

## Assignment : C3

**TITLE :** Bigmart Sales Analysis

**Problem :** Bigmart Sales Analysis: For data statement comparison of transactions / records of a sales store. The data has 8,523 rows of 12 variables. Predict the sales of a store.

**Objective :**

- i] To understand the sales analysis and Prediction.
- ii] To understand the concept of supervised learning in Data Analysis.

**Outcome :**

- i] learn the supervised learning Algorithm
- ii] understand sales Prediction.

**Software :** OS: windows / ubuntu distribution  
**and Hardware :** 4GB RAM, 500GB HDD, Python Libraries,  
**Requirement :** Python framework: Bigmart sales Dataset.

**Theory :**

**Bigmart Sales Analysis :**

According to the information provided, Bigmart is a big Supermarket chain, which stores all around the country and its current board set out a challenge to all data scientist out there, to help they create a Model that can predict the sales per product, for each store. Bigmart has



collected sales data from the year 2013 from 1559 products across 10 stores in different cities. with this information the corporation hopes we can identify the products and stores which play a key role in their sales and uses that information to take the correct measures to ensure success of their business.

we will explore the problem in following steps:

### 1] Hypothesis generation:

this is a very pivotal step in the process of analysing the data. this involves understanding the problem and making some hypothesis about what could potentially have a good impact on the outcome. this is done before looking at the data and we end up creating a laundry list of the different analysis which we can potentially perform if data is available.

### 2] Data Exploration:

we will be performing some basic data exploration here and come up with some inferences about the data. we will figure out some irregularities and address them in the next session. the first step is to look at the data and try to identify the information which we hypothesized vs the data available. A comparison between the data dictionaries on the competition page and our hypothesis. we will invariably find features which are hypothesized but data doesn't carry and vice versa.



```
train = pd.read_csv("train.csv")
test = pd.read_csv("test.csv")
train['source'] = train
test['source'] = test
data = pd.concat([train, test], ignore_index=True)
some observations:
```

1. Item visibility: has a min value of zero. This makes no practical sense because when a product is being sold in a store, the visibility cannot be 0.
2. Outlet-Establishment-Years: vary from 1985 to 2009. The values might not be apt in this form.

3. Data cleaning: This step typically involves imputing missing values and treating outliers. Though outlier removal is very important in regression techniques, advanced tree based algorithms are insensitive to outliers. Imputing missing values: we have found two values with missing values: Item-weight, Outlet-size.

4. Feature Engineering: we explored some nuances in the data in the data exploration section. We will create some new variables using existing ones in this section:

Step 1: Consider combining outlet-type:  
During exploration we decided to consider combining the Supermarket Type 2 and Type 3 variables.

Step 2: Modify item-visibility: we noticed that the minimum value here is 0 which makes



no Practical sense. we can use the Visibility-avg-item variable to achieve this.

Step 3: Create a broad category of Type of item. Earlier we saw that the Item-Type variable has 16 categories which might prove to be very useful in analysis. So it's good idea to combine them. If you look at the Item-identifier i.e. the unique ID of each item it starts either FD, NC, or DR. If you see the categories these look like being Food, Drinks and Non-consumables.

Step 4: Determine the years of operation of a store:

we wanted to make a new column depicting the year of operation of store

Step 5: Modify Categories of Item-Fat-Content. we found typos and difference in representation in categories of Item-Fat-Content variables.

Step 6: Numerical and one hot coding of categorical values: I have converted all categories of nominal variables into numeric types. Also I wanted Outlet-Identifier as a variable, as well. One-hot coding refers to creating dummy variables, one for each category of a categorical variable. For example the item-fat-content.



Test Cases :

Description	Result Values	Pass/Fail
i) Imputing Missing Values : Item-weight outlet-size	Item-weight : original missing : 2439 Final missing : 0  outlet-size : original missing : 4016 Final missing : 0	Pass     Pass
ii) Item-visibility	original zeros : 829 Final zeros : 0	Pass
iii) Model Training	linear regression : accuracy score : 56 Decision tree : accuracy score : 62 XGB Regressor : accuracy score : 67	Pass

Conclusion : Thus we successfully done bigmart sales analysis and predicted the sales of a store.