# problem Statement:To predict How Best the data fits

## 1)Data Collection

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

In [2]:
```python
df=pd.read_csv(r"C:\Users\chinta pavani\Documents\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

## 2)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]:
```
<bound method NDFrame.describe of        age     sex     bmi   children smo
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
...    ...     ...     ...       ...    ...        ...          ...
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]:
```python
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 |

In [12]:
```python
convert={"smoker":{"yes":1, "no":0}}
df=df.replace(convert)
df
```

Out[12]:

| | 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 [13]:  ▶| 
```python
convert={"region":{"southwest":1, "southeast":2,"northwest":3,"northeast":4
df=df.replace(convert)
df
```
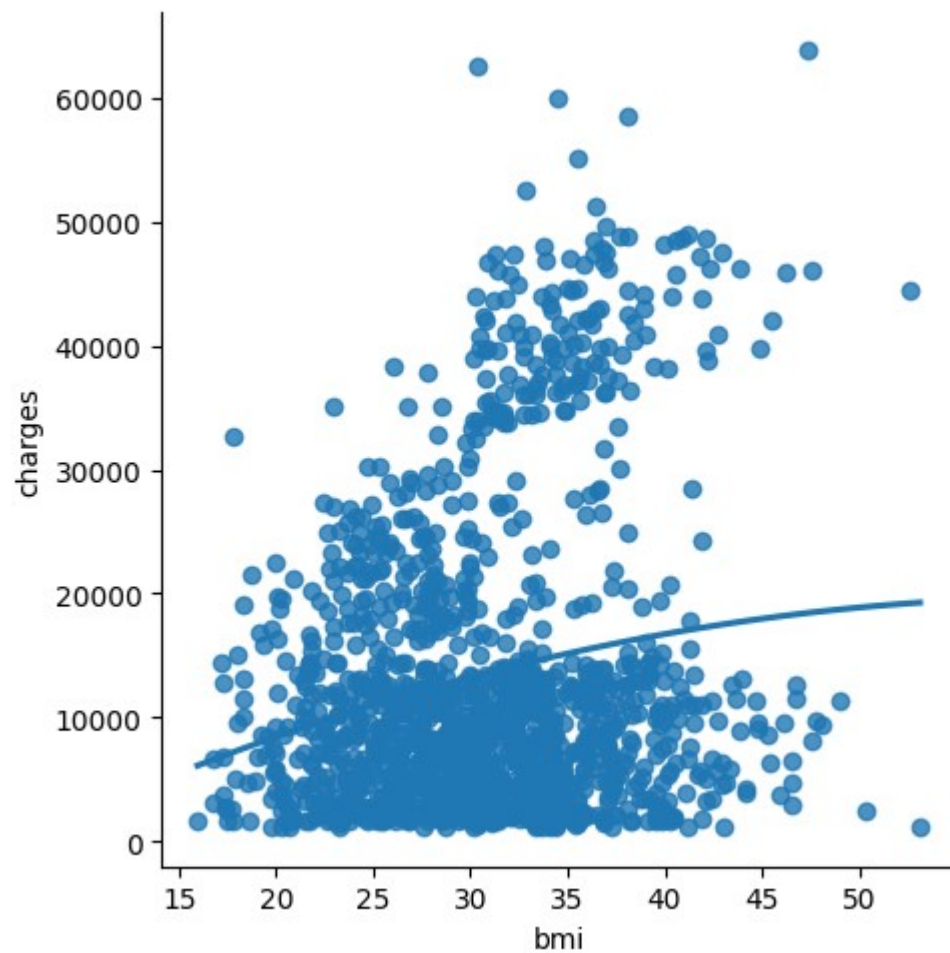
Out[13]:

|  | 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 [14]: 

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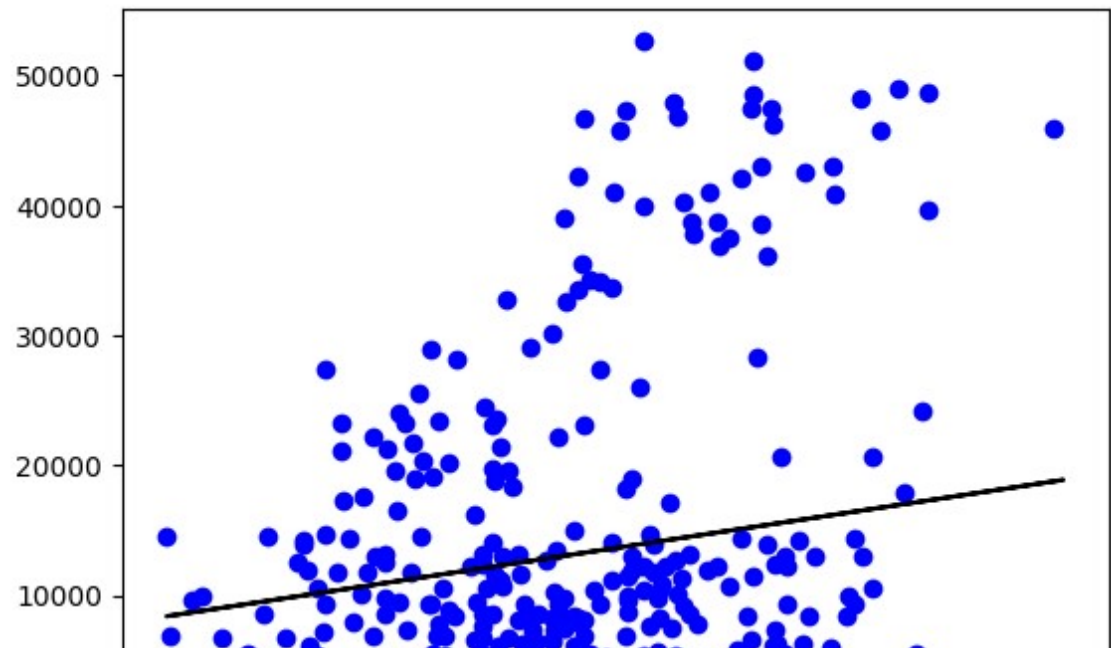


In [15]: 

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

In [16]: 

```python
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25,random_s
lr=LinearRegression()
lr.fit(x_train,y_train)
print(lr.score(x_test,y_test))
```
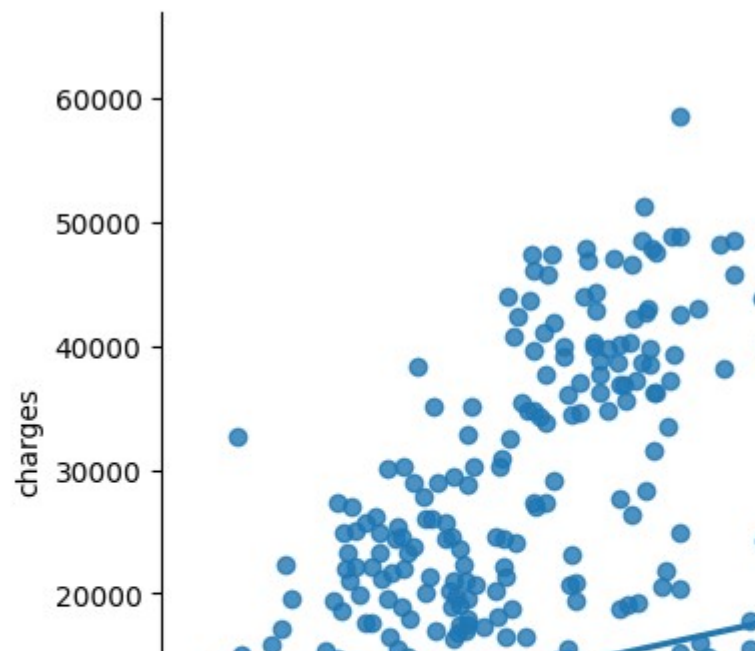
```
0.050136213258239914
```

In [17]: ▶|
```python
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 [18]: ▶|
```python
df700=df[:][:700]
sns.lmplot(x="bmi",y="charges",order=2,ci=None,data=df700)
plt.show()
```



In [19]: ▶|
```python
df700.fillna(method='ffill',inplace=True)
```
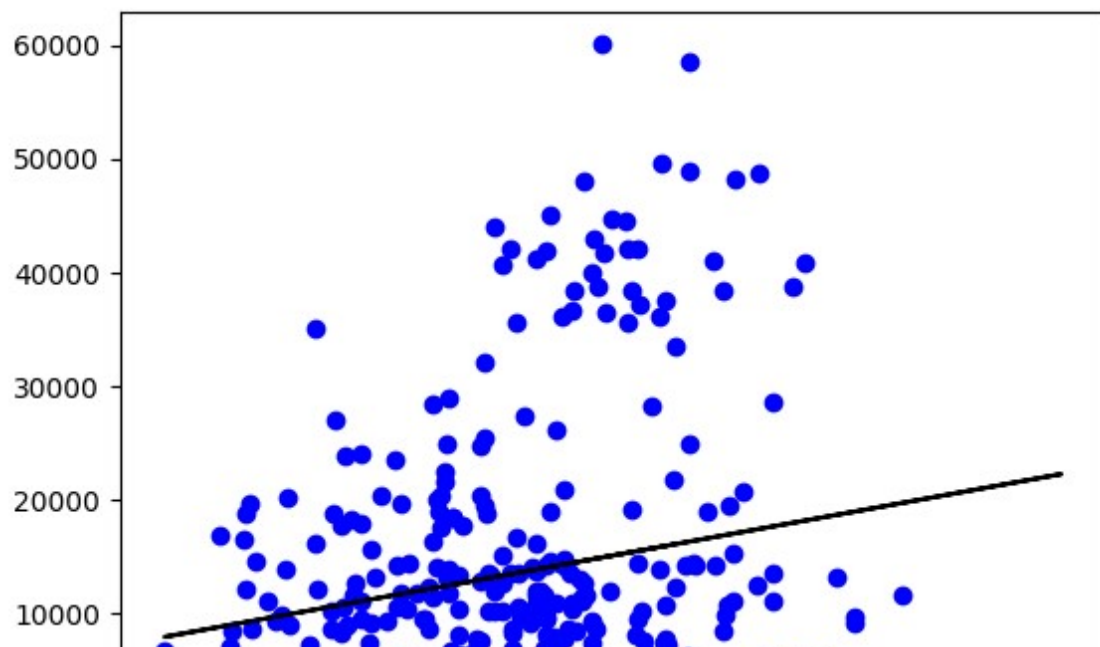
In [20]:
```python
x=np.array(df['bmi']).reshape(-1,1)
y=np.array(df['charges']).reshape(-1,1)
```

In [21]:
```python
df700.dropna(inplace=True)
```

In [22]:
```python
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))
```

0.03834819458441918

In [23]:
```python
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 [24]:
```python
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)
```

0.03834819458441918

# Ridge Regression

In [25]:
```python
from sklearn.linear_model import Lasso,Ridge
from sklearn.preprocessing import StandardScaler
```

In [26]:
```python
plt.figure(figsize= (10,10))
sns.heatmap(df700.corr(),annot = True)
plt.show()
```



In [27]:
```python
features =df.columns[0:1]
target = df.columns[-1]
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,randol
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)
```

```python
In [28]:    lr = LinearRegression()
            #Fit model
            lr.fit(x_train,y_train)
            #Predict
            #Prediction =Lr.predict(x_test)
            #actual
            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
```

```python
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))
             plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',mar
             plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersiz
             plt.xticks(rotation =90)
             plt.legend()
             plt.show()
```



## Lasso Regression

In [31]: 
```python
lasso=Lasso(alpha=10)
lasso.fit(x_train,y_train)
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 ls model is {}".format(train_score_ls))
print("The test score for ls model is {}".format(test_score_ls))
```

```
Ridge Model:

The train score for ls model is 0.09109639395809055
The test score for ls model is 0.08490704421828055
```
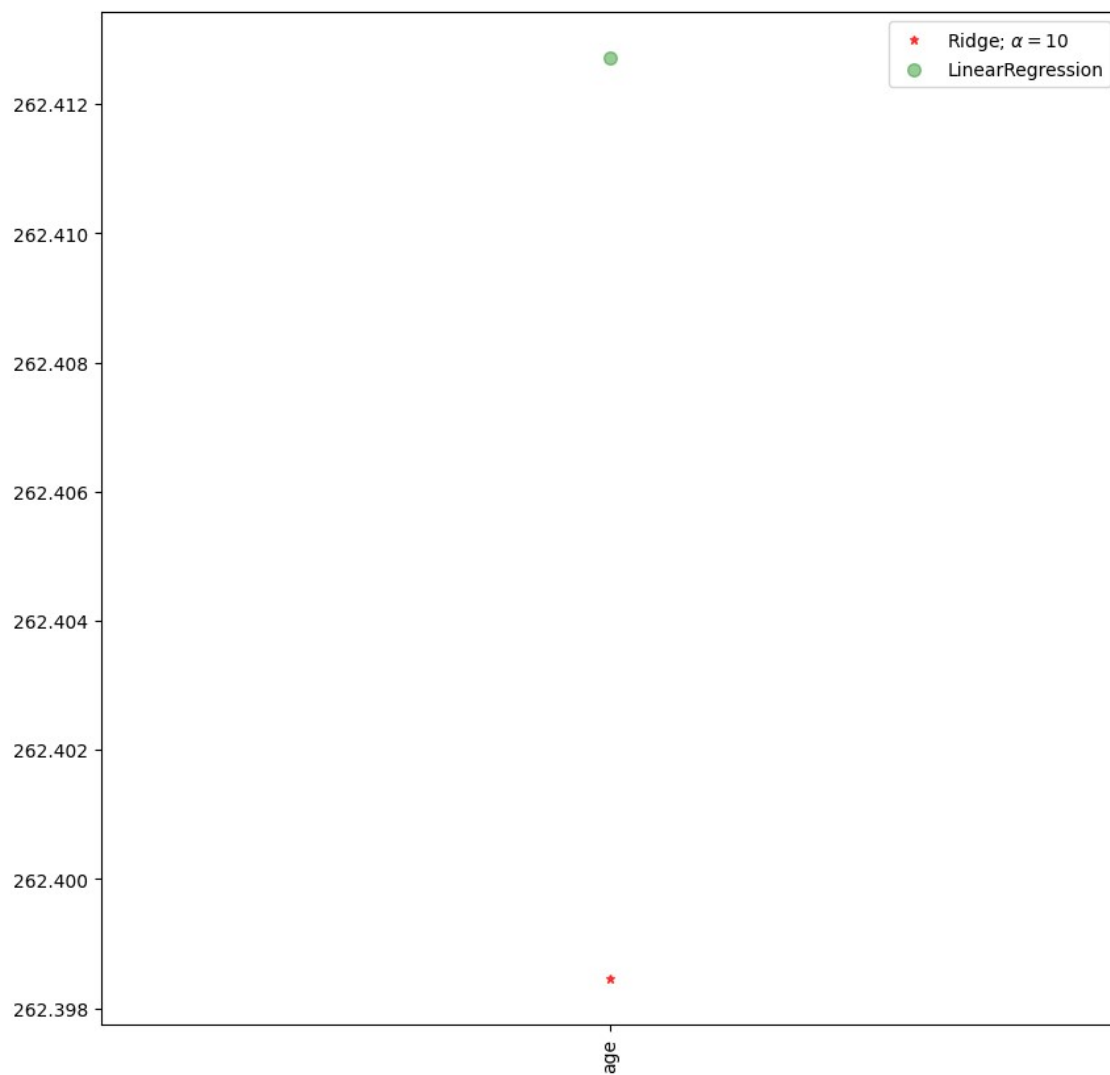
In [32]: 
```python
plt.figure(figsize=(10,10))
pd.Series(lasso.coef_, features).sort_values(ascending =True).plot(kind =
plt.show()
```



In [33]: 
```python
from sklearn.linear_model import LassoCV
```

In [34]: ▶| 
```python
#using the Linear CV model
from sklearn.linear_model import RidgeCV
#Lasso Cross validation
ridge_cv = RidgeCV(alphas = [0.0001,0.001,0.01,0.1,1,10]).fit(x_train,y_tra
#score
print(ridge_cv.score(x_train,y_train))
print(ridge_cv.score(x_test,y_test))
```

```
0.09109639711159612
0.08490538609884779
```

In [47]: ▶| 
```python
#using the linear CV model
from sklearn.linear_model import LassoCV
#Ridge Cross validation
lasso_cv=LassoCV(alphas = [0.0001,0.001,0.01,1,10]).fit(x_train,y_train)
#score
print("The train score for ls model is {}".format(ridge_cv.score(x_train,y_
print("The test score for ls model is {}".format(ridge_cv.score(x_test,y_t
```

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

In [48]: ▶| 
```python
plt.figure(figsize = (10, 10))
#add plot for ridge regression
plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',mar
#add plot for lasso regression
plt.plot(lasso_cv.coef_,alpha=0.5,linestyle='none',marker='d',markersize=6
#add plot for linear model
plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersiz
#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

# Elastic Net Regression

In [49]:
```python
from sklearn.linear_model import ElasticNet
el=ElasticNet()
el.fit(x,y)
print(el.coef_)
print(el.intercept_)
```

```
[257.0684655]
3191.532406056678
```

In [50]:
```python
y_pred_elastic=el.predict(x_train)
```

In [51]:
```python
mean_squared_error=np.mean((y_pred_elastic-y_train)**2)
print("Mean Squared Error on test",mean_squared_error)
```

```
Mean Squared Error on test 135093713.48214507
```

In [52]:
```python
el=ElasticNet()
el.fit(x_train,y_train)
print(el.score(x_train,y_train))
```

```
0.09109580670592365
```

# Logistic Regression

In [80]:
```python
import numpy as np
import pandas as pd
from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import StandardScaler
```

```
In [81]:   ▶ df=pd.read_csv(r"C:\Users\chinta pavani\Documents\insurance.csv")
             df
```

Out[81]:

|    | 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 [82]:   ▶ df.shape
```

Out[82]:  (1338, 7)

```
In [83]:   ▶ pd.set_option('display.max_rows',10000000000)
             pd.set_option('display.max_columns',10000000000)
             pd.set_option('display.width',95)
```

```
In [84]:   ▶ print('This Dataset has %d rows and %d columns'%(df.shape))
```

This Dataset has 1338 rows and 7 columns

```
In [85]:   ▶ df.head()
```

Out[85]:

|   | 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 [86]:   ▶| df.describe

Out[86]:  <bound method NDFrame.describe of         age      sex      bmi   children  smo
          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    23    male  23.845         0     no  northeast   2205.171550

In [87]:   ▶| 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 [88]:   ▶| df.isnull().sum()

Out[88]:  age        0
          sex        0
          bmi        0
          children   0
          smoker     0
          region     0
          charges    0
          dtype: int64

In [89]: ▶| 
```python
convert={"smoker":{"yes":1,"no":0}}
df=df.replace(convert)
df
```

Out[89]:

|    | 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 [90]: ▶| 
```python
convert={"sex":{"female":1,"male":0}}
df=df.replace(convert)
df
```

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

Out[91]:

| | 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 [92]: ▶
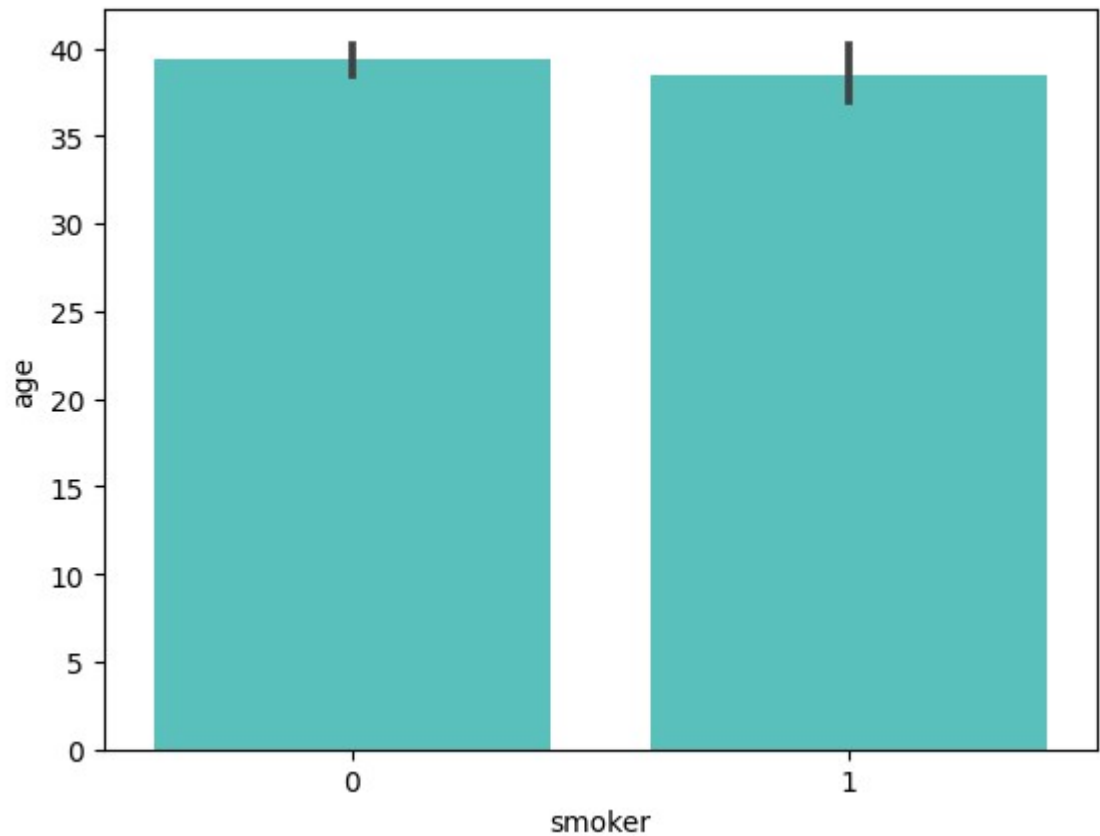```python
features_matrix=df.iloc[:,0:4]
```

In [93]: ▶
```python
target_vector=df.iloc[:,-3]
```

In [94]: ▶
```python
print('The Feature Matrix has %d Rows and %d columns(s)'%(features_matrix.
print('The Target Matrix has %d Rows and %d columns(s)'%(np.array(target_v
```

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

In [95]: ▶|
```python
import matplotlib.pyplot as plt
import seaborn as sns
sns.barplot(x='smoker', y='age', data=df, color="mediumturquoise")
plt.show()
```



In [136]: ▶|
```python
features_matrix_standardized=StandardScaler().fit_transform(features_matri
```

In [137]: ▶|
```python
algorithm=LogisticRegression(max_iter=10000)
```

In [138]: ▶|
```python
Logistic_Regression_Model=algorithm.fit(features_matrix_standardized,targe
```

In [139]: ▶|
```python
observation=[[1,0,0.99539,-0.0588]]
```

In [140]: ▶|
```python
predictions=Logistic_Regression_Model.predict(observation)
print('The model predicted the observation to belong to class %s'%(predict
```

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

In [141]: ▶|
```python
print('The algoritham was trained to predict one of the two classes:%s'%(a
```

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

In [142]: ▶ 
```python
print(" " "The Model says the probability of the observation we passed bel
print()
```

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

In [148]: ▶ 
```python
print(" " "The Model says the probability of the observation we passed bel
print()
```

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

In [146]: ▶ 
```python
x=np.array(df['age']).reshape(-1,1)
y=np.array(df['smoker']).reshape(-1,1)
```

In [102]: ▶ 
```python
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))
```

0.7910447761194029

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

## Decision Tree

In [103]: ▶ 
```python
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 [104]: ▶| 
```python
df=pd.read_csv(r"C:\Users\chinta pavani\Documents\insurance.csv")
df
```

Out[104]:

| | 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 [105]: ▶| 
```python
df.shape
```

Out[105]: (1338, 7)

In [106]: ▶| 
```python
df.isnull().any()
```

Out[106]: 
```
age        False
sex        False
bmi        False
children   False
smoker     False
region     False
charges    False
dtype: bool
```

In [107]: ▶| 
```python
df['region'].value_counts()
```

Out[107]: 
```
region
southeast    364
southwest    325
northwest    325
northeast    324
Name: count, dtype: int64
```

In [108]:
```python
convert={"sex":{"female":1,"male":0}}
df=df.replace(convert)
df
```

Out[108]:

|  | 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 [109]:
```python
convert={"smoker":{"yes":1,"no":0}}
df=df.replace(convert)
df
```

Out[109]:

|  | 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 [110]:
```python
x=["bmi","children"]
y=["Yes","No"]
all_inputs=df[x]
all_classes=df["sex"]
```

In [111]:
```python
(x_train,x_test,y_train,y_test)=train_test_split(all_inputs,all_classes,te
```

In [112]: ▶| `clf=DecisionTreeClassifier(random_state=0)`

In [113]: ▶| `clf.fit(x_train,y_train)`

Out[113]:
```
▾            DecisionTreeClassifier

DecisionTreeClassifier(random_state=0)
```

In [114]: ▶|
```
DecisionTreeClassifier(random_state=0)
score=clf.score(x_test,y_test)
print(score)
```

0.4878048780487805

# Random Forest

In [115]: ▶|
```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt ,seaborn as sns
```

In [116]: ▶|
```python
df=pd.read_csv(r"C:\Users\chinta pavani\Documents\insurance.csv")
df
```

Out[116]:

|    | 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 [117]: ▶| `df.shape`

Out[117]: `(1338, 7)`

```
In [118]:    ▶  df['region'].value_counts()
```

```
Out[118]:   region
            southeast    364
            southwest    325
            northwest    325
            northeast    324
            Name: count, dtype: int64
```

```
In [119]:    ▶  df['bmi'].value_counts()
```

```
Out[119]:   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
            24.330     7
```

```
In [121]:    ▶  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   56    0  40.300         0     no  southwest  10602.385000
```

```
In [122]:    ▶ 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 [123]:    ▶ from sklearn.ensemble import RandomForestClassifier
               rfc=RandomForestClassifier()
               rfc.fit(x_train,y_train)
```

Out[123]:
```
▾ RandomForestClassifier

RandomForestClassifier()
```

```
In [124]:    ▶ 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 [126]:    ▶ from sklearn.model_selection import GridSearchCV
               grid_search=GridSearchCV(estimator=rf,param_grid=params,cv=2,scoring="accu
               grid_search.fit(x_train,y_train)
```

Out[126]:
```
▸          GridSearchCV

▸ estimator: RandomForestClassifier

    ▸ RandomForestClassifier
```
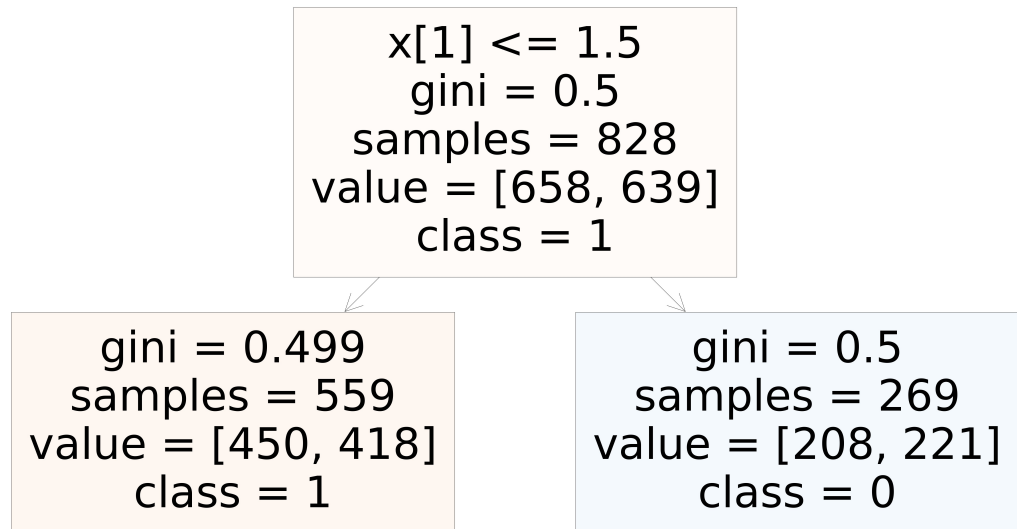
```
In [127]:    ▶ grid_search.best_score_
```
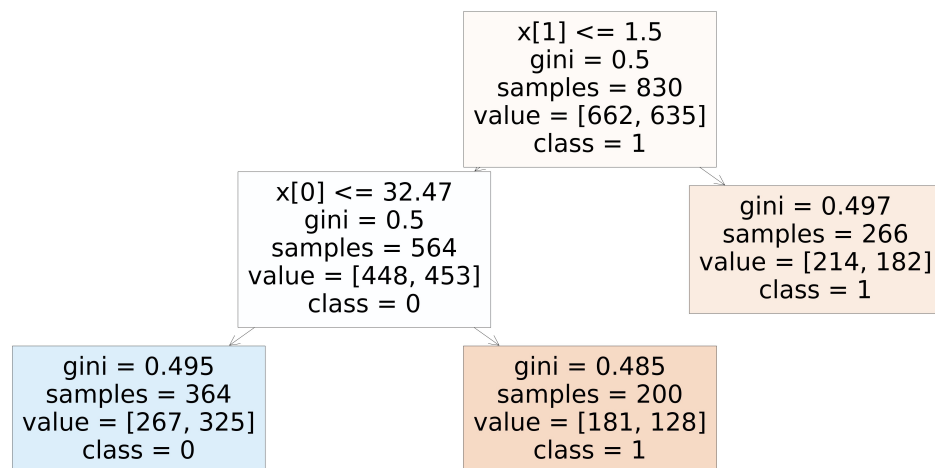
Out[127]:  0.5304492666780802

In [128]: ▶|
```python
rf_best=grid_search.best_estimator_
print(rf_best)
```

RandomForestClassifier(max_depth=3, min_samples_leaf=200, n_estimators=1
0)

In [129]: ▶|
```python
from sklearn.tree import plot_tree
plt.figure(figsize=(80,40))
plot_tree(rf_best.estimators_[4],class_names=['1','0'],filled=True);
```

```
x[1] <= 1.5
gini = 0.5
samples = 828
value = [658, 639]
class = 1
```

```
gini = 0.499
samples = 559
value = [450, 418]
class = 1
```

```
gini = 0.5
samples = 269
value = [208, 221]
class = 0
```

In [131]: ▶|
```python
from sklearn.tree import plot_tree
plt.figure(figsize=(70,30))
plot_tree(rf_best.estimators_[6],class_names=["1","0"],filled=True);
```

```
x[1] <= 1.5
gini = 0.5
samples = 830
value = [662, 635]
class = 1
```

```
x[0] <= 32.47
gini = 0.5
samples = 564
value = [448, 453]
class = 0
```

```
gini = 0.497
samples = 266
value = [214, 182]
class = 1
```

```
gini = 0.495
samples = 364
value = [267, 325]
class = 0
```

```
gini = 0.485
samples = 200
value = [181, 128]
class = 1
```

In [132]: ▶|
```python
rf_best.feature_importances_
```

Out[132]: array([0.53541719, 0.46458281])

In [133]: ▶|
```python
rf=RandomForestClassifier(random_state=0)
```

In [134]: ▶| `rf.fit(x_train,y_train)`

Out[134]:
```
▼              RandomForestClassifier

RandomForestClassifier(random_state=0)
```

In [135]: ▶|
```python
score=rf.score(x_test,y_test)
print(score)
```

```
0.4878048780487805
```

In [ ]: ▶|