

Medical Cost Prediction

The aim of this analysis is to predict the medical expense based on the patients' information. The dataset used for this analysis is Insurance dataset from [Kaggle](#). The dataset contains 1338 observations and 7 variables. The variables are as follows:

Variable	Description
age	age of primary beneficiary
bmi	body mass index
children	number of children covered by health insurance
smoker	smoking
region	the beneficiary's residential area in the US
charges	individual medical costs billed by health insurance

```
In [ ]: #importing the Libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [ ]: df = pd.read_csv('insurance.csv')
df.head()
```

```
Out[ ]:
```

	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

Data Preprocessing

```
In [ ]: #number of rows and columns
df.shape
```

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

```
In [ ]: #checking for missing values
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 [ ]: #checking discriptive statistics
df.describe()
```

```
Out[ ]:
```

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

```
In [ ]: #value counts for categorical variables
print(df.sex.value_counts(),'\n',df.smoker.value_counts(),'\n',df.region.value_c
```

```
sex
male      676
female    662
Name: count, dtype: int64
smoker
no        1064
yes        274
Name: count, dtype: int64
region
southeast  364
southwest  325
northwest  325
northeast  324
Name: count, dtype: int64
```

Replacing the categorical variables with numerical values.

- sex : 1 - male, 0 - female
- smoker : 1 - yes, 0 - no
- region : 0 - northeast, 1 - northwest, 2 - southeast, 3 - southwest

```
In [ ]: #changing categorical variables to numerical
df['sex'] = df['sex'].map({'male':1,'female':0})
df['smoker'] = df['smoker'].map({'yes':1,'no':0})
df['region'] = df['region'].map({'southwest':0,'southeast':1,'northwest':2,'north':3})
```

```
In [ ]: df.head(10)
```

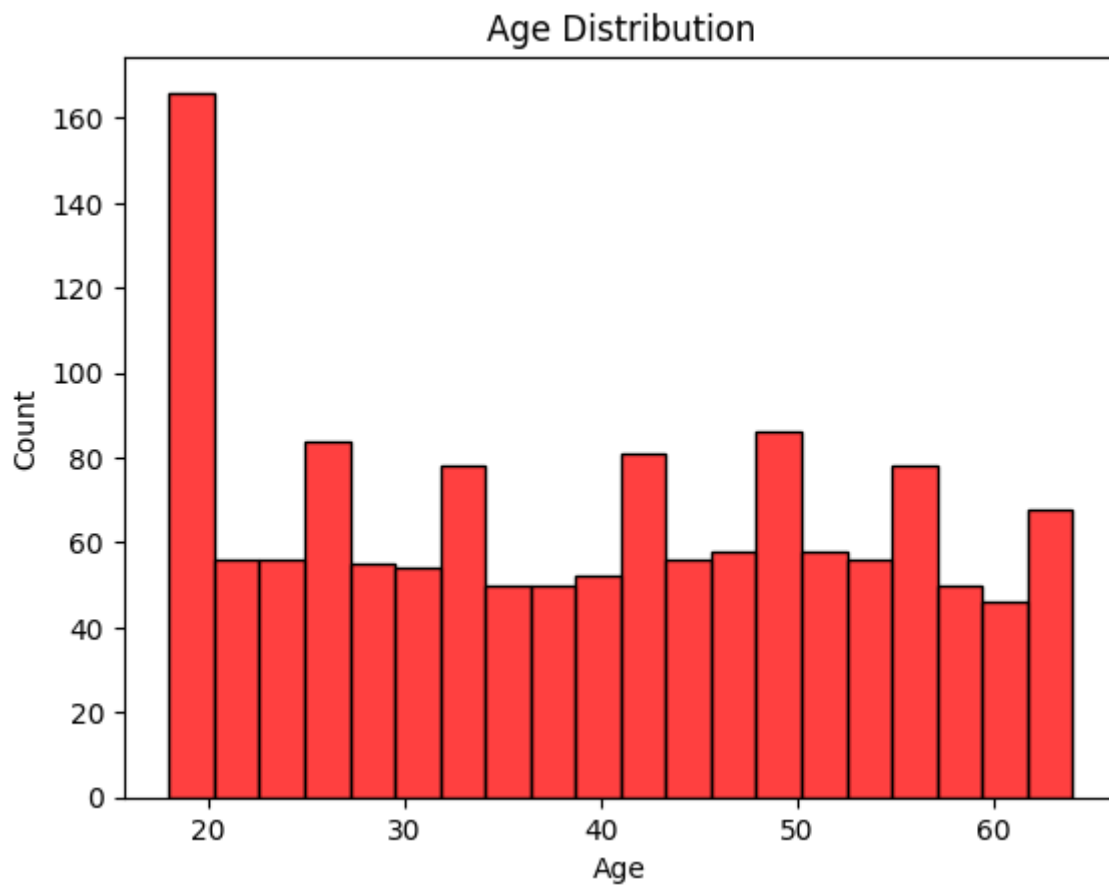
```
Out[ ]:
```

	age	sex	bmi	children	smoker	region	charges
0	19	0	27.900	0	1	0	16884.92400
1	18	1	33.770	1	0	1	1725.55230
2	28	1	33.000	3	0	1	4449.46200
3	33	1	22.705	0	0	2	21984.47061
4	32	1	28.880	0	0	2	3866.85520
5	31	0	25.740	0	0	1	3756.62160
6	46	0	33.440	1	0	1	8240.58960
7	37	0	27.740	3	0	2	7281.50560
8	37	1	29.830	2	0	3	6406.41070
9	60	0	25.840	0	0	2	28923.13692

Exploratory Data Analysis

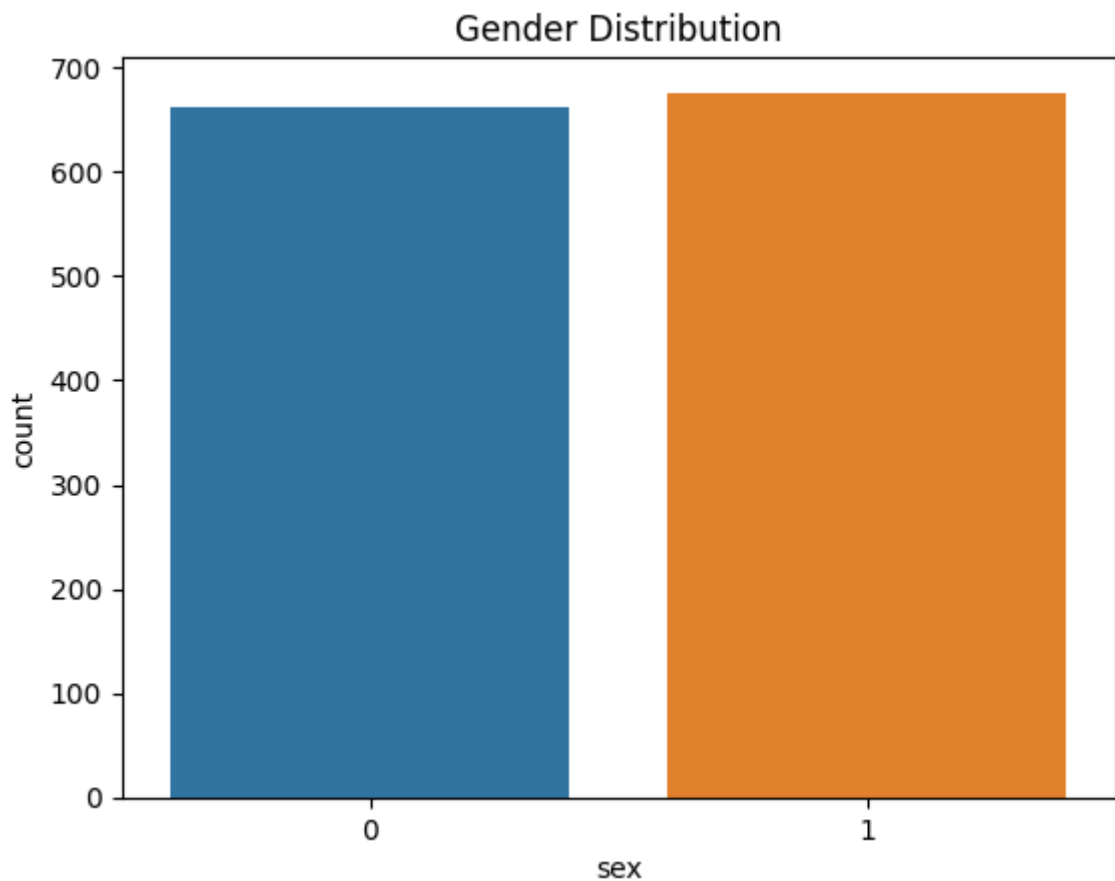
Visualization of the data is a good way to understand the data. In this section, I will plot the distribution of each variable to get an overview about their counts and distributions.

```
In [ ]: #age distribution
sns.histplot(df.age,bins=20, kde=False,color='red')
plt.title('Age Distribution')
plt.xlabel('Age')
plt.ylabel('Count')
plt.show()
```



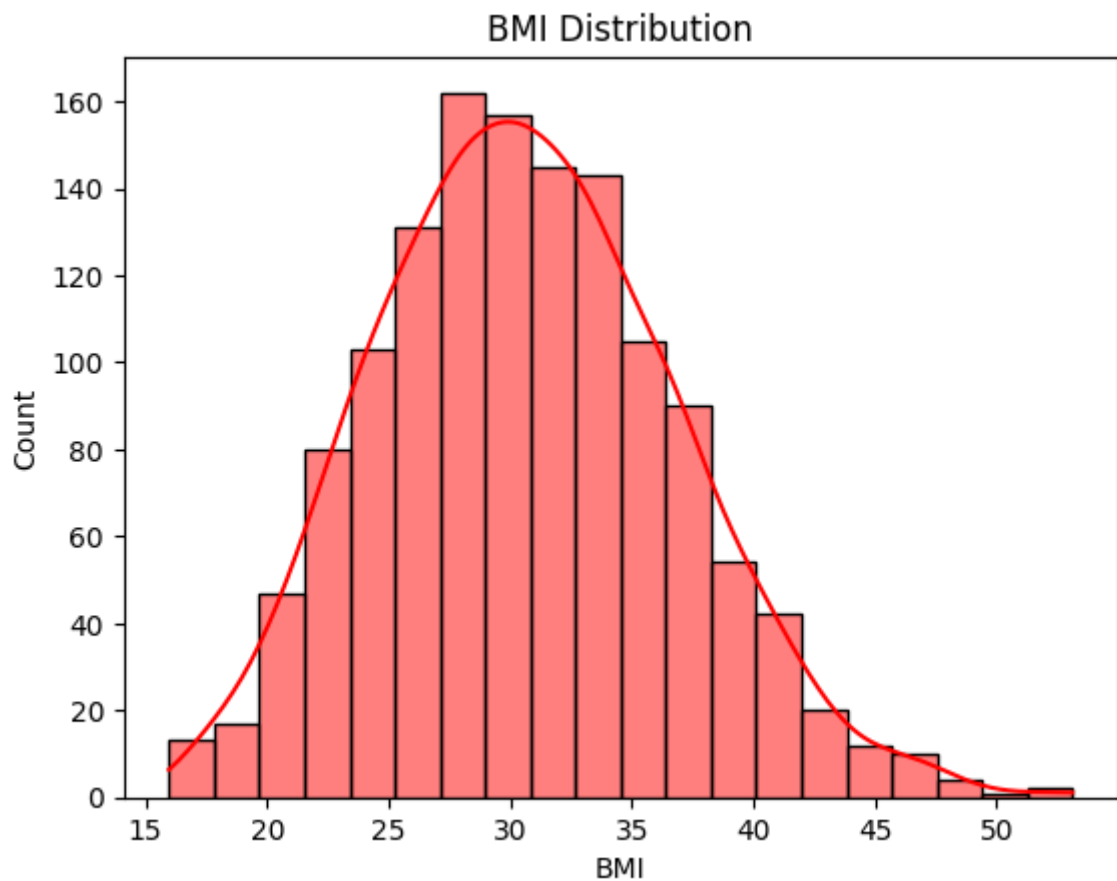
```
In [ ]: #gender plot
sns.countplot(x = 'sex', data = df)
plt.title('Gender Distribution')
```

```
Out[ ]: Text(0.5, 1.0, 'Gender Distribution')
```



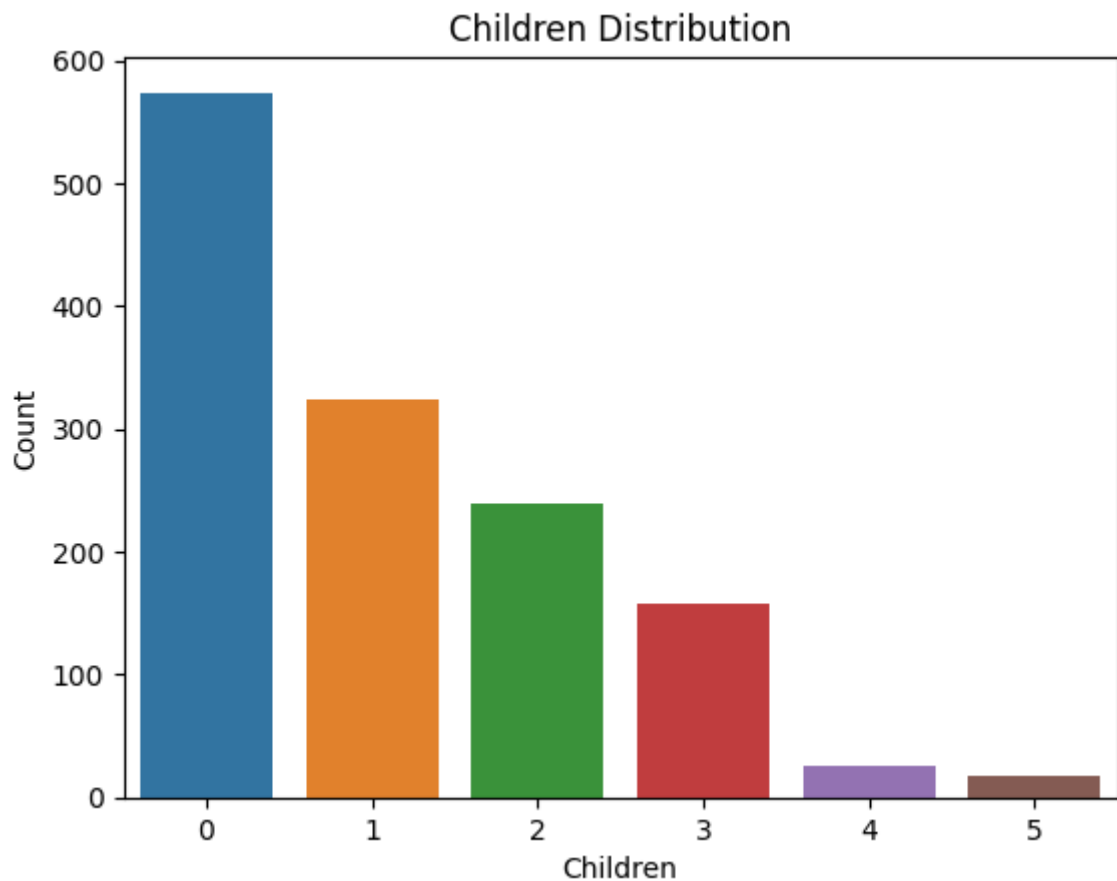
It is clear that number of males and females are almost equal in the dataset.

```
In [ ]: #bmi distribution
sns.histplot(df.bmi,bins=20, kde=True,color='red')
plt.title('BMI Distribution')
plt.xlabel('BMI')
plt.ylabel('Count')
plt.show()
```



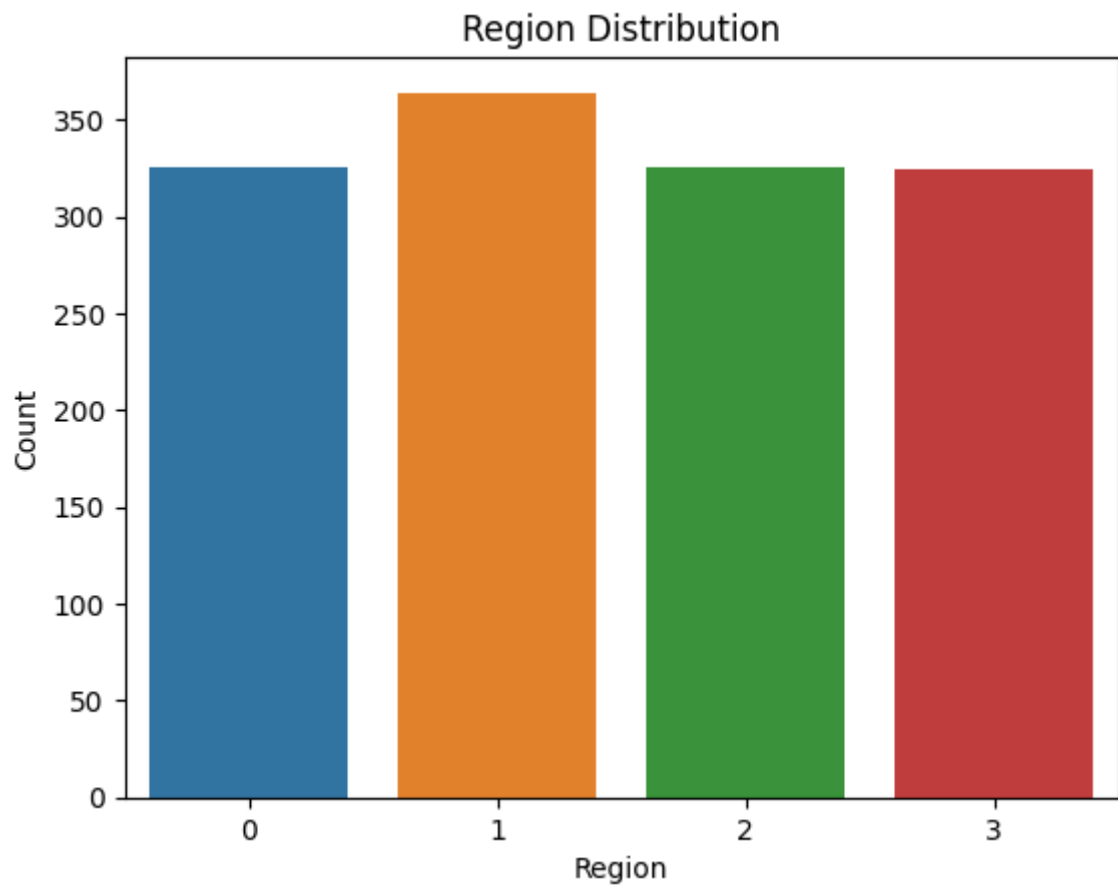
The majority of the patients have BMI between 25 and 40 which is considered as overweight and could be a major factor in increasing the medical cost.

```
In [ ]: #child count distribution
sns.countplot(x = 'children', data = df)
plt.title('Children Distribution')
plt.xlabel('Children')
plt.ylabel('Count')
plt.show()
```



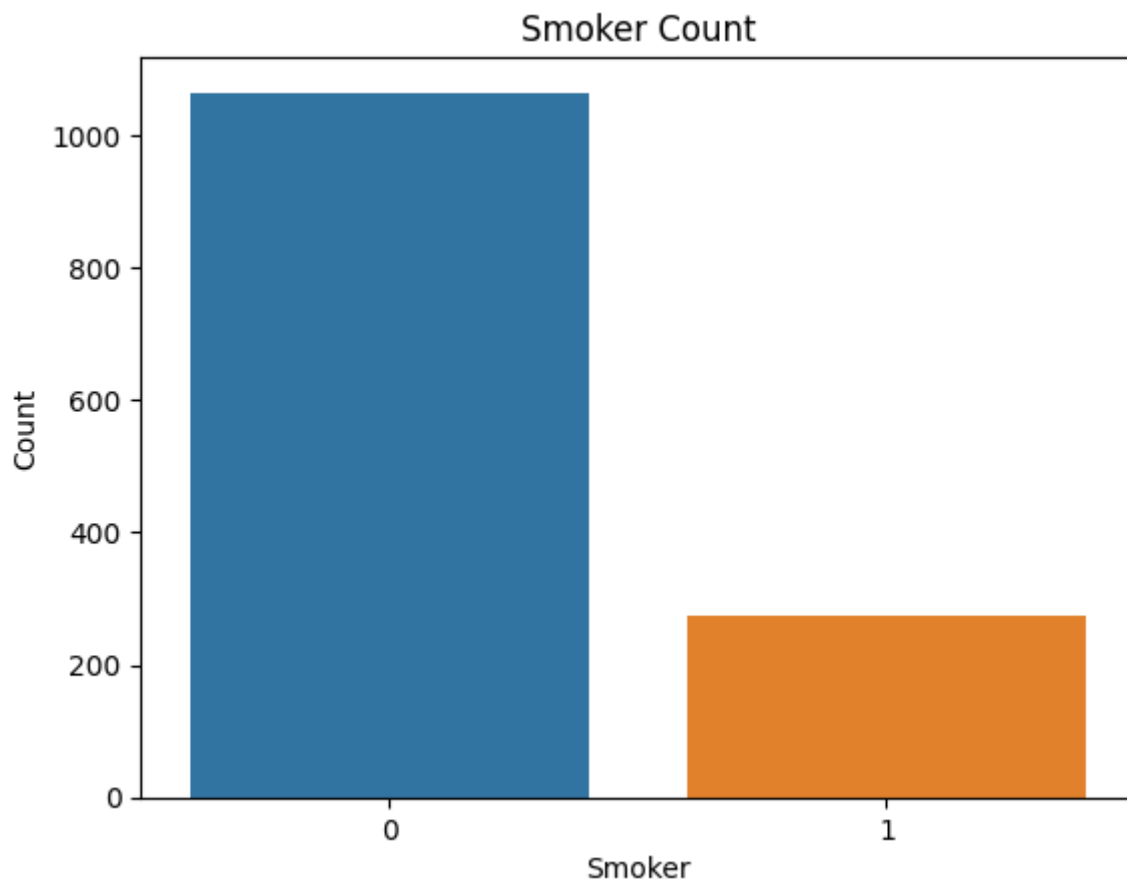
The graph clearly shows that most of the patients have no children and very few patients have more than 3 children.

```
In [ ]: #regionwise plot
sns.countplot(x = 'region', data = df)
plt.title('Region Distribution')
plt.xlabel('Region')
plt.ylabel('Count')
plt.show()
```



The count of patient from northwest is slightly higher than the other regions, but the number of patients from other regions are almost equal.

```
In [ ]: #count of smokers
sns.countplot(x = 'smoker', data = df)
plt.title('Smoker Count')
plt.xlabel('Smoker')
plt.ylabel('Count')
plt.show()
```

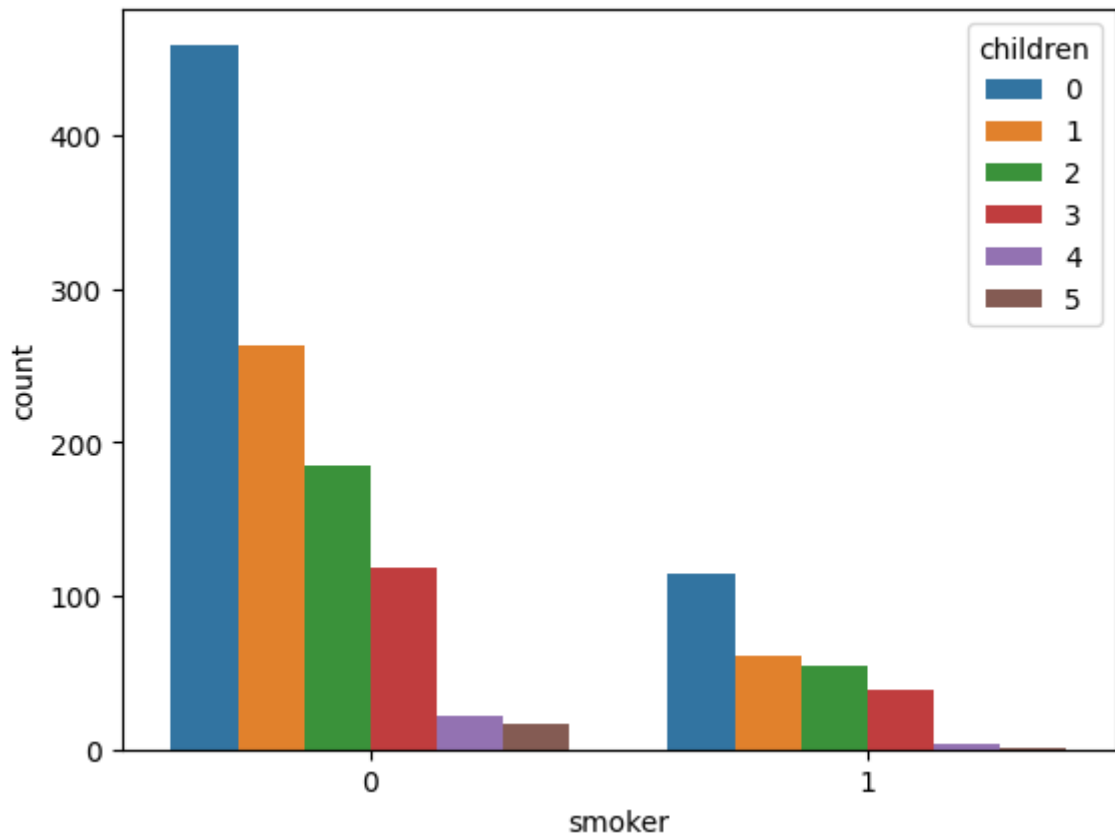



smokers are very few in the dataset. Nearly 80% of the patients are non-smokers.

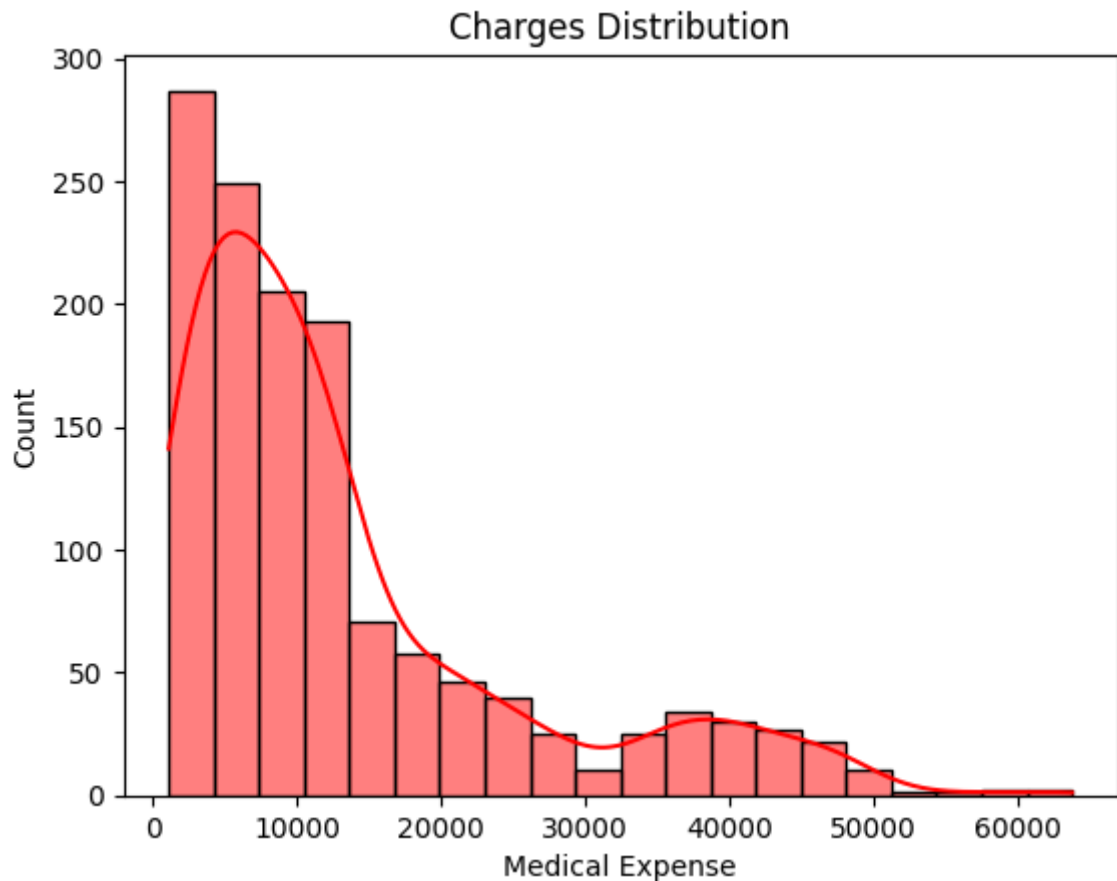
Smoker count with respect to the children count.

```
In [ ]: sns.countplot(x = df.smoker, hue = df.children)
```

```
Out[ ]: <Axes: xlabel='smoker', ylabel='count'>
```



```
In [ ]: #charges distribution
sns.histplot(df.charges,bins=20, kde=True,color='red')
plt.title('Charges Distribution')
plt.xlabel('Medical Expense')
plt.ylabel('Count')
plt.show()
```



Most of the medical expenses are below 20000, with negligible number of patients having medical expenses above 50000.

From all the above plots, we have a clear understanding about the count of patients under each category of the variables. Now I will look into the coorelation between the variables.

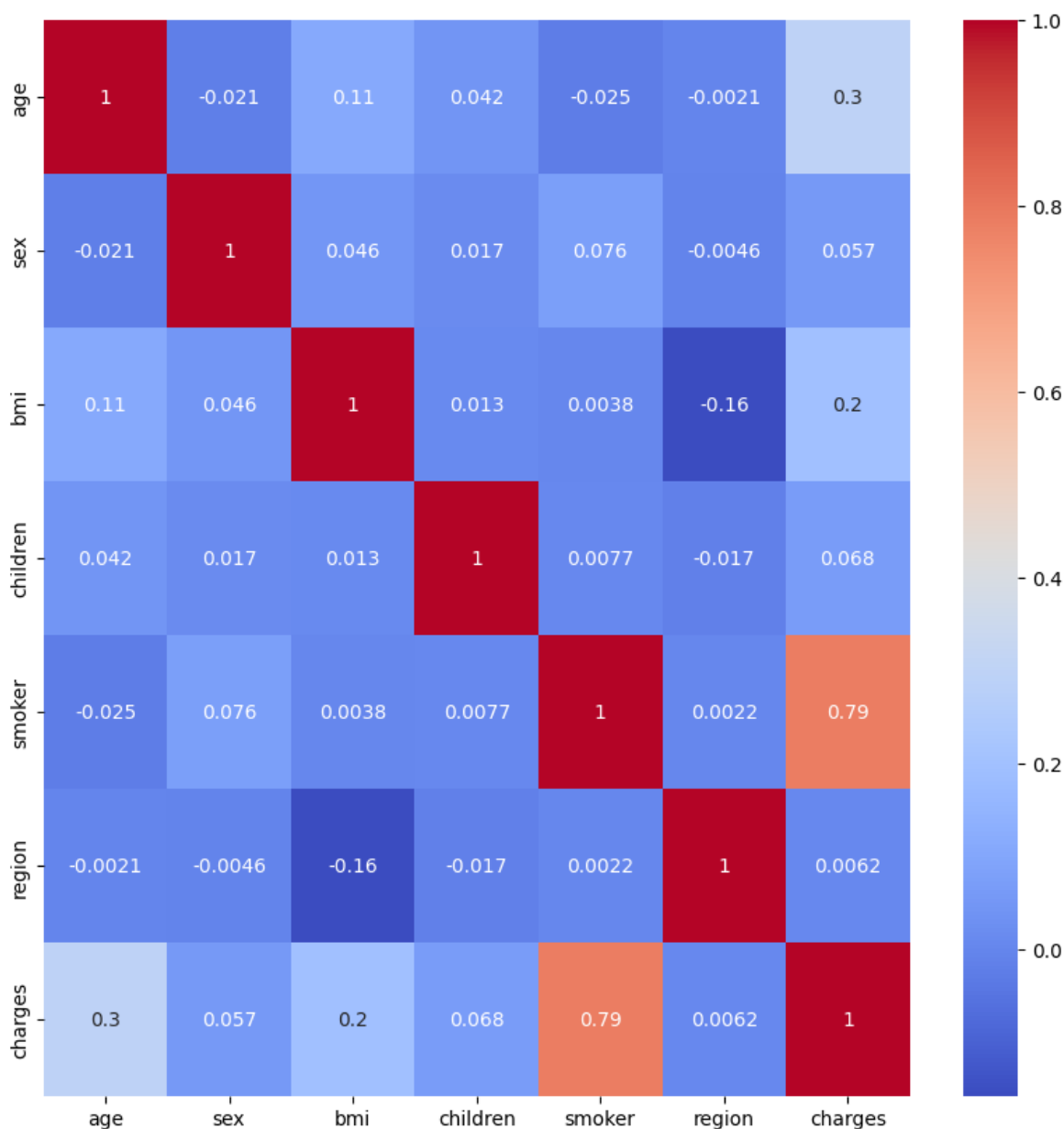
Coorelation

```
In [ ]: #coorelation matrix
df.corr()
```

```
Out[ ]:
```

	age	sex	bmi	children	smoker	region	charges
age	1.000000	-0.020856	0.109272	0.042469	-0.025019	-0.002127	0.299008
sex	-0.020856	1.000000	0.046371	0.017163	0.076185	-0.004588	0.057292
bmi	0.109272	0.046371	1.000000	0.012759	0.003750	-0.157566	0.198341
children	0.042469	0.017163	0.012759	1.000000	0.007673	-0.016569	0.067998
smoker	-0.025019	0.076185	0.003750	0.007673	1.000000	0.002181	0.787251
region	-0.002127	-0.004588	-0.157566	-0.016569	0.002181	1.000000	0.006208
charges	0.299008	0.057292	0.198341	0.067998	0.787251	0.006208	1.000000

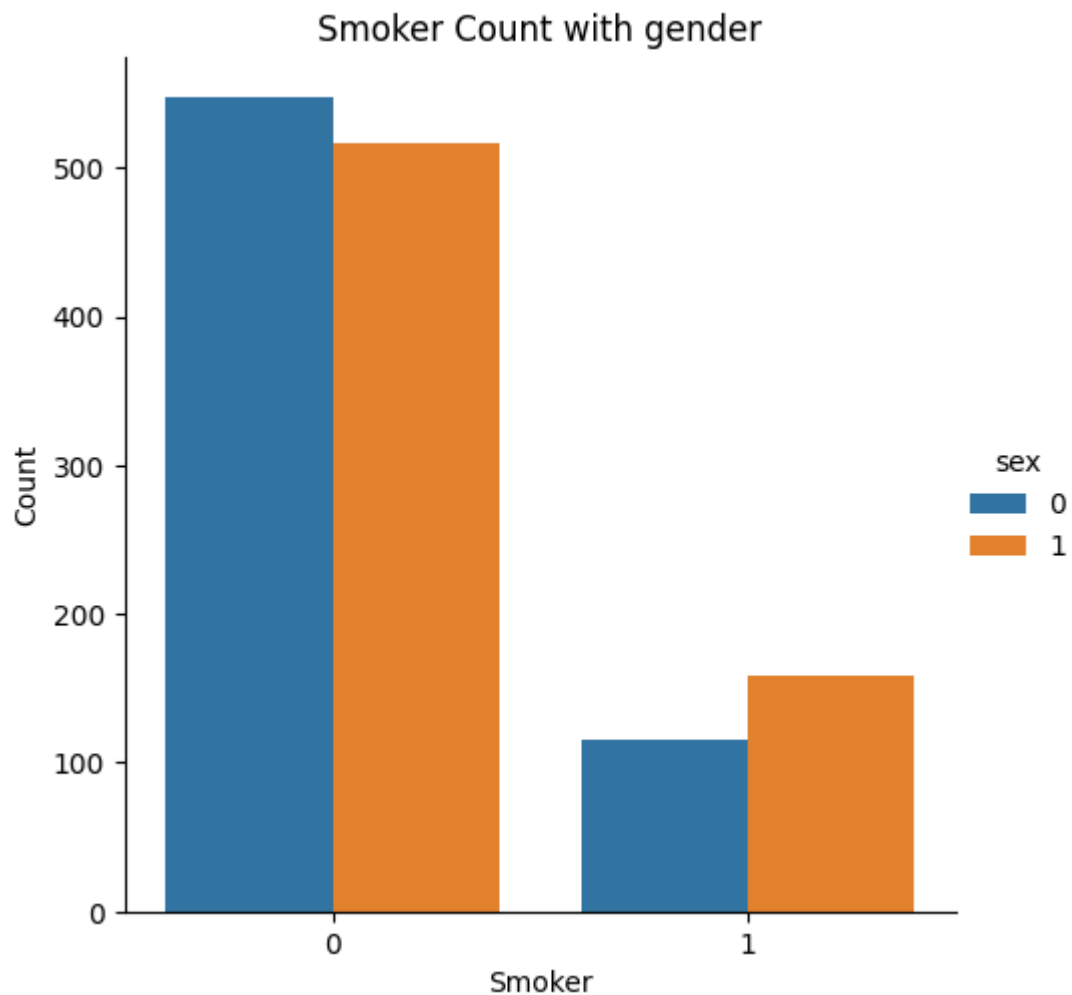
```
In [ ]: #plotting the coorelation heatmap
plt.figure(figsize=(10,10))
sns.heatmap(df.corr(),annot=True,cmap='coolwarm')
plt.show()
```



The variable smoker shows a significant coorelation with the medical expenses. Now I will explore more into patients' smoking habits and their relationa with other factors.

Plotting the smoker count with patient's gender

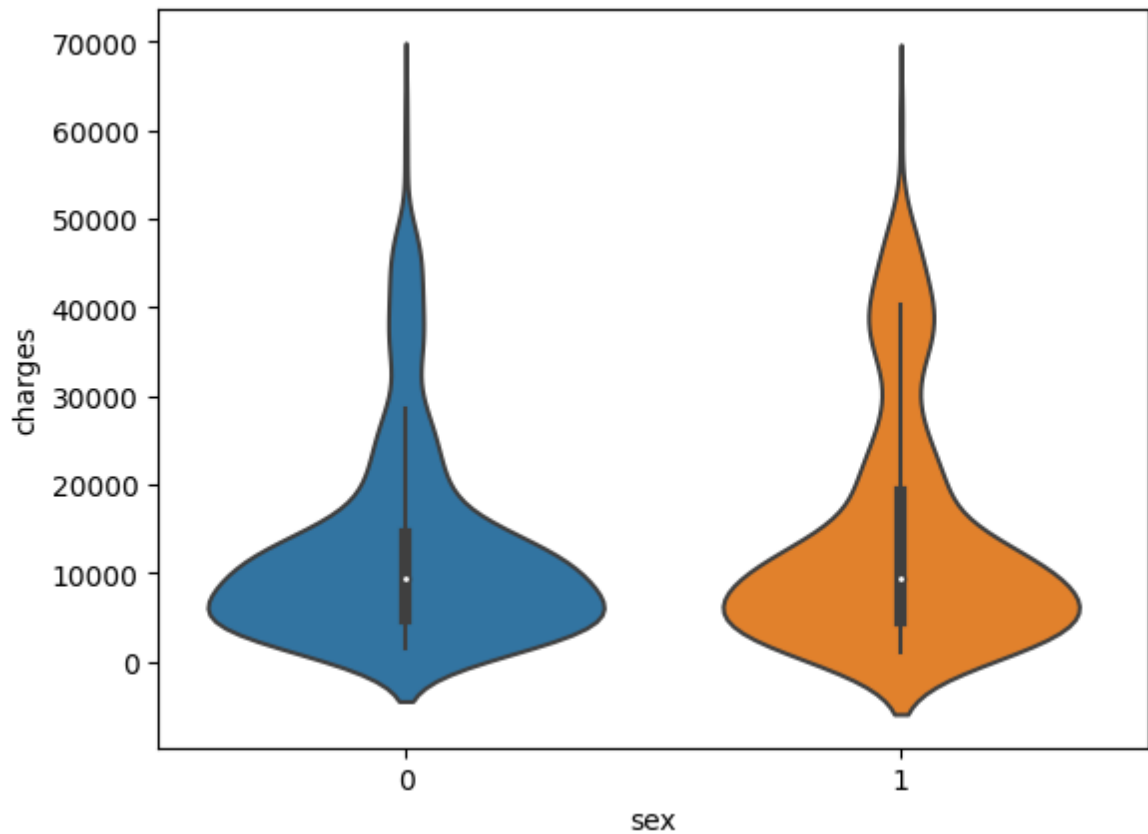
```
In [ ]: sns.catplot(x="smoker", kind="count",hue = 'sex', data=df)
plt.title('Smoker Count with gender')
plt.xlabel('Smoker')
plt.ylabel('Count')
plt.show()
```



We can notice more male smokers than female smokers. So, I will assume that medical treatment expense for males would be more than females, given the impact of smoking on the medical expenses.

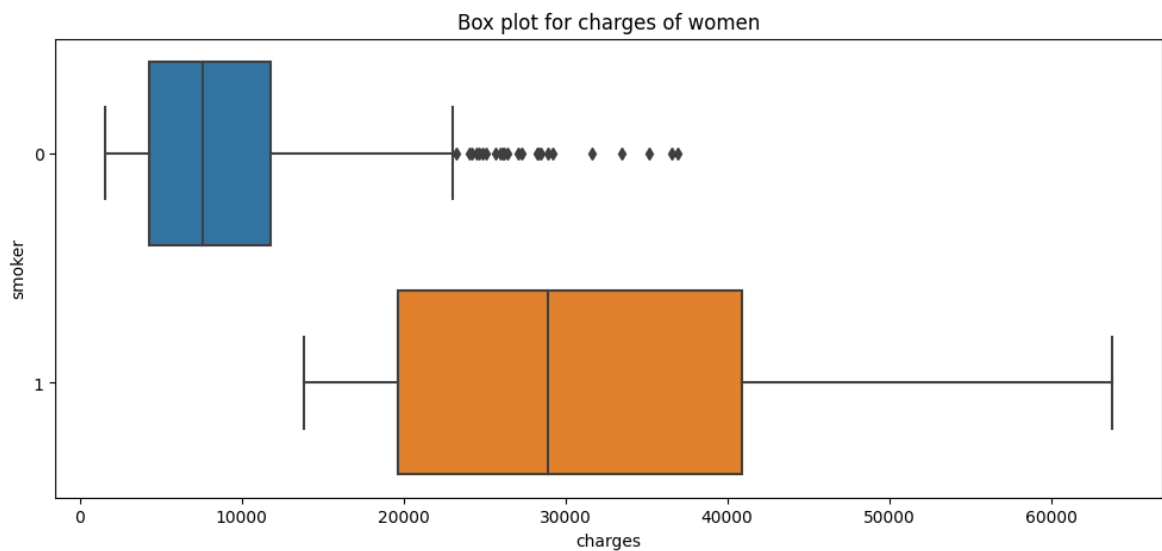
```
In [ ]: sns.violinplot(x = 'sex', y = 'charges', data = df)
```

```
Out[ ]: <Axes: xlabel='sex', ylabel='charges'>
```



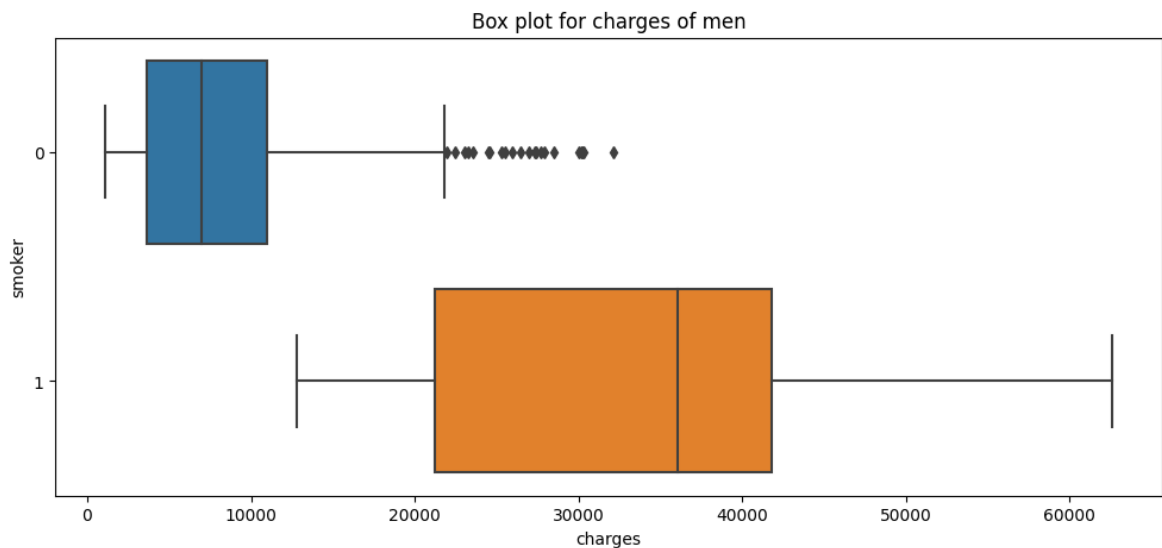
```
In [ ]: plt.figure(figsize=(12,5))
plt.title("Box plot for charges of women")
sns.boxplot(y="smoker", x="charges", data = df[(df.sex == 0)] , orient="h")
```

```
Out[ ]: <Axes: title={'center': 'Box plot for charges of women'}, xlabel='charges', yla
bel='smoker'>
```



```
In [ ]: plt.figure(figsize=(12,5))
plt.title("Box plot for charges of men")
sns.boxplot(y="smoker", x="charges", data = df[(df.sex == 1)] , orient="h")
```

```
Out[ ]: <Axes: title={'center': 'Box plot for charges of men'}, xlabel='charges', yla
bel='smoker'>
```

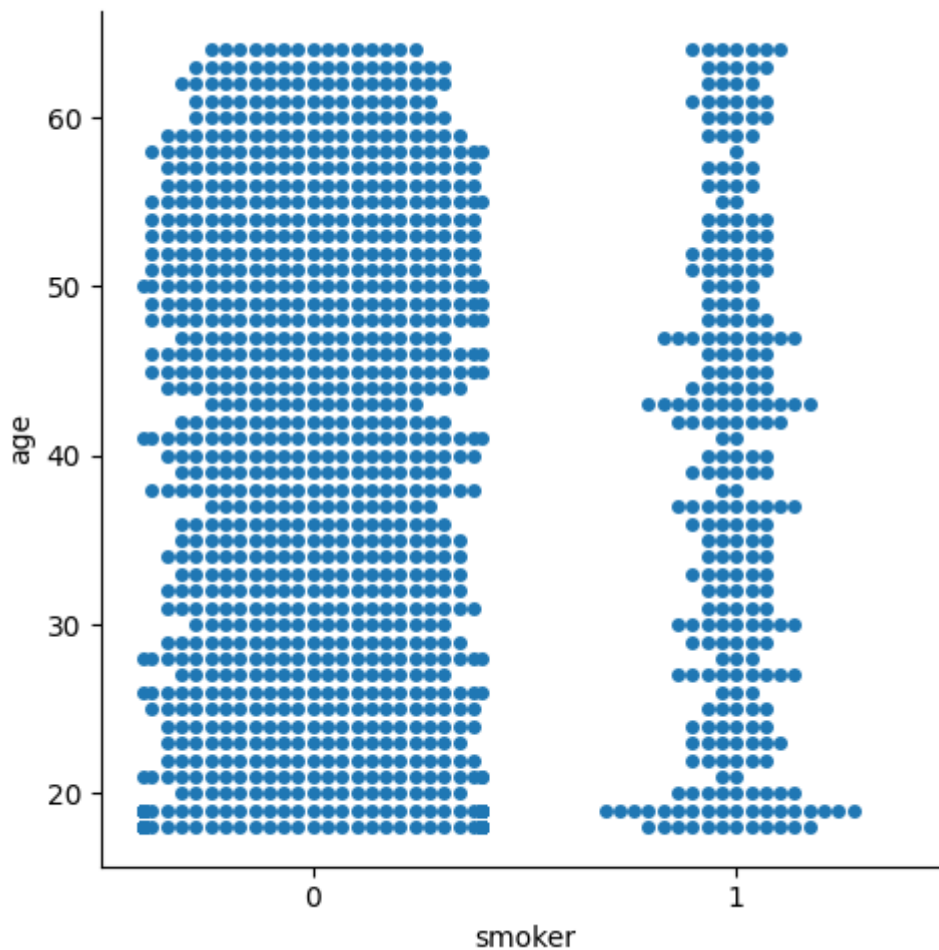


The assumption is true, that the medical expense of males is greater than that of females. In addition to that medical expense of smokers is greater than that of non-smokers.

Smokers and age distribution

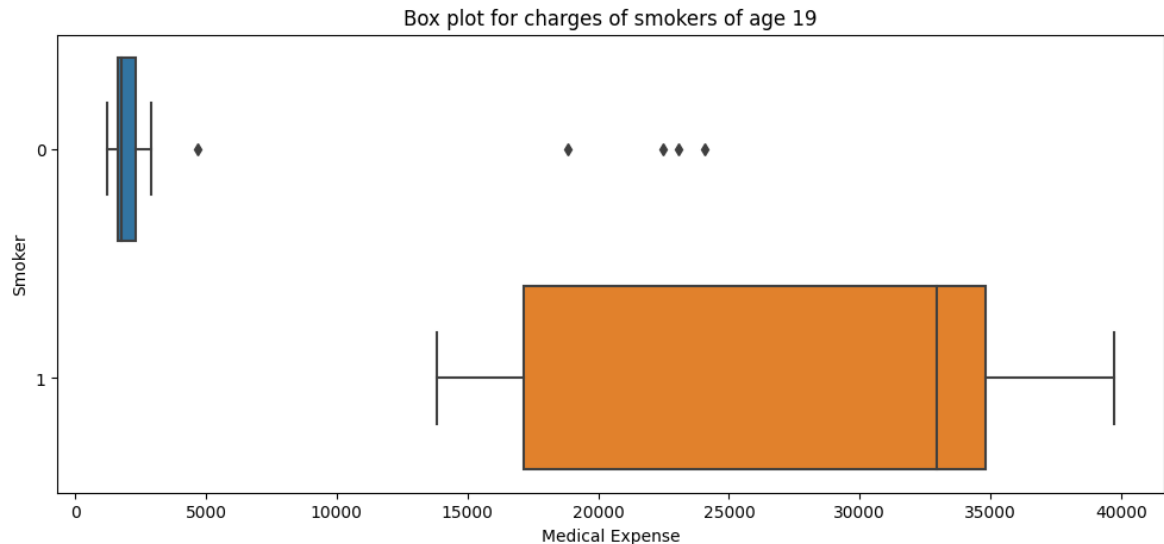
```
In [ ]: #smokers and age distribution
sns.catplot(x="smoker", y="age", kind="swarm", data=df)
```

```
Out[ ]: <seaborn.axisgrid.FacetGrid at 0x1443d53d690>
```



From the graph, we can see that there significant number of smokers of age 19. Now I will study the medical expense of smokers of age 19.

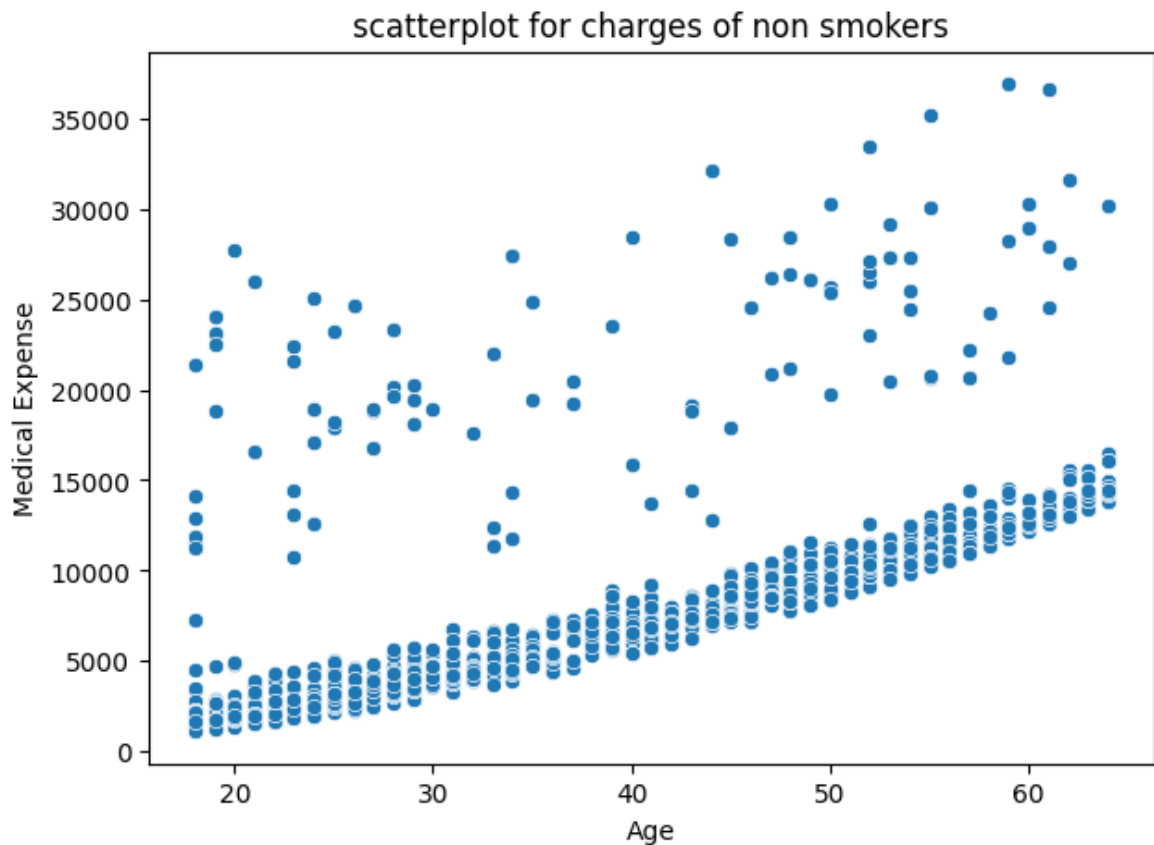
```
In [ ]: #smokers of age 19
plt.figure(figsize=(12,5))
plt.title("Box plot for charges of smokers of age 19")
sns.boxplot(y="smoker", x="charges", data = df[(df.age == 19)] , orient="h")
plt.xlabel('Medical Expense')
plt.ylabel('Smoker')
plt.show()
```



Surprisingly the medical expense of smokers of age 19 is very high in comparison to non smokers. In non smokers we can see some outliers, which may be due to illness or accidents.

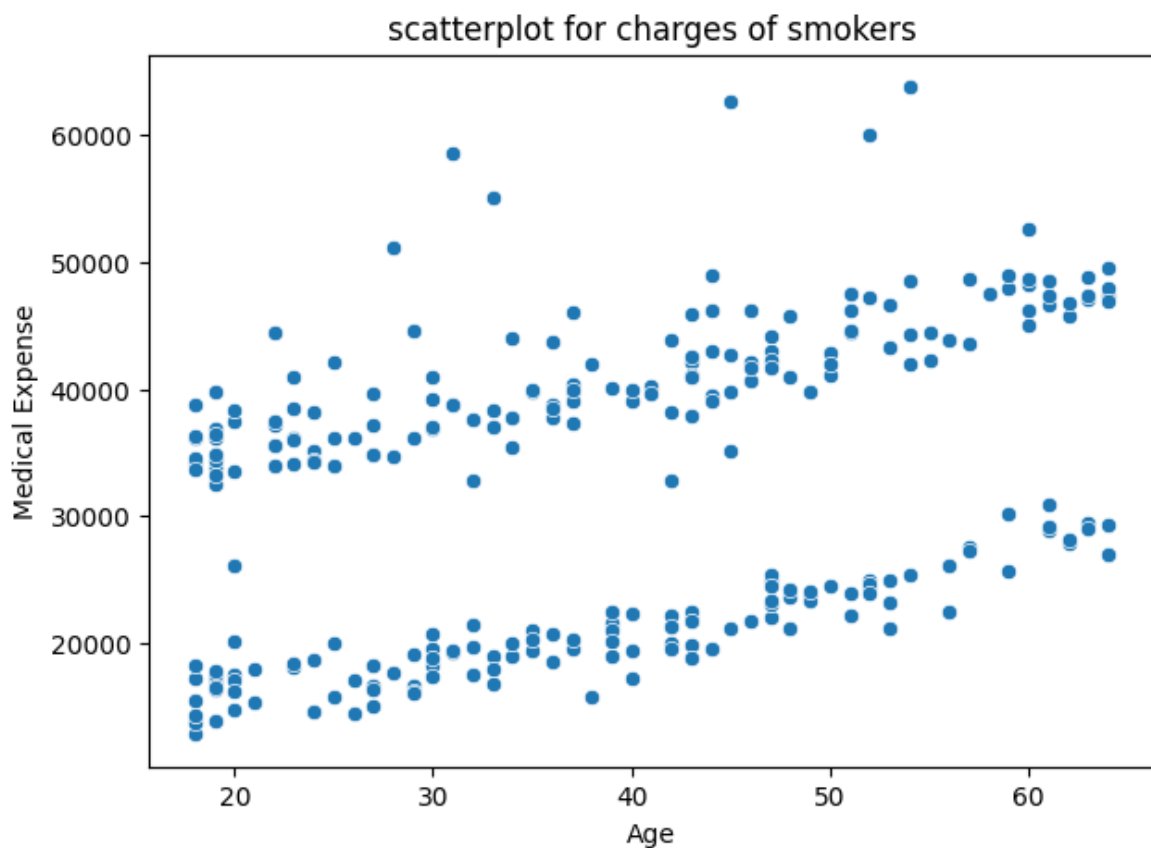
It is clear that the medical expense of smokers is higher than that of non-smokers. Now I will plot the charges distribution with respect to patients age of smokers and non-smokers.

```
In [ ]: #non smokers charge distribution
plt.figure(figsize=(7,5))
plt.title("scatterplot for charges of non smokers")
sns.scatterplot(x="age", y="charges", data = df[(df.smoker == 0)])
plt.xlabel('Age')
plt.ylabel('Medical Expense')
plt.show()
```

Majority of the points shows that medical expense increases with age which may be due to the fact that older people are more prone to illness. But there are some outliers which shows that there are other illness or accidents which may increase the medical expense.

```
In [ ]: #smokers charge distribution
plt.figure(figsize=(7,5))
plt.title("scatterplot for charges of smokers")
sns.scatterplot(x="age", y="charges", data = df[(df.smoker == 1)])
plt.xlabel('Age')
plt.ylabel('Medical Expense')
plt.show()
```

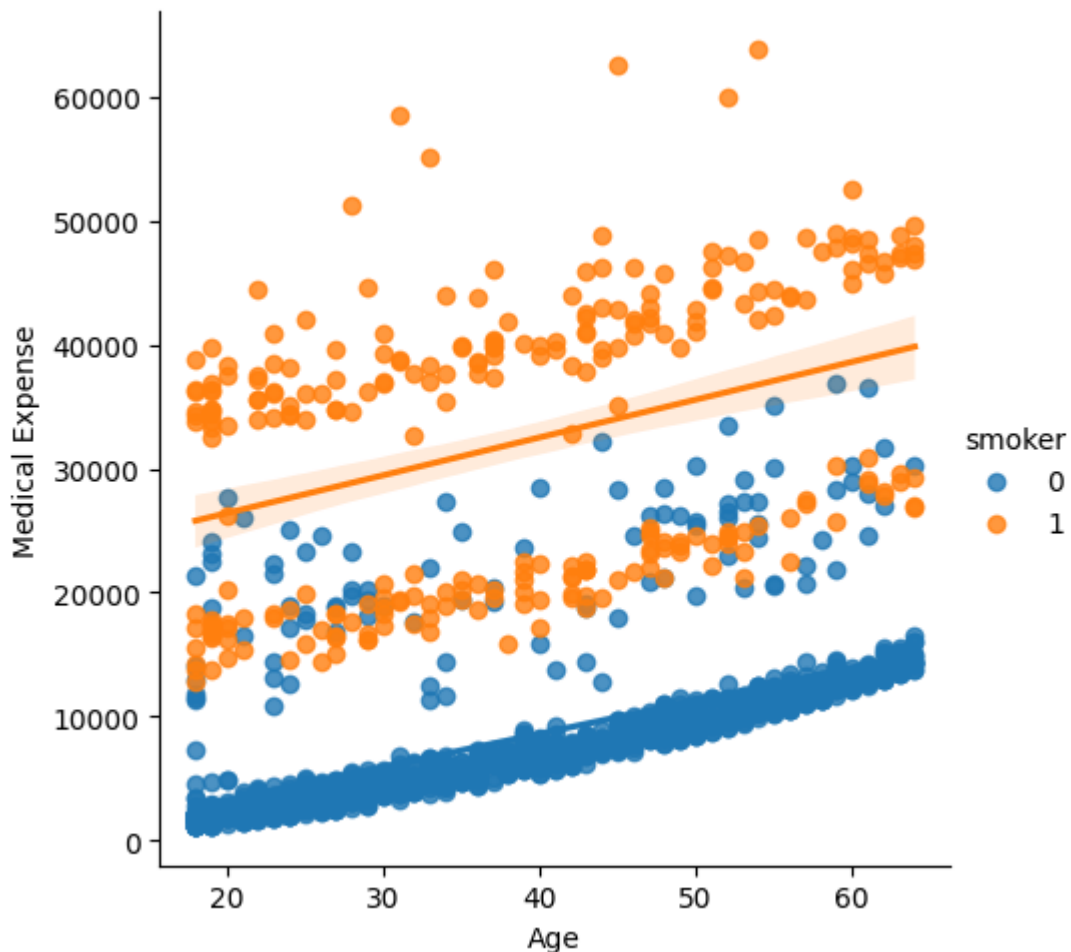


Here we see peculiarity in the graph. In the graph there are two segments, one with high medical expense which may be due to smoking related illness and the other with low medical expense which may be due age related illness.

Now, in order to get a more clear picture, I will combine these two graphs.

```
In [ ]: #age charges distribution

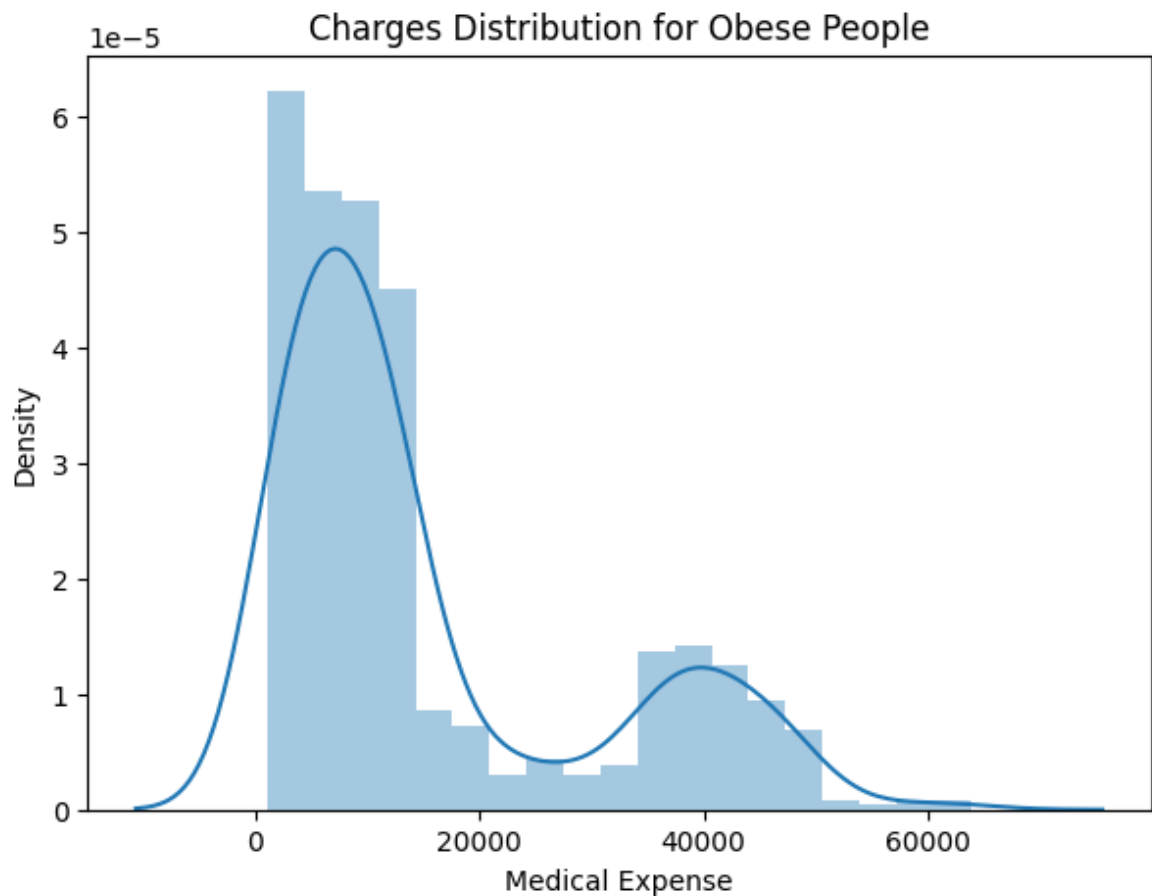
sns.lmplot(x="age", y="charges", data = df, hue = 'smoker')
plt.xlabel('Age')
plt.ylabel('Medical Expense')
plt.show()
```



Now, we clearly understand the variation in charges with respect to age and smoking habits. The medical expense of smokers is higher than that of non-smokers. In non-smokers, the cost of treatment increases with age which is obvious. But in smokers, the cost of treatment is high even for younger patients, which means the smoking patients are spending upon their smoking related illness as well as age related illness.

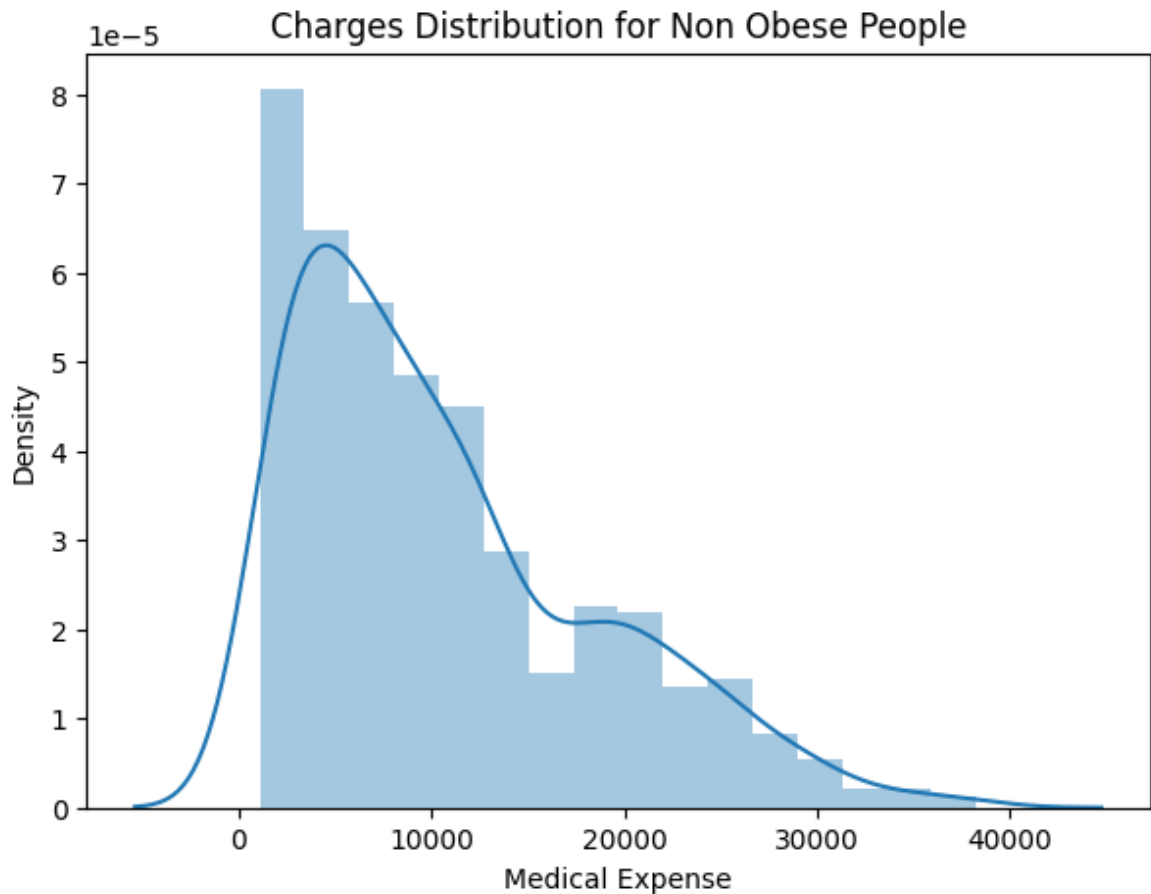
Charges distribution for patients with BMI greater than 30 i.e. obese patients

```
In [ ]: #bmi charges distribution for obese people
plt.figure(figsize=(7,5))
sns.distplot(df[(df.bmi >= 30)][['charges']])
plt.title('Charges Distribution for Obese People')
plt.xlabel('Medical Expense')
plt.show()
```



Charges distribution for patients with BMI less than 30 i.e. healthy patients

```
In [ ]: plt.figure(figsize=(7,5))
sns.distplot(df[(df.bmi < 30)]['charges'])
plt.title('Charges Distribution for Non Obese People')
plt.xlabel('Medical Expense')
plt.show()
```



Therefore, patients with BMI less than 30 are spending less on medical treatment than those with BMI greater than 30.

Through the EDA, we have a clear understanding about the data and the correlation between the variables. Now, I will build a model to predict the medical expense of patients.

Train Test Split

```
In [ ]: from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(df.drop('charges',axis=1), c
```

Model Building

Linear Regression

```
In [ ]: #Linear Regression
from sklearn.linear_model import LinearRegression
lr = LinearRegression()
lr
```

```
Out[ ]: ▼ LinearRegression
LinearRegression()
```

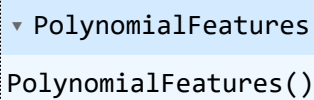
```
In [ ]: #model training
lr.fit(x_train,y_train)
#model accuracy
lr.score(x_train,y_train)
```

Out[]: 0.7368306228430945

```
In [ ]: #model prediction
y_pred = lr.predict(x_test)
```

Polynomial Regression

```
In [ ]: from sklearn.preprocessing import PolynomialFeatures
poly_reg = PolynomialFeatures(degree=2)
poly_reg
```

Out[]: 
PolynomialFeatures()

```
In [ ]: #transforming the features to higher degree
x_train_poly = poly_reg.fit_transform(x_train)
#splitting the data
x_train, x_test, y_train, y_test = train_test_split(x_train_poly, y_train, test_
```

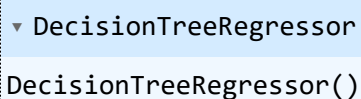
```
In [ ]: plr = LinearRegression()
#model training
plr.fit(x_train,y_train)
#model accuracy
plr.score(x_train,y_train)
```

Out[]: 0.8372892262994722

```
In [ ]: #model prediction
y_pred = plr.predict(x_test)
```

Decision Tree Regressor

```
In [ ]: #decision tree regressor
from sklearn.tree import DecisionTreeRegressor
dtree = DecisionTreeRegressor()
dtree
```

Out[]: 
DecisionTreeRegressor()

```
In [ ]: #model training
dtree.fit(x_train,y_train)
#model accuracy
dtree.score(x_train,y_train)
```

Out[]: 0.9993688476658964

```
In [ ]: #model prediction
dtree_pred = dtree.predict(x_test)
```

Random Forest Regressor

```
In [ ]: #random forest regressor
from sklearn.ensemble import RandomForestRegressor
rf = RandomForestRegressor(n_estimators=100)
rf
```

```
Out[ ]: ▼ RandomForestRegressor
RandomForestRegressor()
```

```
In [ ]: #model training
rf.fit(x_train,y_train)
#model accuracy
rf.score(x_train,y_train)
```

```
Out[ ]: 0.9753382671595934
```

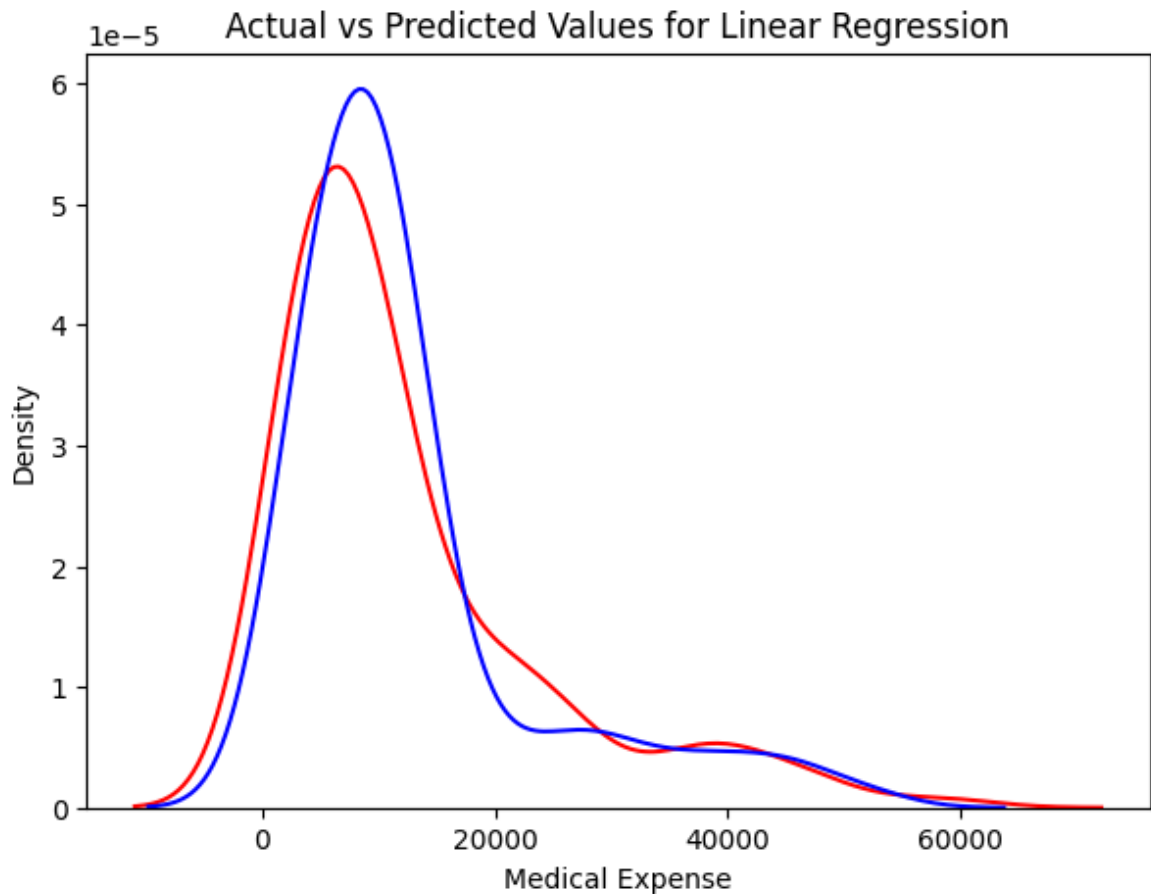
```
In [ ]: #model prediction
rf_pred = rf.predict(x_test)
```

Model Evaluation

```
In [ ]: from sklearn.metrics import mean_squared_error,mean_absolute_error,r2_score
```

Linear Regression

```
In [ ]: #distribution of actual and predicted values
plt.figure(figsize=(7,5))
ax1 = sns.distplot(y_test,hist=False,color='r',label='Actual Value')
sns.distplot(y_pred,hist=False,color='b',label='Predicted Value',ax=ax1)
plt.title('Actual vs Predicted Values for Linear Regression')
plt.xlabel('Medical Expense')
plt.show()
```

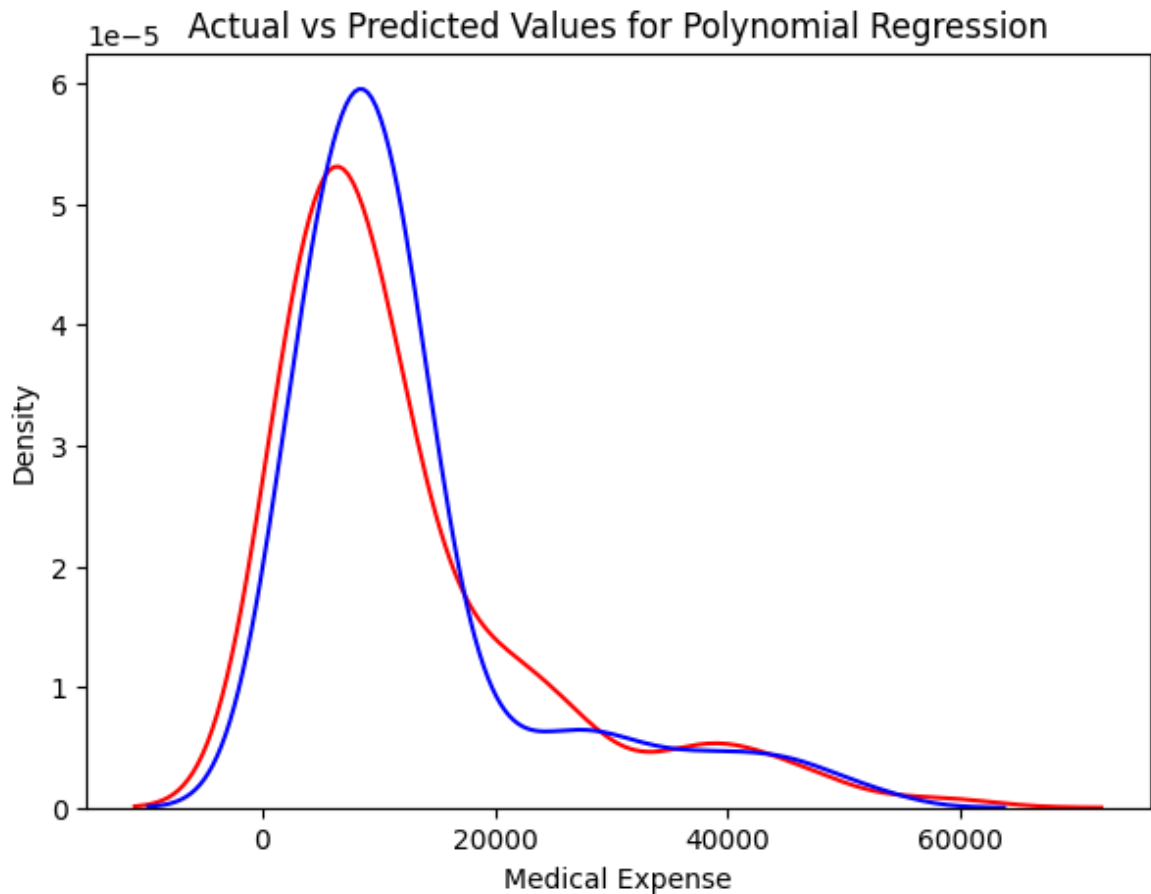


```
In [ ]: print('MAE:', mean_absolute_error(y_test, y_pred))
        print('MSE:', mean_squared_error(y_test, y_pred))
        print('RMSE:', np.sqrt(mean_squared_error(y_test, y_pred)))
        print('R2 Score:', r2_score(y_test, y_pred))
```

MAE: 2988.626627897196
MSE: 24512834.56541676
RMSE: 4951.043785447344
R2 Score: 0.8221477010678055

Polynomial Regression

```
In [ ]: #actual vs predicted values for polynomial regression
plt.figure(figsize=(7,5))
ax1 = sns.distplot(y_test,hist=False,color='r',label='Actual Value')
sns.distplot(y_pred,hist=False,color='b',label='Predicted Value',ax=ax1)
plt.title('Actual vs Predicted Values for Polynomial Regression')
plt.xlabel('Medical Expense')
plt.show()
```

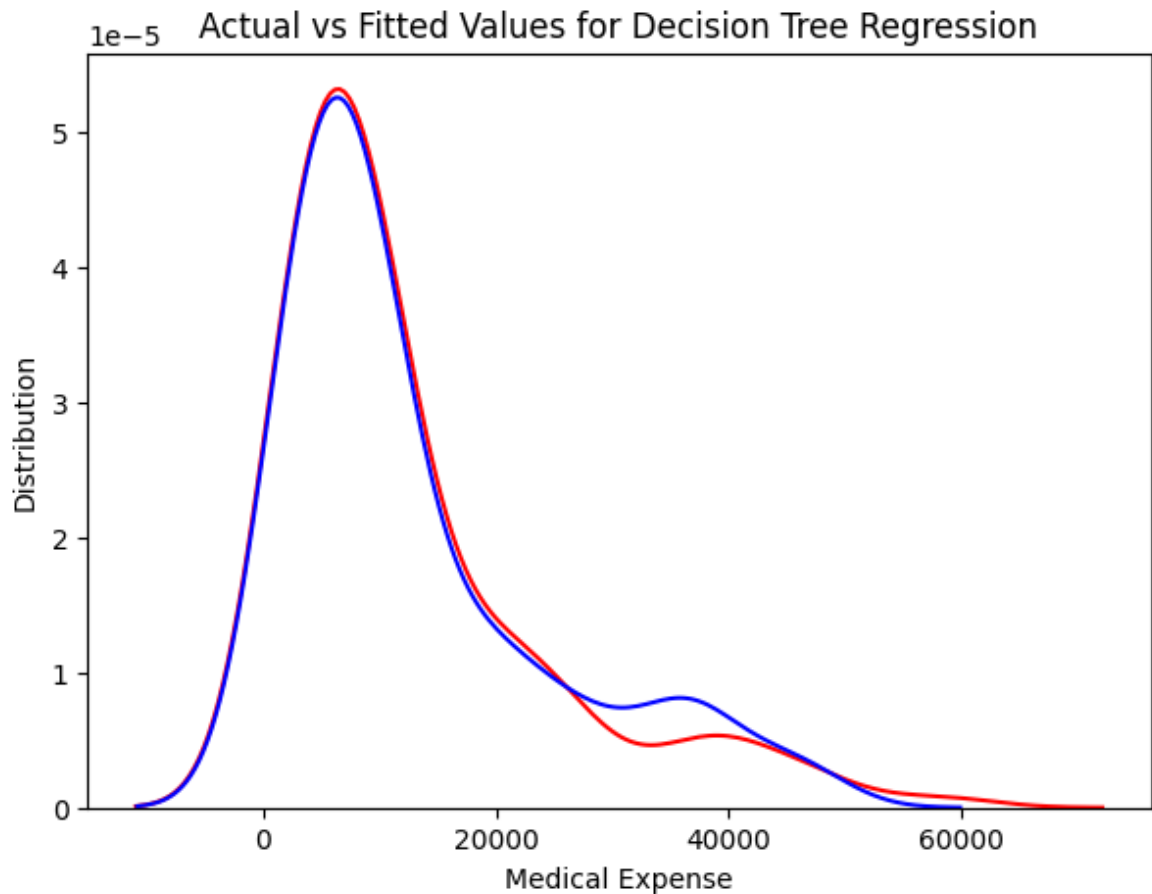



```
In [ ]: print('MAE:', mean_absolute_error(y_test, y_pred))
        print('MSE:', mean_squared_error(y_test, y_pred))
        print('RMSE:', np.sqrt(mean_squared_error(y_test, y_pred)))
        print('R2 Score:', r2_score(y_test, y_pred))
```

MAE: 2988.626627897196
MSE: 24512834.56541676
RMSE: 4951.043785447344
R2 Score: 0.8221477010678055

Decision Tree Regressor

```
In [ ]: #distribution plot of actual and predicted values
plt.figure(figsize=(7,5))
ax = sns.distplot(y_test, hist=False, color="r", label="Actual Value")
sns.distplot(dtrees_pred, hist=False, color="b", label="Fitted Values" , ax=ax)
plt.title('Actual vs Fitted Values for Decision Tree Regression')
plt.xlabel('Medical Expense')
plt.ylabel('Distribution')
plt.show()
```

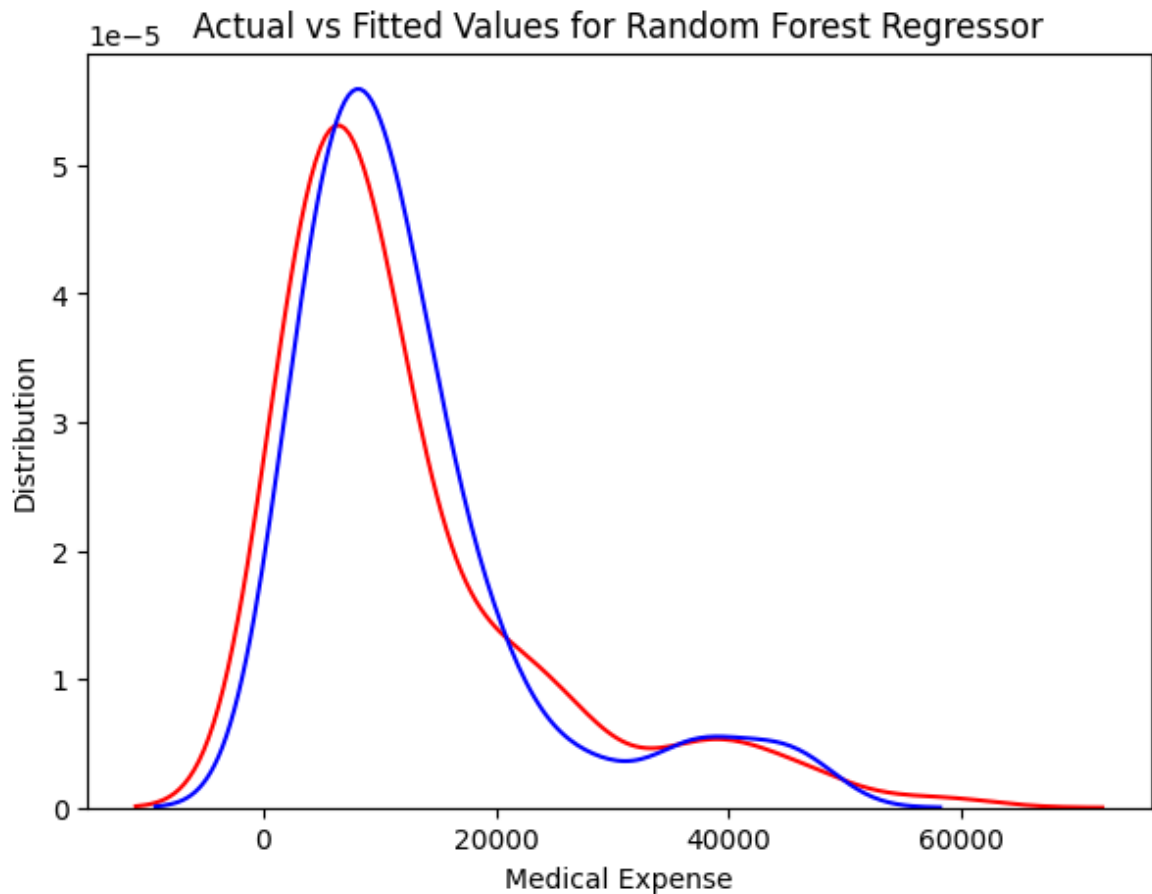


```
In [ ]: print('MAE:', mean_absolute_error(y_test, dtree_pred))
        print('MSE:', mean_squared_error(y_test, dtree_pred))
        print('RMSE:', np.sqrt(mean_squared_error(y_test, dtree_pred)))
        print('Accuracy:', dtree.score(x_test, y_test))
```

MAE: 3432.357628878505
MSE: 51680664.19095652
RMSE: 7188.926497812905
Accuracy: 0.6250321474582967

Random Forest Regressor

```
In [ ]: #distribution plot of actual and predicted values
plt.figure(figsize=(7,5))
ax = sns.distplot(y_test, hist=False, color="r", label="Actual Value")
sns.distplot(rf_pred, hist=False, color="b", label="Fitted Values" , ax=ax)
plt.title('Actual vs Fitted Values for Random Forest Regressor')
plt.xlabel('Medical Expense')
plt.ylabel('Distribution')
plt.show()
```



```
In [ ]: print('MAE:', mean_absolute_error(y_test, rf_pred))
        print('MSE:', mean_squared_error(y_test, rf_pred))
        print('RMSE:', np.sqrt(mean_squared_error(y_test, rf_pred)))
        print('Accuracy:', rf.score(x_test, y_test))
```

```
MAE: 2937.5177587331
MSE: 27234125.722924933
RMSE: 5218.632552970647
Accuracy: 0.8024034366036092
```

Conclusion

From the above models, we can see that Decision Tree Regressor and Random Forest Regressor are giving the best results. But, Random Forest Regressor is giving the best results with the least RMSE value. Therefore, I will use Random Forest Regressor to predict the medical expense of patients.

Moreover, the medical expense of smokers is higher than that of non-smokers. The medical expense of patients with BMI greater than 30 is higher than that of patients with BMI less than 30. The medical expense of older patients is higher than that of younger patients.

Thus, from the overall analysis, we can conclude that the medical expense of patients depends on their age, BMI, smoking habits.