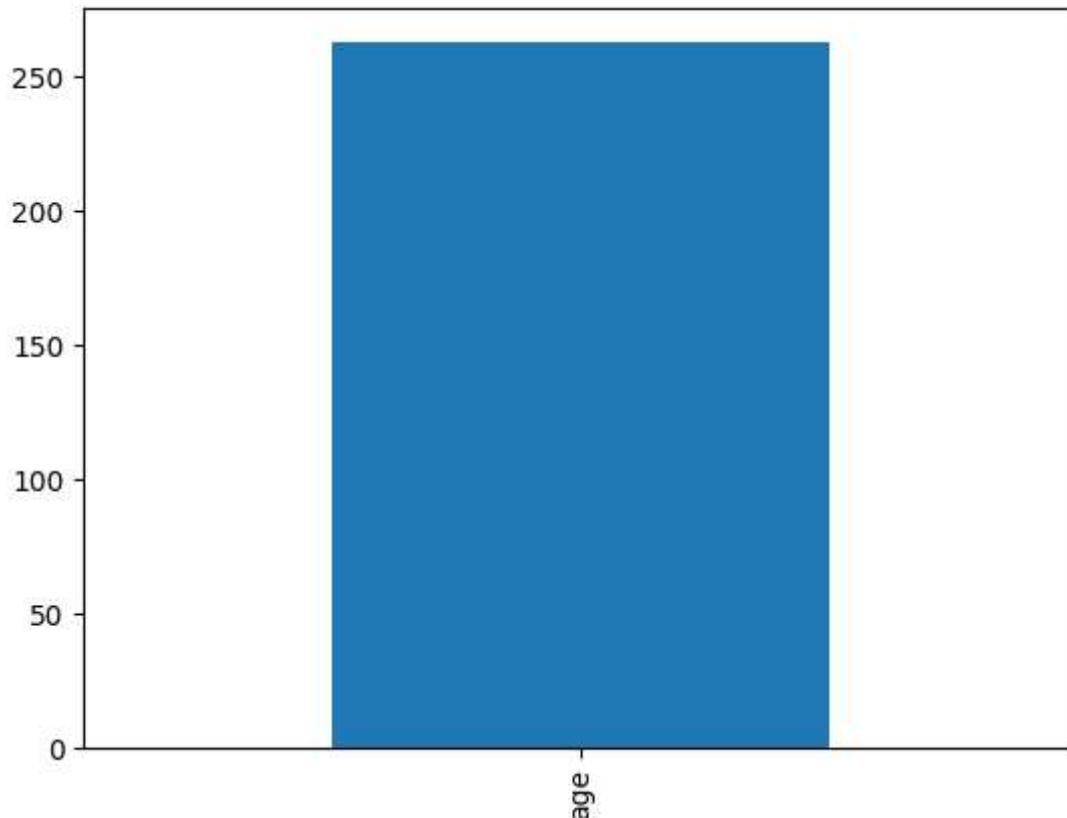
Ridge Model:

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

```
In [35]: plt.figure(figsize=(10,10))
```

```
Out[35]: <Figure size 1000x1000 with 0 Axes>
<Figure size 1000x1000 with 0 Axes>
```

```
In [36]: pd.Series(lasso.coef_, features).sort_values(ascending = True).plot(kind = "bar")
plt.show()
```



```
In [37]: from sklearn.linear_model import LassoCV
```

```
In [38]: #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 [39]: #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
```

## Elasticnet regression

```
In [40]: from sklearn.linear_model import ElasticNet
```

```
In [41]: el=ElasticNet()
el.fit(x_train,y_train)
print(el.coef_)
print(el.intercept_)
```

```
[261.74450967]
3115.083177426244
```

```
In [42]: y_pred_elastic=el.predict(x_train)
```

```
In [43]: mean_squared_error=np.mean((y_pred_elastic-y_train)**2)
print(mean_squared_error)
```

```
135077142.70714515
```

```
In [44]: el=ElasticNet()
el.fit(x_train,y_train)
print(el.score(x_train,y_train))
```

```
0.09109580670592365
```

## Logistic Regression

```
In [45]: import numpy as np
import pandas as pd
from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import StandardScaler
```

```
In [46]: df=pd.read_csv(r"C:\Users\sudheer\Downloads\insurance.csv")
df
```

Out[46]:

	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 [47]: df.shape
```

Out[47]: (1338, 7)

```
In [48]: pd.set_option('display.max_rows',10000000000)
pd.set_option('display.max_columns',10000000000)
pd.set_option('display.width',95)
```

```
In [49]: print('This Dataset has %d rows and %d columns'%(df.shape))
```

This Dataset has 1338 rows and 7 columns

```
In [50]: df.head()
```

Out[50]:

	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 [51]: df.describe
```

```
Out[51]: <bound method NDFrame.describe of  
er      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  
--  --  --  --  --  --  --  --  
          ^  ^  ^  ^  ^  ^  ^  ^
```

```
In [52]: 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 [53]: df.isnull().sum()
```

```
Out[53]: age      0  
sex      0  
bmi      0  
children  0  
smoker    0  
region    0  
charges    0  
dtype: int64
```

```
In [54]: convert={"smoker":{"yes":1,"no":0}}
df=df.replace(convert)
df
```

```
Out[54]:
```

	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
11	62	female	26.290	0	1	southeast	27808.725100

```
In [55]: convert={"sex":{"female":1,"male":0}}
df=df.replace(convert)
df
```

```
Out[55]:
```

	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

```
In [56]: convert={"region":{"southeast":1,"southwest":2,"northeast":3,"northwest":4}}  
df=df.replace(convert)  
df
```

Out[56]:

	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
11	62	1	26.290	0	1	1	27808.725100

```
In [57]: features_matrix=df.iloc[:,0:4]
```

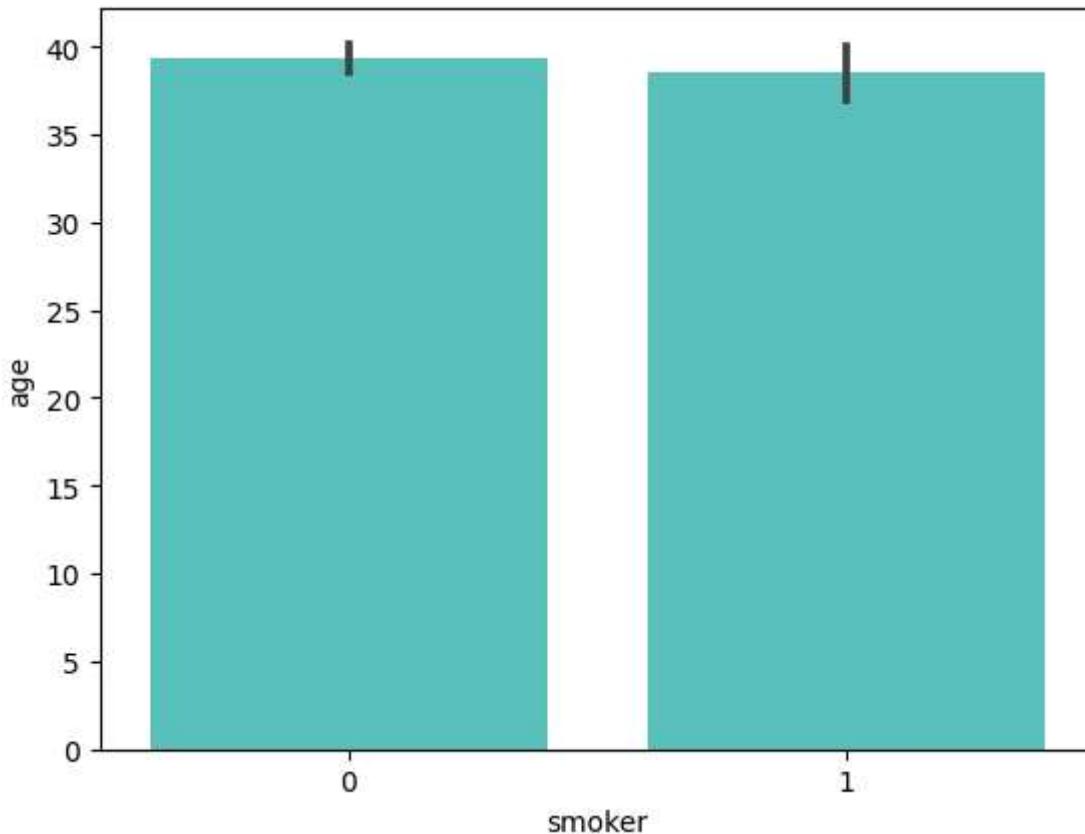
```
In [58]: target_vector=df.iloc[:, -3]
```

```
In [59]: 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 [60]: import matplotlib.pyplot as plt  
import seaborn as sns
```

```
In [61]: sns.barplot(x='smoker', y='age', data=df, color="mediumturquoise")
plt.show()
```



```
In [62]: features_matrix_standardized=StandardScaler().fit_transform(features_matrix)
```

```
In [63]: algorithm=LogisticRegression(max_iter=10000)
```

```
In [64]: Logistic_Regression_Model=algorithm.fit(features_matrix_standardized,target_ve
```

```
In [65]: observation=[[1,0,0.99539,-0.0588]]
```

```
In [66]: 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 [67]: print('The algoritham was trained to predict one of the two classes:%s'%(algor
```

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

```
In [68]: x=np.array(df['age']).reshape(-1,1)
y=np.array(df['smoker']).reshape(-1,1)
```

```
In [69]: 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))
```

```
0.7611940298507462
```

```
C:\Users\sudheer\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 [70]: 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 [71]: df=pd.read_csv(r"C:\Users\sudheer\Downloads\insurance.csv")
df
```

Out[71]:

	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
11	62	female	26.290	0	yes	southeast	27808.725100

```
In [72]: df.shape
```

```
Out[72]: (1338, 7)
```

```
In [73]: df.isnull().any()
```

```
Out[73]: age      False
          sex      False
          bmi      False
         children  False
        smoker   False
       region    False
      charges   False
      dtype: bool
```

```
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  16884.924000
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
```

```
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	16884.924000
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

```
In [77]: x=["bmi","children"]
y=["Yes","No"]
all_inputs=df[x]
all_classes=df["sex"]
```

```
In [78]: (x_train,x_test,y_train,y_test)=train_test_split(all_inputs,all_classes,test_s
```

```
In [79]: clf=DecisionTreeClassifier(random_state=0)
```

```
In [80]: clf.fit(x_train,y_train)
```

```
Out[80]: 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 [81]: score=clf.score(x_test,y_test)
print(score)
```

```
0.4878048780487805
```

## Random Forest

```
In [82]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt ,seaborn as sns
```

```
In [83]: df=pd.read_csv(r"C:\Users\sudheer\Downloads\insurance.csv")
df
```

Out[83]:

	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
11	62	female	26.290	0	yes	southeast	27808.725100

```
In [84]: df.shape
```

Out[84]: (1338, 7)

```
In [85]: df['region'].value_counts()
```

Out[85]:

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

```
In [86]: df['bmi'].value_counts()
```

```
Out[86]: 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  
31.320     7
```

```
In [87]: 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
18	55	0	32.200	2	no	southeast	16662.225200

```
In [88]: 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
18	55	0	30.260	0	0	northeast	10600.205000

```
In [89]: from sklearn.ensemble import RandomForestClassifier
rfc=RandomForestClassifier()
rfc.fit(x_train,y_train)
```

Out[89]: 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](http://nbviewer.org).

```
In [90]: 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 [91]: 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[91]: 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](#).

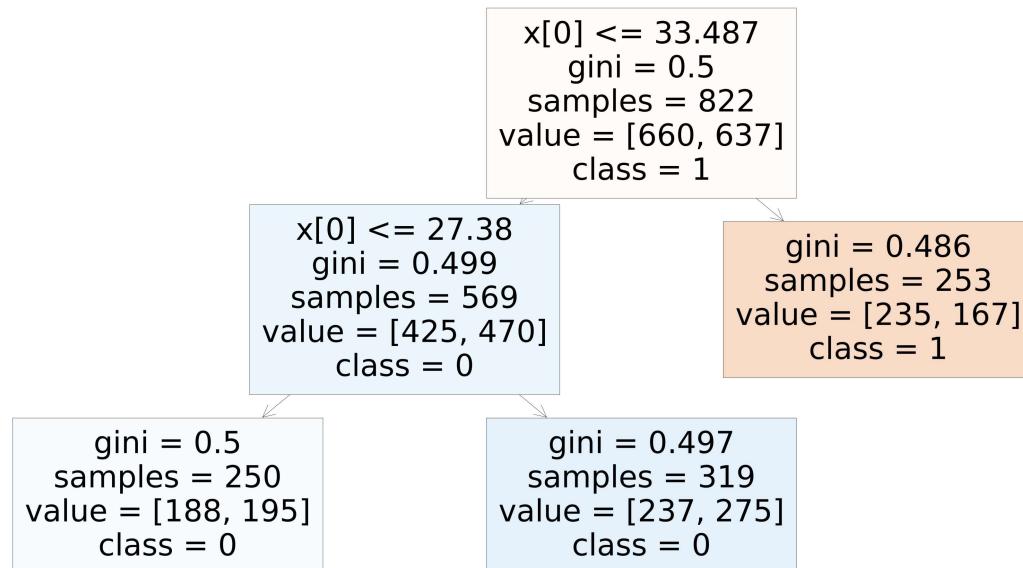
```
In [92]: grid_search.best_score_
```

```
Out[92]: 0.5219770682341304
```

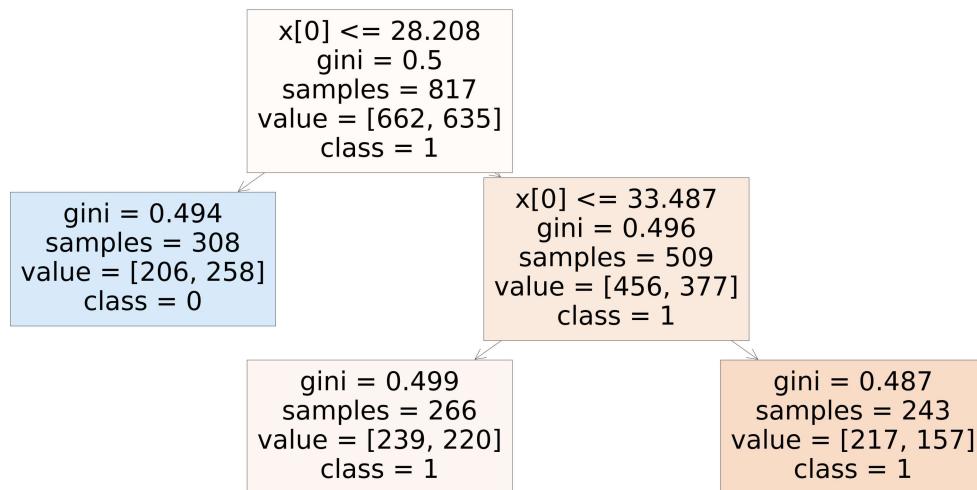
```
In [93]: rf_best=grid_search.best_estimator_
print(rf_best)
```

```
RandomForestClassifier(max_depth=2, min_samples_leaf=200, n_estimators=10)
```

```
In [94]: from sklearn.tree import plot_tree
plt.figure(figsize=(80,40))
plot_tree(rf_best.estimators_[4],class_names=['1','0'],filled=True);
```



```
In [95]: from sklearn.tree import plot_tree
plt.figure(figsize=(70,30))
plot_tree(rf_best.estimators_[6], class_names=["1","0"], filled=True);
```



```
In [96]: rf_best.feature_importances_
```

```
Out[96]: array([0.78802287, 0.21197713])
```

```
In [97]: rf=RandomForestClassifier(random_state=0)
```

```
In [98]: rf.fit(x_train,y_train)
```

```
Out[98]: 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 [99]: score=rf.score(x_test,y_test)
print(score)
```

```
0.5365853658536586
```

**Conclusion: When compared to the above accuracies Logistic regression is getting**