# BIG MART SALES REPORT
## ASSIGNMENT 1

**PREPARED BY : BERNARD, LIONEL, JOEL, ETHAN**

# Contents

# 1.1 FRAMING THE PROBLEM

We must first understand the objective of working on this BigMart Sales Data. The first objective is to determine the products that has the potential to increase their sales based on their current sales value. The second objective is to determine the store that has the potential to increase their sales based on their current sales output on the products.

The problem of the missing data in for the products in some stores due to technical glitches has the impact of BigMart being unable to determine the store or products with the highest sales accurately. This will affect BigMart as they are not able to increase sales, so a good starting point would be to use the existing perfect data and work from there to determine the sales value of the product and the stores.
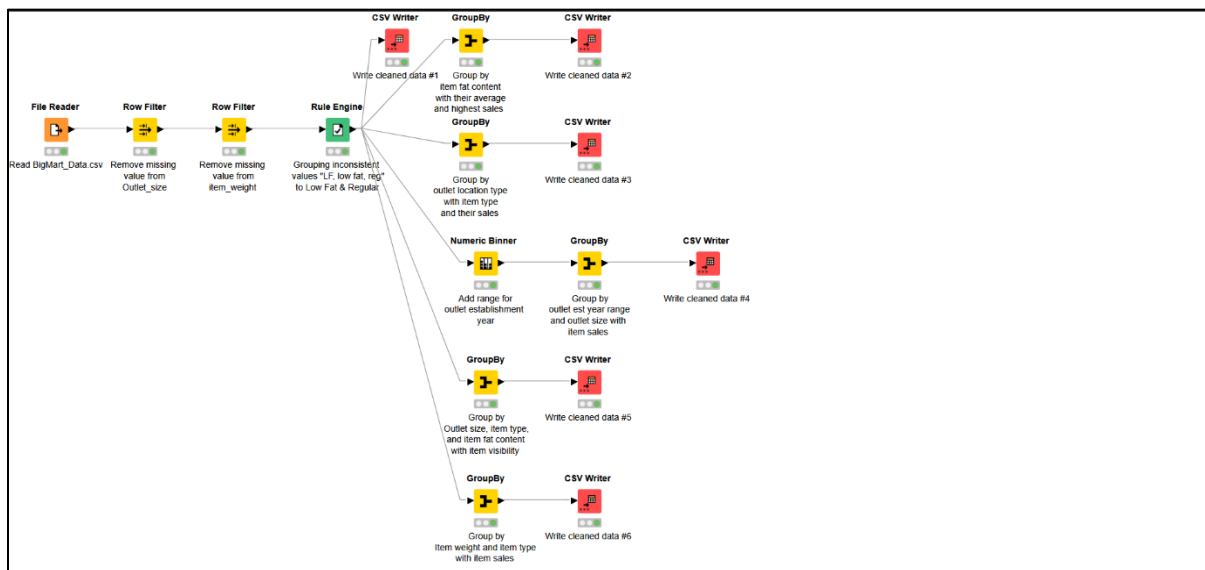
Based on the case study, because of the impacts are directed towards the revenue, we can assume that the stakeholder could be a leadership-level person on the marketing team, or it could also be someone from the executive level, which prioritizes revenue over the maintainability of the product.

In order to achieve the objectives, we must first determine the factors that affect the sales of the product. This can be ranged from the size and location of the store down to the product itself. Next, we can prepare the data by identifying the trends and patterns based on the cleaned data. The model will be built and trained to identify how much of an error can the business tolerate. Finally, we can make the initial predictions and to further refine them later. The expected outcome would be to successfully predict the products and stores that have potential to increase their sales.

    a. Clean the dataset and record all the changes.
        i.      Read BigMart_Data.csv –
             File Reader → Insert BigMart_Data.csv

        ii.     Remove the missing values from Outlet_size –
             Row Filter → Exclude rows by attribute values, Column = Outlet_size, select "only missing values match"

        iii.    Remove the missing values from item_weight –
             Row Filter → Exclude rows by attribute values, Column = item_size, select "only missing values match"

        iv.    Groping inconsistent values "LF, low fat and reg" to the new values "Low Fat and Regular" –
             Rule Engine → $Item_Fat_Content$  LIKE "LF" => "Low Fat"
                        $Item_Fat_Content$  LIKE "low fat" => "Low Fat"
                        $Item_Fat_Content$  LIKE "reg" => "Regular"
                        $Item_Fat_Content$  LIKE "Low Fat" => "Low Fat"
                        $Item_Fat_Content$  LIKE "Regular" => "Regular"
             Replace column "Item_Fat_Content"

        v.      Group by item fat content with their average and highest sales –
             GroupBy → Include Item_Fat_Content, Aggregation → Item_Outlet_Sales (Mean), Item_Outlet_Sales (Max)

vi.       Group by outlet location type with their item types and sales –
GroupBy → Include Item_Type + Outlet_Location_Type, Aggregation → Item_Outlet_Sales (Sum)

vii.       Group by outlet establishment year and outlet size with item sales –
Numeric Binner → Select Outlet_Establishment_Year, Add bin (∞ … 1990), (1990 … 2000), (2000 … ∞)
GroupBy → Include Outlet_Establishment_Year + Outlet_Size, Aggregation → Item_Outlet_Sales (Mean), Item_Outlet_Sales (Sum)

viii.       Group by outlet size, item type and item fat content with item visibility –
GroupBy → Item_Fat_Content + Item_Type + Outlet_Size, Aggregation → Item_Visibility (Max)

ix.       Group by item weight and item type with item sales –
GroupBy → Item_Weight + Item_Type, Aggregation → Item_Outlet_Sales (Max)

Screenshot of Workflow



b.    Identify dependent and independent variables + determine data type.

| Dependent variable | Independent variable |
|---|---|
| Item_Outlet_Sales (Continuous) | Item_Identifier (Categorical) |
| | Item_Weight (Continuous) |
| | Item_Fat_Content (Categorical) |
| | Item_Type(Categorical) |
| | Item_MRP (Continuous) |
| | Outlet_Identifier (Categorical) |
| | Outlet_Establishment_Year (Continous) |
| | Outlet_Size (Categorical) |
| | Outlet_Location_Type (Categorical) |
| | Outlet_Type (Categorical) |

| | Item_Visibility (Continuous) |
|---|---|

c. Find relevance of each independent variable for prediction of dependent variable.

    i. Outlet_Location_Type + Item_Type → Item_Outlet_Sales

    The sales can be affected by the outlet location type and the item type because we can find out the item types that sell best at every outlet location.

    ii. Outlet_Establishment_Year + Outlet_Size → Item_Outlet_Sales

    The sales can be affected by the year the outlet is established as well as the outlet size. This is because customers might prefer new and larger stores rather than older and smaller variants.

    iii. Item_Fat_Content → Item_Outlet_Sales

    The sales of a product may depend on its fat content. However, this does not mean that more is better or vice versa.

    iv. Outlet_Size + Item_Type + Item_Fat_Content → Item_Visibility

    The visibility of a product can be affected by the size of the outlet as well as the product type and its fat content. This is because larger outlets store more products, making them more available to the customers. However, the visibility alone will not determine the sales of a product as higher product visibility might be a sign that that product did not leave the shelf often.

    v. Item_Weight + Item_Type → Item_Outlet_Sales

    The sales of a product can be affected by its weight and type. This is because we can determine the weight of a given product that sells best.

d. Develop (5) store/product-level hypothesis.
- Stores which are very big in size should have higher sales as they act like one-stop-shops and people would prefer getting everything from one place.
- Low fat products have higher sales compared to product with regular fat content.
- Stores located in urban, or Tier 1 cities should have higher sales because of the higher income levels of people there.
- Medium sized stores have higher visibility of household items because the low demand for them.
- Products that weigh heavier have higher sales because they pack more content.

## 2.1 **CONSTRUCTING A PREDICTION DATA SET**

### Step 1 – File reading

1. KNIME was used to execute this task. The "File Reader" node was used to read the pre-processed.csv file.
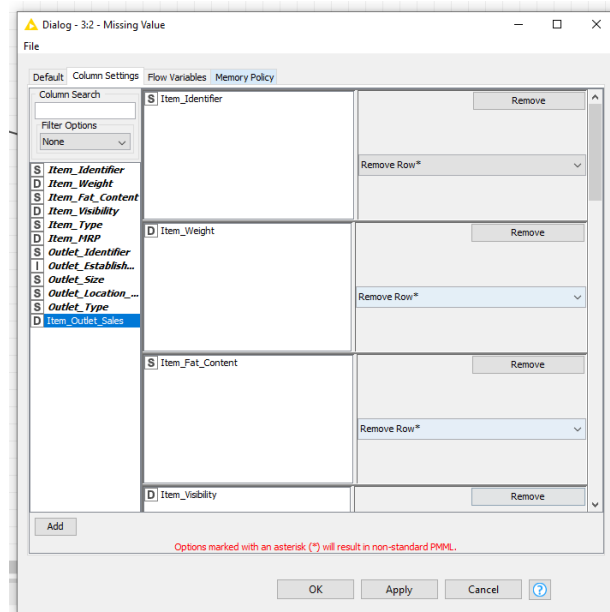
Result of read file



**Input strategies:**

- KNIME new workflow - pre-processed.csv was loaded into KNIME by importing it
- File Reader – This node needed some configuration, such as setting the file location so that the reader can locate it
- After the file is read, certain columns had missing values indicated with a ?
- To remove the missing values from the file, step 2 is used and explained below

### Step 2- Removing missing values

2. Using the "Missing Value" node, missing values in all the rows were filtered out. This filter needed some parameters to be inputted such as (Remove Row) for all column. After setting the parameters, the node would automatically remove the missing values row

**After filtering the missing values**



**Input strategies:**
- Missing value – The parameters to remove the missing values in the file is explained above.
- Results – In the 2$^{nd}$ picture, we can see that there is no "?" anymore, which confirms that missing values have been removed
- To help with the visualization process, visual representations are used to assist in interpreting the data because reading from a csv file may be confusing for certain people
- The visualization representations are shown in step 3 such as using charts

## Step 3 – Visual representations
3. Both the bar & pie chart was used to represent the data visually. There simply is too many data to show so 3 category types were chosen. The 3 categories are as below
Outlet type, Location type, and Item type followed by the graphs below.

Outlet_Location_Type

● Tier 1 ● Tier 2 ● Tier 3

1,860.00    1,860.00

930.00



Item_Type

● Occurrence Count

## Explanations

These explanations help to describe the transformation of data from numerical to categorical data better. For the **outlet type,** bar charts are used to represent this data type

Outlet Type

Based on this data type, there are only 2 types of supermarkets which consists of (Type 1 & Type 2). From the graph, we can see that **Supermarket Type 1** has higher occurrences compared to Type 2 from various aspects such as the item type, outlet establishment year, outlet size, and many other attributed that can be found from the file. However, this graph gives an overview of the data.

Outlet Location Type
Based on this data type, there are 3 main outlet location types (Tier 1 , Tier 2, Tier 3). Both Tier 1 and Tier 3 has almost an equal amount of 1860 occurrences. Tier 2 has the least occurrences which is 930. Based on this data, we can conclude that this company has more Tier 1 and Tier 3 supermarkets in the particular region probably due to the amount of customers around the area. In addition, Tier 2 supermarkets consist of a smaller demographic which can explain why Tier 2 has less supermarkets

Item Type
Based on this data type, there are 3 mostly bought items from the store consisting of (Fruits & Vegetables , Snack foods, Household).The highest bought is fruits and vegetables followed by snack foods, and household items. Based on this data, we can conclude that there's a high number of healthy eaters and unhealthy eaters at the same time. Healthy eaters are slightly more than unhealthy eaters.

From the household items, we can conclude there the area is a matured area, with more people living in the area which is why household items is the 3rd most bought items.

## Step 4. Partitioning the data
The following step was to begin partitioning the data set. Why partition the data set ? According to *(Data Partition - Statistics.com: Data Science, Analytics & Statistics Courses, 2021)*, partitioning a data is used when the data scientist is trying to choose a type of model from a broad list of models. The main concept of partitioning the dataset, is to keep a small percentage of the data to be used for verification purposes after the model has been developed.
Moreover, data partitioning is normally used in supervised learning in data mining which is an overarching concept under the **Data Preparation** phase. Normally, the performance of the model will be evaluated using the training data set.

To put this into context of the project, the partitioning step for the data is explained below.
4.  The data is partitioned into 70% training data and 30% test data. The "Partition" node then links to the "Shuffle" node to randomize the data

    **Normalization** node was used to normalize data.
    •The upper gate puts standardized data into the training set.
    •The light blue square gate below is to repeat the same normalized data for another node "Normalization (Apply)" node applies normalization from "Normalization" and get 30% of original file data from the lower node of "Partitioning".

    **CSV Writer**: 2 nodes
    • Top node is connected to the arrow gate of Normalization to get 70% of data has been shuffled and normalized -> train.csv.
    • Bottom node is connected to the arrow gate of "Normalization (Apply)" to get 30% of normalized data ->test.csv file

This image shows the entire workflow from reading the file to partitioning the dataset.

## 3.1 PROS AND CONS OF TREATING THE OUTPUT VARIABLE AS EITHER CATEGORICAL OR CONTINUOUS

Categorical Data.

**Pros**
Among the pros of treating the output variable as a Categorical Variable is that it will be easier to present it to the stakeholders. If the output variable is treated as a categorical variable the most probable way to analyse it will be through grouping. By grouping these data we get to compare them to the predictor variables(independent variables) which will allow for an easier understanding or visualization of the data. As data visualization with grouped data can be done in the form of pie charts and bar charts. Hence, it will be easier for the stakeholders to understand the data that is being presented to them. Below are examples.

**Total Sales Based On Item Type**
Training Data

Snack Foods
Sum(Item_Outlet_Sales)  100.70753545519895



**Pie Chart For Outlet Location Type and Number of Sales**
Training Data Set

Tier 1 285.92

Furthermore, it is easier to identify relationships between the independent variable and output variable when we make the output variable as a categorical variable. This is also because of how the data has been reduced from a large amount of data into small group which makes analysis of the relationship easier.

## Cons

Large amounts of data goes missing. This is because categorical variables usually show us the data/value in categories. In this case it only shows us the total sum of each item type based on the total number of occurrences of the said item type. For example, let us say that in the year 2005 due to economic problems in the country the value of bread was lower than its usual. This value will not be shown as all the values would be added up together when the data is presented. This could lead to predicting future sales of these kind of item to be affected.

Besides that, most prediction models require a numeric form of data to conduct analysis. For example, most of the algorithms (or ML libraries) produce better result with numerical variable. In python, library

"sklearn" requires features in numerical arrays. If a categorical variable is added there will be an error and the programme would fail. Since, it is harder to get categorical variables to fit in a regression equation and regression equations are essential to machine learning models(prediction models). In this case all the categorical variables have very large gaps between each other and have very little plots on the regression charts therefore it would be hard to plot a line of best fit or so on to later have the model use to analyse and predict an outcome from. Even if we did not categorized(group) the output variable there will still be no pattern in the regression equation to analysis to identify a relationship from. Below are the examples.



Plots are too random and far away from each other and does not show a definitive relationship even though there



Barely any plots available for analysis and plots are too far away. It also does not show a relationship. Even though there is one.

Sales vs Location Type
Training Data Set

No identifiable

Furthermore, the task of converting categorical data to numeric data through encoding is too tedious and time consuming which in turn may increase operational cost for creating a prediction model as the data science team may need to put in more manpower and time to convert all the categorical variables to numeric variables(binary for example) via an encoding method such one hot encoding since there are some categorical variables that have ordinal attributes.

Continuous Variable

**Pros**

Among the pros of treating the output variable as a continuous variable is that it can be fitted directly into a regression equation which can later be used for prediction modelling. Since regression equations are best or more efficient when using continuous variables(numeric data).This is because the regression models are mainly used to compare two numeric variables and these two variables will produce a pattern which can later be analysed to identify the relationship unlike categorical variables. Below are the examples.



Scatter plot for Sales vs MRP
Training Data Set

There is a pattern that corelates to a positive linear regression.

**Sales vs Item_Weight**
Training Data Set

Pattern is present however to deduce relationship a line of best fit should be drawn.



**Sales vs Item_Visibility**
Training Data Set

There is also an identifiable pattern in this regression equation.

In conjunction with that, we can say that using a continuous variable will save time as we would not have to encode the variables manually into a format that the machine would understand. This would direct improve the time taken to make a prediction model. Which in turn would lead to increased productivity reduced cost of production as the time taken to create the prediction model would be significantly increased.

Not only that but continuous variable data can also be visualized in a simpler form to allow stakeholders with no technical knowledge in data science to understand the data. Such forms of visualizations are histograms and box-plots.

**Cons**

Not many cons when compared to categorical data. However, the only disadvantage that cam up during the analysis is that it requires a lot of preparation and have to be really thorough as continuous variables/data are highly sensitive. For example, in the data there might be values that do not carry any significance for example Item_Visibility = 0 (This cannot be the case as if the item is not visible then nobody will know that it can be bought). In categorical variables this is not a problem as when the values are grouped(summed up/averaged/mean) the 0 would gradually disappear as it would be added to another value but this is not the case in continuous variables as it is really sensitive to the input that is given to it therefore even the smallest change it can detect it. Therefore, if not careful then it can effect the prediction model and give the incorrect output. Below is an example.



4.Justify several business values to be gained from the ability to automatically predict the expected sales of a product [15 marks]

Thanks to the rising popularity of Data Science these days, businesses are starting to get more and more interested in being able to use data to predict expected sales on products they sell and with good reason. Being able to predict the expected sales of a product brings numerous benefits for a business. Below are just a few business values that businesses can gained from the ability to automatically predict the expected sales of a product:

**Better Allocation of Resources**

With the ability to be able to predict the expected sales of a product, a business can work around this knowledge and maximize the resources available to them. All businesses have a limited amount of resources and capital that is available to them. By having the ability to predict the expected sales of a product, the company can better utilize what finite resources they have to improve their profits.

Example if a business is putting the same amount of resources into two products, but one product is expected to do have better sales, the business can now choose to put more resources into the more profitable product as that will be a better use of its resources or choose to reduce resources spent on hiring people. Example of these resources could be Man power, as in the number of workers working to produce that product, and machinery that is used to create that product.

**Improve Manufacturing and Production**
If a businesses is able to predict the amount of sales a product will have, they can plan their manufacturing to suit accordingly. Appropriate planning in resources in both Man Power and Raw materials can be forecasted in line with the predicted expected sales. Once they know how much a product is expected to sell, they have no use in overbuying raw materials for manufacturing of the product, because that would just be a waste in capital. The company then can able to capitalize the "just in time" inventory to improve their cash flow.

By being able to predict the expected sales businesses can also calculate how much man power and machinery that they need to meet said sales. That means businesses can allocate the excess manpower and machinery to other parts of the business to improve efficiency.

An example we can use is with a furniture company. The furniture company makes assembled furniture and sells it to customers. To make the furniture we need workers and machinery to cut and design the raw material into furniture to sell. If the furniture company has the ability to predict expected sales of each furniture, They can allocate workers/machinery to help with assembling the expected more popular furniture.

**Improve Negotiations with Vendors**
Once a business knows the expected sales of a product, they are able to have a more informed decision when negotiating with vendors for supplies. With the ability to better forecast their Sales from , there can better negotiate with their respective Vendors for both The Cost ( Volume Discount ) and also the delivery of their raw materials.

This will enable the Business to have upper hand when negotiating with vendors. By being able to get a grasp on the expected sales of a product, a business can schedule with vendors and work together with them to get a lower cost and supplies they need on time.

**Better Utilization of Capital**
Capital in business terms is the amount of money they have available to use for daily operations and funding for future growths. With being able to know expected sales of a product, businesses can better utilize capital available to them to maximize their profits and planned for future growth.

Example with knowing the expected sales of a product, a business can estimate the amount of capital needed to be allocated to manufacture that product and choose to allocate excess capital elsewhere. In addition, the company can also use the forecast Sales to plan for future expanding accordingly.

With McDonalds for example, if they are able to predict the expected sales from each branch, they allocate any capital that some branches don't need and use that capital in future projects example creating more branches.

**Improve on Cashflow and Forecast Revenue**
Cash flow in layman terms is the amount of money going in and coming out from a business. By knowing expected sales, a business can get a better understanding on its cashflow specifically the amount of profit that is going to come back to the business. With this knowledge businesses can more confidently make decisions on putting more money in to the business knowing that it will be able to get its money back. Cash flow is harder to give example because it is more of a business concept.

**Marketing**

With the knowledge on expected sales, a business can make a better decision regarding marketing of a product, both in advertising and promotion. Once a business knows the expected sales, the business will now know which product to focus its marketing on. If a product is expected to have higher sales compared to other products, a business will naturally focus more on advertising and promoting the product as it will be overall more profitable for the business.

An example of this could be with the car brand Toyota. Toyota like every other company has a budget on the marketing they can do on their cars. With the ability to predict expected sales of each car, they can focus their budget on marketing the less popular models to enhance sales.

# 4.1 REFERENCES:

Statistics.com: Data Science, Analytics & Statistics Courses. 2021. *Data Partition - Statistics.com: Data Science, Analytics & Statistics Courses.* [online] Available at: <https://www.statistics.com/glossary/data-partition/#:~:text=Partitioning%20is%20normally%20used%20when,for%20verification%20of%20the%20model.> [Accessed 4 May 2021].

| Contribution | Members |
|---|---|
| Part I | Ethan |
| Part II | Joel |
| Part III | Bernard |
| Part IV | Lionel |
| Part V | Joel, Bernard, Lionel |