

```
# Importing libraries
import pandas as pd
from sklearn.preprocessing import LabelEncoder, OneHotEncoder
from sklearn.metrics import r2_score, mean_squared_error
from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score
from sklearn.model_selection import train_test_split
from sklearn.metrics import r2_score, mean_squared_error
```

In [2]:
train=pd.read_csv("bigdatamart_train.csv")
test=pd.read_csv("bigdatamart_test.csv")

In [3]:
train.head()

	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	FDAIS1	9.30	Low Fat	0.010047	Dairy	249.8092	OUT049	1999	Medium	Tier 1	Supermarket Type1	3735.1390
1	DRCOL1	5.92	Regular	0.010278	Soft Drinks	48.2692	OUT018	2009	Medium	Tier 3	Supermarket Type2	443.4228
2	FDAIS1	17.50	Low Fat	0.016790	Meat	141.6180	OUT049	1999	Medium	Tier 1	Supermarket Type1	2097.2700
3	FDAIS1	19.20	Regular	0.000000	Fruits and Vegetables	182.0950	OUT010	1998	NAN	Tier 3	Grocery Store	732.3800
4	KCDI9	8.93	Low Fat	0.000000	Household	53.8614	OUT013	1987	High	Tier 3	Supermarket Type1	994.7052

In [4]:
train.shape
(6223, 13)

In [5]:
test.head()

	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	FDAIS6	20.750	Low Fat	0.007065	Snack Foods	107.8022	OUT049	1999	Medium	Tier 1	Supermarket Type1	
1	FDAIS4	8.300	Reg	0.038428	Others	87.3198	OUT017	2007	NAN	Tier 2	Supermarket Type2	
2	NCHV5	14.600	Low Fat	0.099575	Others	241.7538	OUT010	1998	NAN	Tier 3	Grocery Store	
3	FQDF9	7.315	Low Fat	0.015388	Snack Foods	135.0340	OUT017	2007	NAN	Tier 2	Supermarket Type1	
4	PQDS1	NAN	Regular	0.115999	Dairy	234.2300	OUT027	1995	Medium	Tier 3	Supermarket Type1	

In [6]:
test.shape
(5681, 13)

In [7]:
train.dtypes, test.dtypes

```
Item_Identifier    object
Item_Weight       float64
Item_Fat_Content   object
Item_Visibility    float64
Item_Type         object
Item_MRP          float64
Outlet_Identifier  object
Outlet_Establishment_Year    int64
Outlet_Size        object
Outlet_Location_Type    object
Outlet_Type         object
Item_Outlet_Sales    float64
dtype: object
```

In [9]:
train['source']='train'
test['source']='test'
both = pd.concat([train, test], ignore_index=True)
both.head()

	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	source
0	FDAIS1	9.30	Low Fat	0.010047	Dairy	249.8092	OUT049	1999	Medium	Tier 1	Supermarket Type1	3735.1390	train
1	DRCOL1	5.92	Regular	0.010278	Soft Drinks	48.2692	OUT018	2009	Medium	Tier 3	Supermarket Type2	443.4228	train
2	FDAIS1	17.50	Low Fat	0.016790	Meat	141.6180	OUT049	1999	Medium	Tier 1	Supermarket Type1	2097.2700	train
3	FDAIS1	19.20	Regular	0.000000	Fruits and Vegetables	182.0950	OUT010	1998	NAN	Tier 3	Grocery Store	732.3800	train
4	KCDI9	8.93	Low Fat	0.000000	Household	53.8614	OUT013	1987	High	Tier 3	Supermarket Type1	994.7052	train

In [10]:
both.describe()

	Item_Weight	Item_Visibility	Item_MRP	Outlet_Establishment_Year	Item_Outlet_Sales
count	11765.000000	14204.000000	14204.000000	14204.000000	8623.000000
mean	12.730664	0.066963	142.000977	1997.890661	2383.288814
std	4.662502	0.051499	62.069938	8.371664	1706.499168
min	4.555000	0.000000	31.290000	1986.000000	33.290000
25%	8.710000	0.027038	94.012000	1987.000000	634.247400
50%	12.000000	0.054021	142.247000	1989.000000	1794.231000
75%	16.350000	0.094037	180.000000	2004.000000	3111.296400
max	21.300000	0.328393	260.880000	2009.000000	13206.954800

In [11]:
sns.scatterplot(train["Item_Visibility"], train["Item_Outlet_Sales"])

C:\ProgramData\Anaconda3\lib\site-packages\matplotlib\decoders.py:361 FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be "data", and passing other arguments without an explicit keyword will result in an error or misinterpretation.

In [16]:
#reference - values having the most sales have visibility less than 0.2
fig, ax = plt.subplots(figsize=(10, 1))
plt.bar(train["Item_Type"].train["Item_Outlet_Sales"], width=0.5, color='pink')
plt.show()

In [21]:
#plot of sales vs type of product
fig, ax = plt.subplots(figsize=(10, 1))
plt.bar(train["Item_Type"].train["Item_Outlet_Sales"], width=0.5, color='grey')
plt.show()

In [24]:
data = pd.concat([train["Item_Type"], train["Item_Outlet_Sales"]], axis=1)
f, ax = plt.subplots(figsize=(10, 6))
fig = sns.boxplot(x="Item_Type", y="Item_Outlet_Sales", data=data)
fig.set_ylim(0, ymax=20000)
plt.show()

In [24]:
#BOXPLOT ANALYSIS
sns.boxplot(x="Outlet_Size", y="Item_Outlet_Sales", data=combi)
plt.show()

In [25]:
sns.boxplot(x="Outlet_Location_Type", y="Item_Outlet_Sales", data=combi)
plt