

Programmed By : Rithik Tripathi

[Connect with me on LinkedIn \(https://www.linkedin.com/in/rithik-tripathi-data-scientist/\)](https://www.linkedin.com/in/rithik-tripathi-data-scientist/)

## Linear Regression

In [1]:

```
1 #importing Libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
```

In [2]:

```
1 # Importing Data
2 df = pd.read_csv('Dataset/train_cleaned.csv')
3 df.head()
```

Out[2]:

	Item_Weight	Item_Visibility	Item_MRP	Outlet_Establishment_Year	Item_Outlet_Sales	Item_Fat_Content_LF	Item_Fat_Content_Low Fat	Item_Fat_Content_Re
0	9.30	0.016047	249.8092	1999	3735.1380	0	1	
1	5.92	0.019278	48.2692	2009	443.4228	0	0	
2	17.50	0.016760	141.6180	1999	2097.2700	0	1	
3	19.20	0.000000	182.0950	1998	732.3800	0	0	
4	8.93	0.000000	53.8614	1987	994.7052	0	1	

5 rows × 46 columns

### Separating Target Variables

In [3]:

```
1 x = df.drop(['Item_Outlet_Sales'], axis=1)
2 y = df['Item_Outlet_Sales']
3 df.shape, x.shape, y.shape
```

Out[3]:

((8523, 46), (8523, 45), (8523,))

### Train Test Split

In [4]:

```
1 from sklearn.model_selection import train_test_split
2 train_x, test_x, train_y, test_y = train_test_split(x,y, random_state = 9)
3
4 train_x.shape, test_x.shape, train_y.shape, test_y.shape
```

Out[4]:

((6392, 45), (2131, 45), (6392,), (2131,))

### Implementing Linear Regression

In [5]:

```
1 from sklearn.linear_model import LinearRegression as LR
2 from sklearn.metrics import mean_absolute_error as mae
```

In [6]:

```
1 # creating instance of Linear Regression Class
2 lr = LR()
3
4 # Training the model with Gradient Descent
5 lr.fit(train_x, train_y)
```

Out[6]:

LinearRegression()

### predicting over Training Set

In [7]:

```
1 train_pred = lr.predict(train_x)
2 train_mae = mae(train_pred, train_y)
3 print('Training Mean Absolute Error :', train_mae)
```

Training Mean Absolute Error : 827.3664175128886

predicting over Test Set

In [8]:

```
1 test_pred = lr.predict(test_x)
2 test_mae = mae(test_pred, test_y)
3 print('Test Mean Absolute Error :', test_mae)
```

Test Mean Absolute Error : 861.806594100068

Parameters of Linear Regression

In [9]:

```
1 # coefficients of the equation
2 lr.coef_
```

Out[9]:

array([ -1.13512756, -179.79233388, 15.65668767, -18.6574453 ,
 4.26675415, 4.2154459 , 43.9152877 , 20.49224006,
 -72.8897278 , -22.38153707, -1.17784843, -91.0598699 ,
 -15.79938771, -39.47904273, -22.75735283, 8.6798061 ,
 -52.75260518, -20.67316491, -20.45574445, 19.09990758,
 29.94212896, 268.66029346, -41.52897269, -46.52031544,
 48.20370526, -495.50373377, -46.06491435, 158.18888447,
 -38.63257523, -468.75193706, 595.48427957, 143.00628086,
 -89.44056606, 309.29119287, -67.57691129, -46.06491435,
 489.27479304, -16.45446332, -227.03765548, 211.75459927,
 15.28305621, -964.25567083, 407.4039665 , -38.63257523,
 595.48427957])

In [20]:

```
1 feature_importance_df = pd.DataFrame({
2     'Feature' : train_x.columns,
3     'Importance': lr.coef_
4 })
5
6 feature_importance_df['Importance_magnitude'] = feature_importance_df.Importance.abs()
7
8 feature_importance_df.sort_values(by='Importance_magnitude', inplace=True, ascending=False)
9 # sorting so we can plot only top n important features for better visibility
10
11 feature_importance_df.head()
```

Out[20]:

	Feature	Importance	Importance_magnitude
41	Outlet_Type_Grocery Store	-964.255671	964.255671
44	Outlet_Type_Supermarket Type3	595.484280	595.484280
30	Outlet_Identifier_OUT027	595.484280	595.484280
25	Outlet_Identifier_OUT010	-495.503734	495.503734
36	Outlet_Size_Medium	489.274793	489.274793

In [23]:

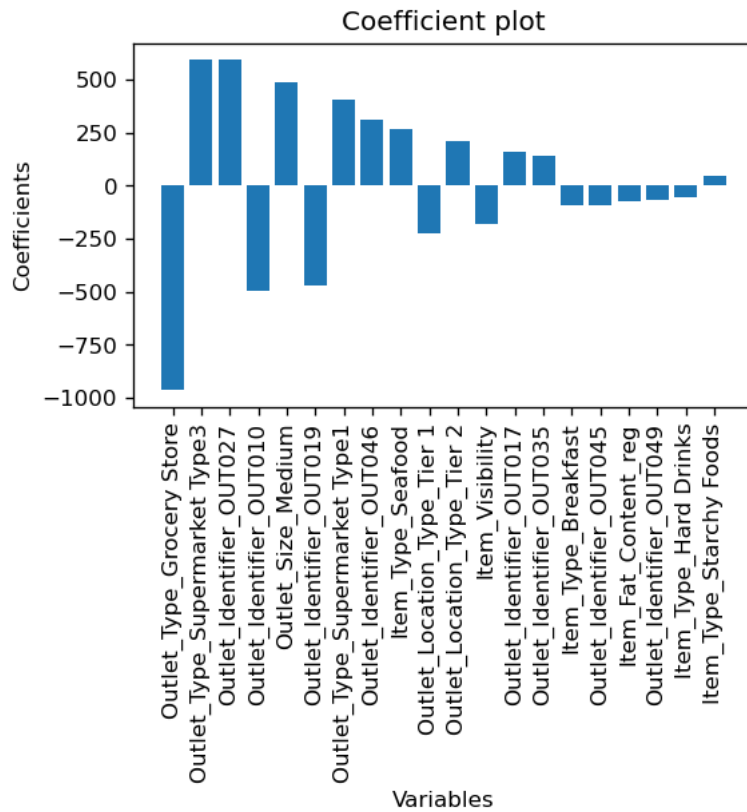
```

1 # plotting the coefficients
2
3
4 plt.figure(figsize=(5,3), dpi=120, facecolor='w', edgecolor='b')
5 x = feature_importance_df['Feature'].head(20)
6 y = feature_importance_df['Importance'].head(20)
7 plt.bar( x, y )
8 plt.xlabel( "Variables")
9 plt.ylabel('Coefficients')
10 plt.xticks(rotation=90)
11 plt.title('Coefficient plot')

```

Out[23]:

Text(0.5, 1.0, 'Coefficient plot')



Note : Here we can see that the model depends upon some Independent variables too much, But these coefficients are not suitable for interpretation because these are not scaled, we will look into this later on in details.

## Validating Assumptions of Linear Regression

In [11]:

```

1 # Calculating Error Terms
2
3 error_df = pd.DataFrame({
4     'original_values': test_y,
5     'predicted_values': test_pred
6 })
7
8 error_df['error'] = error_df['original_values'] - error_df['predicted_values']
9
10 error_df.head()

```

Out[11]:

	original_values	predicted_values	error
2344	1095.2410	3460.817769	-2365.576769
4005	5070.7328	3000.322073	2070.410727
2897	892.1720	712.864464	179.307536
6252	2976.1260	3200.235119	-224.109119
6414	365.5242	1956.201769	-1590.677569

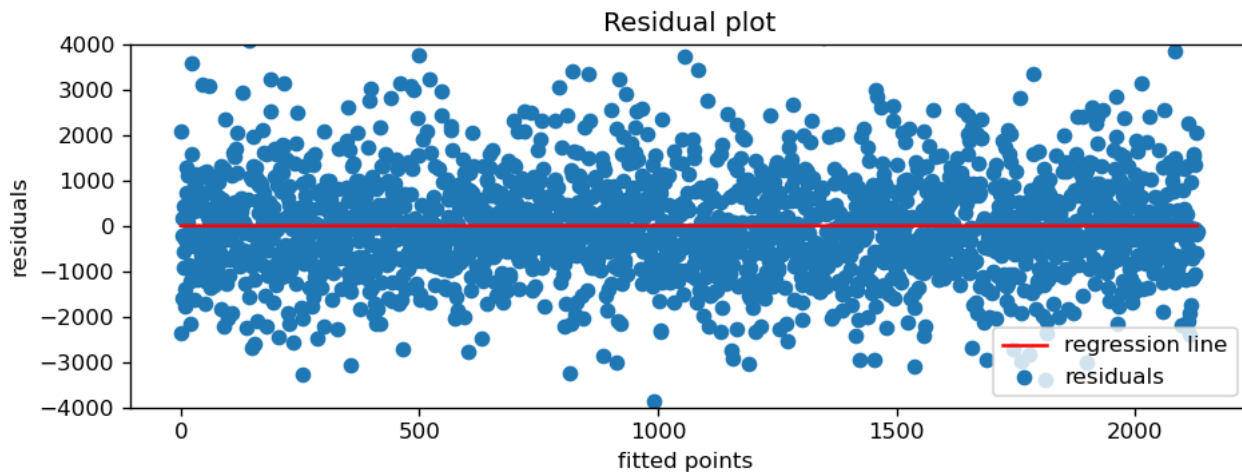
## Homoscedasticity : Constant Vairance of Error Terms

In [12]:

```
1 plt.figure(figsize=(9,3), dpi=120, facecolor='w', edgecolor='b')
2 f = range(0,2131)
3 k = [0 for i in range(0,2131)]
4 plt.scatter( f, error_df.error, label = 'residuals')
5 plt.plot( f, k , color = 'red', label = 'regression line' )
6 plt.xlabel('fitted points ')
7 plt.ylabel('residuals')
8 plt.title('Residual plot')
9 plt.ylim(-4000, 4000)
10 plt.legend()
```

Out[12]:

<matplotlib.legend.Legend at 0x1a7cd95b400>

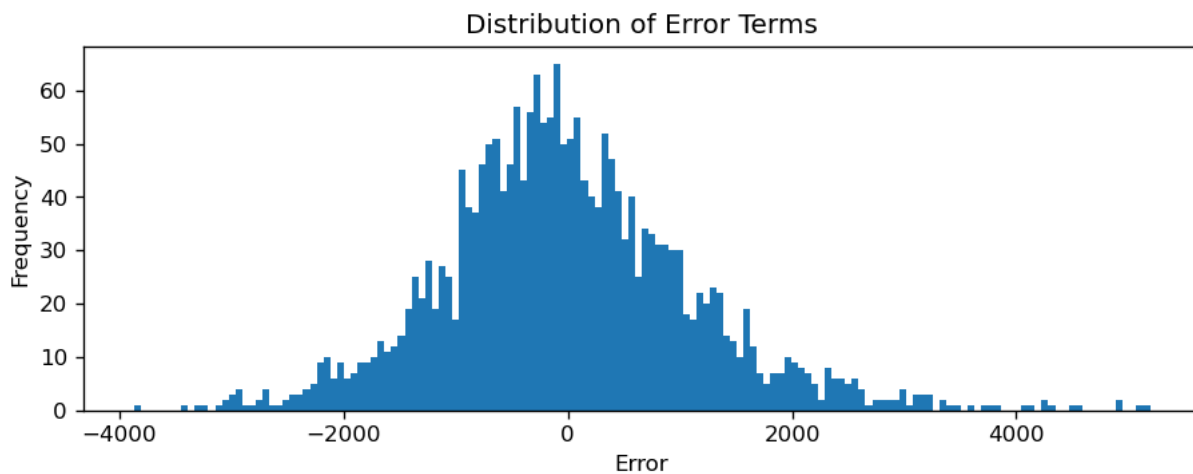


The Residual plot clearly Looks Homoscedastic, i.e. the the variance of the error across the dataset is nearly constant.

## Normally Distributed Error Terms

In [13]:

```
1 # Histogram for distribution of error terms
2 plt.figure(figsize=(9,3), dpi=120, facecolor='w', edgecolor='b')
3 plt.hist(error_df.error, bins = 150)
4 plt.xlabel('Error')
5 plt.ylabel('Frequency')
6 plt.title('Distribution of Error Terms')
7 plt.show()
```



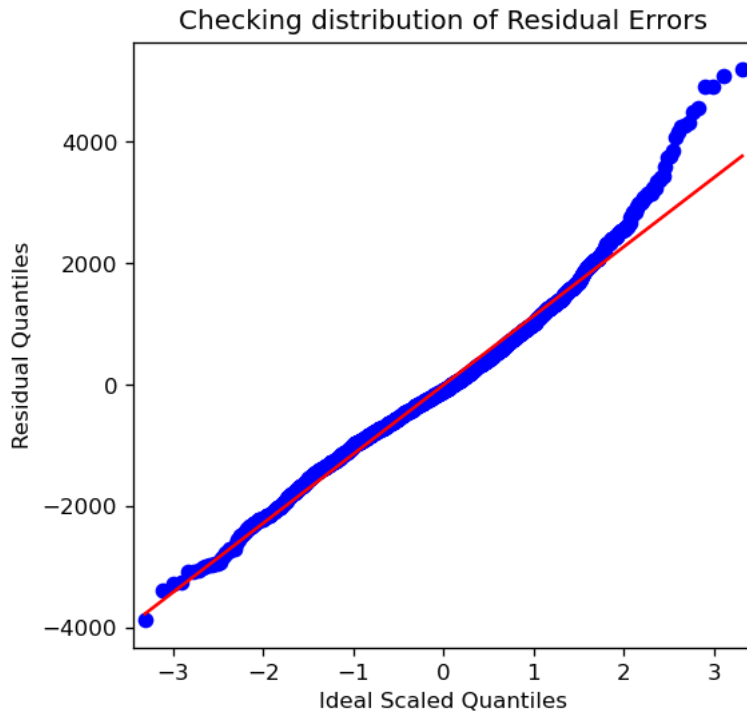
According to the Histogram, the distribution of error is nearly normal, But there are some outliers on the Higher end of the errors which could be handled easily.

Another way to check normal distribution is Q-Q Plot

In [14]:

```
1 # importing the QQ-plot from the statsmodels
2 import numpy as np
3 from statsmodels.graphics.gofplots import qqplot
4
5 ## Plotting the QQ plot
6 fig, ax = plt.subplots(figsize=(5,5) , dpi = 120)
7 qqplot(error_df.error, line = 's' , ax = ax)
8 plt.ylabel('Residual Quantiles')
9 plt.xlabel('Ideal Scaled Quantiles')
10 plt.title('Checking distribution of Residual Errors')
11 plt.show()
```

C:\Users\rkt7k\anaconda3\lib\site-packages\statsmodels\graphics\gofplots.py:993: UserWarning: marker is redundantly defined by the 'marker' keyword argument and the fmt string "bo" (-> marker='o'). The keyword argument will take precedence.  
ax.plot(x, y, fmt, \*\*plot\_style)



The QQ-plot clearly verifies our findings from the the histogram of the residuals, the data is mostly normal in nature, but there sre some outliers on the higher end of the Residues.

In [15]:

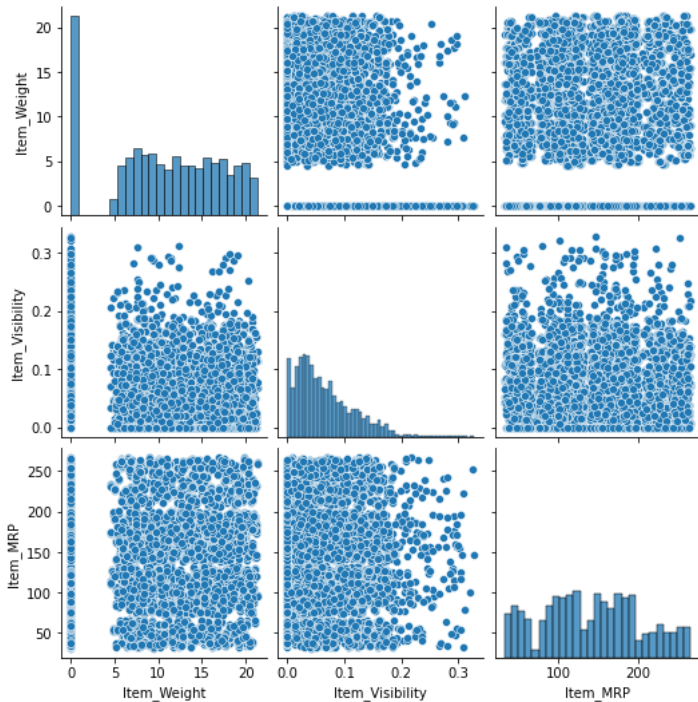
```

1 import seaborn as sns
2
3 sns.pairplot(df[['Item_Weight', 'Item_Visibility', 'Item_MRP']])

```

Out[15]:

&lt;seaborn.axisgrid.PairGrid at 0x1a7d0480700&gt;



Since there are lot of features and feature selection is not in the scope of this notebook, displayed pairplot only between some columns.

### Variance Inflation Factor (VIF) (Checking for multi collinearity)

In [25]:

```

1 # Importing Variance_inflation_Factor funtion from the Statsmodels
2 from statsmodels.stats.outliers_influence import variance_inflation_factor
3 from statsmodels.tools.tools import add_constant
4

```

```

1 # note : VIF by default does not calculate for the constant b in the linear reg equation : y = Bx+b
2     so we add it separately
3
4 Also, this new constang is TREATED AS COEFFICIENT OF A NEW FEATURE WHOSE VALUE IS 1 , so there is no effect of its value on data and
5 we can get a constant as well
6
7 => y = Bx + 1b

```

In [32]:

```

1 trainx_with_constant = add_constant(train_x.values)
2 trainx_with_constant[1]
3 # we can observe, one has beed added as the first value

```

Out[32]:

```

array([1.0000000e+00, 1.6000000e+01, 7.2655379e-02, 2.2986680e+02,
       2.0040000e+03, 0.0000000e+00, 1.0000000e+00, 0.0000000e+00,
       0.0000000e+00, 0.0000000e+00, 0.0000000e+00, 0.0000000e+00,
       0.0000000e+00, 0.0000000e+00, 0.0000000e+00, 0.0000000e+00,
       0.0000000e+00, 0.0000000e+00, 0.0000000e+00, 1.0000000e+00,
       0.0000000e+00, 0.0000000e+00, 0.0000000e+00, 0.0000000e+00,
       0.0000000e+00, 0.0000000e+00, 0.0000000e+00, 0.0000000e+00,
       0.0000000e+00, 0.0000000e+00, 0.0000000e+00, 0.0000000e+00,
       1.0000000e+00, 0.0000000e+00, 0.0000000e+00, 0.0000000e+00,
       0.0000000e+00, 0.0000000e+00, 1.0000000e+00, 0.0000000e+00,
       1.0000000e+00, 0.0000000e+00, 0.0000000e+00, 1.0000000e+00,
       0.0000000e+00, 0.0000000e+00])

```

In [34]:

```

1 # Calculating VIF for every column (only works for the not Catagorical)
2 VIF= pd.Series([variance_inflation_factor(trainx_with_constant, i) for i in range(1,trainx_with_constant.shape[1])],
3                 index =train_x.columns)
4 VIF

```

C:\Users\rkt7k\anaconda3\lib\site-packages\statsmodels\stats\outliers\_influence.py:193: RuntimeWarning: divide by zero encountered in double\_scalars  
 vif = 1. / (1. - r\_squared\_i)

Out[34]:

```

Item_Weight                2.335412
Item_Visibility            1.099468
Item_MRP                   1.015434
Outlet_Establishment_Year      inf
Item_Fat_Content_LF         inf
Item_Fat_Content_Low Fat    inf
Item_Fat_Content_Regular    inf
Item_Fat_Content_low fat    inf
Item_Fat_Content_reg        inf
Item_Type_Baking Goods      inf
Item_Type_Breads            inf
Item_Type_Breakfast         inf
Item_Type_Canned            inf
Item_Type_Dairy             inf
Item_Type_Frozen Foods      inf
Item_Type_Fruits and Vegetables inf
Item_Type_Hard Drinks       inf
Item_Type_Health and Hygiene inf
Item_Type_Household         inf
Item_Type_Meat              inf
Item_Type_Others            inf
Item_Type_Seafood           inf
Item_Type_Snack Foods       inf
Item_Type_Soft Drinks       inf
Item_Type_Starchy Foods     inf
Outlet_Identifier_OUT010     inf
Outlet_Identifier_OUT013     inf
Outlet_Identifier_OUT017     inf
Outlet_Identifier_OUT018     inf
Outlet_Identifier_OUT019     inf
Outlet_Identifier_OUT027     inf
Outlet_Identifier_OUT035     inf
Outlet_Identifier_OUT045     inf
Outlet_Identifier_OUT046     inf
Outlet_Identifier_OUT049     inf
Outlet_Size_High            inf
Outlet_Size_Medium          inf
Outlet_Size_Small           inf
Outlet_Location_Type_Tier 1  inf
Outlet_Location_Type_Tier 2  inf
Outlet_Location_Type_Tier 3  inf
Outlet_Type_Grocery Store    inf
Outlet_Type_Supermarket Type1 inf
Outlet_Type_Supermarket Type2 inf
Outlet_Type_Supermarket Type3 inf
dtype: float64

```

```

1 note :
2 vif = 1 / (1-R2)
3 so if R2 value is coming 1, the vif will be infinity
4
5 but why VIF here is infinity, if we see closely, INF is only for the ONE HOT ENCODED variables
6 so what happens is, when we encode all n categories, then if (n-1) categories is 0, it means that that last nth category is 1. so
  actually we dont need to encode all n categories and rather encode only (n-1) categories.
7
8 Here as all n categories are encoded, the (n-1) features could be easily used to predict that last nth feature and hence R2 for that
  is 1

```

## Model Inpretability

So far we have simply been predicting the values using the linear regression, But in order to Interpret the model, the normalising of the data is essential.

In [40]:

```

1 lr = LR(normalize=True)
2 lr.fit(train_x, train_y)

```

Out[40]:

LinearRegression(normalize=True)

In [41]:

```
1 train_pred = lr.predict(train_x)
2 train_mae = mae(train_pred, train_y)
3 print('Training Mean Absolute Error :', train_mae)
```

Training Mean Absolute Error : 850.426924968711

In [42]:

```
1 test_pred = lr.predict(test_x)
2 test_mae = mae(test_pred, test_y)
3 print('Training Mean Absolute Error :', test_mae)
```

Training Mean Absolute Error : 880.1328247770999

In [43]:

```
1 feature_importance_df = pd.DataFrame({
2     'Feature' : train_x.columns,
3     'Importance': lr.coef_
4 })
5
6 feature_importance_df['Importance_magnitude'] = feature_importance_df.Importance.abs()
7
8 feature_importance_df.sort_values(by='Importance_magnitude', inplace=True, ascending=False)
9 # sorting so we can plot only top n important features for better visibility
10
11 feature_importance_df.head()
```

Out[43]:

	Feature	Importance	Importance_magnitude
29	Outlet_Identifier_OUT019	-9.311582e+15	9.311582e+15
26	Outlet_Identifier_OUT013	-7.190359e+15	7.190359e+15
35	Outlet_Size_High	-6.300353e+15	6.300353e+15
44	Outlet_Type_Supermarket Type3	-6.171283e+15	6.171283e+15
43	Outlet_Type_Supermarket Type2	6.060654e+15	6.060654e+15



In [44]:

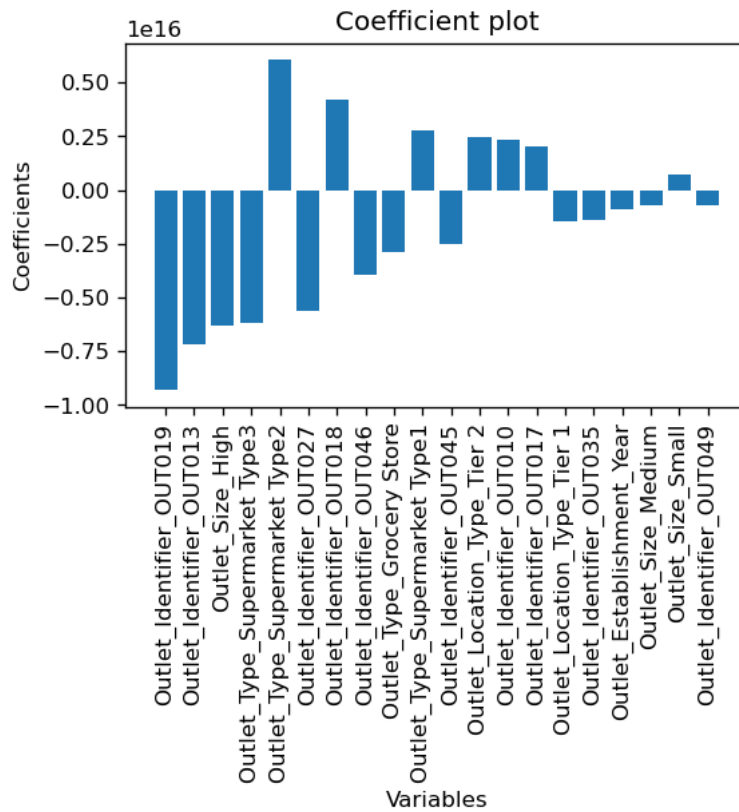
```

1 # plotting the coefficients
2
3
4 plt.figure(figsize=(5,3), dpi=120, facecolor='w', edgecolor='b')
5 x = feature_importance_df['Feature'].head(20)
6 y = feature_importance_df['Importance'].head(20)
7 plt.bar( x, y )
8 plt.xlabel( "Variables")
9 plt.ylabel('Coefficients')
10 plt.xticks(rotation=90)
11 plt.title('Coefficient plot')

```

Out[44]:

Text(0.5, 1.0, 'Coefficient plot')



Now the coefficients we see are normalised and we can easily make final inferences out of it.

In [ ]:

1