

# About Dataset

The data were collected in 2013 for 1559 products across 10 stores in different cities. Also, certain attributes of each product and store have been defined.

## Attributes

- **Item\_Identifier** -> Unique Product ID
- **Item\_Weight** -> Weight of the product
- **Item\_Fat\_Content** -> Whether the product is low fat or not
- **Item\_Visibility** -> The % of the total display area of all products in a store allocated to the particular product
- **Item\_Type** -> The category to which the product belongs
- **Item\_MRP** -> Maximum Retail Price (list price) of the product
- **Outlet\_Identifier** -> Unique store ID
- **Outlet\_Establishment\_Year** -> The year in which the store was established
- **Outlet\_Size** -> The size of the store in terms of ground area covered
- **Outlet\_Location\_Type** -> The type of city in which the store is located
- **Outlet\_Type** -> Whether the outlet is just a grocery store or some sort of supermarket
- **Item\_Outlet\_Sales** -> Sales of the product in the particular store.

# Data cleaning

1. Using the count blank function to count the null values for each column, there are 1463 missing values in the item weight column and 2410 in the outlet size column.
2. Use the countif function to count the number of rows in the item visibility column that has a 0 value.
3. Fill null values in the item weight column with the item weight value for each item identifier.
4. Fill null values in the item visibility column with the mean.
5. Remove the outlet size column because it has many null values, and the column is less important for analysis.
6. Replace the “lf” value with ‘Low Fat’ in item fat content.
7. Replace the “reg” value with “Regular” in the item fat content column.
8. Convert the type of item\_mrp and item\_outlet\_sales columns to currency type.