

# Problem Statement:

To predict whether the given dataset is best fit or not.

In [1]:

```

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

```

In [2]:

```

1 #Step-2:Reading the dataset
2 df=pd.read_csv(r"C:\Users\91955\Downloads\insurance.csv")
3 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

In [3]:

```

1 df=df[['age','charges']]
2 #Taking only the selected two attributes from the dataset
3 df.columns=['Age','Charges']
4 #Renaming the columns for easier writing of the code

```

In [4]:

```
1 df.head(10)
2 #Displaying only the 1st 10 rows
```

Out[4]:

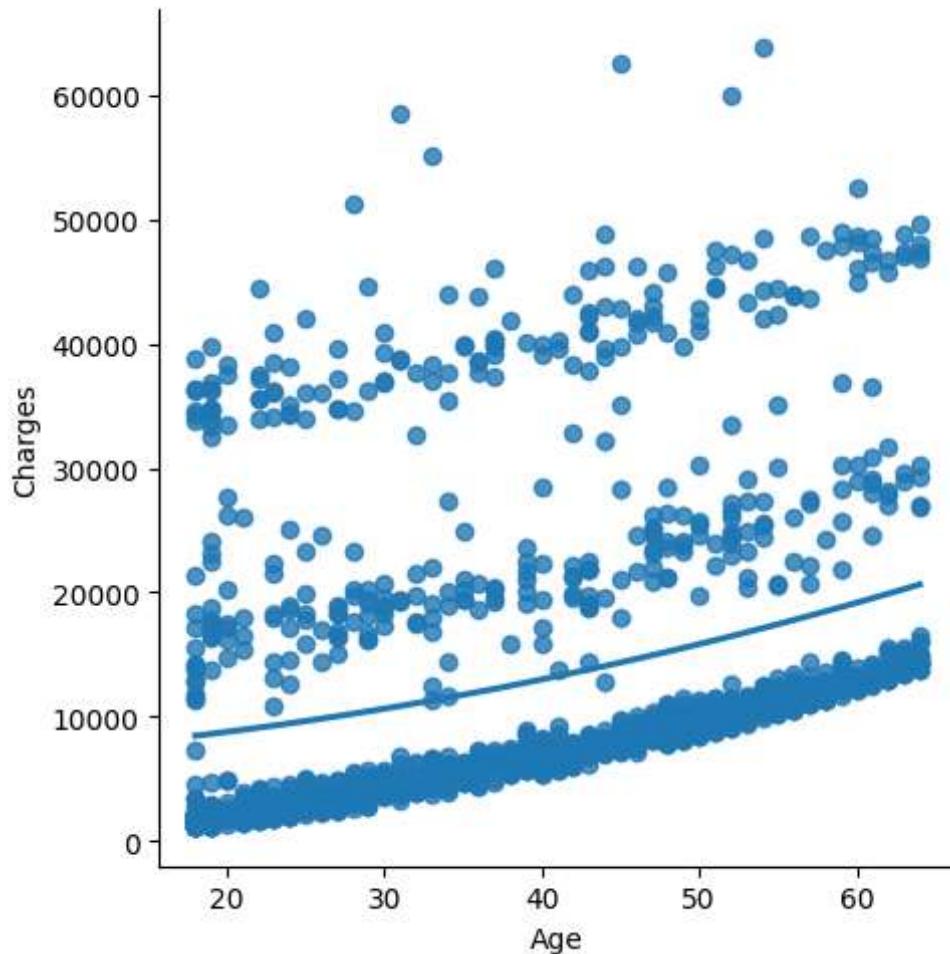
	Age	Charges
0	19	16884.92400
1	18	1725.55230
2	28	4449.46200
3	33	21984.47061
4	32	3866.85520
5	31	3756.62160
6	46	8240.58960
7	37	7281.50560
8	37	6406.41070
9	60	28923.13692

In [5]:

```
1 sns.lmplot(x='Age',y='Charges',data=df,order=2,ci=None)
2
```

Out[5]:

&lt;seaborn.axisgrid.FacetGrid at 0x16ec2c85e40&gt;



In [6]:

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

Out[6]:

	Age	Charges
count	1338.000000	1338.000000
mean	39.207025	13270.422265
std	14.049960	12110.011237
min	18.000000	1121.873900
25%	27.000000	4740.287150
50%	39.000000	9382.033000
75%	51.000000	16639.912515
max	64.000000	63770.428010

In [7]:

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

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1338 entries, 0 to 1337
Data columns (total 2 columns):
 #   Column   Non-Null Count  Dtype  
--- 
 0   Age       1338 non-null   int64  
 1   Charges   1338 non-null   float64 
dtypes: float64(1), int64(1)
memory usage: 21.0 KB
```

In [8]:

```
1 df.fillna(method='ffill')
```

Out[8]:

	Age	Charges
0	19	16884.92400
1	18	1725.55230
2	28	4449.46200
3	33	21984.47061
4	32	3866.85520
...	...	...
1333	50	10600.54830
1334	18	2205.98080
1335	18	1629.83350
1336	21	2007.94500
1337	61	29141.36030

1338 rows × 2 columns

In [9]:

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

Out[9]:

```
Age      0
Charges  0
dtype: int64
```

In [10]:

```
1 x=np.array(df['Age']).reshape(-1,1)
```

In [11]:

```
1 y=np.array(df['Charges']).reshape(-1,1)
```

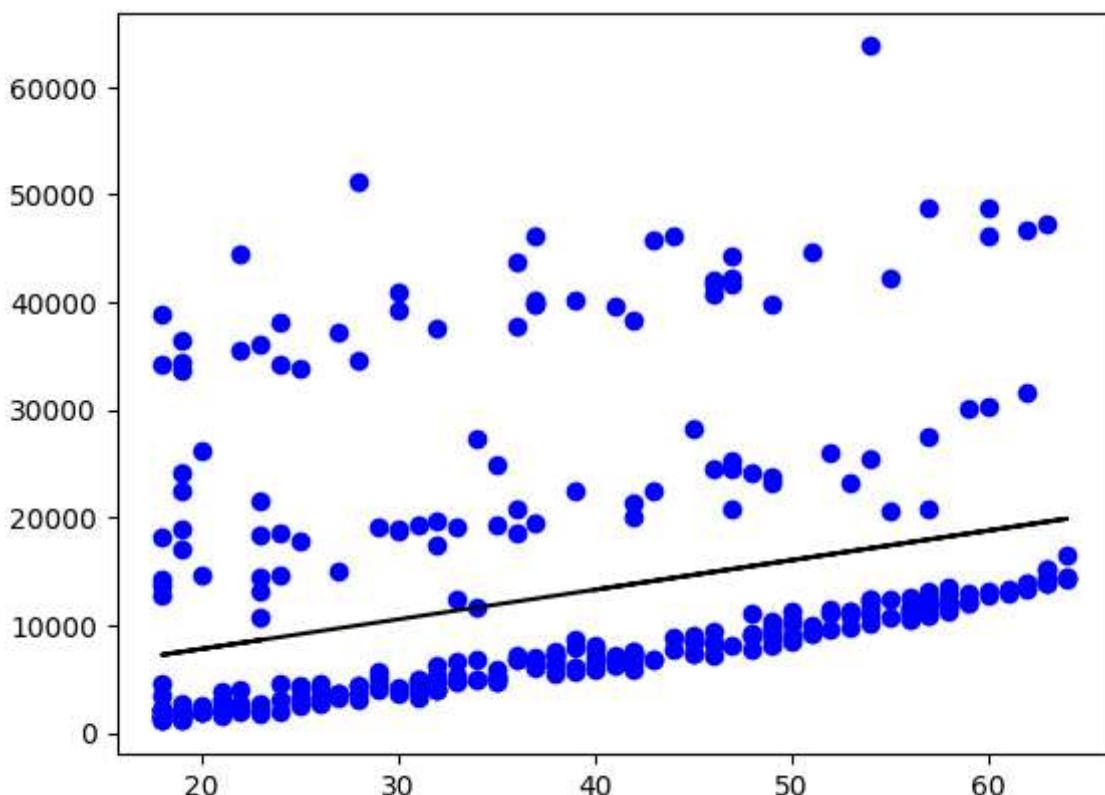
In [12]:

```
1 X_train,X_test,y_train,y_test=train_test_split(x,y,test_size=0.25)
2 reg=LinearRegression()
3 reg.fit(X_train,y_train)
4 print(reg.score(X_test,y_test))
```

0.04503038004350668

In [13]:

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

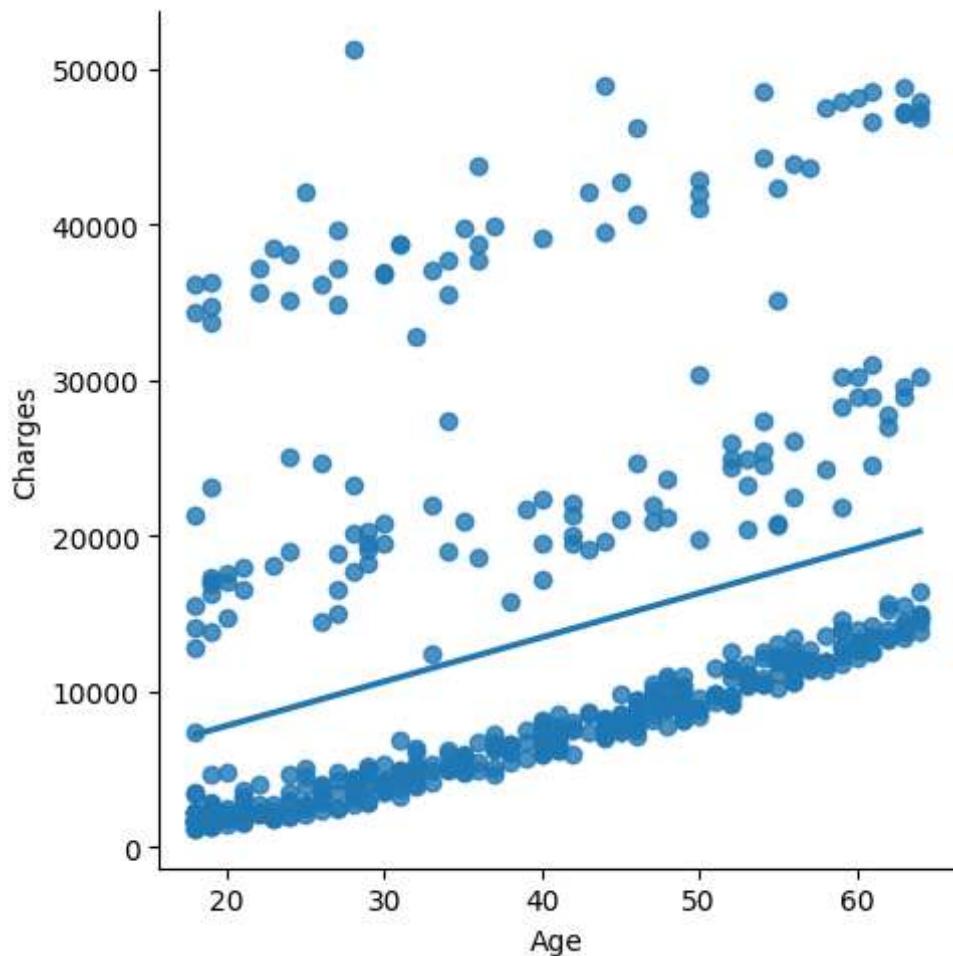


In [14]:

```
1 df500=df[:][:500]
2 sns.lmplot(x='Age',y='Charges',data=df500,order=1,ci=None)
```

Out[14]:

&lt;seaborn.axisgrid.FacetGrid at 0x16ec827b2b0&gt;



In [15]:

```

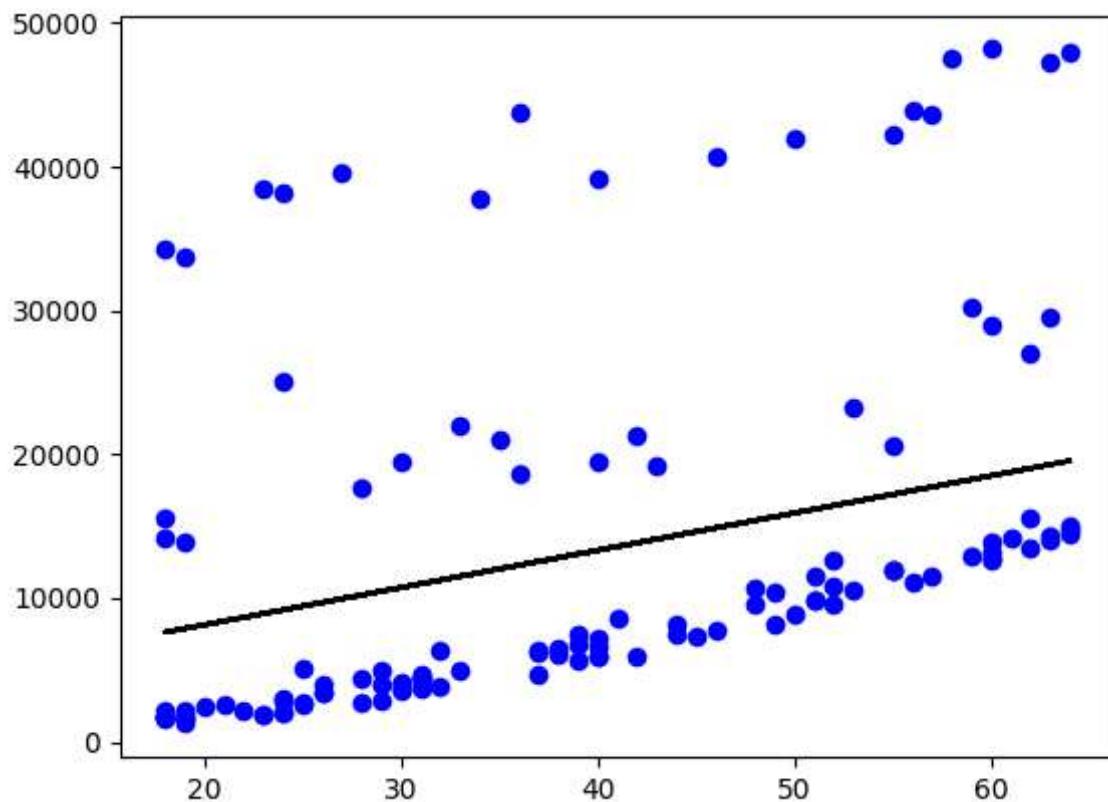
1 df500.fillna(method='ffill',inplace=True)
2 X=np.array(df500['Age']).reshape(-1,1)
3 y=np.array(df500['Charges']).reshape(-1,1)
4 df500.dropna(inplace=True)
5 X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.25)
6 reg=LinearRegression()
7 reg.fit(X_train,y_train)
8 print("Regression:",reg.score(X_test,y_test))
9 y_pred=reg.predict(X_test)
10 plt.scatter(X_test,y_test,color='b')
11 plt.plot(X_test,y_pred,color='k')
12 plt.show

```

Regression: 0.1525537912566497

Out[15]:

&lt;function matplotlib.pyplot.show(close=None, block=None)&gt;



In [16]:

```

1 from sklearn.linear_model import LinearRegression
2 from sklearn.metrics import r2_score
3 mode1=LinearRegression()
4 mode1.fit(X_train,y_train)
5 y_pred=mode1.predict(X_test)
6 r2=r2_score(y_test,y_pred)
7 print("R2 score: ",r2)

```

R2 score: 0.1525537912566497

In [17]:

```
1 from sklearn import metrics
2 print('MAE:',metrics.mean_absolute_error(y_test,y_pred))
3 print('MSE:',metrics.mean_squared_error(y_test,y_pred))
4 print('RMSE:',np.sqrt(metrics.mean_absolute_error(y_test,y_pred)))
```

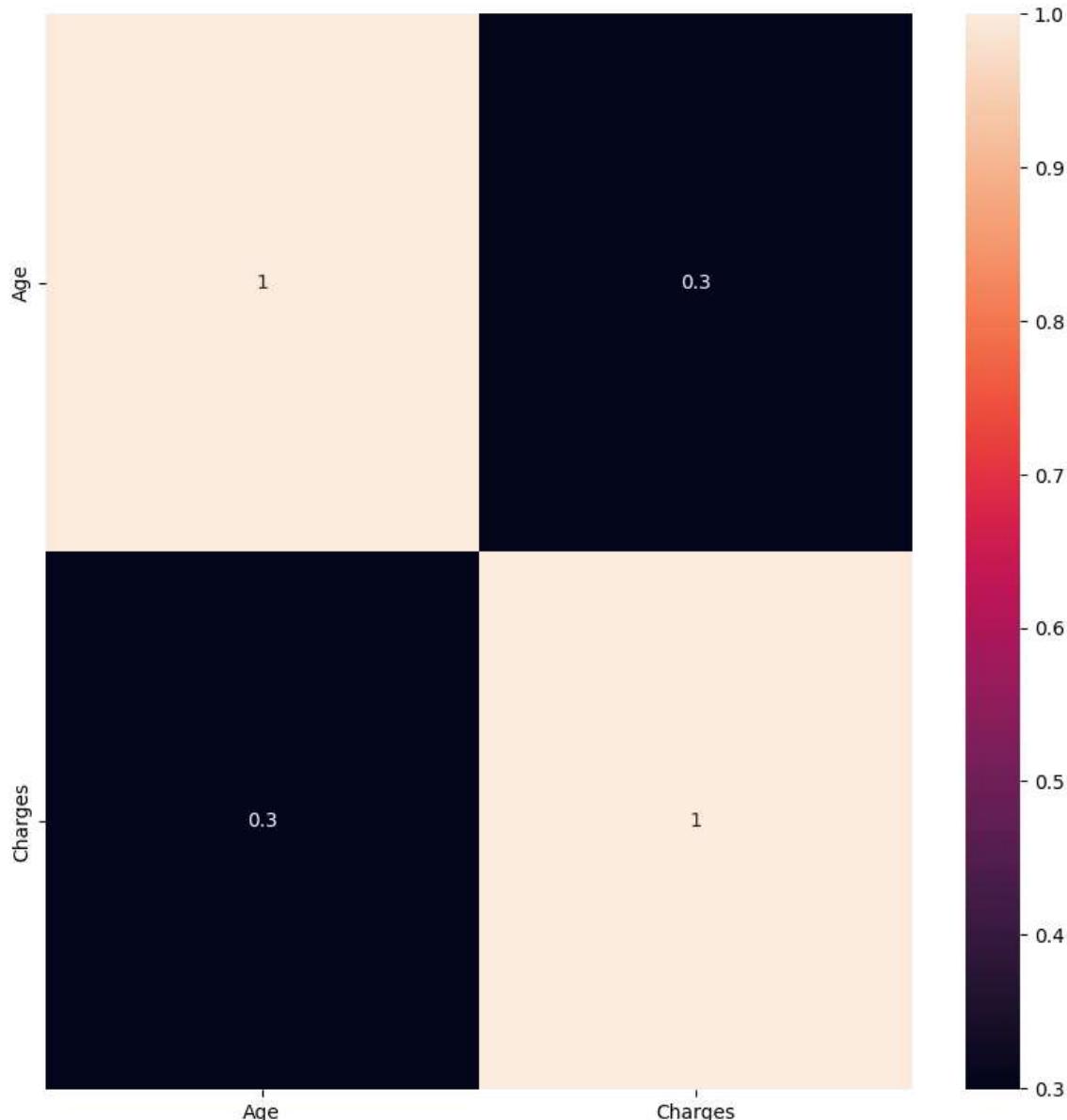
MAE: 9346.906572307227  
MSE: 142424837.01171198  
RMSE: 96.67940097201279

In [18]:

```
1 plt.figure(figsize = (10, 10))
2 sns.heatmap(df.corr(), annot = True)
```

Out[18]:

&lt;Axes: &gt;



## Ridge Regression

In [19]:

```

1 features = df.columns[0:2]
2 target = df.columns[-1]
3 #X and y values
4 X = df[features].values
5 y = df[target].values
6 #splot
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
8 print("The dimension of X_train is {}".format(X_train.shape))
9 print("The dimension of X_test is {}".format(X_test.shape))
10 #Scale features
11 from sklearn.preprocessing import StandardScaler
12 scaler = StandardScaler()
13 X_train = scaler.fit_transform(X_train)
14 X_test = scaler.transform(X_test)

```

The dimension of X\_train is (936, 2)

The dimension of X\_test is (402, 2)

In [20]:

```

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

```

Linear Regression Model:

The train score for lr model is 1.0

The test score for lr model is 1.0

In [21]:

```

1 #Using the Linear CV model
2 from sklearn.linear_model import RidgeCV
3 #Ridge Cross validation
4 ridge_cv = RidgeCV(alphas = [0.0001, 0.001, 0.01, 0.1, 1, 10]).fit(X_train, y_train)
5 #score
6 print("The train score for ridge model is {}".format(ridge_cv.score(X_train, y_train)))
7 print("The test score for ridge model is {}".format(ridge_cv.score(X_test, y_test)))

```

The train score for ridge model is 0.999999999999869

The train score for ridge model is 0.9999999999999875

In [22]:

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

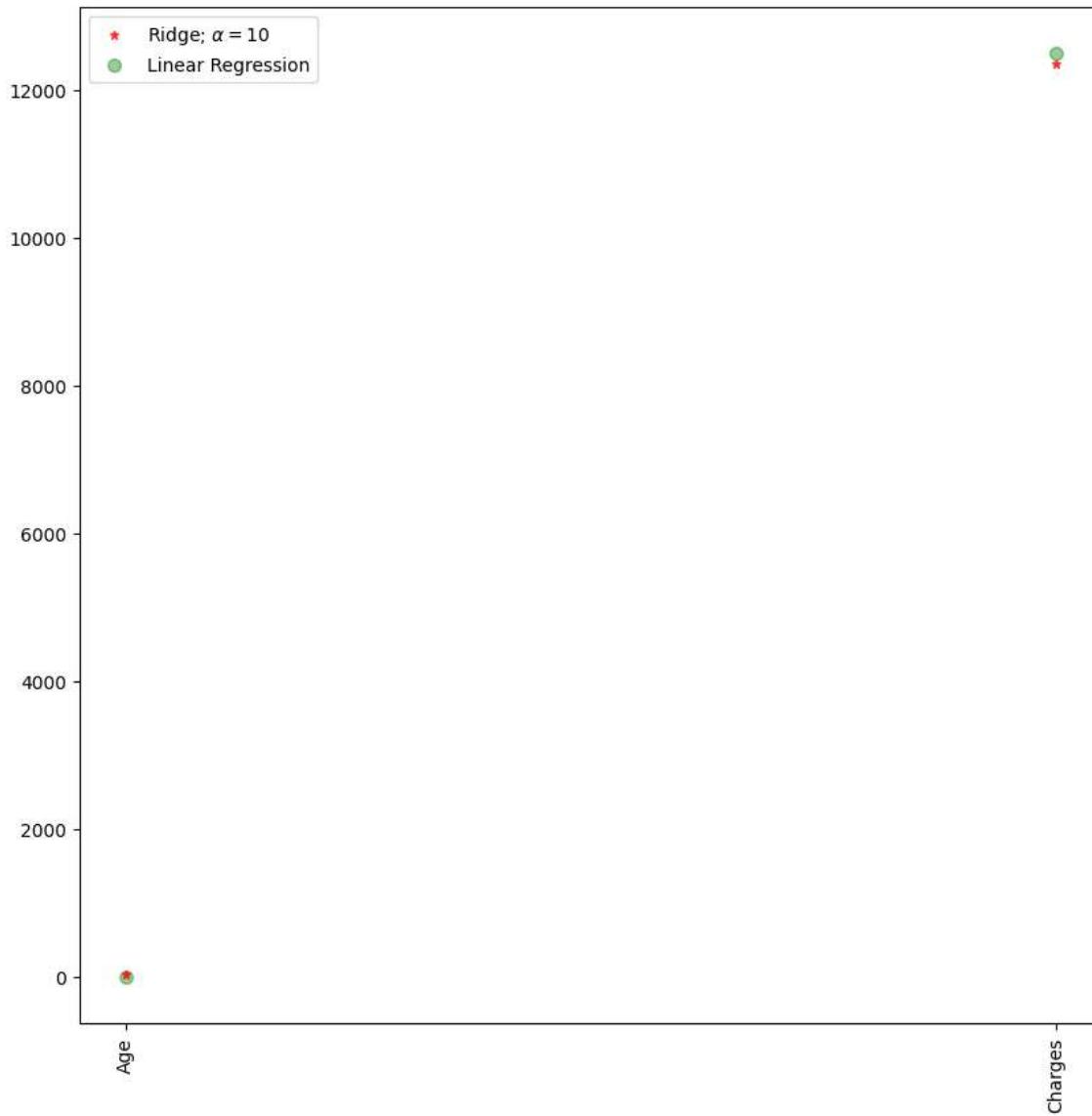
Ridge Model:

The train score for ridge model is 0.9998796730521575

The test score for ridge model is 0.9998841048873327

In [23]:

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



## Lasso Regression

In [24]:

```
1 #Using the linear CV model
2 from sklearn.linear_model import LassoCV
3 #Lasso Cross validation
4 lasso_cv = LassoCV(alphas = [0.0001, 0.001, 0.01, 0.1, 1, 10], random_state=0).fit(X_
5 #score
6 print(lasso_cv.score(X_train, y_train))
7 print(lasso_cv.score(X_test, y_test))
```

0.9999999999996456

0.99999999999531

In [25]:

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

Lasso Model:

The train score for ls model is 0.9999993594011218

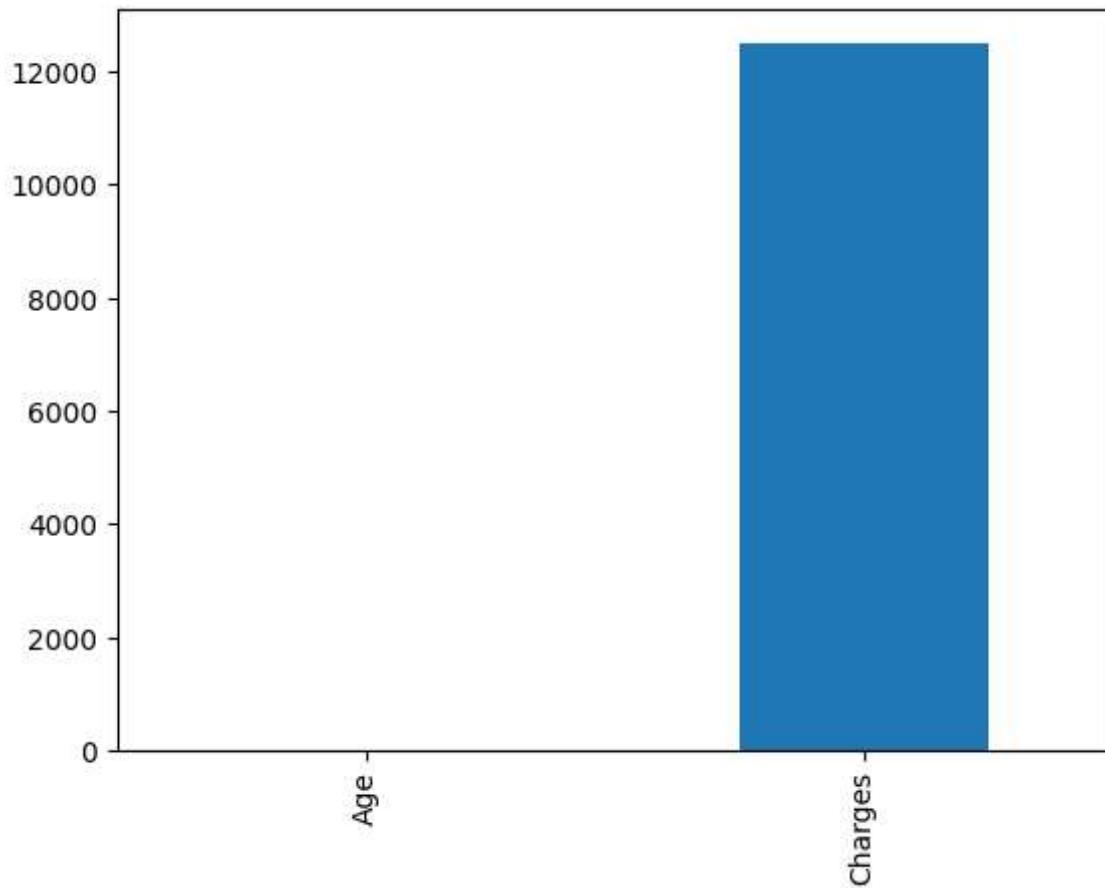
The test score for ls model is 0.9999993415784849

In [26]:

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

Out[26]:

&lt;Axes: &gt;

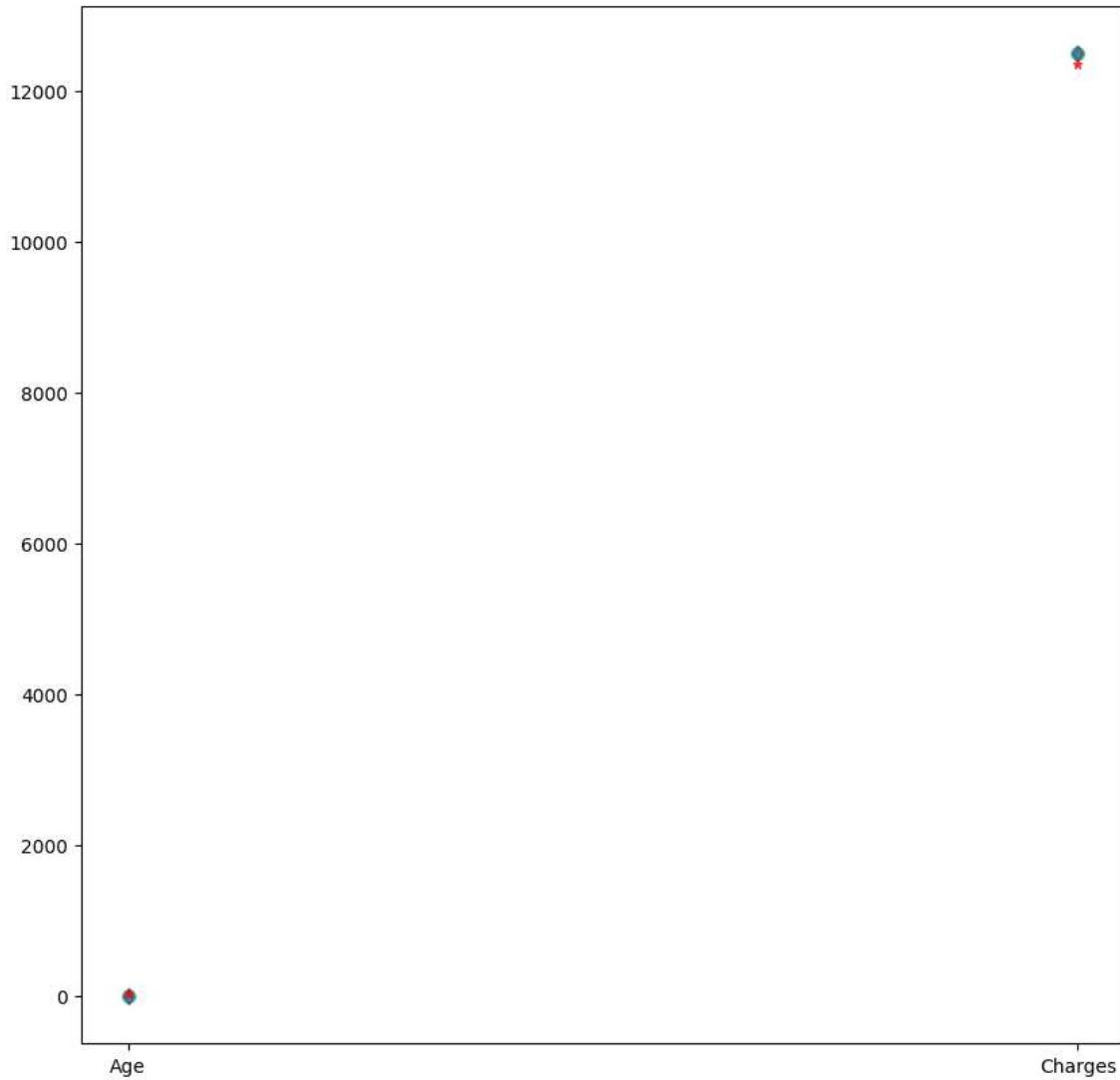


In [27]:

```
1 #plot size
2 plt.figure(figsize = (10, 10))
3 #add plot for ridge regression
4 plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,
5 #add plot for Lasso regression
6 plt.plot(lasso_cv.coef_,alpha=0.5,linestyle='none',marker='d',markersize=6,color='blue')
7 #add plot for Linear model
8 plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersize=7,color='red')
```

Out[27]:

[&lt;matplotlib.lines.Line2D at 0x16ec8fdd990&gt;]



## Elastic Regression

In [28]:

```
1 from sklearn.linear_model import ElasticNet
2 regr=ElasticNet()
3 regr.fit(X,y)
4 print(regr.coef_)
5 print(regr.intercept_)
```

```
[0.          0.99999999]
9.055665395862889e-05
```

In [29]:

```
1 y_pred_elastic=regr.predict(X_train)
2 mean_squared_error=np.mean((y_pred_elastic-y_train)**2)
3 print("Mean squared error on test set",mean_squared_error)
```

```
Mean squared error on test set 347174894.85874784
```

## Logistic Regression

In [30]:

```
1 import pandas as pd
2 import numpy as np
3 from sklearn import preprocessing
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import StandardScaler
7 import matplotlib.pyplot as plt
8 import seaborn as sns
9 sns.set(style="white")
10 sns.set(style="whitegrid",color_codes=True)
11 import warnings
12 warnings.simplefilter(action='ignore')
```

In [31]:

```
1 db=pd.read_csv(r"C:\Users\91955\Downloads\insurance.csv")
2 db
```

Out[31]:

	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 [32]:

```
1 db.head()
```

Out[32]:

	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 [33]:

```
1 db.tail()
```

Out[33]:

	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 [34]:

```
1 db.shape
```

Out[34]:

(1338, 7)

In [35]:

```
1 db.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 [36]:

```
1 db.describe()
```

Out[36]:

	age	bmi	children	charges
<b>count</b>	1338.000000	1338.000000	1338.000000	1338.000000
<b>mean</b>	39.207025	30.663397	1.094918	13270.422265
<b>std</b>	14.049960	6.098187	1.205493	12110.011237
<b>min</b>	18.000000	15.960000	0.000000	1121.873900
<b>25%</b>	27.000000	26.296250	0.000000	4740.287150
<b>50%</b>	39.000000	30.400000	1.000000	9382.033000
<b>75%</b>	51.000000	34.693750	2.000000	16639.912515
<b>max</b>	64.000000	53.130000	5.000000	63770.428010

In [37]:

```
1 db.isnull().sum()
```

Out[37]:

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

In [38]:

```

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

```

Out[38]:

	age	sex	bmi	children	smoker	region	charges
0	19	1	27.900	0	yes	southwest	16884.92400
1	18	2	33.770	1	no	southeast	1725.55230
2	28	2	33.000	3	no	southeast	4449.46200
3	33	2	22.705	0	no	northwest	21984.47061
4	32	2	28.880	0	no	northwest	3866.85520
...	...	...	...	...	...	...	...
1333	50	2	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 [39]:

```

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

```

Out[39]:

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

1338 rows × 7 columns

In [40]:

```

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

```

Out[40]:

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

1338 rows × 7 columns

In [41]:

```

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

```

In [42]:

```
1 features_matrix=db.iloc[:,0:6]
```

In [43]:

```
1 target_vector=db.iloc[:, -2]
```

In [44]:

```

1 print("The features matrix has %d rows and %d columns"%(features_matrix.shape))
2 print("The target matrix has %d rows and %d columns"%(np.array(target_vector).reshape(1338,1)))

```

The features matrix has 1338 rows and 6 columns

The target matrix has 1338 rows and 1 columns

In [45]:

```

1 x=np.array(db['sex']).reshape(-1,1)
2 y=np.array(db['smoker']).reshape(-1,1)

```

In [46]:

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

In [47]:

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

In [48]:

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

In [49]:

```
1 Observation=[[1,0,0.99539,-0.085889,0.8524299999999999,0.02306]]
```

In [50]:

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

The model predicted the observation to belong to class [2]

In [51]:

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

The algorithm was trained to predict one of the two classes:[1 2 3 4]

In [52]:

```
1 print("""The model says the probability of the observation we passed belonging to cl
2 print("""The model says the probability of the observation we passed belonging to cl
3 print("""The model says the probability of the observation we passed belonging to cl
4 print("""The model says the probability of the observation we passed belonging to cl
```

The model says the probability of the observation we passed belonging to class['0'] is 0.0006464896816406193

The model says the probability of the observation we passed belonging to class['1'] is 0.5652740537171539

The model says the probability of the observation we passed belonging to class['2'] is 0.43340043770162473

The model says the probability of the observation we passed belonging to class['3'] is 0.0006790188995806598

In [53]:

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

0.835820895522388

In [54]:

```
1 from sklearn import metrics
2 print('MAE:',metrics.mean_absolute_error(X_test,y_test))
3 print('MSE:',metrics.mean_squared_error(X_test,y_test))
4 print('RMSE:',np.sqrt(metrics.mean_absolute_error(X_test,y_test)))
```

MAE: 0.44776119402985076

MSE: 0.44776119402985076

RMSE: 0.6691496051182058

## Decision Tree

In [55]:

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

In [56]:

```
1 dt=pd.read_csv(r"C:\Users\91955\Downloads\insurance.csv")
2 dt
```

Out[56]:

	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

In [57]:

```
1 dt.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 [58]:

```
1 dt['smoker'].value_counts()
```

Out[58]:

```
smoker
no      1064
yes     274
Name: count, dtype: int64
```

In [59]:

```
1 dt['children'].value_counts()
```

Out[59]:

```
children
0    574
1    324
2    240
3    157
4     25
5     18
Name: count, dtype: int64
```

In [60]:

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

Out[60]:

	age	sex	bmi	children	smoker	region	charges
0	19	1	27.900	0	yes	southwest	16884.924000
1	18	2	33.770	1	no	southeast	1725.552300
2	28	2	33.000	3	no	southeast	4449.462000
3	33	2	22.705	0	no	northwest	21984.470610
4	32	2	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	2	29.830	2	no	northeast	6406.410700
9	60	1	25.840	0	no	northwest	28923.136920

In [61]:

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

Out[61]:

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

In [62]:

```
1 x=["sex", "children", "region"]
2 y=["yes", "no"]
3 all_inputs=dt[x]
4 all_classes=dt["smoker"]
```

In [63]:

```
1 (x_train,x_test,y_train,y_test)=train_test_split(all_inputs,all_classes,test_size=0.
```

In [64]:

```
1 clf=DecisionTreeClassifier(random_state=0)
```

In [65]:

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

Out[65]:

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

In [66]:

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

0.817910447761194

## Random Forest

In [67]:

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

In [68]:

```
1 ds=pd.read_csv(r"C:\Users\91955\Downloads\insurance.csv")
2 ds
```

47	28	female	34.770	0	no	northwest	3556.922300
48	60	female	24.530	0	no	southeast	12629.896700
49	36	male	35.200	1	yes	southeast	38709.176000
50	18	female	35.625	0	no	northeast	2211.130750
51	21	female	33.630	2	no	northwest	3579.828700
52	48	male	28.000	1	yes	southwest	23568.272000
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

In [69]:

```
1 ds.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 [70]:

```
1 x=ds.drop('smoker',axis=1)
2 y=ds['smoker']
```

In [71]:

```
1 ds['smoker'].value_counts()
```

Out[71]:

```
smoker
no      1064
yes     274
Name: count, dtype: int64
```

In [72]:

```

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

```

10	25	2	26.220	0	no	northeast	2721.320800
11	62	1	26.290	0	yes	southeast	27808.725100
12	23	2	34.400	0	no	southwest	1826.843000
13	56	1	39.820	0	no	southeast	11090.717800
14	27	2	42.130	0	yes	southeast	39611.757700
15	19	2	24.600	1	no	southwest	1837.237000
16	52	1	30.780	1	no	northeast	10797.336200
17	23	2	23.845	0	no	northeast	2395.171550
18	56	2	40.300	0	no	southwest	10602.385000
19	30	2	35.300	0	yes	southwest	36837.467000
20	60	1	36.005	0	no	northeast	13228.846950
21	30	1	32.400	1	no	southwest	4149.736000
22	18	2	34.100	0	no	southeast	1137.011000

In [73]:

```

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

```

Out[73]:

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

In [74]:

```
1 convert={"smoker": {"yes": 1, "no": 2}}
```

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

```
3 ds
```

Out[74]:

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

In [75]:

```
1 x=ds.drop('smoker',axis=1)
```

```
2 y=ds['smoker']
```

In [76]:

```
1 x
```

Out[76]:

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

In [77]:

```
1 y
```

Out[77]:

```
0      1
1      2
2      2
3      2
4      2
5      2
6      2
7      2
8      2
9      2
10     2
11     1
12     2
13     2
14     1
15     2
16     2
17     2
```

In [78]:

```
1 from sklearn.model_selection import train_test_split
2 x_train,x_test,y_train,y_test=train_test_split(x,y,train_size=0.7,random_state=42)
3 x_train.shape,x_test.shape
```

Out[78]:

```
((936, 6), (402, 6))
```

In [79]:

```
1 from sklearn.ensemble import RandomForestClassifier
2 rfc=RandomForestClassifier()
3
```

In [80]:

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

Out[80]:

```
▼ RandomForestClassifier
  RandomForestClassifier()
```

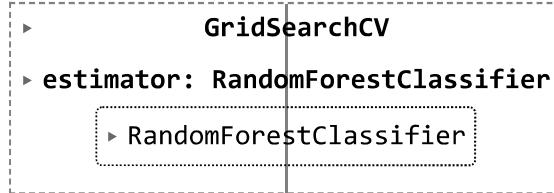
In [81]:

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

In [82]:

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

Out[82]:



In [83]:

```
1 grid_search.best_score_
```

Out[83]:

0.9508547008547008

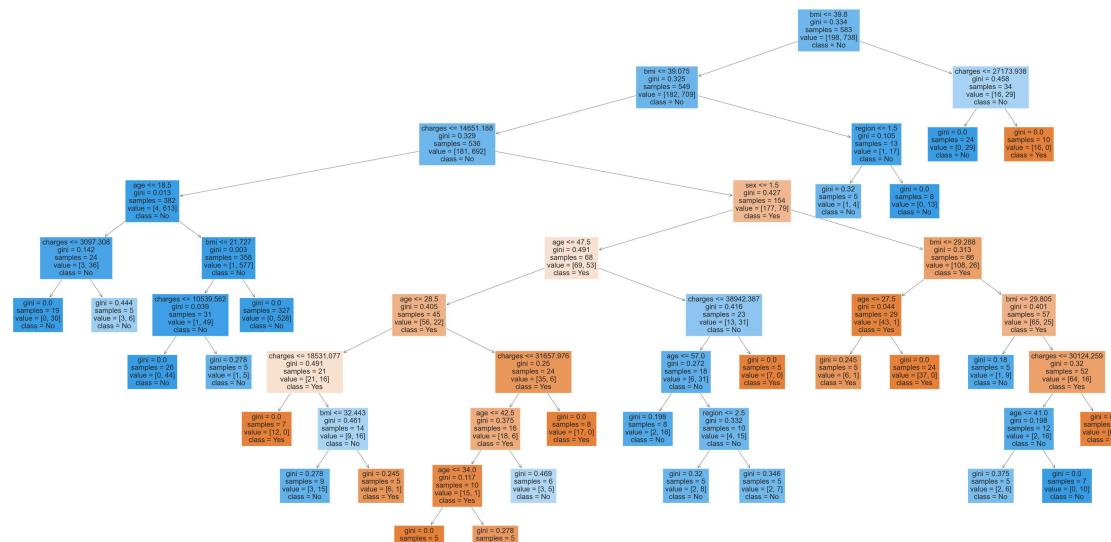
In [84]:

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

RandomForestClassifier(max\_depth=20, min\_samples\_leaf=5, n\_estimators=25)

In [85]:

```
1 from sklearn.tree import plot_tree
2 plt.figure(figsize=(80,40))
3 plot_tree(rf_best.estimators_[4],feature_names=x.columns,class_names=["Yes","No"],fi
```



In [86]:

```
1 from sklearn.tree import plot_tree
2 plt.figure(figsize=(80,40))
3 plot_tree(rf_best.estimators_[7], feature_names=x.columns, class_names=["Yes", "No"], fi
= 62\nvalue = [68, 31]\nnclass = Yes'),
Text(0.6666666666666666, 0.4375, 'age <= 26.0\nngini = 0.26\nnsamples =
33\nvalue = [44, 8]\nnclass = Yes'),
Text(0.6296296296296297, 0.3125, 'gini = 0.444\nnsamples = 5\nvalue =
[2, 4]\nnclass = No'),
Text(0.7037037037037037, 0.3125, 'age <= 49.5\nngini = 0.159\nnsamples =
28\nvalue = [42, 4]\nnclass = Yes'),
Text(0.6666666666666666, 0.1875, 'age <= 32.0\nngini = 0.056\nnsamples =
20\nvalue = [34, 1]\nnclass = Yes'),
Text(0.6296296296296297, 0.0625, 'gini = 0.245\nnsamples = 5\nvalue =
[6, 1]\nnclass = Yes'),
Text(0.7037037037037037, 0.0625, 'gini = 0.0\nnsamples = 15\nvalue = [2
8, 0]\nnclass = Yes'),
Text(0.7407407407407407, 0.1875, 'gini = 0.397\nnsamples = 8\nvalue =
[8, 3]\nnclass = Yes'),
Text(0.8888888888888888, 0.4375, 'bmi <= 39.225\nngini = 0.5\nnsamples =
29\nvalue = [24, 23]\nnclass = Yes'),
Text(0.8518518518518519, 0.3125, 'charges <= 33094.077\nngini = 0.488\nn
samples = 24\nvalue = [16, 22]\nnclass = No'),
Text(0.8148148148148148. 0.1875. 'gini = 0.0\nnsamples = 12\nvalue =
```

**Conclusion:**For the given dataset,we have performed linear regression,ridge regression,lasso regression,elastic regression,logistic regression,decision tree,random forest classification.Among all the models,we observed that ,in random forest the accuracy is 0.94, and in the lasso regression we observed,the accuracy is 0.99 where as lasso regression got the highest accuracy than the random forest.So, the best model that suits for the given dataset is linear regression and lasso regression.

In [ ]:

1