

```
In [1]: import pandas as pd

In [2]: import numpy as np

In [3]: df = pd.read_csv('https://raw.githubusercontent.com/Y81-Foundation/dataset/main/Big20Sales20data.csv')

In [4]: df.head()

Out[4]:
   Item_Identifier  Item_Weight  Item_Fat_Content  Item_Visibility  Item_Type  Item_MRP  Outlet_Identifier  Outlet_Establishment_Year  Outlet_Size  Outlet_Location_Type  Outlet_Type  Item_Outlet_Sales
0      FDT36         12.3          Low Fat         0.111448      Baking Goods   33.4874      OUT049              1999             Medium      Tier 1  Supermarket Type1      436.608721
1      FDT36         12.3          Low Fat         0.111904      Baking Goods   33.9874      OUT017              2007             Medium      Tier 2  Supermarket Type1      443.127721
2      FDT36         12.3           LF         0.111728      Baking Goods   33.9874      OUT018              2009             Medium      Tier 3  Supermarket Type2      564.598400
3      FDT36         12.3          Low Fat         0.000000      Baking Goods   34.3874      OUT019              1985             Small      Tier 1  Grocery Store      1719.370000
4      FDP12          9.8          Regular         0.045523      Baking Goods   35.0874      OUT017              2007             Medium      Tier 2  Supermarket Type1      352.874000

In [5]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14204 entries, 0 to 14203
Data columns (total 12 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Item_Identifier        14204 non-null object
 1   Item_Weight            14204 non-null float64
 2   Item_Fat_Content       14204 non-null object
 3   Item_Visibility        14204 non-null float64
 4   Item_Type              14204 non-null object
 5   Item_MRP               14204 non-null float64
 6   Outlet_Identifier      14204 non-null object
 7   Outlet_Establishment_Year 14204 non-null int64
 8   Outlet_Size            14204 non-null object
 9   Outlet_Location_Type   14204 non-null object
10   Outlet_Type            14204 non-null object
11   Item_Outlet_Sales      14204 non-null float64
dtypes: float64(4), int64(1), object(7)
memory usage: 1.3+ MB

In [6]: df.columns

Out[6]:
Index(['Item_Identifier', 'Item_Weight', 'Item_Fat_Content', 'Item_Visibility',
       'Item_Type', 'Item_MRP', 'Outlet_Identifier',
       'Outlet_Establishment_Year', 'Outlet_Size', 'Outlet_Location_Type',
       'Outlet_Type', 'Item_Outlet_Sales'],
      dtype='object')

In [7]: df.describe()

Out[7]:
   Item_Weight  Item_Visibility  Item_MRP  Outlet_Establishment_Year  Item_Outlet_Sales
count  11815.000000    14204.000000    14204.000000    14204.000000    14204.000000
mean      12.789355     0.069593    141.004977    1997.830681     2185.836320
std       4.654126     0.051459     62.086938     8.371664     1827.479550
min       4.555000     0.000000     31.290000    1985.000000     33.290000
25%      8.710000     0.027036    94.012000    1987.000000     922.135101
50%     12.500000     0.054021   142.247000    1999.000000    1768.287600
75%     16.750000     0.094037   185.855600    2004.000000    2988.119400
max     30.000000     0.328391   266.888400    2009.000000   31224.726950

In [8]: df['Item_Weight'].fillna(df.groupby(['Item_Type'])['Item_Weight'].transform('mean'), inplace=True)

In [9]: df.info()


<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14204 entries, 0 to 14203
Data columns (total 12 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Item_Identifier        14204 non-null object
 1   Item_Weight            14204 non-null float64
 2   Item_Fat_Content       14204 non-null object
 3   Item_Visibility        14204 non-null float64
 4   Item_Type              14204 non-null object
 5   Item_MRP               14204 non-null float64
 6   Outlet_Identifier      14204 non-null object
 7   Outlet_Establishment_Year 14204 non-null int64
 8   Outlet_Size            14204 non-null object
 9   Outlet_Location_Type   14204 non-null object
10   Outlet_Type            14204 non-null object
11   Item_Outlet_Sales      14204 non-null float64
dtypes: float64(4), int64(1), object(7)
memory usage: 1.3+ MB

In [10]: df.describe()

Out[10]:
   Item_Weight  Item_Visibility  Item_MRP  Outlet_Establishment_Year  Item_Outlet_Sales
count  14204.000000    14204.000000    14204.000000    14204.000000    14204.000000
mean     12.790542     0.069593    141.004977    1997.830681     2185.836320
std       4.251186     0.051459     62.086938     8.371664     1827.479550
min       4.555000     0.000000     31.290000    1985.000000     33.290000
25%      9.300000     0.027036    94.012000    1987.000000     922.135101
50%     12.800000     0.054021   142.247000    1999.000000    1768.287600
75%     16.000000     0.094037   185.855600    2004.000000    2988.119400
max     30.000000     0.328391   266.888400    2009.000000   31224.726950

In [11]: import seaborn as sns
sns.pairplot(df)

<seaborn.axisgrid.PairGrid at 0x1c0a56539>

Out[11]:


In [12]: df[['Item_Identifier']].value_counts()

Out[12]:
Item_Identifier
FDD08      18
FDD24      18
FDD19      18
FDD25      18
FDD21      18
FDD52       1
FDD50       7
FDD56       7
FDD10       7
FDD31       7
Length: 1559, dtype: int64

In [13]: df[['Item_Fat_Content']].value_counts()

Out[13]:
Item_Fat_Content
Low Fat      8485
Regular      4824
LF           522
reg          195
Low Fat      178
dtype: int64

In [14]: df.replace({'Item_Fat_Content': {'LF': 'Low Fat', 'reg': 'Regular', 'Low Fat': 'Low Fat'}}, inplace=True)

In [15]: df[['Item_Fat_Content']].value_counts()

Out[15]:
Item_Fat_Content
Low Fat      9185
Regular      5019
dtype: int64

In [16]: df.replace({'Item_Fat_Content': {'Low Fat': 0, 'Regular': 1}}, inplace=True)

In [17]: df[['Item_Type']].value_counts()

Out[17]:
Item_Type
Fruits and Vegetables  2813
Snack Foods           2809
Household             1548
Frozen Foods          1426
Dairy                1136
Baking Goods          1886
Canned               1884
Health and Hygiene    858
Meat                 736
Soft Drinks          728
Breads               416
Hard Drinks          382
Others               288
Starchy Foods        259
Breakfast            186
Seafood              89
dtype: int64

In [18]: df.replace({'Item_Type': {'Fruits and Vegetables': 0, 'Snack Foods': 0, 'Household': 1,
                                   'Frozen Foods': 0, 'Dairy': 0, 'Baking Goods': 0,
                                   'Canned': 0, 'Health and Hygiene': 1,
                                   'Meat': 0, 'Soft Drinks': 0,
                                   'Breads': 0, 'Hard Drinks': 0,
                                   'Others': 2, 'Starchy Foods': 0, 'Breakfast': 0, 'Seafood': 0 }}, inplace=True)

In [19]: df[['Item_Type']].value_counts()

Out[19]:
Item_Type
0      11518
1       2488
2        280
dtype: int64

In [20]: df[['Outlet_Identifier']].value_counts()

Out[20]:
Outlet_Identifier
OUT027      1559
OUT013      1553
OUT035      1558
OUT045      1558
OUT049      1559
OUT045      1548
OUT018      1546
OUT017      1543
OUT019      925
OUT019      880
dtype: int64

In [21]: df.replace({'Outlet_Identifier': {'OUT027': 0, 'OUT013': 1,
                                           'OUT049': 2, 'OUT040': 3, 'OUT035': 4,
                                           'OUT045': 5, 'OUT018': 6,
                                           'OUT017': 7, 'OUT019': 8, 'OUT019': 9
                                           }}, inplace=True)

In [22]: df[['Outlet_Identifier']].value_counts()

Out[22]:
Outlet_Identifier
0      1559
1      1553
2      1558
3      1559
4      1558
5      1548
6      1548
7      1543
8      925
9      880
dtype: int64

In [23]: df[['Outlet_Size']].value_counts()

Out[23]:
Outlet_Size
Medium      7122
Small       5529
High        1553
dtype: int64

In [24]: df.replace({'Outlet_Size': {'Small': 0, 'Medium': 1, 'High': 2}}, inplace=True)

In [25]: df[['Outlet_Size']].value_counts()

Out[25]:
Outlet_Size
1      7122
0      5529
2      1553
dtype: int64

In [26]: df[['Outlet_Location_Type']].value_counts()

Out[26]:
Outlet_Location_Type
Tier 3      5583
Tier 2      4641
Tier 1      3980
dtype: int64

In [27]: df.replace({'Outlet_Location_Type': {'Tier 1': 0, 'Tier 2': 1, 'Tier 3': 2}}, inplace=True)

In [28]: df[['Outlet_Location_Type']].value_counts()

Out[28]:
Outlet_Location_Type
2      5583
1      4641
0      3980
dtype: int64

In [29]: df[['Outlet_Type']].value_counts()

Out[29]:
Outlet_Type
Supermarket Type1  9294
Grocery Store     1805
Supermarket Type3  1559
Supermarket Type2  1546
dtype: int64

In [30]: df.replace({'Outlet_Type': {'Grocery Store': 0, 'Supermarket Type1': 1, 'Supermarket Type2': 2, 'Supermarket Type3': 3}}, inplace=True)

In [31]: df[['Outlet_Type']].value_counts()

Out[31]:
Outlet_Type
1      9294
0      1805
3      1559
2      1546
dtype: int64

In [32]: df.head()

Out[32]:
   Item_Identifier  Item_Weight  Item_Fat_Content  Item_Visibility  Item_Type  Item_MRP  Outlet_Identifier  Outlet_Establishment_Year  Outlet_Size  Outlet_Location_Type  Outlet_Type  Item_Outlet_Sales
0      FDT36         12.3          0         0.111448      0      33.4874      2      1999      1      0      1      436.608721
1      FDT36         12.3          0         0.111904      0      33.9874      7      2007      1      1      1      443.127721
2      FDT36         12.3          0         0.111728      0      33.9874      6      2009      1      2      2      564.598400
3      FDT36         12.3          0         0.000000      0      34.3874      9      1985      0      0      0      1719.370000
4      FDP12          9.8          1         0.045523      0      35.0874      7      2007      1      1      1      352.874000

In [33]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14204 entries, 0 to 14203
Data columns (total 12 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Item_Identifier        14204 non-null object
 1   Item_Weight            14204 non-null float64
 2   Item_Fat_Content       14204 non-null int64
 3   Item_Visibility        14204 non-null float64
 4   Item_Type              14204 non-null int64
 5   Item_MRP               14204 non-null float64
 6   Outlet_Identifier      14204 non-null int64
 7   Outlet_Establishment_Year 14204 non-null int64
 8   Outlet_Size            14204 non-null int64
 9   Outlet_Location_Type   14204 non-null int64
10   Outlet_Type            14204 non-null int64
11   Item_Outlet_Sales      14204 non-null float64
dtypes: float64(4), int64(7), object(1)
memory usage: 1.3+ MB

In [34]: y=df['Item_Outlet_Sales']

In [35]: y

Out[35]:
0      436.608721
1      443.127721
2      564.598400
3      1719.370000
4      352.874000
...
14199  4984.178800
14200  2885.577200
14201  2885.577200
14202  3863.676454
14203  3644.354765
Name: Item_Outlet_Sales, Length: 14204, dtype: float64

In [36]: X=df[['Item_Weight', 'Item_Fat_Content', 'Item_Visibility',
               'Item_Type', 'Item_MRP', 'Outlet_Identifier',
               'Outlet_Establishment_Year', 'Outlet_Size', 'Outlet_Location_Type',
               'Outlet_Type']]

In [37]: X

Out[37]:
   Item_Weight  Item_Fat_Content  Item_Visibility  Item_Type  Item_MRP  Outlet_Identifier  Outlet_Establishment_Year  Outlet_Size  Outlet_Location_Type  Outlet_Type
0      12.300000      0         0.111448      0      33.4874      2      1999      1      0      1
1      12.300000      0         0.111904      0      33.9874      7      2007      1      1      1
2      12.300000      0         0.111728      0      33.9874      6      2009      1      2      2
3      12.300000      0         0.000000      0      34.3874      9      1985      0      0      0
4      9.800000      1         0.045523      0      35.0874      7      2007      1      1      1
...
14199  12.800000      0         0.069006      0      261.9252      4      2004      0      1      1
14200  12.800000      0         0.070013      0      262.8252      7      2007      1      1      1
14201  12.800000      0         0.069561      0      263.0252      1      1987      2      2      1
14202  12.697578      0         0.069282      0      263.0252      0      1985      1      2      3
14203  12.800000      0         0.069727      0      263.0252      2      1999      1      0      1
14204 rows x 10 columns

In [38]: from sklearn.preprocessing import StandardScaler

In [39]: sc = StandardScaler()

In [40]: X_std = sc.fit_transform(X_std)

In [41]: X_std = sc.fit_transform(X_std)

In [42]: X_std

Out[42]:
array([[ -0.11541705,  0.88413635, -1.73278716,  0.13968866],
       [ -0.11541705,  0.89308616, -1.72773366,  1.09531886],
       [ -0.11541705,  0.88998031, -1.72773366,  1.3342284 ],
       ...,
       [  0.90220132,  0.07611952,  1.96538148, -1.28377659],
       [  0.20444792,  0.06469366,  1.97343459, -1.5326861 ],
       [  0.90220132,  0.07334891,  1.97504569,  0.13968866]])

In [43]: X[['Item_Weight', 'Item_Visibility', 'Item_MRP', 'Outlet_Establishment_Year']] = pd.DataFrame(X_std, columns=['Item_Weight', 'Item_Visibility', 'Item_MRP', 'Outlet_Establishment_Year'])

C:\Users\shale\AppData\Local\Temp\ipykernel_24632\3276057488.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
X[['Item_Weight', 'Item_Visibility', 'Item_MRP', 'Outlet_Establishment_Year']] = pd.DataFrame(X_std, columns=['Item_Weight', 'Item_Visibility', 'Item_MRP', 'Outlet_Establishment_Year'])

In [47]: X

Out[47]:
   Item_Weight  Item_Fat_Content  Item_Visibility  Item_Type  Item_MRP  Outlet_Identifier  Outlet_Establishment_Year  Outlet_Size  Outlet_Location_Type  Outlet_Type
0      12.300000      0         0.111448      0      33.4874      2      1999      1      0      1
1      12.300000      0         0.894306      0      33.9874      7      2007      1      1      1
2      12.300000      0         0.889893      0      33.9874      6      2009      1      2      2
3      12.300000      0      -1.281712      0      34.3874      9      1985      0      0      0
4      9.703509      1      -0.397031      0      35.0874      7      2007      1      1      1
...
14199  12.800000      0         0.069006      0      261.9252      4      2004      0      1      1
14200  0.002201      0         0.070898      0      1.842664      7      2007      1      1      1
14201  0.002201      0         0.070120      0      1.965381      1      1987      2      2      1
14202  0.204448      0         0.064694      0      1.973435      0      1985      1      2      3
14203  0.002201      0         0.073349      0      1.975046      2      0199081      1      0      1
14204 rows x 10 columns

In [48]: from sklearn.model_selection import train_test_split

In [49]: X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.1, random_state=22529)

In [50]: from sklearn.ensemble import RandomForestRegressor

In [51]: rfr = RandomForestRegressor(random_state=22529)

In [52]: rfr.fit(X_train, y_train)

Out[52]:
RandomForestRegressor
RandomForestRegressor(random_state=22529)

In [53]: y_pred = rfr.predict(X_test)

In [54]: y_pred

Out[54]:
array([[ 842.17961988,  989.97031772, 2975.19528623, ...,  818.18826824,
        2954.7788205 , 1676.8202906]])

In [55]: from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score

In [56]: mean_squared_error(y_test, y_pred)

Out[56]:
1779270.7898065657

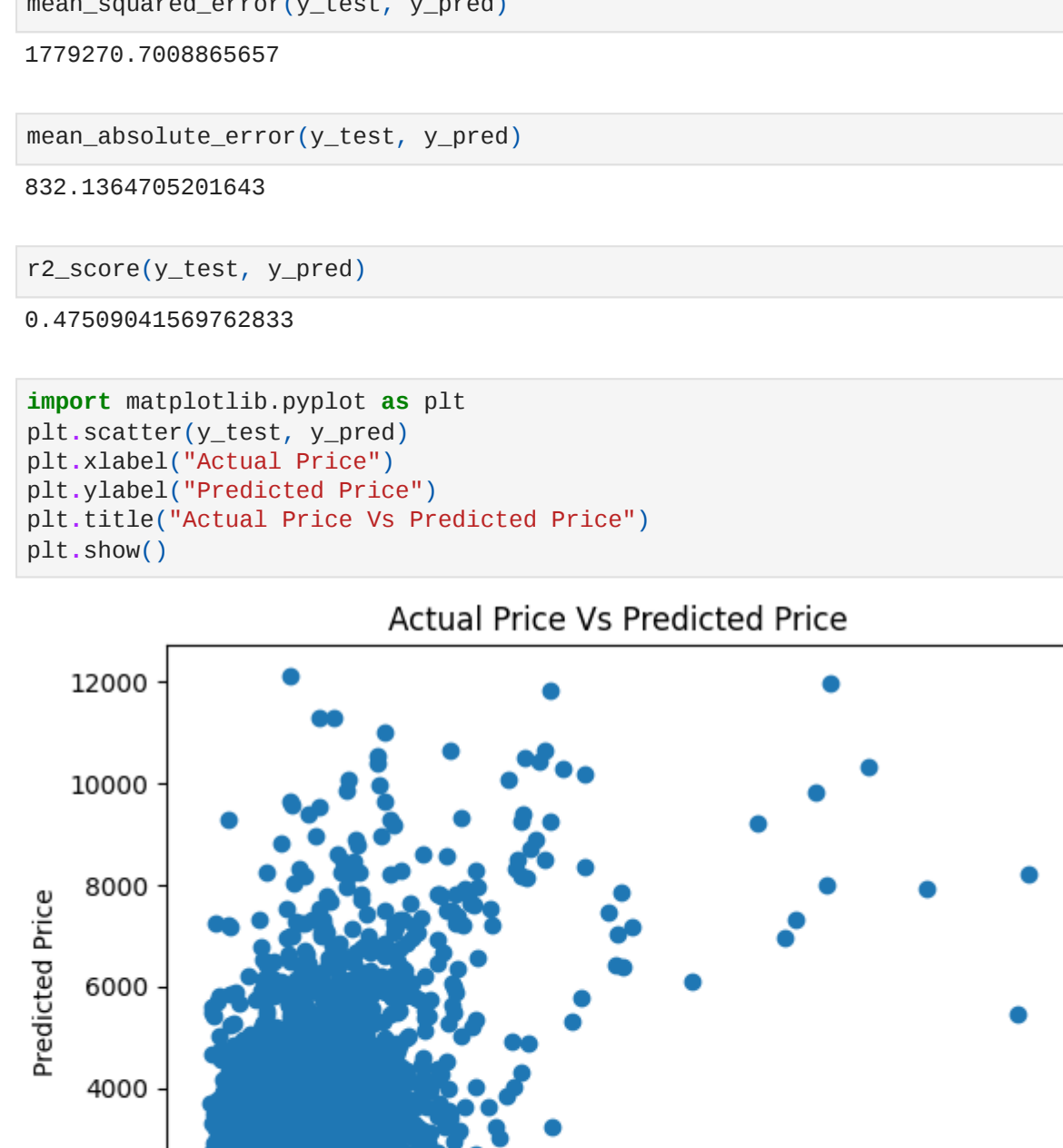
In [57]: mean_absolute_error(y_test, y_pred)

Out[57]:
832.1564785261643

In [58]: r2_score(y_test, y_pred)

Out[58]:
0.47589941569762833

In [59]: import matplotlib.pyplot as plt
plt.scatter(y_test, y_pred)
plt.xlabel("Actual Price")
plt.ylabel("Predicted Price")
plt.title("Actual Price Vs Predicted Price")
plt.show()



In [ ]:
```