

Context

Machine Learning with R by Brett Lantz is a book that provides an introduction to machine learning using R. As far as I can tell, Packt Publishing does not make its datasets available online unless you buy the book and create a user account which can be a problem if you are checking the book out from the library or borrowing the book from a friend. All of these datasets are in the public domain but simply needed some cleaning up and recoding to match the format in the book.

Content Columns

age: age of primary beneficiary

sex: insurance contractor gender, female, male

bmi: Body mass index, providing an understanding of body, weights that are relatively high or low relative to height, objective index of body weight (kg / m ²) using the ratio of height to weight, ideally 18.5 to 24.9

children: Number of children covered by health insurance / Number of dependents

smoker: Smoking

region: the beneficiary's residential area in the US, northeast, southeast, southwest, northwest.

charges: Individual medical costs billed by health insurance

Acknowledgements The dataset is available on [Kaggle](#) here.

Inspiration Can you accurately predict insurance costs?

```
In [1]: #We will try to analyse and visualize the dataset and try to predict the insurance costs of different individuals  
checking the regression models and try to understand different variable selection techniques like  
Backward elimination and forward selection
```

Importing Libraries

```
In [2]: import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns  
import plotly.express as px  
import matplotlib inline
```

Getting Dataset

```
In [76]: data_file=r'insurance.csv'
```

```
In [77]: dataset=pd.read_csv(data_file)
```

Getting insights from Dataset

```
In [78]: dataset
```

```
Out [78]:
```

	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	2007.84500
1337	61	female	29.070	0	yes	northwest	29141.36030

1338 rows x 7 columns

```
In [6]: dataset.head()
```

```
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

```
In [7]: dataset.tail()
```

```
Out [7]:
```

	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	2007.8450
1337	61	female	29.07	0	yes	northwest	29141.3603

```
In [8]: dataset.info()
```

```
Out [8]:
```

```
<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 [9]: dataset.describe(include='all')
```

```
Out [9]:
```

	age	sex	bmi	children	smoker	region	charges
count	1338.000000	1338	1338.000000	1338.000000	1338	1338	1338.000000
unique	NaN	2	NaN	NaN	2	4	NaN
top	NaN	male	NaN	NaN	NaN	southwest	NaN
freq	NaN	676	NaN	NaN	1064	364	NaN
mean	39.207025	NaN	30.863397	1.094918	NaN	NaN	13270.422265
std	14.049960	NaN	6.098187	1.205493	NaN	NaN	12110.011237
min	18.000000	NaN	15.960000	0.000000	NaN	NaN	1121.873900
25%	37.000000	NaN	26.296250	0.000000	NaN	NaN	4740.287150
50%	29.000000	NaN	30.400000	1.000000	NaN	NaN	9382.030000
75%	51.000000	NaN	34.693750	2.000000	NaN	NaN	16639.912515
max	64.000000	NaN	53.130000	5.000000	NaN	NaN	63770.428010

```
In [10]: dataset.isnull().sum()
```

```
Out [10]:
```

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

```
In [11]: dataset.columns
```

```
Out [11]: Index(['age', 'sex', 'bmi', 'children', 'smoker', 'region', 'charges'], dtype='object')
```

Visualization

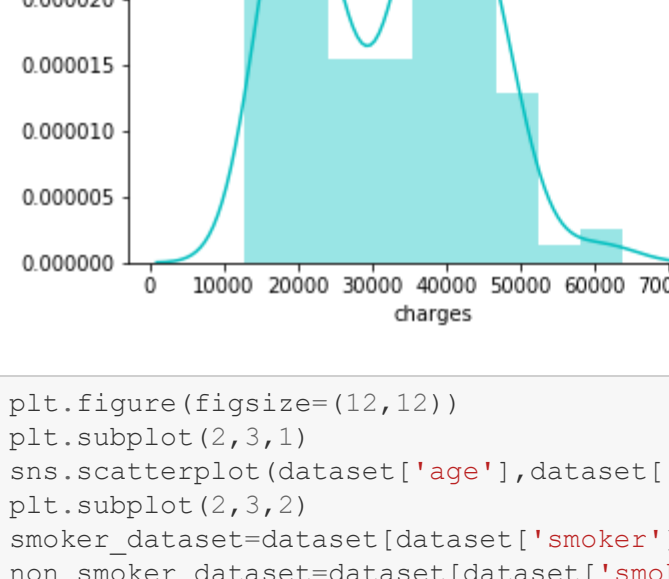
```
In [12]: correlation_plot=dataset.corr()  
mask = np.triu(np.ones_like(correlation_plot, dtype=np.bool))  
sns.heatmap(correlation_plot,mask=mask,annot=True,fmt='0.2f',linewidth=0.5)
```

```
Out [12]: <matplotlib.axes._subplots.AxesSubplot at 0x1a47aa4643b>
```



```
In [79]: correlation_plot=dataset.corr()  
mask = np.triu(np.ones_like(correlation_plot, dtype=np.bool))  
sns.heatmap(correlation_plot,mask=mask,annot=True,fmt='0.2f',linewidth=0.5)
```

```
Out [79]: <matplotlib.axes._subplots.AxesSubplot at 0x1a47b676780>
```



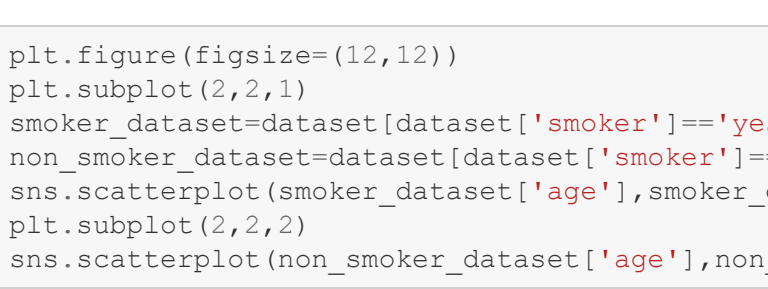
```
In [14]: correlation_plot['charges'].sort_values()
```

```
Out [14]:
```

```
children    0.067998  
bmi         0.198341  
age         0.299008  
charges     1.000000  
Name: charges, dtype: float64
```

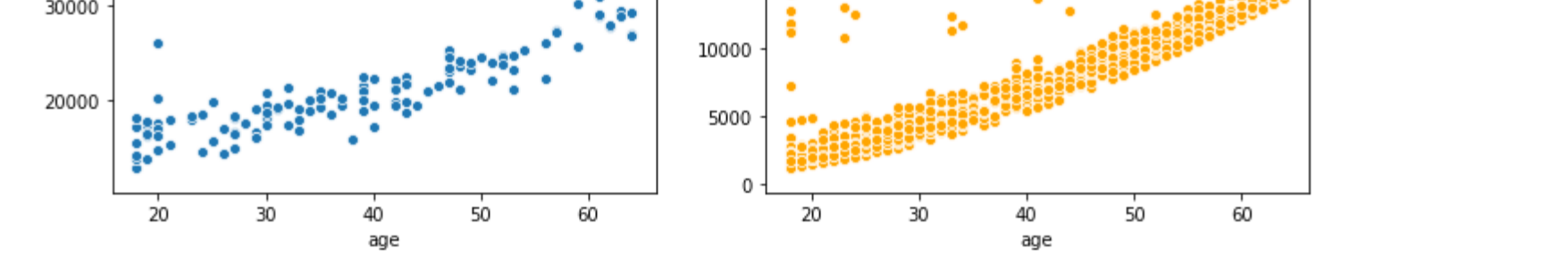
```
In [15]: sns.distplot(dataset['charges'])
```

```
Out [15]: <matplotlib.axes._subplots.AxesSubplot at 0x1a47b2a2cf8>
```



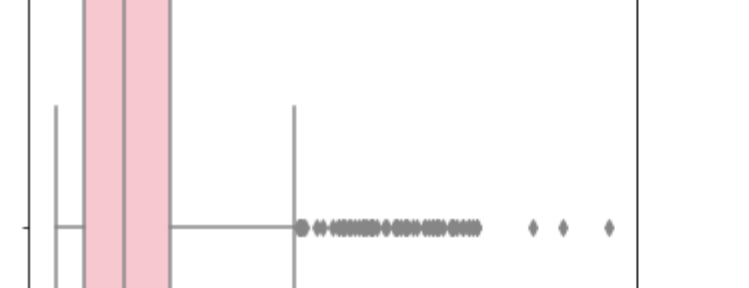
```
In [16]: plt.figure(figsize=(12,12))  
plt.subplot(2,2,1)  
explode=(0.1,0)  
color='blue', 'orange')  
dataset['smoker'].value_counts().plot.pie(autopct='%2.2f%%',shadow=True,explode=explode,color=color,lab  
el=labels)
```

```
Out [16]: Text(0.5, 1.0, 'Smokers vs Non-Smokers')
```



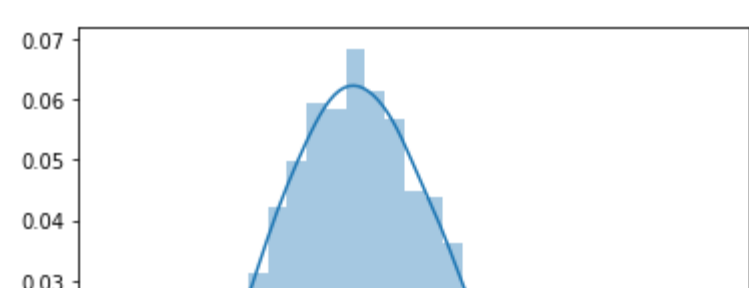
```
In [17]: sns.scatterplot(dataset['charges'],dataset['smoker'])
```

```
Out [17]: <matplotlib.axes._subplots.AxesSubplot at 0x1a47b6b24a8>
```



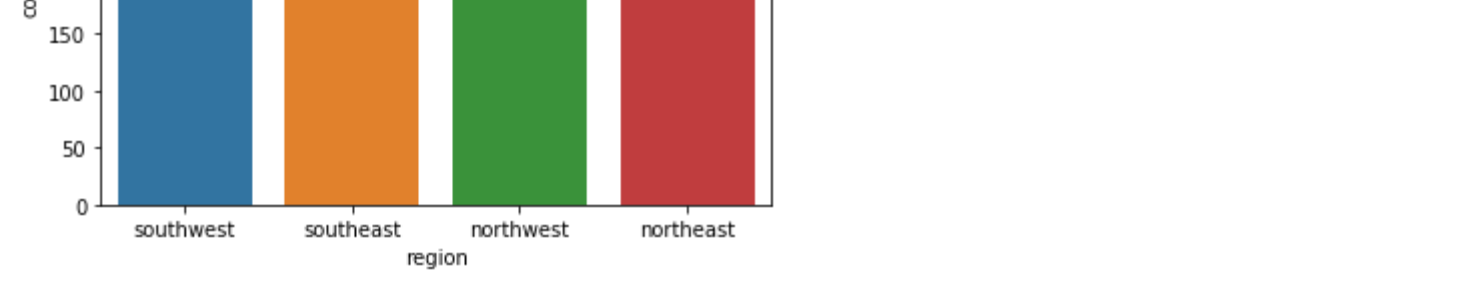
```
In [18]: sns.violinplot(dataset['charges'],dataset['smoker'],hue=dataset['sex'])
```

```
Out [18]: <matplotlib.axes._subplots.AxesSubplot at 0x1a47b4c76a0>
```



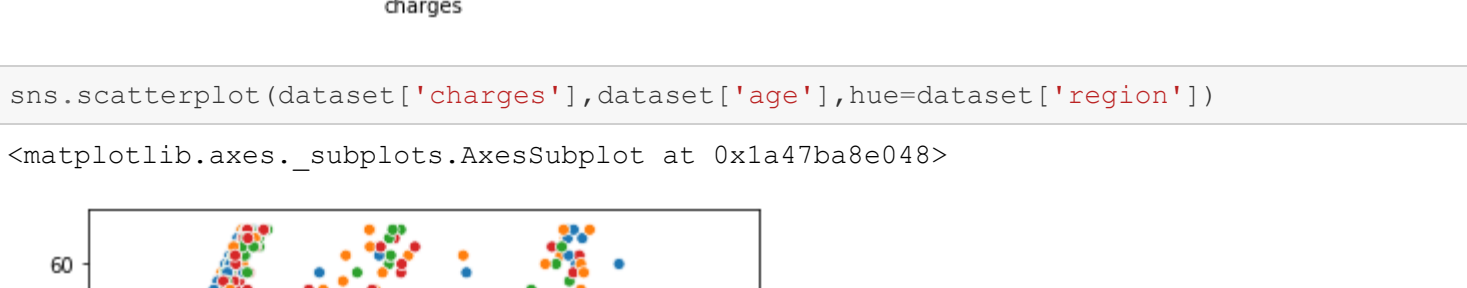
```
In [19]: f=plt.figure(figsize=(12,5))  
ax=f.add_subplot(1,2,1)  
sns.distplot(dataset[(dataset.smoker=='yes')]['charges'],color='c',ax=ax)  
ax=f.add_subplot(1,2,2)  
sns.distplot(dataset[(dataset.smoker=='no')]['charges'],color='b',ax=ax)
```

```
Out [19]: <matplotlib.axes._subplots.AxesSubplot at 0x1a47b539ec0>
```



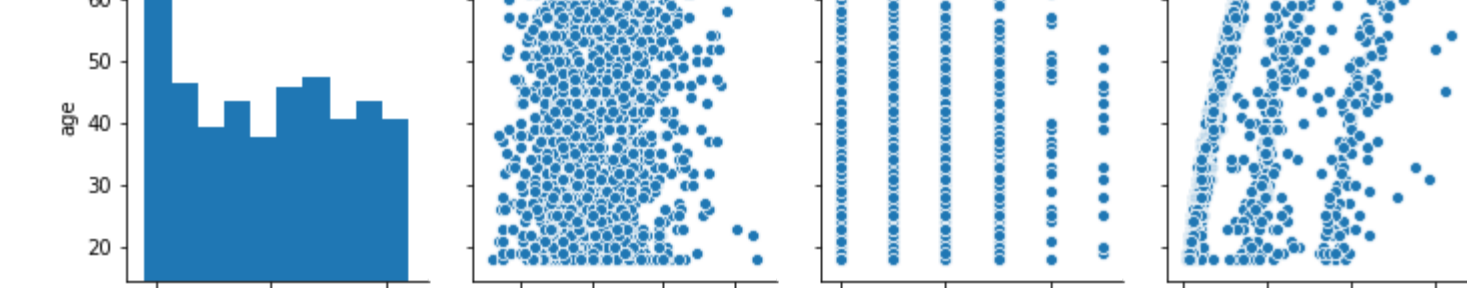
```
In [20]: plt.figure(figsize=(12,12))  
plt.subplot(2,2,1)  
sns.scatterplot(dataset['age'],dataset['charges'],hue=dataset['smoker'])  
plt.subplot(2,2,2)  
smoker_dataset=dataset[(dataset['smoker']=='yes')]  
non_smoker_dataset=dataset[(dataset['smoker']=='no')]  
sns.scatterplot(smoker_dataset['age'],smoker_dataset['charges'])  
plt.subplot(2,2,3)  
sns.scatterplot(non_smoker_dataset['age'],non_smoker_dataset['charges'],color='orange')
```

```
Out [20]: <matplotlib.axes._subplots.AxesSubplot at 0x1a47b6b24a8>
```



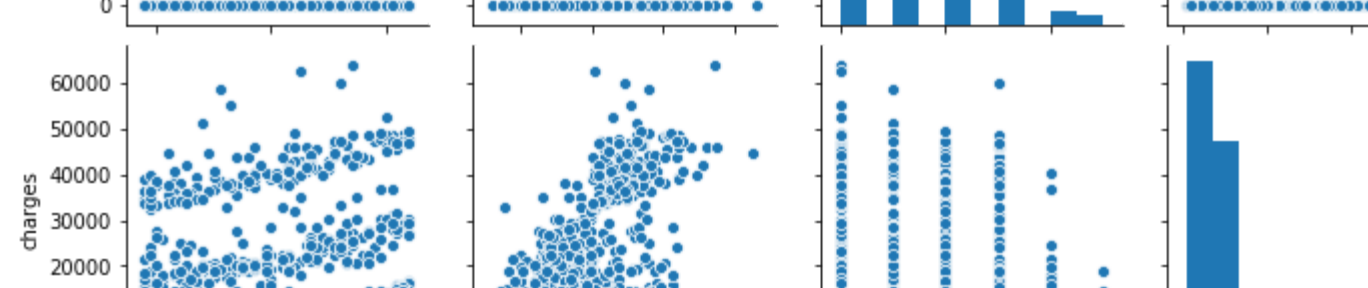
```
In [21]: plt.figure(figsize=(12,12))  
plt.subplot(2,2,1)  
smoker_dataset=dataset[(dataset['smoker']=='yes')]  
non_smoker_dataset=dataset[(dataset['smoker']=='no')]  
sns.scatterplot(smoker_dataset['age'],smoker_dataset['charges'])  
plt.subplot(2,2,2)  
sns.scatterplot(non_smoker_dataset['age'],non_smoker_dataset['charges'],color='orange')
```

```
Out [21]: <matplotlib.axes._subplots.AxesSubplot at 0x1a47b73b128>
```



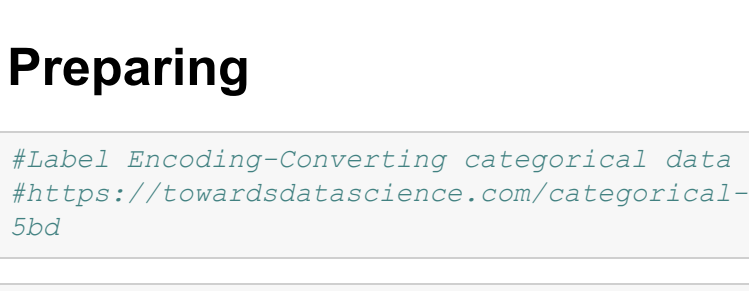
```
In [22]: female_dataset=dataset[(dataset['sex']=='female')]  
male_dataset=dataset[(dataset['sex']=='male')]  
plt.figure(figsize=(12,12))  
plt.subplot(2,2,1)  
sns.boxplot(female_dataset['charges'],color='pink')
```

```
Out [22]: <matplotlib.axes._subplots.AxesSubplot at 0x1a47ba0b5e0>
```



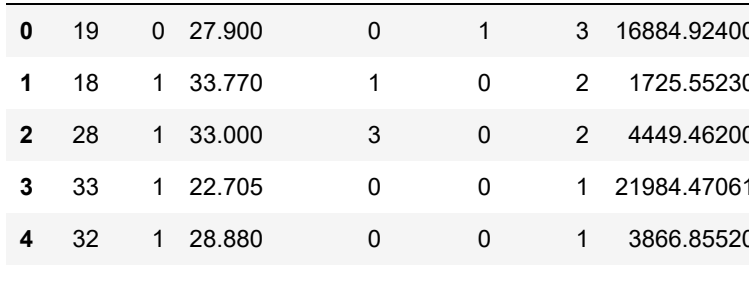
```
In [23]: sns.distplot(dataset['bmi'])
```

```
Out [23]: <matplotlib.axes._subplots.AxesSubplot at 0x1a47b8d7e48>
```



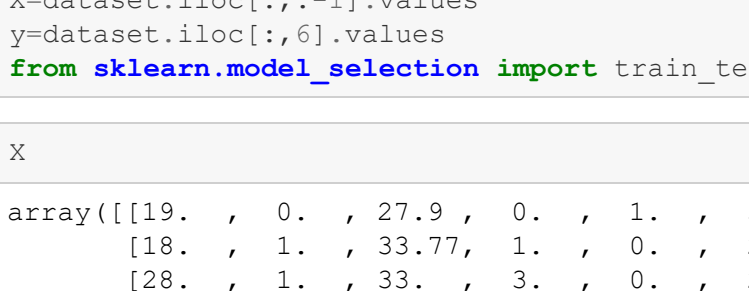
```
In [24]: sns.countplot(dataset['region'])
```

```
Out [24]: <matplotlib.axes._subplots.AxesSubplot at 0x1a47b845f8>
```



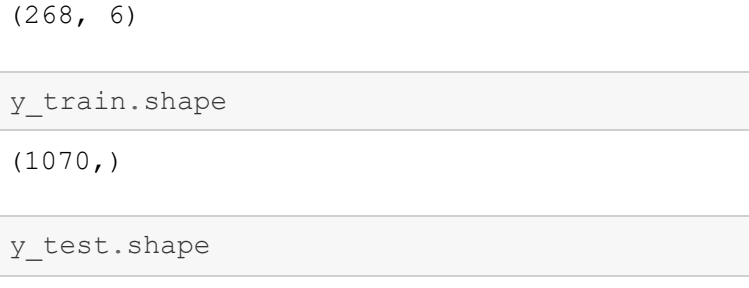
```
In [25]: sns.scatterplot(dataset['charges'],dataset['bmi'],hue=dataset['region'])
```

```
Out [25]: <matplotlib.axes._subplots.AxesSubplot at 0x1a47ba0b9e8>
```



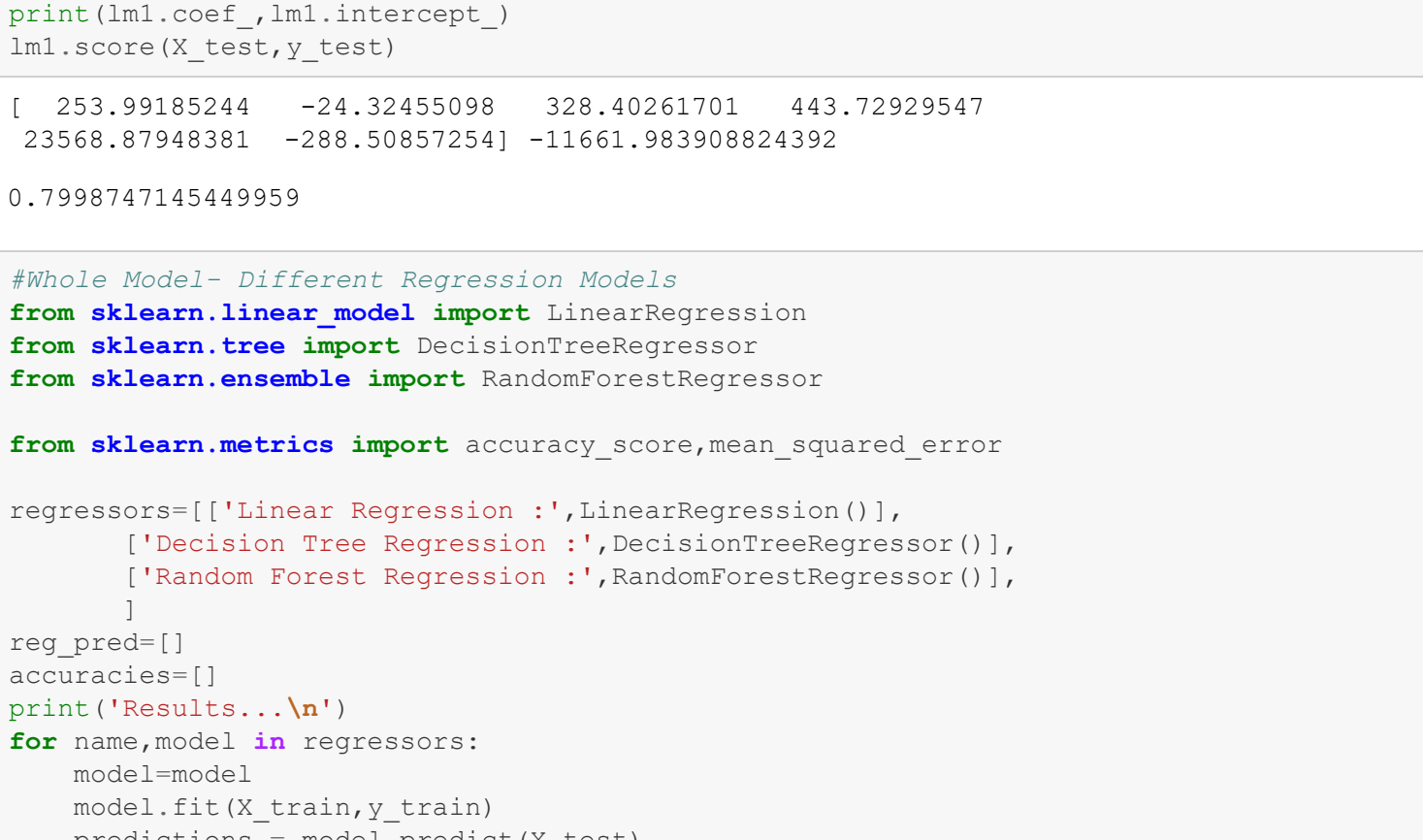
```
In [26]: sns.scatterplot(dataset['charges'],dataset['age'],hue=dataset['region'])
```

```
Out [26]: <matplotlib.axes._subplots.AxesSubplot at 0x1a47ba0b8e0>
```



```
In [27]: sns.pairplot(dataset)
```

```
Out [27]: <seaborn.axisgrid.PairGrid at 0x1a47bb0ba8>
```



Analysis

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

Preparing

```
In [28]: #Label Encoding-Converting categorical data to numerical data  
#https://towardsdatascience.com/categorical-encoding-using-label-encoding-and-one-hot-encoder-91ef77fb5b9a
```

```
In [29]: from sklearn.preprocessing import LabelEncoder  
le=LabelEncoder()  
le.fit(dataset.sex)  
dataset.sex=le.transform(dataset.sex)  
le.fit(dataset.smoker)  
dataset.smoker=le.transform(dataset.smoker)  
le.fit(dataset.region)  
dataset.region=le.transform(dataset.region)
```

```
In [30]: dataset.head()
```

```
Out [30]:
```

	age	sex	bmi	children	smoker	region	charges
0	19	0	27.900	0	1	3	16884.92400
1	18	1	33.770	1	0	2	1725.55230
2	28	1	33.000	3	0	2	4449.46200
3	33	1	22.705	0	0	1	21984.47061
4	32	1	28.880	0	0	1	3866.85520

Modelling

```
In [31]: #whole model(all 6 independent variables) used to predict the independent variable
```

```
In [32]: #backward elimination
```

```
In [33]: #forward selection
```

```
In [34]: #splitting data into training and testing dataset  
X=dataset.iloc[:,1:].values  
y=dataset.iloc[:,0].values  
from sklearn.model_selection import train_test_split
```

```
In [35]: X
```

```
Out [35]: array([[19, 0, 27.9, 0, 1, 3, ],  
 [18, 1, 33.77, 1, 0, 2, ],  
 [28, 1, 33, 3, 0, 2, ],  
 ...,  
 [18, 0, 36.85, 0, 0, 2, ],  
 [21, 0, 25.8, 0, 0, 3, ],  
 [61, 0, 29.07, 0, 1, 1, ]])
```

```
In [36]: y
```

```
Out [36]: array([16884.924, 1725.5523, 4449.462, ..., 1629.8335, 2007.945,  
 29141.3603])
```

```
In [37]: X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.20,random_state=0)
```

```
In [38]: X_train.shape
```

```
Out [38]: (1070, 6)
```

```
In [39]: X_test.shape
```

```
Out [39]: (268, 6)
```

```
In [40]: y_train.shape
```

```
Out [40]: (1070,)
```

```
In [41]: y_test.shape
```

```
Out [41]: (268,)
```

```
In [42]: #from sklearn.linear model import LinearRegression  
lm1=LinearRegression()  
lm1.fit(X_train,y_train)  
lm1.predict(X_test)  
print(lm1.coef_,lm1.intercept_)
```

```
Out [42]: [ 253.99185244 -24.32455098 328.40261701 443.72929547  
23568.87948381 -288.50857254] -11661.983908824392
```

```
Out [43]: 0.799877415449959
```

```
In [43]: #Whole Model- Different Regression Models  
from sklearn.linear_model import LinearRegression  
from sklearn.tree import DecisionTreeRegressor  
from sklearn.ensemble import RandomForestRegressor  
from sklearn.metrics import accuracy_score,mean_squared_error
```

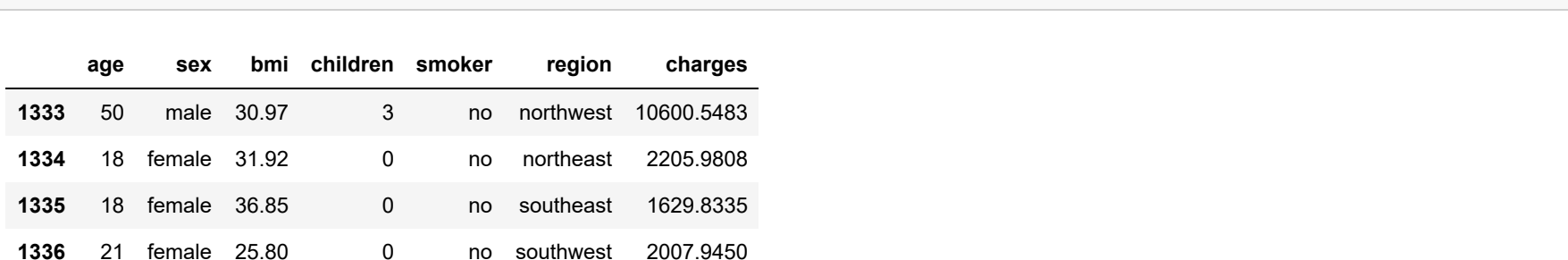
```
regressors=[('Linear Regression ',LinearRegression()),  
 ['Decision Tree Regression ',DecisionTreeRegressor()],  
 ['Random Forest Regression ',RandomForestRegressor()],  
 ]
```

```
reg_pred=[]  
accuracies=[]  
print('Results...\n')
```

```
for name,model in regressors:  
    model=model  
    model.fit(X_train,y_train)  
    predictions = model.predict(X_test)  
    rmse=np.sqrt(mean_squared_error(y_test, predictions))  
    reg_pred.append(rmse)  
    accuracy = accuracy_score(X_test,y_test)  
    accuracies.append(accuracy)  
    print(name,rmse,accuracy)
```

```
Out [43]: Linear Regression : 5643.219748880902 0.799877415449959  
Decision Tree Regression : 7132.62962421751 0.6802965034567527  
Random Forest Regression : 4462.589466666236 0.8748525767929204
```

```
Out [44]:
```



```
In [44]: #-----#
```

```
In [45]: #Backward Elimination  
import statsmodels.api as sm  
X_be=X[:,0:5,4,3]  
regressor_OLS=sm.OLS(endog=y,exog=X_be).fit()  
regressor_OLS.summary()
```

```
Out [45]:
```

Dep. Variable:	y	R-squared (uncentered):	0.874
Model:	OLS	Adj. R-squared (uncentered):	0.873
Method:	Least Squares	F-statistic:	1537.
Date:	Sat, 04 Dec 2021	Prob (F-statistic):	0.00
Time:	19:18:20	Log-Likelihood:	-13621.
No. Observations:	1338	AIC:	2.725e+04
Df Residuals:	1332	BIC:	2.725e+04
Df Model:	6		
Covariance Type:	nonrobust		

unique	NaN	2	NaN	NaN	2
top	NaN	male	NaN	NaN	no sou
freq	NaN	676	NaN	NaN	1064
mean	39.207025	NaN	30.663397	1.054918	NaN
std	14.049960	NaN	6.098187	1.205493	NaN
min	18.000000	NaN	15.000000	0.000000	NaN


```
In [47]: X_be3=X[:,[0,1,4,5]]
regressor_OLS=sm.OLS(endog=y,exog=X_be3).fit()
regressor_OLS.summary()
```

Model:

OLS

Adj. R-squared (uncentered):

0.872

Method:

Least Squares

F-statistic:

2273.

Date:

Sat, 04 Dec 2021

Prob (F-statistic):

0.00

Time:

19:18:20

Log Likelihood:

-13630.

No. Observations:

1338

AIC:

2.727e+04

DF Residuals:

1334

BIC:

2.729e+04

DF Model:

4

Covariance Type:

nonrobust

	coef	std err	t	P> t	[0.025	0.975]
x1	235.7394	7.296	32.313	0.000	221.427	250.051
x2	-317.3367	337.286	-0.941	0.347	-979.006	344.332
x3	2.356e+04	432.714	54.549	0.000	2.28e+04	2.45e+04
x4	-297.2862	146.476	-2.030	0.043	-584.635	-9.937
Omnibus:	264.736	Durbin-Watson:	2.077			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	616.743			
Skew:	1.084	Prob(JB):	1.19e-134			
Kurtosis:	5.523	Cond. No.	104.			

Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

In [48]:

X = be4X[[0,4,5]]

regressor_OLS = sm.OLS(endog=y,exog=X_be4).fit()

regressor_OLS.summary()

Out[48]:

OLS Regression Results

Dep. Variable:

y

R-squared (uncentered):

0.872

Model:

OLS

Adj. R-squared (uncentered):

0.872

Method:

Least Squares

F-statistic:

2031.

Date:

Sat, 04 Dec 2021

Prob (F-statistic):

0.00

Time:

19:18:21

Log Likelihood:

-13630.

No. Observations:

1338

AIC:

2.727e+04

DF Residuals:

1335

BIC:

2.729e+04

DF Model:

3

Covariance Type:

nonrobust

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```
In [48]: X_be4=X[:,[0,4,5]]
regressor_OLS=sm.OLS(endog=y,exog=X_be4).fit()
regressor_OLS.summary()
```

In [49]:	#Best Features are Age, Smoker, Regions
In [50]:	X=dataset.iloc[:,[0,4,5]].values y=dataset.iloc[:,6].values
In [51]:	X
Out[51]:	array([[19, 1, 3], [18, 0, 2], [28, 0, 2], [18, 0, 2], [21, 0, 3], [61, 1, 1]], dtype=int64)
In [52]:	y
Out[52]:	array([16884.924 , 1725.5523, 4449.462 , ..., 1629.8335, 29141.3603])
In [53]:	X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0
In [54]:	X_train.shape
Out[54]:	(1070, 3)
In [55]:	X_test.shape
Out[55]:	(268, 3)
In [56]:	y_train.shape
Out[56]:	(1070,)
In [57]:	y_test.shape
Out[57]:	(268,)
In [58]:	from sklearn.linear_model import LinearRegression from sklearn.tree import DecisionTreeRegressor from sklearn.ensemble import RandomForestRegressor from sklearn.metrics import accuracy_score,mean_squared_error regressor=[('Linear Regression ',LinearRegression()), [('Decision Tree Regression ',DecisionTreeRegressor()), [('Random Forest Regression ',RandomForestRegressor()),]] reg_pred=[]

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```
In [49]: #Best Features are Age, Smoker, Regions
```

```
In [50]: X=dataset.iloc[:,[0,4,5]].values
y=dataset.iloc[:,4].values
```

```
In [51]: X
```

```
Out [51]: array([[19,  1,  31,
        [18,  0,  21,
        [28,  0,  21,
        ...,
        [18,  0,  21,
        [21,  0,  31,
        [61,  1,  31], dtype=int64])
```

```
In [52]: y
```

```
Out [52]: array([16884.924 , 1725.5523, 4449.462 , ..., 1629.8335, 2007.945 ,
        [29141.3603])
```

```
In [53]: X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.20,random_state=0)
```

```
In [54]: X_train.shape
```

```
Out [54]: (1070, 3)
```

```
In [55]: X_test.shape
```

```
Out [55]: (268, 3)
```

```
In [56]: y_train.shape
```

```
Out [56]: (1070,)
```

```
In [57]: y_test.shape
```

```
Out [57]: (268,)
```

```
In [58]: from sklearn.linear_model import LinearRegression
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor

from sklearn.metrics import accuracy_score,mean_squared_error

regressors=[('Linear Regression ',LinearRegression()),
             ('Decision Tree Regression ',DecisionTreeRegressor()),
             ('Random Forest Regression ',RandomForestRegressor()),
             ]

reg_pred=[]
accuracies=[]
print("Results...\n")
for name,model in regressors:
    model=model
    model.fit(X_train,y_train)
    predictions = model.predict(X_test)
    rms=np.sqrt(mean_squared_error(y_test, predictions))
    reg_pred.append(rms)
    accuracy= model.score(X_test,y_test)
    accuracies.append(accuracy)
    print(name,rms,accuracy)

y_ax=['Linear Regression' , 'Decision Tree Regression' , 'Random Forest Regression' ]
x_ax=reg_pred

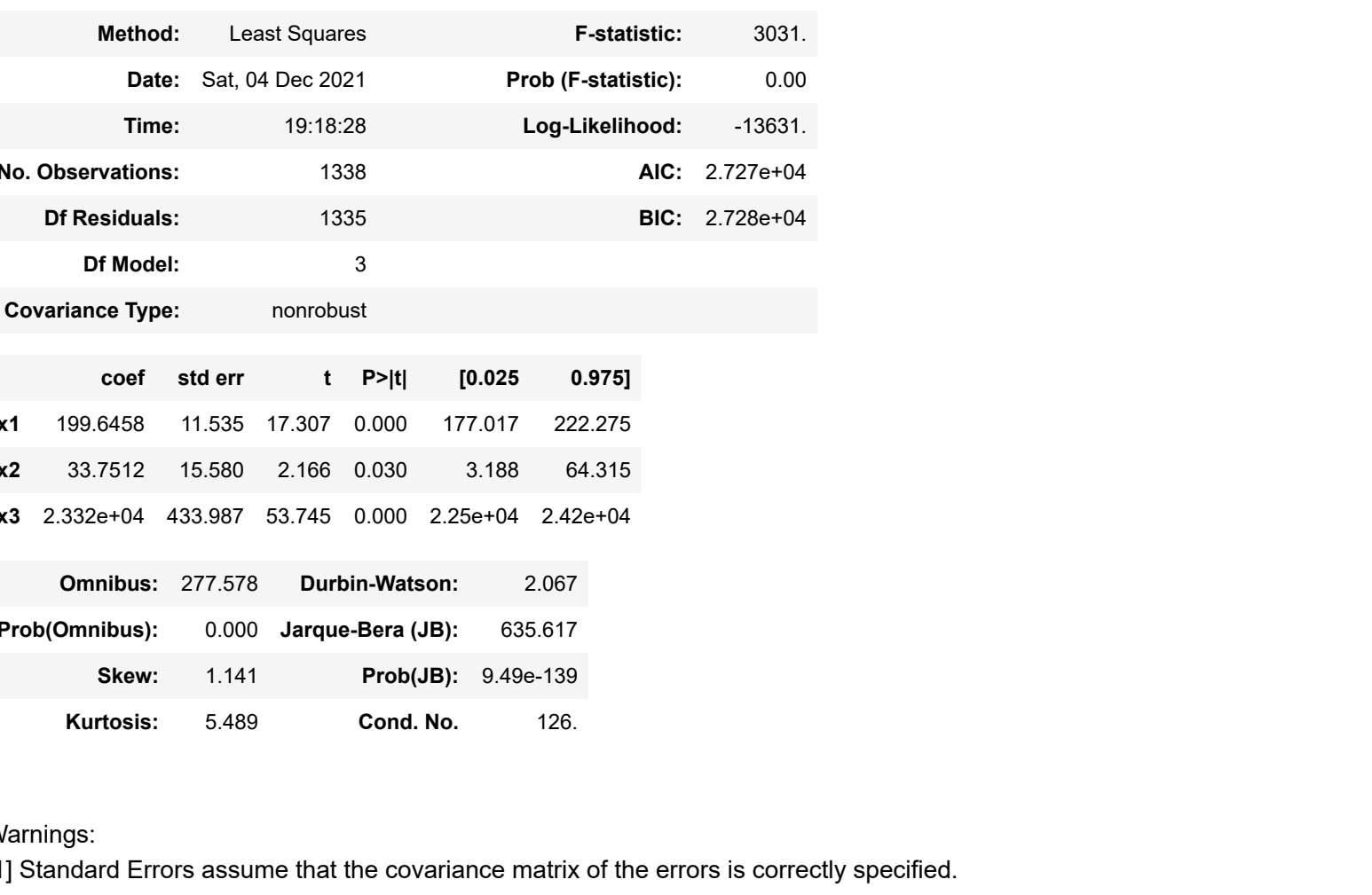
plt.figure(figsize=(12,5))
plt.subplot(1,2,1)
sns.barplot(x=x_ax,y=y_ax,linewidth=1.5,edgecolor="0.1")
plt.title('RMSE Scores')
plt.plot()

plt.subplot(1,2,2)
sns.barplot(x=accuracies,y=y_ax,linewidth=1.5,edgecolor="0.1")
plt.title('Accuracies')
plt.plot()

Results...
```

Linear Regression : 6016.21369593682 0.7725454923581263
Decision Tree Regression : 7543.001912826166 0.6424502524438279
Random Forest Regression : 6938.857439698769 0.6974313205014286

```
Out [58]: []
```



```
In [59]: #-----#
```

```
In [60]: #Forward Selection
```

```
In [61]: X_fsl=dataset.iloc[:,0].values
y=dataset.iloc[:,4].values
```

```
In [62]: X_fsl
```

```
Out [62]: array([19, 18, 28, ..., 18, 21, 61], dtype=int64)
```

```
In [63]: y
```

```
Out [63]: array([16884.924 , 1725.5523, 4449.462 , ..., 1629.8335, 2007.945 ,
        [29141.3603])
```

```
In [64]: X_fsl=dataset.iloc[:,0].values
regressor_OLS=sm.OLS(endog=y,exog=X_fsl).fit()
regressor_OLS.summary()
```

```

Random Forest Regression : RandomForestRegressor()

reg_pred=[]
accuracies=[]
print('Results...\\n')
for name,model in regressors:
    model=model
    model.fit(X_train,y_train)
    predictions = model.predict(X_test)
    rms=np.sqrt(mean_squared_error(y_test, predictions))
    reg_pred.append(rms)
    accuracy= model.score(X_test,y_test)
    accuracies.append(accuracy)
    print(name,rms,accuracy)

y_ax='Linear Regression' , 'Decision Tree Regression' , 'Random
x_ax=req_pred

plt.figure(figsize=(12,9))
plt.subplot(1,2,1)
sns.barplot(x=x_ax,y=y_ax,linewidth=1.5,edgecolor="0.1")
sns.barplot(x='RMS Scores'
plt.plot()
plt.subplot(1,2,2)
sns.barplot(x=accuracies,y=y_ax,linewidth=1.5,edgecolor="0.1")
plt.title('Accuracies')
plt.plot()

Results...

Linear Regression : 5698.626766322956 , 0.7966412232221931
Decision Tree Regression : 7371.159031214616 , 0.65856278771606
Random Forest Regression : 4777.263907022525 , 0.856581109523612

Out[4]: []

```

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```
In [65]: X_fsl=dataset.iloc[:,[0,1]].values
regressor_OLS=sm.OLS(endog=y,exog=X_fsl2).fit()
regressor_OLS.summary()
```

OLS Regression Results						
Dep. Variable:	y	R-squared (uncentered):	0.587			
Model:	OLS	Adj. R-squared (uncentered):	0.586			
Method:	Least Squares	F-statistic:	947.6			
Date:	Sat, 04 Dec 2021	Prob (F-statistic):	6.01e-257			
Time:	19:18:27	Log-Likelihood:	-14415.			
No. Observations:	1338	AIC:	2.883e+04			
DF Residuals:	1336	BIC:	2.884e+04			
DF Model:	2					
Covariance Type:	nonrobust					
	coef	std err	t	P> t		
				[0.025	0.975]	
x1	306.1242	10.150	30.160	0.000	286.213	326.036
x2	2043.2674	594.697	3.436	0.001	876.626	3209.908
Omnibus:	387.299	Durbin-Watson:	2.052			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	819.051			
Skew:	1.691	Prob(JB):	1.40e-178			
Kurtosis:	4.803	Cond. No.	78.4			

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```
In [66]: X_fsl=dataset.iloc[:,[0,1,2]].values
regressor_OLS=sm.OLS(endog=y,exog=X_fsl3).fit()
regressor_OLS.summary()
```

OLS Regression Results						
Dep. Variable:	y	R-squared (uncentered):	0.596			
Model:	OLS	Adj. R-squared (uncentered):	0.595			
Method:	Least Squares	F-statistic:	655.7			
Date:	Sat, 04 Dec 2021	Prob (F-statistic):	6.56e-262			
Time:	19:18:27	Log-Likelihood:	-14400.			
No. Observations:	1338	AIC:	2.881e+04			
DF Residuals:	1335	BIC:	2.882e+04			
DF Model:	3					
Covariance Type:	nonrobust					
	coef	std err	t	P> t		
				[0.025	0.975]	
x1	207.6709	20.510	10.125	0.000	167.436	247.906
x2	979.3428	619.212	1.582	0.114	-235.392	2194.078
x3	158.2079	28.739	5.505	0.000	101.830	214.586
Omnibus:	364.363	Durbin-Watson:	2.024			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	731.946			
Skew:	1.631	Prob(JB):	1.21e-159			
Kurtosis:	4.578	Cond. No.	102.			

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```
In [67]: X_fsl=dataset.iloc[:,[0,2]].values
regressor_OLS=sm.OLS(endog=y,exog=X_fsl4).fit()
regressor_OLS.summary()
```

OLS Regression Results						
Dep. Variable:	y	R-squared (uncentered):	0.595			
Model:	OLS	Adj. R-squared (uncentered):	0.594			
Method:	Least Squares	F-statistic:	981.2			
Date:	Sat, 04 Dec 2021	Prob (F-statistic):	6.45e-263			
Time:	19:18:28	Log-Likelihood:	-14401.			
No. Observations:	1338	AIC:	2.881e+04			
DF Residuals:	1336	BIC:	2.882e+04			
DF Model:	2					
Covariance Type:	nonrobust					
	coef	std err	t	P> t		
				[0.025	0.975]	
x1	208.8632	20.507	10.185	0.000	168.633	249.094
x2	172.3943	27.318	6.311	0.000	118.803	225.986
Omnibus:	364.702	Durbin-Watson:	2.016			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	732.185			
Skew:	1.634	Prob(JB):	1.02e-159			
Kurtosis:	4.567	Cond. No.	5.51			

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```
In [68]: X_fsl=dataset.iloc[:,[0,2,4]].values
regressor_OLS=sm.OLS(endog=y,exog=X_fsl5).fit()
regressor_OLS.summary()
```

OLS Regression Results						
Dep. Variable:	y	R-squared (uncentered):	0.872			
Model:	OLS	Adj. R-squared (uncentered):	0.872			
Method:	Least Squares	F-statistic:	3031.			
Date:	Sat, 04 Dec 2021	Prob (F-statistic):	0.00			
Time:	19:18:28	Log-Likelihood:	-13631.			
No. Observations:	1338	AIC:	2.727e+04			
DF Residuals:	1335	BIC:	2.728e+04			
DF Model:	3					
Covariance Type:	nonrobust					
	coef	std err	t	P> t		
				[0.025	0.975]	
x1	199.6458	11.535	17.307	0.000	177.017	222.275
x2	33.7512	15.580	2.166	0.030	3.188	64.315
x3	2.332e+04	433.987	53.745	0.000	2.25e+04	2.42e+04
Omnibus:	277.578	Durbin-Watson:	2.067			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	635.617			
Skew:	1.141	Prob(JB):	9.49e-139			
Kurtosis:	5.489	Cond. No.	126.			

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```
In [69]: X_fsl=dataset.iloc[:,[0,2,4,5]].values
regressor_OLS=sm.OLS(endog=y,exog=X_fsl6).fit()
regressor_OLS.summary()
```

OLS Regression Results						
Dep. Variable:	y	R-squared (uncentered):	0.873			
Model:	OLS	Adj. R-squared (uncentered):	0.873			
Method:	Least Squares	F-statistic:	2234.			
Date:	Sat, 04 Dec 2021	Prob (F-statistic):	0.00			
Time:	19:18:29	Log-Likelihood:	-13625.			
No. Observations:	1338	AIC:	2.726e+04			
DF Residuals:	1334	BIC:	2.728e+04			
DF Model:	4					
Covariance Type:	nonrobust					
	coef	std err	t	P> t		
				[0.025	0.975]	
x1	200.8988	11.495	17.478	0.000	178.349	223.448
x2	58.8403	17.145	3.432	0.001	25.205	92.475
x3	2.334e+04	432.249	53.887	0.000	2.25e+04	2.42e+04
x4	-549.2686	159.673	-3.440	0.001	-862.507	-236.031
Omnibus:	273.748	Durbin-Watson:	2.071			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	626.116			
Skew:	1.126	Prob(JB):	1.10e-136			
Kurtosis:	5.482	Cond. No.	126.			

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```
In [70]: X=
```