

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

In [2]: import numpy as np

In [3]: df= pd.read_csv('https://raw.githubusercontent.com/Lorddhaval/Dataset/main/B1gh28Sales%20Data.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          11815 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.describe()

Out[6]:
   Item_Weight  Item_Visibility  Item_MRP  Outlet_Establishment_Year  Item_Outlet_Sales
count  11815.000000    14204.000000    14204.000000    14204.000000    14204.000000
mean      12.788355         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        1768.287880
75%      16.750000         0.094037   185.855600        2004.000000        2988.110400
max      30.000000         0.328391   266.888400        2009.000000        31224.726950

In [7]: df.columns

Out[7]:
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 [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      12.790642         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        1768.287880
75%      16.000000         0.094037   185.855600        2004.000000        2988.110400
max      30.000000         0.328391   266.888400        2009.000000        31224.726950

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

Out[11]:
<seaborn.axisgrid.PairGrid at 0x1a5df3cf5b8>

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

Out[12]:
Item_Identifier
FDD08      10
FDD24      10
FDD19      10
FDD28      10
FDD31      10
FDM52       7
FDM50       7
FDD50       7
FDR51       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     5619
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      2013
Snack Foods                1989
Household                  1548
Frozen Foods               1426
Dairy                      1136
Baking Goods               1086
Canned                     1084
Health and Hygiene         858
Meat                       736
Soft Drinks                726
Breads                     416
Hard Drinks                362
Others                     288
Starchy Foods             269
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       2408
2        288
dtype: int64

In [20]: df[['Item_Type']].info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14204 entries, 0 to 14203
Data columns (total 1 columns):
#   Column              Non-Null Count  Dtype
---  ---
0   Item_Type            14204 non-null  int64
dtypes: int64(1)
memory usage: 111.1 KB

In [21]: df[['Outlet_Identifier']].value_counts()

Out[21]:
Outlet_Identifier
OUT827      1559
OUT813      1553
OUT835      1550
OUT846      1550
OUT849      1550
OUT845      1548
OUT818      1546
OUT827      1543
OUT810      925
OUT819      880
dtype: int64

In [22]: df.replace({'Outlet_Identifier': { 'OUT827': 0, 'OUT813': 1,
'OUT849': 2, 'OUT846': 3, 'OUT835': 4,
'OUT845':5, 'OUT818' : 6,
'OUT817': 7, 'OUT810': 8, 'OUT819': 9,
}},inplace=True)

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

Out[23]:
Outlet_Identifier
0      1559
1      1553
2      1550
3      1550
4      1550
5      1548
6      1546
7      1543
8       925
9       880
dtype: int64

In [24]: df[['Outlet_Size']].value_counts()

Out[24]:
Outlet_Size
Medium      7122
Small       5529
High        1553
dtype: int64

In [25]: df.replace({'Outlet_Size':{'Small':0, 'Medium':1, 'High':2}},inplace=True)

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

Out[26]:
Outlet_Location_Type
Tier 3       5583
Tier 2       4641
Tier 1       3989
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      3989
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]: df.shape

Out[34]:
(14204, 12)

In [35]: y = df['Item_Outlet_Sales']

In [36]: y.shape

Out[36]:
(14204,)

In [37]: Y

Out[37]:
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    3893.676404
14203    3644.354765
Name: Item_Outlet_Sales, Length: 14204, dtype: float64

In [38]: 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 [39]: X = df.drop(['Item_Identifier','Item_Outlet_Sales'],axis = 1)

In [40]: X.shape

Out[40]:
(14204, 10)

In [41]: X

Out[41]:
   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      0.069606         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              1985              2              2              1
14202    13.659758              0      0.069262         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 [42]: from sklearn.preprocessing import StandardScaler

In [43]: sc = StandardScaler()

In [44]: X_std = df[['Item_Weight','Item_Visibility','Item_MRP','Outlet_Establishment_Year']]

In [45]: X_std = sc.fit_transform(X_std)

Out[46]:
array([[ -0.11541705,  0.88413635, -1.73178716,  0.13968968],
       [ -0.11541705,  0.89300616, -1.72373366,  1.09531886],
       [ -0.11541705,  0.88958331, -1.72373366,  1.3342284 ],
       ...,
       [  0.88228132,  0.07011952,  1.96538148, -1.29377659],
       [  0.26444792,  0.06469365,  1.97343489, -1.53286614],
       [  0.88228132,  0.07334891,  1.97504569,  0.13968968]])

In [47]: 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 [48]: X

Out[48]:
   Item_Weight  Item_Fat_Content  Item_Visibility  Item_Type  Item_MRP  Outlet_Identifier  Outlet_Establishment_Year  Outlet_Size  Outlet_Location_Type  Outlet_Type
0   -0.115417              0      0.884136         0      33.4874              2              1999              1              0              1
1   -0.115417              0      0.893006         0      33.9874              7              2007              1              1              1
2   -0.115417              0      0.889583         0      33.9874              6              2009              1              2              2
3   -0.115417              0      0.000000         0      34.3874              9              1985              0              0              0
4   -0.703509              1   -0.397031         0      35.0874              7              2007              1              1              1
...
14199    0.002201              0      0.070990         0      1.947664              4              0.736955              0              1              1
14200    0.002201              0      0.078988         0      1.962160              7              1.095319              1              1              1
14201    0.002201              0      0.070120         0      1.965381              1             -1.293777              2              2              1
14202    0.044448              0      0.064694         0      1.973435              0             -1.532866              1              2              3
14203    0.002201              0      0.073349         0      1.975046              2              0.139681              1              0              1
14204 rows x 10 columns

In [49]: from sklearn.model_selection import train_test_split

In [50]: X_train, X_test, y_train, y_test = train_test_split(X,y , test_size =0.1, random_state= 72529 )

In [51]: X_train.shape, y_train.shape, y_train.shape, y_test.shape

Out[51]:
((12783, 10), (1421, 10), (12783, ), (1421, ))

In [52]: from sklearn.ensemble import RandomForestRegressor

In [53]: rfr = RandomForestRegressor(random_state=72529)

In [54]: rfr.fit(X_train, y_train)

Out[54]:
RandomForestRegressor
RandomForestRegressor(random_state=72529)

In [55]: y_pred = rfr.predict(X_test)

In [56]: y_pred.shape

Out[56]:
(1421,)

In [57]: y_pred

Out[57]:
array([1515.89681418, 1796.47737123, 1896.58147176, ..., 3678.89298551,
       745.11959868, 3592.32874773])

In [58]: from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score

In [59]: mean_squared_error(y_test, y_pred)

Out[59]:
1480354.1533997788

In [60]: mean_absolute_error(y_test, y_pred)

Out[60]:
810.2512914488754

In [61]: r2_score(y_test, y_pred)

Out[61]:
0.5950364101543593

In [62]: 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()

Actual Price Vs Predicted Price

In [ ]:
```