

In (84): `import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns`

In (85): `train=pd.read_csv('Train.csv')
test=pd.read_csv('Test.csv')`

In (86): `train.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8523 entries, 0 to 8522
Data columns (total 12 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   Item_Identifier        8523 non-null   object
 1   Item_Weight            7060 non-null   float64
 2   Item_Fat_Content       8523 non-null   object
 3   Item_Visibility        8523 non-null   float64
 4   Item_Type              8523 non-null   object
 5   Item_MRP               8523 non-null   float64
 6   Outlet_Identifier      8523 non-null   object
 7   Outlet_Establishment_Year 8523 non-null   int64
 8   Outlet_Size            6113 non-null   object
 9   Outlet_Location_Type   8523 non-null   object
10   Outlet_Type            8523 non-null   object
11   Item_Outlet_Sales      8523 non-null   float64
dtypes: float64(4), int64(1), object(7)
memory usage: 799.2+ KB
```

In (87): `train.shape`

Out (87): (8523, 12)

In (88): `train.head()`

Out (88):

	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
0	FDA15	9.30	Low Fat	0.016047	Dairy	249.8092	OUT049	1999	Medium	Tier 1	Supermarket Type1
1	DRC01	5.92	Regular	0.019278	Soft Drinks	48.2692	OUT018	2009	Small	Tier 3	Supermarket Type2
2	FDM15	17.50	Low Fat	0.016780	Meat	141.6180	OUT049	1999	Medium	Tier 1	Supermarket Type1
3	FDX07	19.20	Regular	0.000000	Fruits and Vegetables	182.0950	OUT010	1998	High	Tier 3	Supermarket Type1
4	NCD19	8.93	Low Fat	0.000000	Household	53.8614	OUT013	1987	Medium	Tier 2	Supermarket Type1

**Missing Value**

In (89): `train.isna().sum().sort_values(ascending=False)`

Out (89):

Outlet_Size	2410
Item_Weight	1463
Item_Outlet_Sales	0
Outlet_Type	0
Outlet_Location_Type	0
Outlet_Establishment_Year	0
Outlet_Identifier	0
Item_MRP	0
Item_Type	0
Item_Visibility	0
Item_Fat_Content	0
Item_Identifier	0
dtype:	int64

**Item\_Weight**

In (90): `train.loc[:,('Item_Weight','Item_Fat_Content','Item_Type')]`

Out (90):

	Item_Weight	Item_Fat_Content	Item_Type
0	9.300	Low Fat	Dairy
1	5.920	Regular	Soft Drinks
2	17.500	Low Fat	Meat
3	19.200	Regular	Fruits and Vegetables
4	8.930	Low Fat	Household
...	...	...	...
8518	6.865	Low Fat	Snack Foods
8519	8.385	Regular	Baking Goods
8520	10.600	Low Fat	Health and Hygiene
8521	7.210	Regular	Snack Foods
8522	14.800	Low Fat	Soft Drinks

8523 rows x 3 columns

In (91): `train.Item_Type.unique()`

Out (91): `array(['Dairy', 'Soft Drinks', 'Meat', 'Fruits and Vegetables', 'Household', 'Baking Goods', 'Snack Foods', 'Frozen Foods', 'Breakfast', 'Health and Hygiene', 'Hard Drinks', 'Canned', 'Breads', 'Starchy Foods', 'Others', 'Seafood'], dtype=object)`

In (92): `med=train.groupby(['Item_Fat_Content','Item_Type'])['Item_Weight'].transform('median')
med=test.groupby(['Item_Fat_Content','Item_Type'])['Item_Weight'].transform('median')`

In (93): `train.Item_Weight=np.where(train.Item_Weight.isnull(), med, train.Item_Weight)`

In (94): `test.Item_Weight=np.where(test.Item_Weight.isnull(), med, test.Item_Weight)`

In (95): `train.info()`

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8523 entries, 0 to 8522
Data columns (total 12 columns):
 # Column Non-Null Count Dtype
--- -
 0 Item\_Identifier 8523 non-null object
 1 Item\_Weight 8523 non-null float64
 2 Item\_Fat\_Content 8523 non-null float64
 3 Item\_Visibility 8523 non-null float64
 4 Item\_Type 8523 non-null object
 5 Item\_MRP 8523 non-null object
 6 Outlet\_Identifier 8523 non-null object
 7 Outlet\_Establishment\_Year 8523 non-null int64
 8 Outlet\_Size 8523 non-null object
 9 Outlet\_Location\_Type 8523 non-null object
10 Outlet\_Type 8523 non-null object
11 Item\_Outlet\_Sales 8523 non-null float64
dtypes: float64(4), int64(1), object(7)
memory usage: 799.2+ KB

Outlet\_Size

	Outlet_Size	Outlet_Location_Type	Outlet_Type
10	Medium	Tier 1	Supermarket Type1
11	Small	Tier 1	Supermarket Type1
12	Medium	Tier 1	Supermarket Type1
13	Small	Tier 1	Supermarket Type1
14	High	Tier 3	Supermarket Type1
...	...	...	...
8518	High	Tier 3	Supermarket Type1
8519	NaN	Tier 2	Supermarket Type1
8520	Small	Tier 3	Supermarket Type1
8521	Medium	Tier 2	Supermarket Type2
8522	Small	Tier 1	Supermarket Type1

8513 rows x 3 columns

In (97): `Na=train.loc[:,('Outlet_Size','Outlet_Location_Type','Outlet_Type')][10:]
Na`

Out (97):

	Outlet_Size	Outlet_Location_Type	Outlet_Type
25	NaN	Tier 2	Supermarket Type1
28	NaN	Tier 3	Grocery Store
30	NaN	Tier 3	Grocery Store
33	NaN	Tier 2	Supermarket Type1
45	NaN	Tier 3	Grocery Store
...	...	...	...
8508	NaN	Tier 2	Supermarket Type1
8509	NaN	Tier 3	Grocery Store
8514	NaN	Tier 2	Supermarket Type1
8519	NaN	Tier 2	Supermarket Type1

2407 rows x 3 columns

In (98): `train.Outlet_Type.unique()`

Out (98): `array(['Supermarket Type1', 'Supermarket Type2', 'Grocery Store', 'Supermarket Type3'], dtype=object)`

In (99): `filter1=train['Outlet_Type']=='Grocery Store'
filter2=train['Outlet_Location_Type']=='Tier 1'
train.where(filter1 & filter2).dropna() # we can clearly see that there is some relationship between Outlet`

Out (99):

	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
23	FDC37	12.300	Low Fat	0.057557	Baking Goods	107.8938	OUT019	1985.0	Medium	Tier 1	Supermarket Type1
29	FDC14	13.350	Regular	0.072222	Canned	43.6454	OUT019	1985.0	Medium	Tier 1	Supermarket Type1
49	FDS02	12.700	Regular	0.255395	Canned	198.8794	OUT019	1985.0	Medium	Tier 1	Supermarket Type1
59	FDI26	10.195	Low Fat	0.061082	Dairy	180.0341	OUT019	1985.0	Medium	Tier 1	Supermarket Type1
63	FDY40	13.350	Regular	0.150286	Frozen Foods	51.0692	OUT019	1985.0	Medium	Tier 1	Supermarket Type1
...	...	...	...	...	...	...	...	...	...	...	...
8454	NCH54	13.150	Low Fat	0.127234	Household	158.3920	OUT019	1985.0	Medium	Tier 1	Supermarket Type1
8458	FDX20	12.500	Low Fat	0.074518	Fruits and Vegetables	227.3720	OUT019	1985.0	Medium	Tier 1	Supermarket Type1
8469	FDQ45	11.500	Regular	0.019114	Snack Foods	182.1608	OUT019	1985.0	Medium	Tier 1	Supermarket Type1
8480	FDG58	14.100	Low Fat	0.000000	Snack Foods	154.5340	OUT019	1985.0	Medium	Tier 1	Supermarket Type1
8490	FDH44	13.800	Regular	0.102296	Fruits and Vegetables	162.3552	OUT019	1985.0	Medium	Tier 1	Supermarket Type1

528 rows x 12 columns

In (100): `train.dropna(inplace=True)`

In (101): `train.info()`

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6113 entries, 0 to 8522
Data columns (total 12 columns):
 # Column Non-Null Count Dtype
--- -
 0 Item\_Identifier 6113 non-null object
 1 Item\_Weight 6113 non-null float64
 2 Item\_Fat\_Content 6113 non-null object
 3 Item\_Visibility 6113 non-null float64
 4 Item\_Type 6113 non-null object
 5 Item\_MRP 6113 non-null object
 6 Outlet\_Identifier 6113 non-null object
 7 Outlet\_Establishment\_Year 6113 non-null int64
 8 Outlet\_Size 6113 non-null object
 9 Outlet\_Location\_Type 6113 non-null object
10 Outlet\_Type 6113 non-null object
11 Item\_Outlet\_Sales 6113 non-null float64
dtypes: float64(4), int64(1), object(7)
memory usage: 620.9+ KB

**EDA**

Outlet size

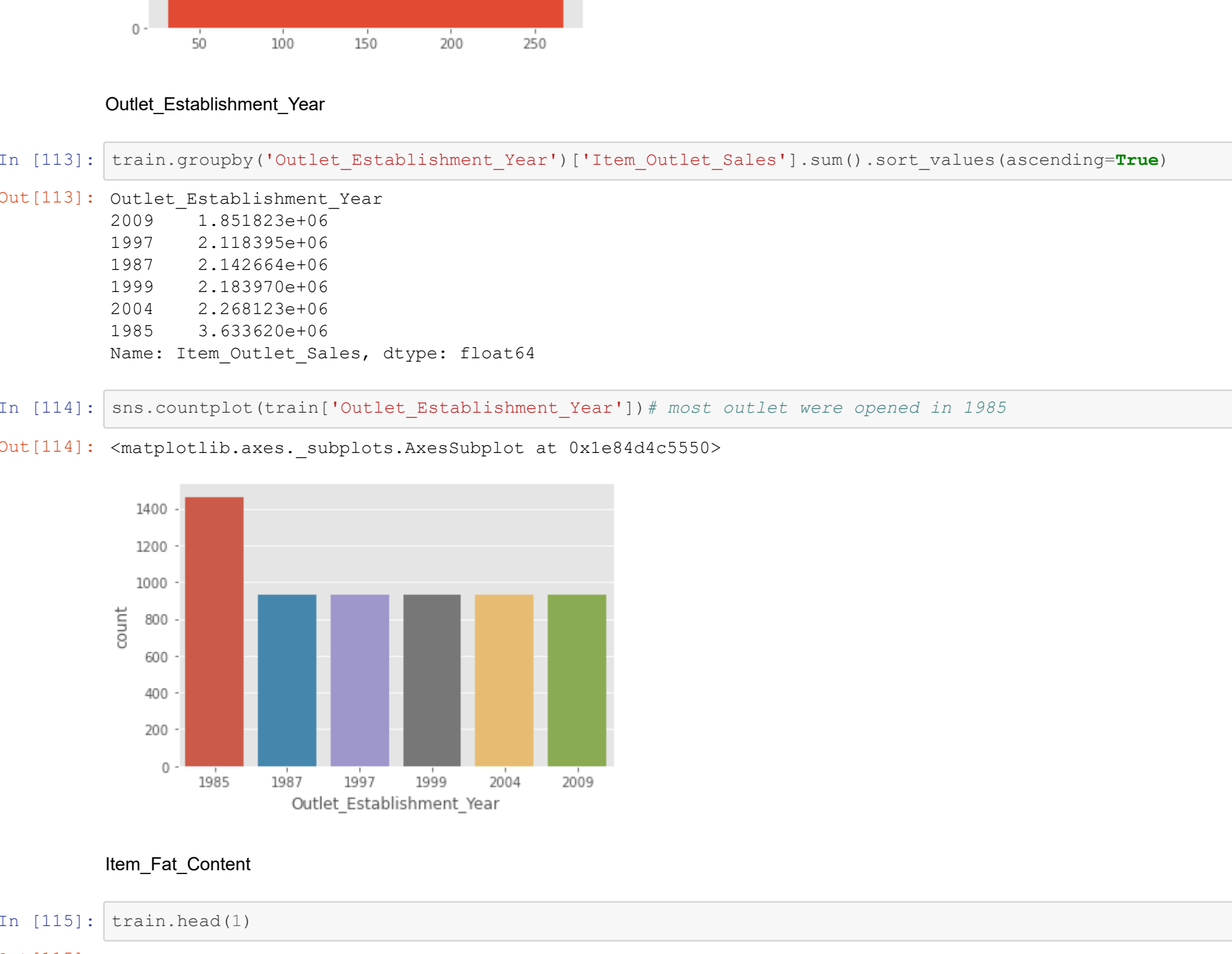
In (102): `len(train)`

Out (102): 6113

In (103): `f,ax=plt.subplots(1,2,figsize=(18,5))
train.groupby('Outlet_Size')['Item_Outlet_Sales'].sum().plot.bar(ax=ax[0],color='red')
ax[0].set_title('Sum')
train.groupby('Outlet_Size')['Item_Outlet_Sales'].mean().plot.bar(ax=ax[1])
ax[1].set_title('Mean')`

Out (103):

Text(0.5, 1.0, 'Mean')



In (104): `sns.countplot(train['Outlet_Size'])`

Out (104): `<matplotlib.axes._subplots.AxesSubplot at 0x1e84b278e0>`

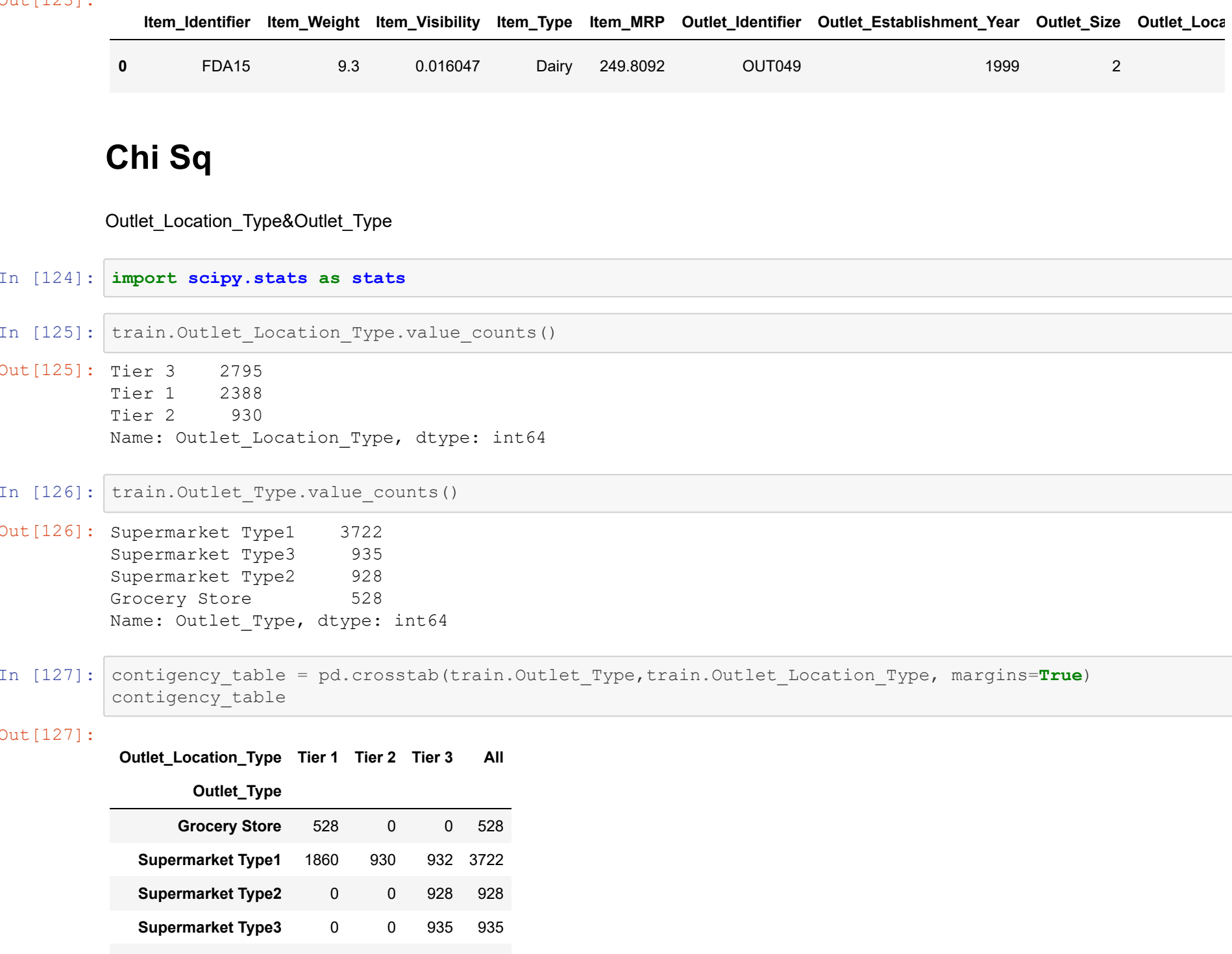


In (105): `plt.figure(figsize=(10,5))
sns.distplot(train[train.Outlet_Size == "1"].Item_Outlet_Sales, hist = False, kde_kws = {'shade': 'True'}, label = 'Small')
sns.distplot(train[train.Outlet_Size == "2"].Item_Outlet_Sales, hist = False, kde_kws = {'shade': 'True'}, label = 'Medium')
sns.distplot(train[train.Outlet_Size == "3"].Item_Outlet_Sales, hist = False, kde_kws = {'shade': 'True'}, label = 'High')`

C:\Users\Rahul\anaconda3\lib\site-packages\seaborn\distributions.py:198: RuntimeWarning: Mean of empty ndarray.

C:\Users\Rahul\anaconda3\lib\site-packages\numpy\core\methods.py:161: RuntimeWarning: invalid value encountered in double\_scalars

Out (105):



In (106): `train.head(3)`

Out (106):

	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
0	FDA15	9.30	Low Fat	0.016047	Dairy	249.8092	OUT049	1999	Medium	Tier 1	Supermarket Type1
1	DRC01	5.92	Regular	0.019278	Soft Drinks	48.2692	OUT018	2009	Small	Tier 3	Supermarket Type2
2	FDM15	17.50	Low Fat	0.016780	Meat	141.6180	OUT049	1999	Medium	Tier 1	Supermarket Type1

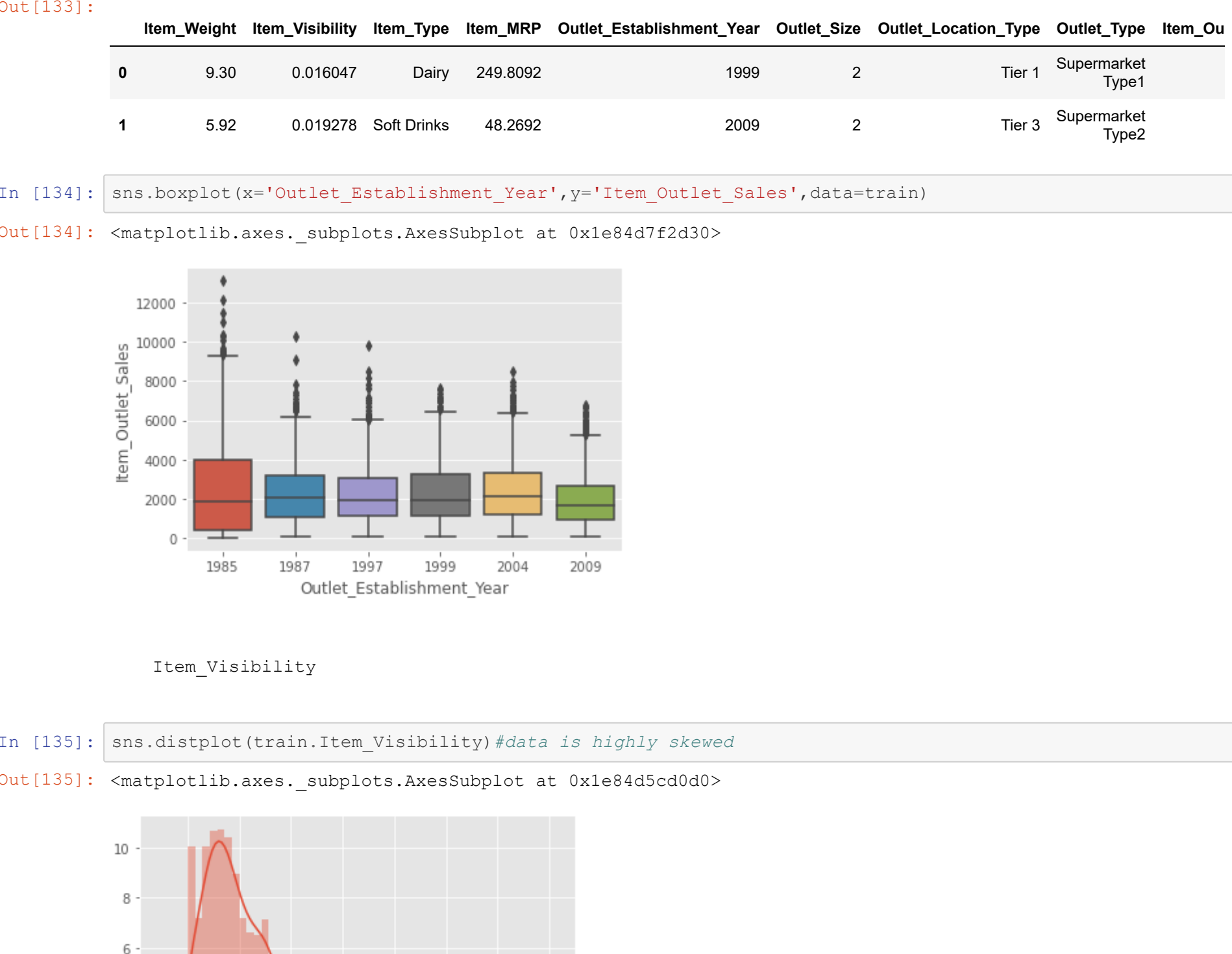
In (107): `train['Outlet_Size']=train.Outlet_Size.replace(['Medium','Small','High'],['2','1','3']) # we replaced t`

Out (107): `train['Outlet_Size'].unique()`

Out (108): `array(['2', '3', '1'], dtype=object)`

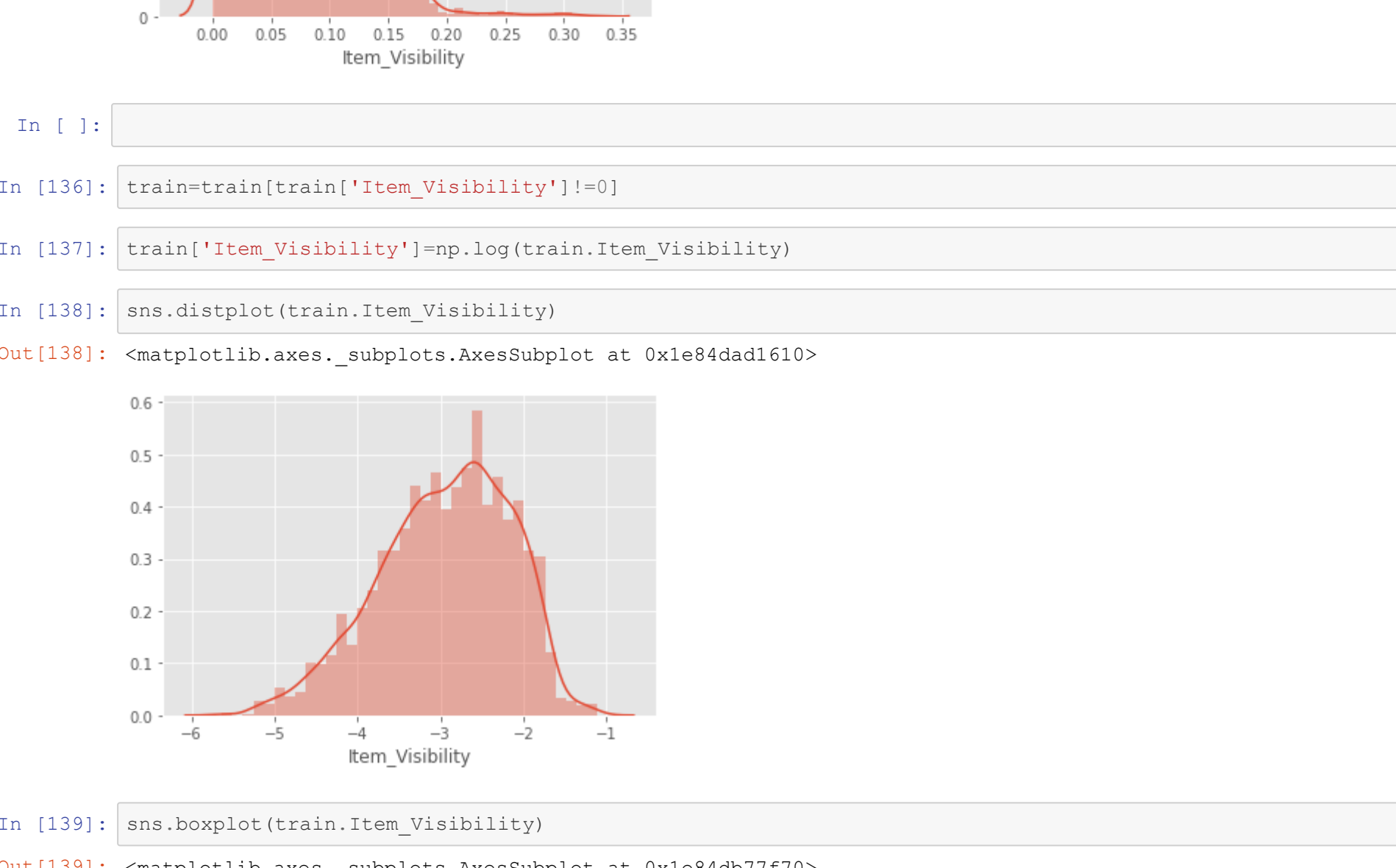
In (109): `plt.style.use('ggplot')
plt.figure(figsize=(7,7))
data=train[['Item_MRP', 'Item_Outlet_Sales']]
plt.xlabel('Item_MRP')
plt.ylabel('Item_Outlet_Sales')
plt.show()`

Out (109):



In (110): `sns.boxplot(train.Item_MRP) # there are no outliers`

Out (110): `<matplotlib.axes._subplots.AxesSubplot at 0x1e84dc2be0>`

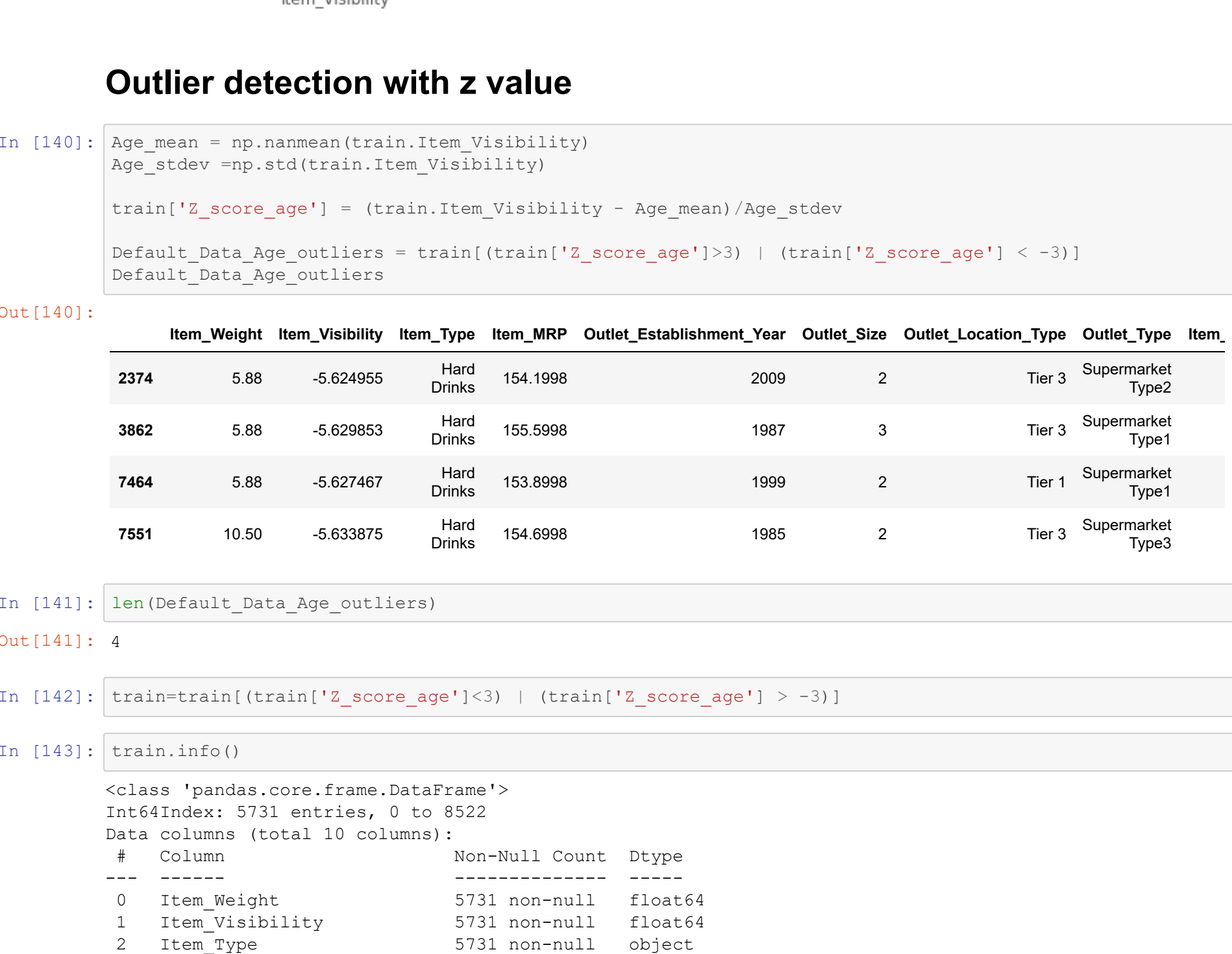


In (111): `print(train.Item_MRP.skew(),train.Item_MRP.kurtosis())`

Out (111): 0.1253530834135324 -0.8926804942521422

In (112): `plt.hist(train.Item_MRP)`

Out (112): `(array([653., 406., 769., 888., 577., 783., 871., 338., 411., 417.]),
array([31.29, 54.84984, 78.40968, 101.96952, 125.52936, 149.0892,
172.64994, 196.20288, 219.76972, 243.32856, 266.8884 ]),
< a list of 10 Patch objects>)`



Outlet\_Establishment\_Year

In (113): `train.groupby('Outlet_Establishment_Year')['Item_Outlet_Sales'].sum().sort_values(ascending=True)`

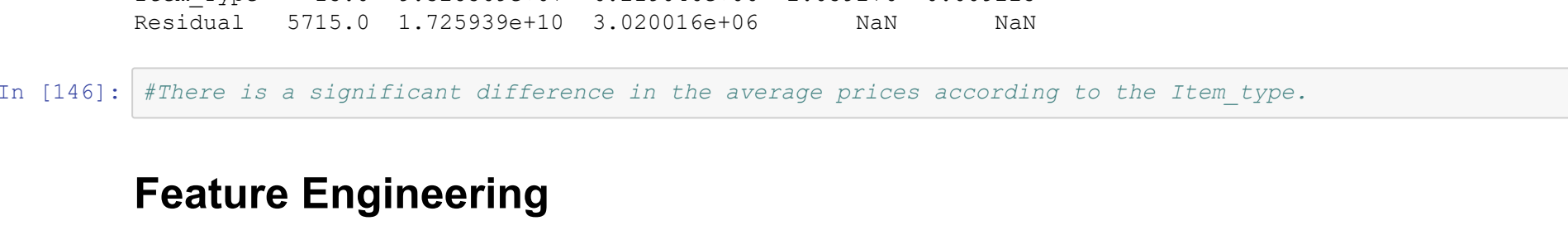
Out (113):

Outlet_Establishment_Year	Item_Outlet_Sales
2009	1.851923e+06
1997	2.118395e+06
1987	2.142664e+06
1999	2.183970e+06
2004	2.268123e+06
1985	3.633620e+06

Name: Item\_Outlet\_Sales, dtype: float64

In (114): `sns.countplot(train['Outlet_Establishment_Year']) # most outlet were opened in 1985`

Out (114): `<matplotlib.axes._subplots.AxesSubplot at 0x1e84dc5550>`



Item\_Fat\_Content

In (115): `train.head(1)`

Out (115):

	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
0	FDA15	9.3	Low Fat	0.016047	Dairy	249.8092	OUT049	1999	Medium	Tier 1	Supermarket Type1

In (116): `train.Item_Fat_Content.unique()`

Out (116): `array(['Low Fat', 'Regular', 'Low Fat', 'reg', 'LP'], dtype=object)`

In (117): `train[['Item_Fat_Content']]=train.Item_Fat_Content.replace(['LP','low fat','reg'],['Low Fat','Low Fat','Regular'])`

In (118): `train.Item_Fat_Content.nunique()`

Out (118): 2

In (119): `test['Item_Fat_Content']=test.Item_Fat_Content.replace(['LP','low fat','reg'],['Low Fat','Low Fat','Regular'])`

In (120): `train.groupby('Item_Fat_Content')['Item_Outlet_Sales'].mean() # so there's no difference so we can drop the column`

Out (120):

Item_Fat_Content	Item_Outlet_Sales
Low Fat	2305.015424
Regular	2355.078065

Name: Item\_Outlet\_Sales, dtype: float64

In (121): `train.Item_Fat_Content.value_counts()` # people tend to prefer low fat things over regular ones

Out (121):

Low Fat	3955
Regular	2158

Name: Item\_Fat\_Content, dtype: int64

In (122): `train.drop('Item_Fat_Content',axis=1,inplace=True)
test.drop('Item_Fat_Content',axis=1,inplace=True)`

In (123): `train.head(1)`

Out (123):

	Item_Identifier	Item_Weight	Item_Visibility	Item_Type	Item_MRP	Outlet_Identifier	Outlet_Establishment_Year	Outlet_Size	Outlet_Location_Type	Outlet_Type
0	FDA15	9.3	0.016047	Dairy	249.8092	OUT049	1999	2	Tier 1	Supermarket Type1

**Chi Sq**

Outlet\_Location\_Type&Outlet\_Type

In (124): `import scipy.stats as stats`

In (125): `train.Outlet_Location_Type.value_counts()`

Out (125):

Tier 3	2795
Tier 1	2388
Tier 2	930

Name: Outlet\_Location\_Type, dtype: int64

In (126): `train.Outlet_Type.value_counts()`

Out (126):

Supermarket Type1	3722
Supermarket Type3	935
Supermarket Type2	528
Grocery Store	528

Name: Outlet\_Type, dtype: int64

In (127): `contingency_table = pd.crosstab(train.Outlet_Type,train.Outlet_Location_Type, margins=True)
contingency_table`

Out (127):

	Outlet_Location_Type	Tier 1	Tier 2	Tier 3	All
Grocery Store	528	0	0	528	
Supermarket Type1	1800	930	932	3722	
Supermarket Type2	0	0	528	928	
Supermarket Type3	0	0	935	935	
Supermarket Type1 All	2388	930	2795	6113	

In (128): `chi_square, p_value, degrees_of_freedom, expected_frequencies=stats.chi2_contingency(contingency_table)
print(expected_frequencies)`

Out (128):

[ 206.25944708  80.3271716  241.41338132  528. 1
[1453.97284476  56.24570587 1701.78144937 3722. 1
[ 362.51660396 141.18109848 424.30230656 928. 1
[ 365.2511042 142.24603504 427.50286278 935. 1
[2388. 930. 2795. 6113. 1]

In (129): `print(chi_square, p_value) # hence we reject the hypothesis there is some relationship between the variables`

Out (129): 3730.483305468889 0.0

In (130): `train.drop(['Item_Identifier','Outlet_Identifier'],inplace=True,axis=1)
test.drop(['Item_Identifier','Outlet_Identifier'],inplace=True,axis=1)`

Out (130):

	Item_Weight	Item_Visibility	Item_Type	Item_MRP	Outlet_Establishment_Year	Outlet_Size	Outlet_Location_Type	Outlet_Type	Item_Outlet_Sales
0	9.300	0.016047	Dairy	249.8092	1999	2	Tier 1	Supermarket Type1	2498.092
1	5.920	0.019278	Soft Drinks	48.2692	2009	2	Tier 3	Supermarket Type2	48.2692
2	17.500	0.016780	Meat	141.6180	1999	2	Tier 1	Supermarket Type1	1416.180
4	8.930	0.000000	Household	53.8614	1987	3			



One Hot Encoding

```
In [150]: train=pd.get_dummies(data=train, columns=['Outlet_Location_Type', 'Outlet_Type'],drop_first=True)
test=pd.get_dummies(data=test, columns=['Outlet_Location_Type', 'Outlet_Type'],drop_first=True)
```

```
In [151]: train.head()
```

```
Out[151]:
```

	Item_Weight	Item_Visibility	Item_Type	Item_MRP	Outlet_Establishment_Year	Outlet_Size	Item_Outlet_Sales	Z_score_age	Outlet_Location_Type
0	9.30	-4.132215	4	249.8092	1999	2	3735.1380	-1.509749	3
1	5.92	-3.948780	14	48.2692	2009	2	443.4228	-1.275974	1
2	17.50	-4.088756	10	141.6180	1999	2	2097.2700	-1.454364	0
6	13.65	-4.362923	13	57.6588	1987	3	343.5528	-1.803770	1
7	14.10	-2.059875	13	107.7622	1985	2	4022.7636	1.131299	1

```
In [152]: train.drop(['Outlet_Establishment_Year','Z_score_age'],inplace=True,axis=1)
```

```
In [153]: train.head()
```

```
Out[153]:
```

	Item_Weight	Item_Visibility	Item_Type	Item_MRP	Outlet_Size	Item_Outlet_Sales	Outlet_Location_Type_Tier_2	Outlet_Location_Type_Tier_3
0	9.30	-4.132215	4	249.8092	2	3735.1380	0	0
1	5.92	-3.948780	14	48.2692	2	443.4228	0	1
2	17.50	-4.088756	10	141.6180	2	2097.2700	0	0
6	13.65	-4.362923	13	57.6588	3	343.5528	0	1
7	14.10	-2.059875	13	107.7622	2	4022.7636	0	1

## SCALING

```
In [154]: from sklearn.preprocessing import StandardScaler
```

```
In [155]: scaler=StandardScaler()
```

```
In [156]: train[['Item_Weight','Item_Visibility','Item_MRP']]=scaler.fit_transform(train[['Item_Weight','Item_Visibility','Item_MRP']])
```

```
In [157]: train.head()
```

	Item_Weight	Item_Visibility	Item_Type	Item_MRP	Outlet_Size	Item_Outlet_Sales	Outlet_Location_Type_Tier_2	Outlet_Location_Type_Tier_3
0	9.30	-1.132215	4	249.8092	2	3735.1380	0	0
1	-1.670407	-1.275974	14	-1.499642	2	443.4228	0	1
2	1.129638	-1.454364	10	0.002105	2	2097.2700	0	0
6	0.198708	-1.803770	13	-1.348587	3	343.5528	0	1
7	0.307518	1.131299	13	-0.542550	2	4022.7636	0	1

```
In [158]: plt.figure(figsize=(10,5))
sns.heatmap(train.corr(),annot=True)
```

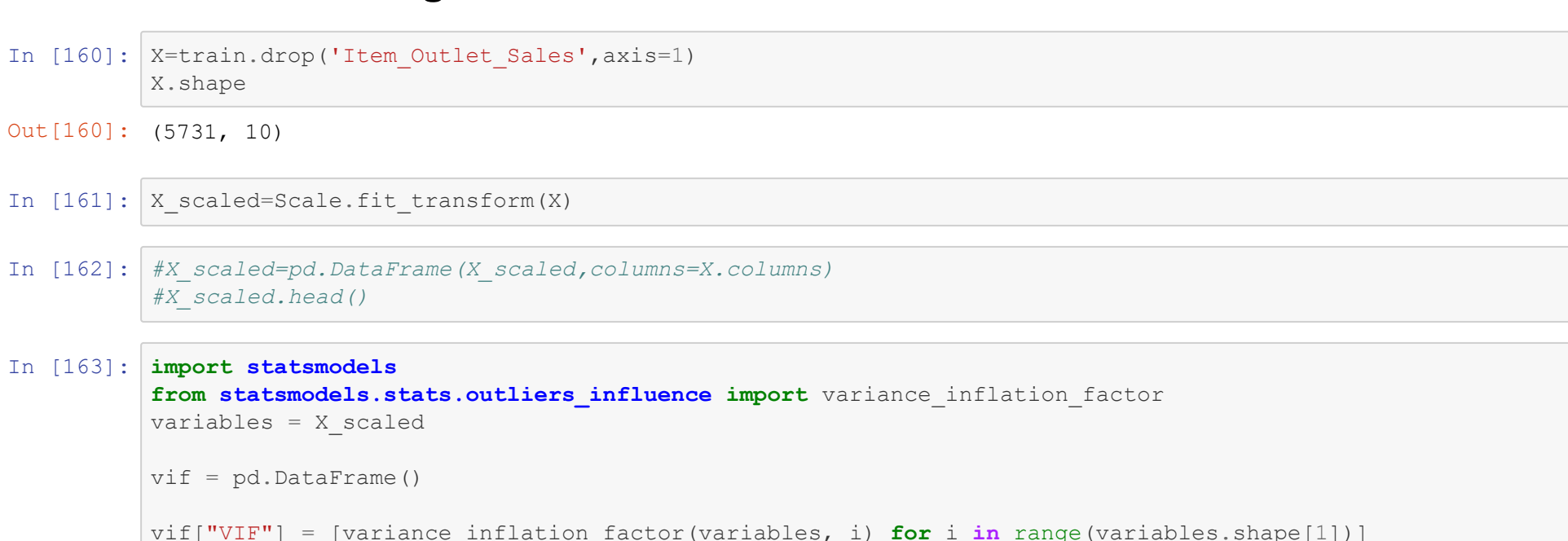
```
Out[158]: <matplotlib.axes._subplots.AxesSubplot at 0x1e84db07580>
```



```
In [159]: plt.figure(figsize=(30,20))
sns.pairplot(train)
```

```
Out[159]: <seaborn.axisgrid.PairGrid at 0x1e84dc41400>
```

<Figure size 2160x1440 with 0 Axes>



## Model Building

```
In [160]: X=train.drop('Item_Outlet_Sales',axis=1)
X.shape
```

```
Out[160]: (5731, 10)
```

```
In [161]: X_scaled=Scale.fit_transform(X)
```

```
In [162]: #X_scaled=pd.DataFrame(X_scaled,columns=X.columns)
#X_scaled.head()
```

```
In [163]: import statsmodels
from statsmodels.stats.outliers_influence import variance_inflation_factor
variables = X_scaled
vif = pd.DataFrame()
vif['VIF'] = [variance_inflation_factor(variables, i) for i in range(variables.shape[1])]
vif['Features'] = X.columns
vif
```

```
Out[163]:
```

	VIF	Features
0	1.004546	Item_Weight
1	1.045574	Item_Visibility
2	1.004702	Item_Type
3	1.002355	Item_MRP
4	6.391509	Outlet_Size
5	1.701447	Outlet_Location_Type_Tier_2
6	9.750995	Outlet_Location_Type_Tier_3
7	4.435168	Outlet_Type_Supermarket Type1
8	0.699739	Outlet_Type_Supermarket Type2
9	4.105065	Outlet_Type_Supermarket Type3

```
In [164]: X_scaled=pd.DataFrame(X_scaled,columns=X.columns)
X_scaled.head()
X_scaled.drop(['Item_Weight','Item_Visibility','Item_Type','Outlet_Size','Outlet_Location_Type_Tier_2',
              'Outlet_Location_Type_Tier_3'],axis=1,inplace=True)
```

```
In [165]: Y=train.loc[:, "Item_Outlet_Sales"]
Y
```

```
Out[165]: 0      3735.1380
1      443.4228
2      2097.2700
6      343.5528
7      4022.7636

8517      3608.0360
8518      2778.3834
8519      1193.1136
8521      1845.5976
8522      765.6700
Name: Item_Outlet_Sales, Length: 5731, dtype: float64
```

```
In [166]: from sklearn.model_selection import train_test_split
```

```
In [167]: X_train, X_test, y_train, y_test = train_test_split(X_scaled,Y, test_size = .3, random_state = 3)
```

```
In [212]: from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from sklearn.preprocessing import PolynomialFeatures
from sklearn.metrics import accuracy_score
from sklearn.svm import SVR
from sklearn.ensemble import AdaBoostRegressor
from sklearn.neighbors import KNeighborsRegressor
```

```
In [169]: model=LinearRegression()
model1=Ridge()
model2=RidgeCV()
model3=Lasso()
model4=LassoCV()
sv=SVR()
```

```
In [217]: model.fit(X_train,y_train)
sv.fit(X_train,y_train)
print(model.intercept_)
print(model.coef_)
print("R2 score of the model on test data SVR")
print(sv.score(X_test,y_test))
print("R2 score of the model on test data")
print(model.score(X_test,y_test))
```

```
2316.3422104849938
[1024.4685246  973.4793457  578.60475165 1213.79075137]
R2 score of the model on test data SVR
0.11901720116178194
R2 score of the model on test data
0.5518203802646215
```

```
In [171]: pd.DataFrame(model.coef_,index=X_scaled.columns)
```

```
Out[171]:
```

	0
Item_MRP	1024.468525
Outlet_Type_Supermarket Type1	973.479346
Outlet_Type_Supermarket Type2	578.604752
Outlet_Type_Supermarket Type3	1213.790751

## AdaBoost

```
In [201]: Ab=AdaBoostRegressor(n_estimators=1000,loss='linear')
```

```
In [202]: Ab.fit(X_train,y_train)
```

```
Out[202]: AdaBoostRegressor(n_estimators=1000)
```

```
In [203]: print(Ab.score(X_test,y_test))
print('=====')
print(Ab.score(X_train,y_train))
```

```
0.47646286731213217
=====
0.49775495077937395
```

```
In [204]: Ab.get_params()
```

```
Out[204]: {'base_estimator': None,
'learning_rate': 1.0,
'loss': 'linear',
'n_estimators': 1000,
'random_state': None}
```

```
In [191]: Features=pd.DataFrame(Ab.feature_importances_,index=X_scaled.columns)
Features
```

```
Out[191]:
```

	0
Item_MRP	0.778474
Outlet_Type_Supermarket Type1	0.054359
Outlet_Type_Supermarket Type2	0.009638
Outlet_Type_Supermarket Type3	0.157529

```
In [192]: from sklearn.model_selection import GridSearchCV
```

```
In [206]: param = {'learning_rate':np.linspace(0.1,1,10)}
```

```
In [208]: grid=GridSearchCV(estimator=Ab,cv=5,param_grid=param,n_jobs=-1)
```

```
In [209]: grid.fit(X_train,y_train)
```

```
Out[209]: GridSearchCV(cv=5, estimator=AdaBoostRegressor(n_estimators=1000), n_jobs=-1,
param_grid={'learning_rate': array([0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1. ])})
```

```
In [210]: grid.best_params_
```

```
Out[210]: {'learning_rate': 0.1}
```

```
In [211]: grid.best_score_
```

```
Out[211]: 0.4776963754777038
```

## Polynomial

```
In [174]: polynomial_features= PolynomialFeatures(degree=2)
```

```
x_poly2 = polynomial_features.fit_transform(X_train)
```

```
In [175]: x_poly2.shape
```

```
Out[175]: (4011, 15)
```

```
In [176]: model.fit(x_poly2 ,y_train)
```

```
Out[176]: LinearRegression()
```

```
In [177]: model.score(x_poly2 ,y_train)
```

```
Out[177]: 0.5981235933727165
```

## Assumption check

```
In [178]: import scipy.stats as stats
import statsmodels.formula.api as sm
```

```
In [ ]:
```

```
In [179]: model_m = sm.ols(formula='Item_Outlet_Sales~ Item_MRP', data = train)
```

```
# fitting the model
model_fit_m = model_m.fit()

# summary of the model
model_fit_m.summary()
```

```
Out[179]:
```

OLS Regression Results

Warnings:

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

Omnibus: 524.572 Durbin-Watson: 1.991

Prob(Omnibus): 0.000 Jarque-Bera (JB): 1174.416

Skew: 0.570 Prob(JB): 9.52e-258

Kurtosis: 4.902 Cond. No. 1.00

Warnings:

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