

## ▼ Big Mart Sales Dataset

- ▼ Objective: To find out the properties of a product, and store which impacts the sales of a product.

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 %matplotlib inline
5
```

```
1 train = pd.read_csv("train.csv")
2 test = pd.read_csv("test.csv")
```

```
1 #Combine test and train into one file
2 train['source']='train'
3 test['source']='test'
4 data = pd.concat([train, test],ignore_index=True)
5 print(train.shape, test.shape, data.shape)
```

```
(8523, 13) (5681, 12) (14204, 13)
```

```
1 data.head()
```



	Item_Fat_Content	Item_Identifier	Item_MRP	Item_Outlet_Sales	Item_Type
0	Low Fat	FDA15	249.8092	3735.1380	Dairy
1	Regular	DRC01	48.2692	443.4228	Soft Drinks
2	Low Fat	FDN15	141.6180	2097.2700	Meat
3	Regular	FDX07	182.0950	732.3800	Fruits and Vegetables
4	Low Fat	NCD19	53.8614	994.7052	Household

```
1 #Numerical data summary:
2 data.describe()
```

	Item_MRP	Item_Outlet_Sales	Item_Visibility	Item_Weight	Outlet_Est
<b>count</b>	14204.000000	8523.000000	14204.000000	11765.000000	
<b>mean</b>	141.004977	2181.288914	0.065953	12.792854	
<b>std</b>	62.086938	1706.499616	0.051459	4.652502	
<b>min</b>	31.290000	33.290000	0.000000	4.555000	
<b>25%</b>	94.012000	834.247400	0.027036	8.710000	
<b>50%</b>	142.247000	1794.331000	0.054021	12.600000	
<b>75%</b>	185.855600	3101.296400	0.094037	16.750000	

## ▼ Data Cleaning

```
1 #Check missing values:
2 data.apply(lambda x: sum(x.isnull()))
```

```
Item_Fat_Content      0
Item_Identifier       0
Item_MRP              0
Item_Outlet_Sales    5681
Item_Type             0
Item_Visibility       0
Item_Weight          2439
Outlet_Establishment_Year  0
Outlet_Identifier     0
Outlet_Location_Type  0
Outlet_Size          4016
Outlet_Type           0
source               0
dtype: int64
```

## ▼ Filling missing values

```
1 data.Item_Outlet_Sales = data.Item_Outlet_Sales.fillna(data.Item_Outlet_Sales.mean())
```

```
1 data.Item_Weight = data.Item_Weight.fillna(data.Item_Weight.mean())
```

```
1 data['Outlet_Size'].value_counts()
```

```
Medium    4655
Small     3980
High      1553
Name: Outlet_Size, dtype: int64
```

```
1 data.Outlet_Size = data.Outlet_Size.fillna('Medium')
```

```
1 data.apply(lambda x: sum(x.isnull()))
```

```

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

```

```
1 data.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14204 entries, 0 to 14203
Data columns (total 13 columns):
Item_Fat_Content      14204 non-null object
Item_Identifier      14204 non-null object
Item_MRP             14204 non-null float64
Item_Outlet_Sales    14204 non-null float64
Item_Type            14204 non-null object
Item_Visibility       14204 non-null float64
Item_Weight          14204 non-null float64
Outlet_Establishment_Year 14204 non-null int64
Outlet_Identifier     14204 non-null object
Outlet_Location_Type  14204 non-null object
Outlet_Size          14204 non-null object
Outlet_Type          14204 non-null object
source              14204 non-null object
dtypes: float64(4), int64(1), object(8)
memory usage: 1.4+ MB

```

```

1 #Item type combine:
2 data['Item_Identifier'].value_counts()
3 data['Item_Type_Combined'] = data['Item_Identifier'].apply(lambda x: x[0:2])
4 data['Item_Type_Combined'] = data['Item_Type_Combined'].map({'FD': 'Food',
5                                                             'NC': 'Non-Consumab',
6                                                             'DR': 'Drinks'})
7 data['Item_Type_Combined'].value_counts()

```

```

Food      10201
Non-Consumable  2686
Drinks     1317
Name: Item_Type_Combined, dtype: int64

```

## ▼ Numerical and One-Hot Coding of Categorical variables

```

1 #Import library:
2 from sklearn.preprocessing import LabelEncoder, OneHotEncoder
3 le = LabelEncoder()

```

```

4 #New variable for outlet
5 data['Outlet'] = le.fit_transform(data['Outlet_Identifier'])
6 var_mod = ['Item_Fat_Content', 'Outlet_Location_Type', 'Outlet_Size', 'Item_Type_Co
7 le = LabelEncoder()
8 for i in var_mod:
9     data[i] = le.fit_transform(data[i])

```

```

1 #One Hot Coding:
2 data = pd.get_dummies(data, columns=['Item_Fat_Content', 'Outlet_Location_Type',

```

```
1 data.head()
```

	Item_Identifier	Item_MRP	Item_Outlet_Sales	Item_Type	Item_Visibility	I
0	FDA15	249.8092	3735.1380	Dairy	0.016047	
1	DRC01	48.2692	443.4228	Soft Drinks	0.019278	
2	FDN15	141.6180	2097.2700	Meat	0.016760	
3	FDX07	182.0950	732.3800	Fruits and Vegetables	0.000000	
4	NCD19	53.8614	994.7052	Household	0.000000	

5 rows × 37 columns

```
1 data.dtypes
```

```

Item_Identifier      object
Item_MRP             float64
Item_Outlet_Sales    float64
Item_Type            object
Item_Visibility      float64
Item_Weight          float64
Outlet_Establishment_Year  int64
Outlet_Identifier    object
source              object
Item_Fat_Content_0   uint8
Item_Fat_Content_1   uint8
Item_Fat_Content_2   uint8
Item_Fat_Content_3   uint8
Item_Fat_Content_4   uint8
Outlet_Location_Type_0  uint8
Outlet_Location_Type_1  uint8
Outlet_Location_Type_2  uint8
Outlet_Size_0         uint8
Outlet_Size_1         uint8
Outlet_Size_2         uint8
Outlet_Type_0         uint8
Outlet_Type_1         uint8
Outlet_Type_2         uint8
Outlet_Type_3         uint8
Item_Type_Combined_0   uint8
Item_Type_Combined_1   uint8
Item_Type_Combined_2   uint8

```

```

Outlet_0      uint8
Outlet_1      uint8
Outlet_2      uint8
Outlet_3      uint8
Outlet_4      uint8
Outlet_5      uint8
Outlet_6      uint8
Outlet_7      uint8
Outlet_8      uint8
Outlet_9      uint8
dtype: object

```

## ▼ Exporting Data

```

1 import warnings
2 warnings.filterwarnings('ignore')
3 #Drop the columns which have been converted to different types:
4 data.drop(['Item_Type', 'Outlet_Establishment_Year'], axis=1, inplace=True)
5
6 #Divide into test and train:
7 train = data.loc[data['source']=="train"]
8 test = data.loc[data['source']=="test"]
9
10 #Drop unnecessary columns:
11 test.drop(['Item_Outlet_Sales', 'source'], axis=1, inplace=True)
12 train.drop(['source'], axis=1, inplace=True)
13
14 #Export files as modified versions:
15 train.to_csv("train_modified.csv", index=False)
16 test.to_csv("test_modified.csv", index=False)

```

## ▼ Model Building

```

1 # Reading modified data
2 train2 = pd.read_csv("train_modified.csv")
3 test2 = pd.read_csv("test_modified.csv")

1 train2.head()

```

Item_Identifier	Item_MRP	Item_Outlet_Sales	Item_Visibility	Item_Weight
-----------------	----------	-------------------	-----------------	-------------

```
1 X_train = train2.drop(['Item_Outlet_Sales', 'Outlet_Identifier', 'Item_Identifier'])
2 y_train = train2.Item_Outlet_Sales
```

0	249.8092	0.016047	9.30	0
1	48.2692	0.019278	5.92	0
2	141.6180	0.016760	17.50	0
3	182.0950	0.000000	19.20	0
4	53.8614	0.000000	8.93	0

```
1 X_test = test2.drop(['Outlet_Identifier', 'Item_Identifier'], axis=1)
```

```
1 X_train.head()
```

	Item_MRP	Item_Visibility	Item_Weight	Item_Fat_Content_0	Item_Fat_Content_1
0	249.8092	0.016047	9.30	0	
1	48.2692	0.019278	5.92	0	
2	141.6180	0.016760	17.50	0	
3	182.0950	0.000000	19.20	0	
4	53.8614	0.000000	8.93	0	

5 rows × 5 columns

```
1 y_train.head()
```

0	3735.1380
1	443.4228
2	2097.2700
3	732.3800
4	994.7052

Name: Item\_Outlet\_Sales, dtype: float64

## ▼ Linear Regression Model:

```
1 # Fitting Multiple Linear Regression to the training set
2 from sklearn.linear_model import LinearRegression
3 regressor = LinearRegression()
4 regressor.fit(X_train, y_train)
```

LinearRegression(copy\_X=True, fit\_intercept=True, n\_jobs=1, normalize=False)

```
1 # Predicting the test set results
2 y_pred = regressor.predict(X_test)
```

```
1 y_pred
```

array([1848.53604783, 1472.81670435, 1875.65285894, ..., 1809.18796433,  
3565.6645235 , 1267.46171871])

```

1 import warnings
2 warnings.filterwarnings('ignore')
3 # Measuring Accuracy
4 from sklearn.metrics import accuracy_score, r2_score, mean_squared_error
5 from sklearn.model_selection import cross_val_score
6 from sklearn import cross_validation, metrics
7

```

```

1 lr_accuracy = round(regressor.score(X_train,y_train) * 100,2)
2 lr_accuracy

```

56.36

```

1 r2_score(y_train, regressor.predict(X_train))

```

0.563589277727048

```

1 import warnings
2 warnings.filterwarnings('ignore')
3 #Perform cross-validation:
4 cv_score = cross_val_score(regressor, X_train, y_train, cv=5, scoring='mean_squared_error')
5

```

```

1 print(np.sqrt(np.abs(cv_score)))

```

[1150.93927648 1118.68414103 1112.89657923 1126.30724065 1140.59735737]

```

1 print("RMSE : %.4g" % np.sqrt(metrics.mean_squared_error(y_train, regressor.predict(X_train))))

```

RMSE : 1127

```

1 submission = pd.DataFrame({
2 'Item_Identifier':test2['Item_Identifier'],
3 'Outlet_Identifier':test2['Outlet_Identifier'],
4 'Item_Outlet_Sales': y_pred
5 },columns=['Item_Identifier','Outlet_Identifier','Item_Outlet_Sales'])

```

```

1 submission.to_csv('submission1.csv',index=False)

```

## ▼ Decision Tree Model:

```

1 # Fitting Decision Tree Regression to the dataset
2 from sklearn.tree import DecisionTreeRegressor
3 regressor = DecisionTreeRegressor(max_depth=15,min_samples_leaf=300)
4 regressor.fit(X_train, y_train)
5

```

DecisionTreeRegressor(criterion='mse', max\_depth=15, max\_features=None,

```
max_leaf_nodes=None, min_impurity_decrease=0.0,
min_impurity_split=None, min_samples_leaf=300,
min_samples_split=2, min_weight_fraction_leaf=0.0,
presort=False, random_state=None, splitter='best')
```

```
1 # Predicting the test set results
2 y_pred = regressor.predict(X_test)
3 y_pred
```

```
array([1673.98398729, 1349.51290433, 471.30684669, ..., 1892.06614452,
       3805.94860417, 1349.51290433])
```

```
1 tree_accuracy = round(regressor.score(X_train,y_train),2)
2 tree_accuracy
```

```
0.59
```

```
1 r2_score(y_train, regressor.predict(X_train))
```

```
0.5884050821570486
```

```
1 import warnings
2 warnings.filterwarnings('ignore')
3 cv_score = cross_val_score(regressor, X_train, y_train, cv=5, scoring='mean_squared_error')
4 print(np.sqrt(np.abs(cv_score)))
```

```
[1138.77137157 1109.42501179 1145.66395939 1113.2648073 1129.0816826 ]
```

```
1 print("RMSE : %.4g" % np.sqrt(metrics.mean_squared_error(y_train, regressor.predict(X_test))))
```

```
RMSE : 1095
```

```
1 submission = pd.DataFrame({
2 'Item_Identifier':test2['Item_Identifier'],
3 'Outlet_Identifier':test2['Outlet_Identifier'],
4 'Item_Outlet_Sales': y_pred
5 },columns=['Item_Identifier','Outlet_Identifier','Item_Outlet_Sales'])
```

```
1 submission.to_csv('submission2.csv',index=False)
```

## ▼ Random Forest Model:

```
1 # Fitting Random Forest Regression to the dataset
2 from sklearn.ensemble import RandomForestRegressor
3 regressor = RandomForestRegressor(n_estimators=100,max_depth=6, min_samples_leaf=1)
4 regressor.fit(X_train, y_train)
```

```
RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=6,
                        max_features='auto', max_leaf_nodes=None,
```



```

min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=50, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=4,
oob_score=False, random_state=None, verbose=0, warm_start=False)

```

```

1 # Predicting the test set results
2 y_pred = regressor.predict(X_test)
3 y_pred

```

```

array([1643.87106725, 1364.24193091, 603.09113992, ..., 1957.62183676,
       3698.60040819, 1290.25320329])

```

```

1 rf_accuracy = round(regressor.score(X_train,y_train),2)
2 rf_accuracy

```

```
0.61
```

```
1 r2_score(y_train, regressor.predict(X_train))
```

```
0.6125814698282157
```

```

1 import warnings
2 warnings.filterwarnings('ignore')
3 cv_score = cross_val_score(regressor, X_train, y_train, cv=5, scoring='mean_squared_error')
4 print(np.sqrt(np.abs(cv_score)))

```

```
[1100.46298396 1077.70836131 1077.65325884 1069.0502564 1083.85364282]
```

```
1 print("RMSE : %.4g" % np.sqrt(metrics.mean_squared_error(y_train, regressor.predict(X_test))))
```

```
RMSE : 1062
```

```

1 submission = pd.DataFrame({
2 'Item_Identifier':test2['Item_Identifier'],
3 'Outlet_Identifier':test2['Outlet_Identifier'],
4 'Item_Outlet_Sales': y_pred
5 },columns=['Item_Identifier','Outlet_Identifier','Item_Outlet_Sales'])

```

```
1 submission.to_csv('submission3.csv',index=False)
```

