

# Business Report

SMDM Project Business Report DSBA

*Capstone project – Supply Chain  
Management*



***Sanjay Srinivasan***

***PGP-DSBA Online***

***JULY' 21 Batch***

***Date: 29-05-2022***

## **INDEX**

<b>S. No</b>	<b>Contents</b>	<b>Page No</b>
<b>1)</b>	<b>Introduction of the business problem</b>	<b>4</b>
	a) Defining problem statement	4
	b) Need of the study/project	4
	c) Understanding business/social opportunity	4
<b>2)</b>	<b>Data Report</b>	<b>6</b>
	a) Understanding how data was collected in terms of time, frequency and methodology	6
	b) Visual inspection of data (rows, columns, descriptive details)	7
	c) Understanding of attributes (variable info, renaming if required)	8
<b>3)</b>	<b>Exploratory data analysis</b>	<b>10</b>
	a) Univariate analysis (distribution and spread for every continuous attribute, distribution of data in categories for categorical ones)	10
	b) Bivariate analysis (relationship between different variables , correlations)	14
	c) Removal of unwanted variables (if applicable)	16
	d) Missing Value treatment (if applicable)	17
	e) Variable transformation (if applicable)	18
	f) Addition of new variables (if required)	18
	g) Addition of new variables (if required)	20
<b>4)</b>	<b>Business insights from EDA</b>	<b>21</b>
	a) Is the data unbalanced? If so, what can be done? Please explain in the context of the business	21
	b) Any business insights using clustering (if applicable)	21
	c) Any other business insights	23

## **List Of Figures**

<b>S.No</b>	<b>Content</b>	<b>Page No</b>
1.1	Sample data of product weight across years	6
1.2	Number of times frequently product refilled across zone	6
1.3	Sample dataset	7
1.4	Shape of the dataset	7
1.5	Description of the data.	8
1.6	Variable info of the dataset.	8
1.7	Missing value count in the dataset.	9
1.8	Missing value count in percentage.	9
1.9	Univariate Analysis	11
1.10	Count plot for approved_wh_govt_certificate	12
1.11	Count plot for Location_type	12
1.12	Count plot for WH_capacity_size	12
1.13	Count plot for Zone	13
1.14	Count plot for Warehouse in regional zone	13
1.15	Count plot for Warehouse owner type	13
1.16	Boxplot for Product weight ton vs. zone across location type.	14
1.17	Stripplot for Product weight ton vs. flood_impacted.	14
1.18	Violinplot for Product weight ton vs. electric_supply with flood_impacted.	14
1.19	Violinplot for Product weight ton vs. wh_regional_zone with wh_capacity_size.	15
1.20	Pairplot for the bivariate analysis.	15
1.21	Sample Correlation data for the bivariate analysis.	16
1.22	Plotting correlation in the heatmap for the bivariate analysis.	16
1.23	Shape of the dataframe after dropping unwanted variables	17
1.24	Info of the dataframe after dropping unwanted variables.	17
1.25	Missing values before treating.	17
1.26	Missing values after treating	18
1.27	Shape of the dataset after variable transformation.	18
1.28	Sample dataset after variable transformation.	18
1.29	Info of the dataset after variable transformation.	19
1.30	Transforming right skewed target variable.	19
1.31	Finding skewness of the data frame.	20
1.32	Boxplot before scaling	20
1.33	Boxplot after scaling	21
1.34	Dendrogram of hierarchial clustering	21
1.35	Dendrogram of hierarchial clustering after Truncating	21
1.36	Finding inertia from K-means clustering	22
1.37	Elbow curve for K – means	22
1.38	Silhouette score	22
1.39	Silhouette width	22
1.40	Silhouette Sample value	23

## Introduction of the business problem

### a) Defining problem statement:

The objective of this report is to find that, how the machine learning model supports the supply chain to overcome the demand and supply mismatch in every zone. A FMCG company has entered into the instant noodles business two years back. The data is gathered based on the FMCG Company's demand and supply mismatch of the product instant noodles. The higher management has noticed that there is a mismatch in the demand and supply of instant noodles.

The demand and supply mismatch can be overcome by following these: first of all, finding the demand and supply mismatch. Secondly, find the optimum weight of the product been shipped to each warehouse at different zone and regions of the country.

Drawback of demand and supply mismatch:

1. Company will lose heavily on logistic movement of goods / products
2. In order to sale the product, goods has to be moved where there is high supply or high demand zone.

### b) Need of the study/project:

- 1) To find the mismatch in demand and supply so that management can optimize the supply quantity in each and every warehouse in entire country.
- 2) The optimum weight of the product shipped every time based on the demand on each zone or region.
- 3) Data analysis will help us to analysis the product quantity sale based on the zone wise demand.
- 4) Loss can be minimized by the management based on data analysis.

To meet demand and supply:

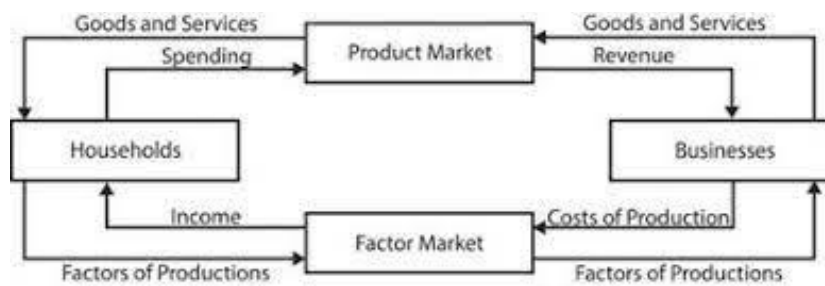
1. Promotion of products through advertisements in low demanding zone.
2. New marketing strategies to be implemented based on low and high demand of product zone. Marketing strategies are:
  - Posters
  - Pamphlets
  - Offers
  - Sign Boards
  - Health benefits of the product

### c) Understanding business/ social opportunity:

1. The demand patterns of the customers across different zones can help the company in understanding the customer behaviour and needs of the customer in every zone.
2. By this, we can forecast the demand of the products in every zone and needed quantity can be shipped to every dealer/ retailer based on the demand in every zone (Resources in demand & shipping of the product).

3. The company can monitor the performance of the goods/product in every zone and offers based on low and high selling zone.
4. The optimum maintenance and cash margin can be given based on the demand of the product in every zone.
5. Based on the demand of the product, the stocks can be refilled in the warehouses in every zone.
6. Rotation of Stocks regularly helps the company to reduce the loss, improves the efficiency and product demand can be improved.
7. This helps the business to optimize the performance.

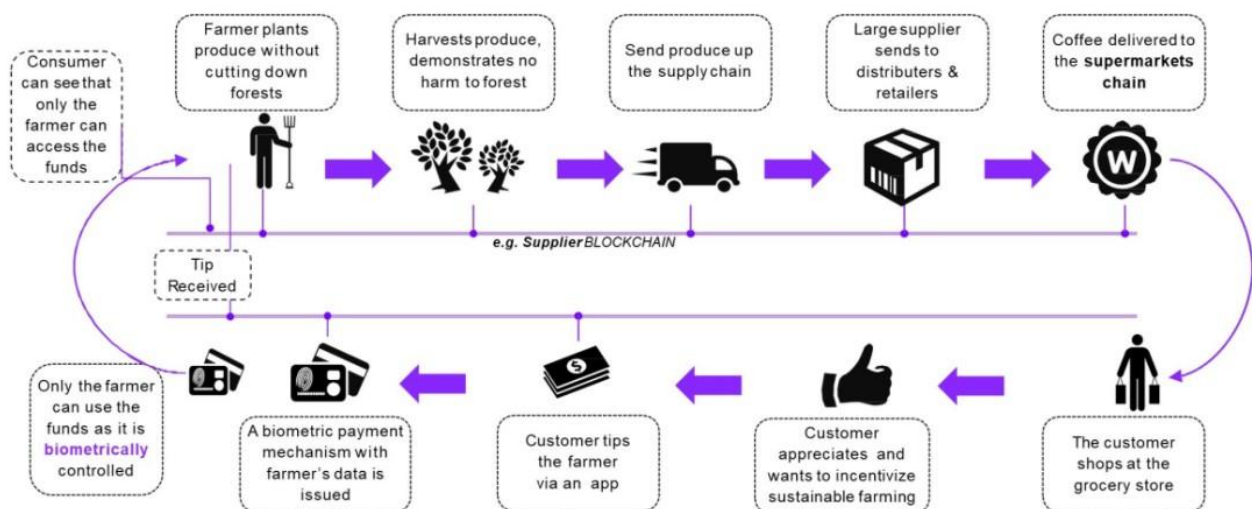
### Without Supply chain and logistics:



### With Supply chain and logistics:

## **CIRCULAR SUPPLY CHAIN – HOW IT WORKS**

Farming Example



## Data Report

### a) Understanding how data was collected in terms of time, frequency and methodology.

#### **Methodology:**

7 Data Collection Methods Used in Business Analytics are :

1. Surveys
2. Transactional Tracking
3. Interviews and Focus Groups
4. Observation
5. Online Tracking
6. Forms
7. Social Media Monitoring

Transactional Tracking and Survey type of method are used in the data collection of this dataset for the FMCG Company.

#### **Time:**

wh_est_year	zone	product_wg_ton
1996.0	East	34614.500000
	North	35020.408451
	South	34025.018182
	West	35422.793651
1997.0	East	38114.000000
	North	33936.573643
	South	35382.709677
	West	35517.509434
1998.0	East	35213.000000
	North	34781.502415
	South	34106.789474
	West	35985.016854

Fig 1.1 Sample data of product weight across years

In the year **1996**, the demand of the product is higher in the **west zone** when compared with the other zone.

#### **Frequency:**

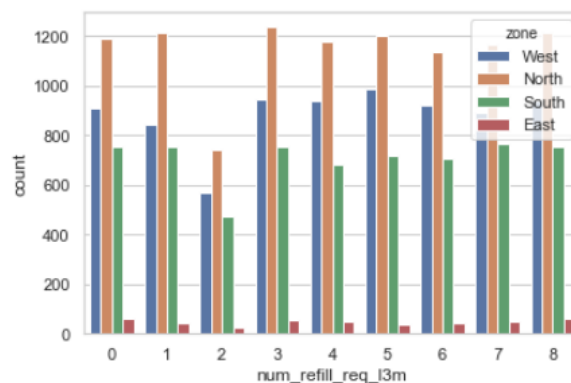


Fig 1.2 Number of times frequently product refilled across zone

North zone has the highest number of times the stock / product has been refilled. From this we can infer that, North zone has high demand of the product.

### b) Visual inspection of data (rows, columns, descriptive details):

Ware_house_ID	WH_Manager_ID	Location_type	WH_capacity_size	zone	WH_regional_zone	num_refill_req_13m	transport_issue_11y	Competitor_in_mkt	retail_
WH_100000	EID_50000	Urban	Small	West	Zone 6	3	1	2	
WH_100001	EID_50001	Rural	Large	North	Zone 5	0	0	4	
WH_100002	EID_50002	Rural	Mid	South	Zone 2	1	0	4	
WH_100003	EID_50003	Rural	Mid	North	Zone 3	7	4	2	
WH_100004	EID_50004	Rural	Large	North	Zone 5	3	1	2	

Fig 1.3 Sample dataset

Dataset consist of 25000 rows and 23 columns. 22 are independent variable and 1 target variable. 'Product\_wg\_ton' is the target column.

Description of each and every variable in the dataset.

Variable	Business Definition
Ware_house_ID	Product warehouse ID
WH_Manager_ID	Employee ID of warehouse manager
Location_type	Location of warehouse like in city or village
WH_capacity_size	Storage capacity size of the warehouse
zone	Zone of the warehouse
WH_regional_zone	Regional zone of the warehouse under each zone
num_refill_req_13m	Number of times refilling has been done in last 3 months
transport_issue_11y	Any transport issue like accident or goods stolen reported in last one year
Competitor_in_mkt	Number of instant noodles competitor in the market
retail_shop_num	Number of retails shop who sell the product under the warehouse area
wh_owner_type	Company is owning the warehouse or they have get the warehouse on rent
distributor_num	Number of distributor works in between warehouse and retail shops
flood_impacted	Warehouse is in the Flood impacted area indicator
flood_proof	Warehouse is flood proof indicators. Like storage is at some height not directly on the ground
electric_supply	Warehouse have electric back up like generator, so they can run the warehouse in load shedding
dist_from_hub	Distance between warehouse to the production hub in Kms
workers_num	Number of workers working in the warehouse
wh_est_year	Warehouse established year
storage_issue_reported_13m	Warehouse reported storage issue to corporate office in last 3 months. Like rat, fungus because of moisture etc.
temp_reg_mach	Warehouse have temperature regulating machine indicator
approved_wh_govt_certificate	What kind of standard certificate has been issued to the warehouse from government regulatory body
wh_breakdown_13m	Number of time warehouse face a breakdown in last 3 months. Like strike from worker, flood, or electrical failure
govt_check_13m	Number of time government Officers have been visited the warehouse to check the quality and expire of stored food in last 3 months
product_wg_ton	Product has been shipped in last 3 months. Weight is in tons

(25000, 23)

Fig 1.4 Shape of the dataset

	num_refill_req_13m	transport_issue_11y	Competitor_in_mkt	retail_shop_num	distributor_num	flood_impacted	flood_proof	electric_supply	dist_from_hut
count	25000.000000	25000.000000	25000.000000	25000.000000	25000.000000	25000.000000	25000.000000	25000.000000	25000.000000
mean	4.089040	0.773680	3.104200	4985.711560	42.418120	0.098160	0.054640	0.656880	163.537320
std	2.606612	1.199449	1.141663	1052.825252	16.064329	0.297537	0.227281	0.474761	62.718600
min	0.000000	0.000000	0.000000	1821.000000	15.000000	0.000000	0.000000	0.000000	55.000000
25%	2.000000	0.000000	2.000000	4313.000000	29.000000	0.000000	0.000000	0.000000	109.000000
50%	4.000000	0.000000	3.000000	4859.000000	42.000000	0.000000	0.000000	1.000000	164.000000
75%	6.000000	1.000000	4.000000	5500.000000	56.000000	0.000000	0.000000	1.000000	218.000000
max	8.000000	5.000000	12.000000	11008.000000	70.000000	1.000000	1.000000	1.000000	271.000000

Fig 1.5 Description of the data.

### c) Understanding of attributes (variable info, renaming if required):

```
<class 'pandas.core.frame.DataFrame'>
Index: 25000 entries, WH_100000 to WH_124999
Data columns (total 23 columns):
WH_Manager_ID                25000 non-null object
Location_type                 25000 non-null object
WH_capacity_size              25000 non-null object
zone                          25000 non-null object
WH_regional_zone              25000 non-null object
num_refill_req_13m            25000 non-null int64
transport_issue_11y           25000 non-null int64
Competitor_in_mkt             25000 non-null int64
retail_shop_num               25000 non-null int64
wh_owner_type                 25000 non-null object
distributor_num               25000 non-null int64
flood_impacted                25000 non-null int64
flood_proof                   25000 non-null int64
electric_supply               25000 non-null int64
dist_from_hub                 25000 non-null int64
workers_num                   24010 non-null float64
wh_est_year                   13119 non-null float64
storage_issue_reported_13m    25000 non-null int64
temp_reg_mach                 25000 non-null int64
approved_wh_govt_certificate  24092 non-null object
wh_breakdown_13m             25000 non-null int64
govt_check_13m               25000 non-null int64
product_wg_ton                25000 non-null int64
dtypes: float64(2), int64(14), object(7)
memory usage: 4.6+ MB
```

Fig 1.6 Variable info of the dataset.

From this we can infer that, 7 variable are object type, 14 variable are integer type and 2 are float type variable. 3 variable has some missing value in the dataset.



WH_Manager_ID	0
Location_type	0
WH_capacity_size	0
zone	0
WH_regional_zone	0
num_refill_req_13m	0
transport_issue_11y	0
Competitor_in_mkt	0
retail_shop_num	0
wh_owner_type	0
distributor_num	0
flood_impacted	0
flood_proof	0
electric_supply	0
dist_from_hub	0
workers_num	990
wh_est_year	11881
storage_issue_reported_13m	0
temp_reg_mach	0
approved_wh_govt_certificate	908
wh_breakdown_13m	0
govt_check_13m	0
product_wg_ton	0

dtype: int64

Fig 1.7 Missing value count in the dataset.

WH_Manager_ID	0.000
Location_type	0.000
WH_capacity_size	0.000
zone	0.000
WH_regional_zone	0.000
num_refill_req_13m	0.000
transport_issue_11y	0.000
Competitor_in_mkt	0.000
retail_shop_num	0.000
wh_owner_type	0.000
distributor_num	0.000
flood_impacted	0.000
flood_proof	0.000
electric_supply	0.000
dist_from_hub	0.000
workers_num	3.960
wh_est_year	47.524
storage_issue_reported_13m	0.000
temp_reg_mach	0.000
approved_wh_govt_certificate	3.632
wh_breakdown_13m	0.000
govt_check_13m	0.000
product_wg_ton	0.000

dtype: float64

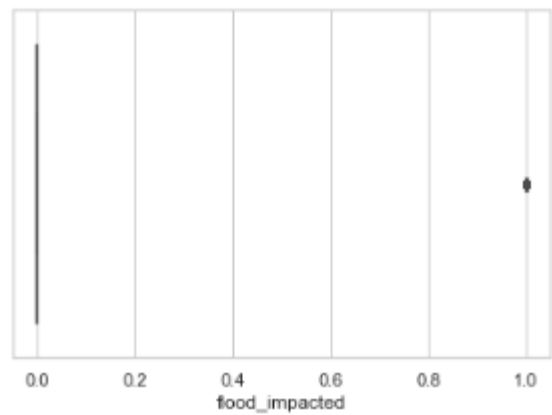
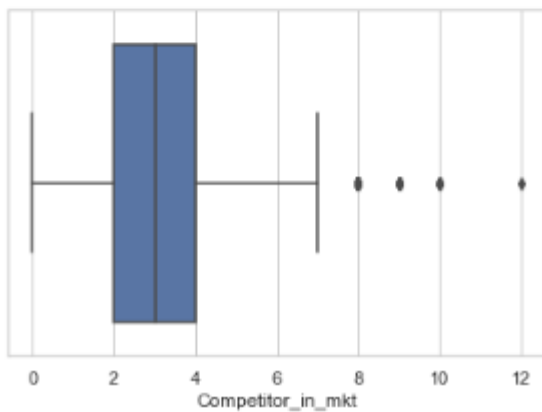
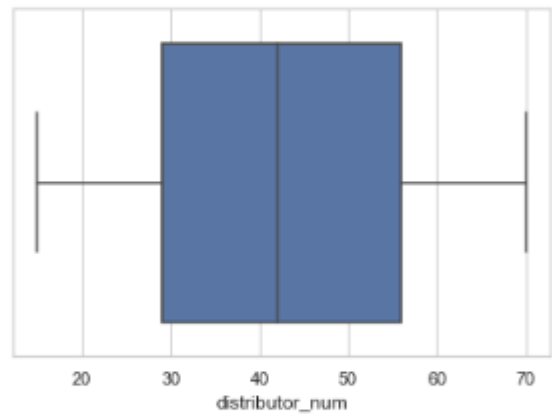
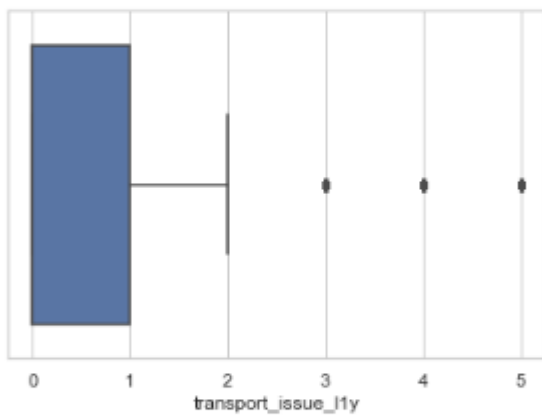
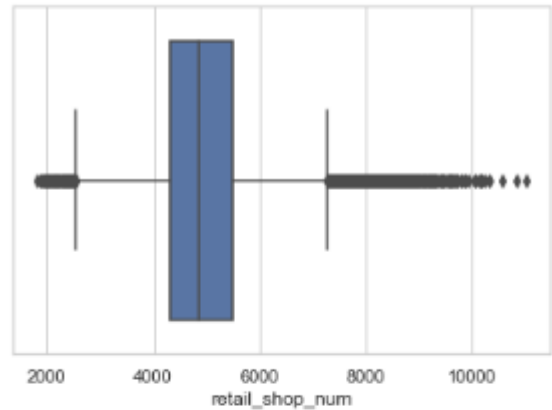
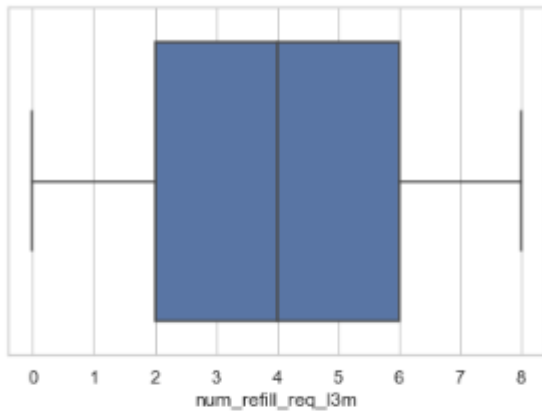
Fig 1.8 Missing value count in percentage.

Around **48 % (11881)** missing values data are present in the “**wh\_est\_year**” variable.

Renaming of the variable is not required for this dataset.

## Exploratory data analysis

a) Univariate analysis (distribution and spread for every continuous attribute, distribution of data in categories for categorical ones):



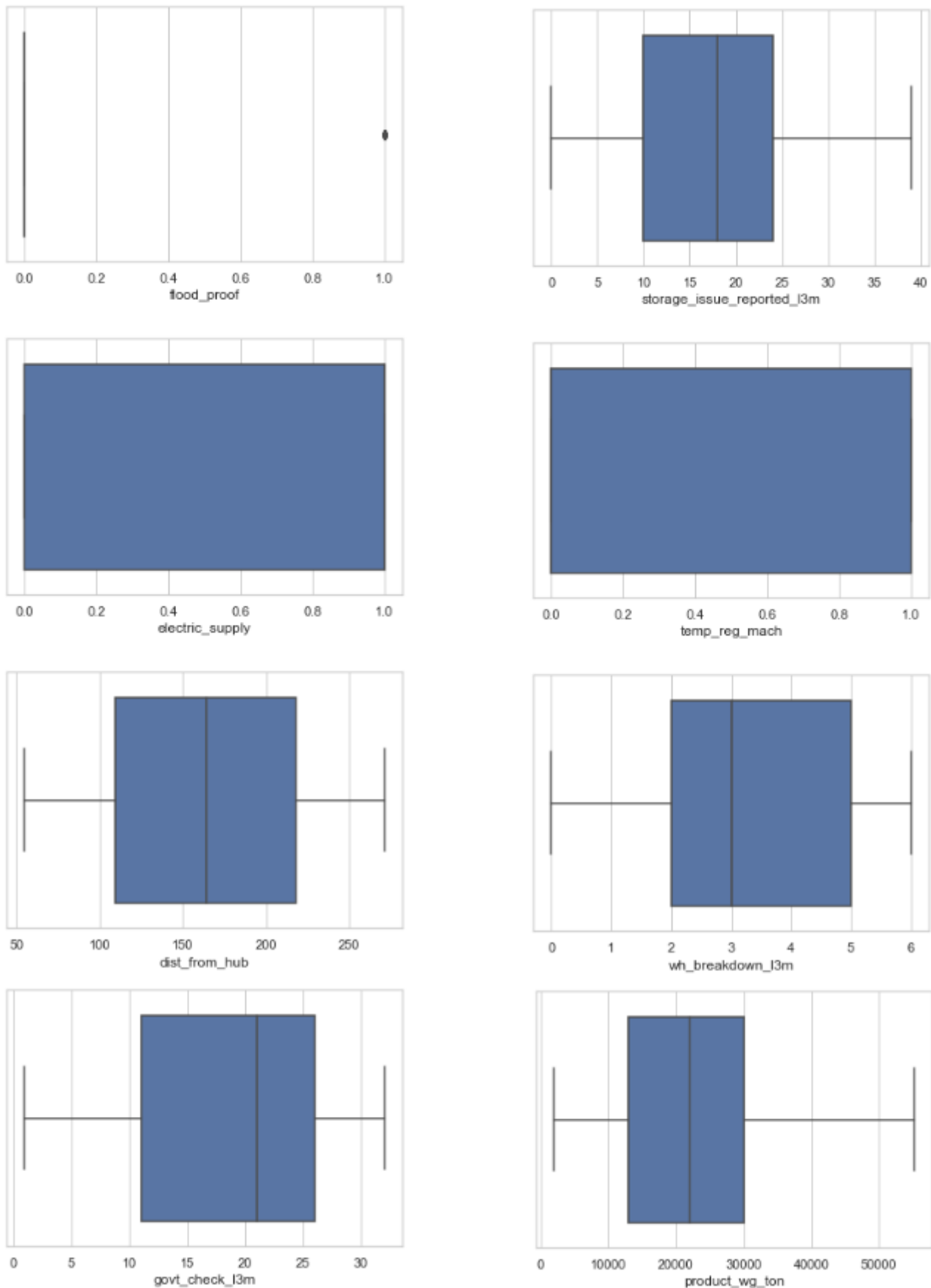


Fig 1.9 Univariate Analysis

**Product\_wg\_ton** – Target column is right skewed and there is no outliers present in the target column.

**Competitor\_in\_mkt** - This independent variable has outlier in the dataset

**transport\_issue\_11y** - This independent variable has outlier in the dataset and right skewed values are present in the dataset.

**retail\_shop\_num** - This independent variable has more outliers and values are slightly right skewed.

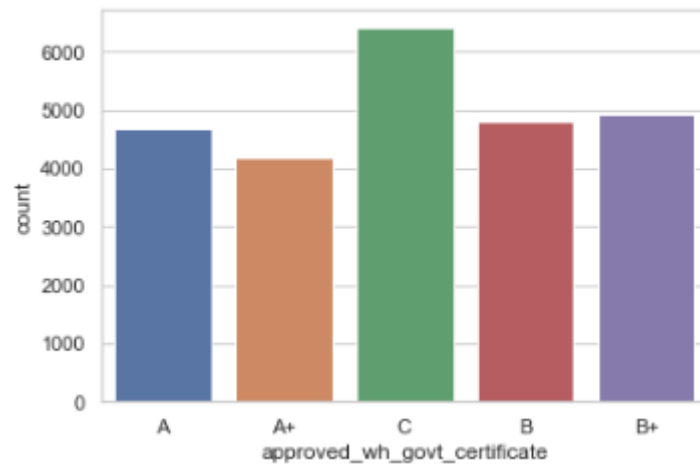


Fig 1.10 Count plot for approved\_wh\_govt\_certificate

From this count plot C certificate has the highest number of warehouse with government certificate A+ has the least number of warehouse government certificate.

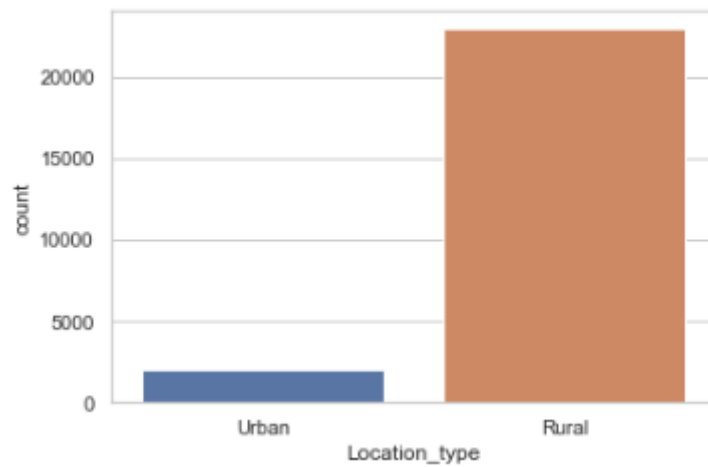


Fig 1.11 Count plot for Location\_type

Most number of ware houses is located in the rural area.

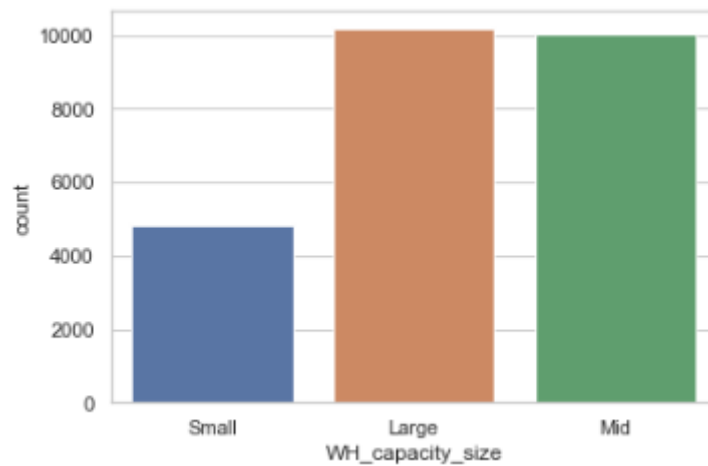


Fig 1.12 Count plot for WH\_capacity\_size

The Large capacity warehouse are having the ware house capacity higher than the other capacity ware houses.

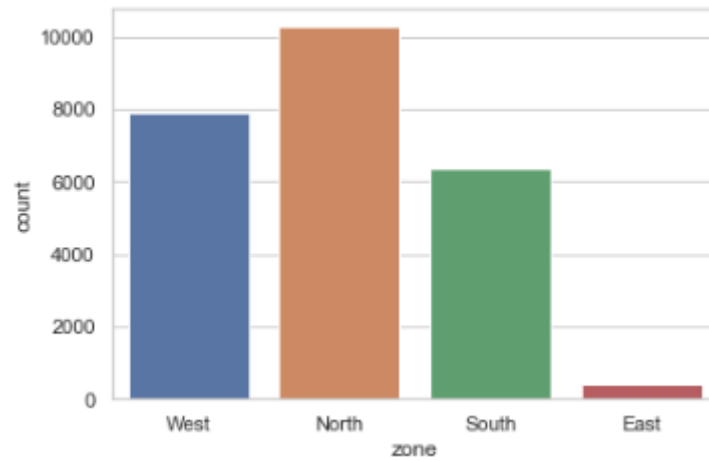


Fig 1.13 Count plot for Zone

North zone has the highest number of warehouse are built and East zone has the least number of ware houses.

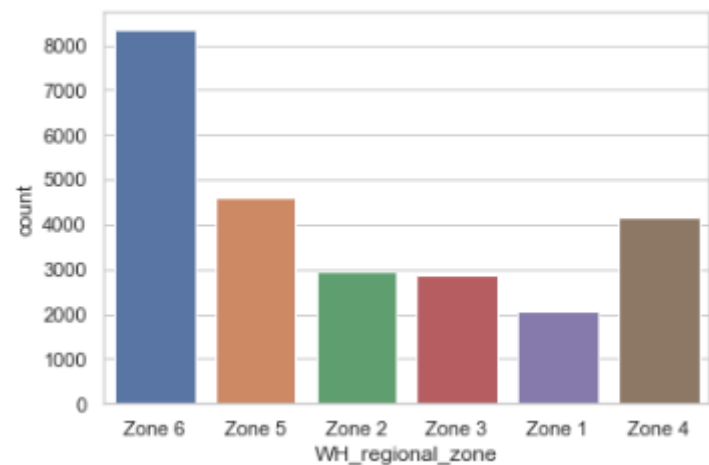


Fig 1.14 Count plot for Warehouse in regional zone

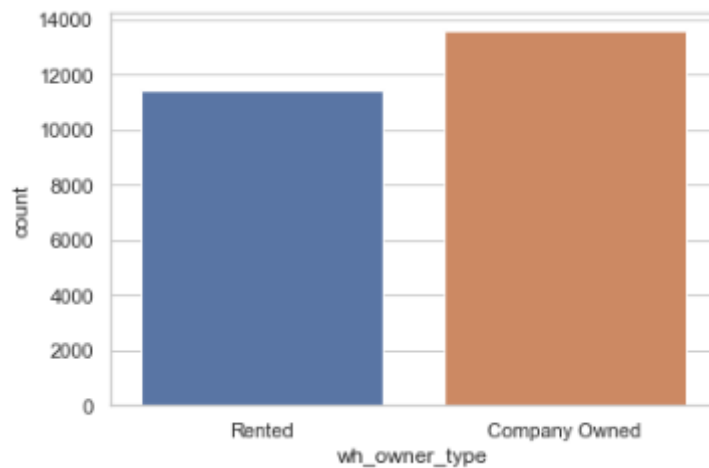


Fig 1.15 Count plot for Warehouse owner type

Most number of warehouse are owned by companies

b) Bivariate analysis (relationship between different variables , correlations):

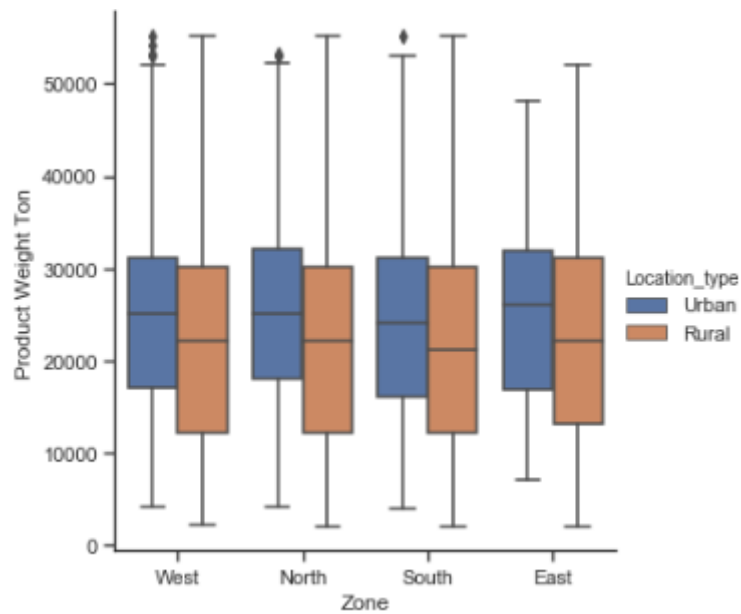


Fig 1.16 Boxplot for Product weight ton vs. zone across location type.

The demand for the product in the urban areas is higher when compared with the rural area. Outliers are present in the urban area, the inference from the above boxplot.

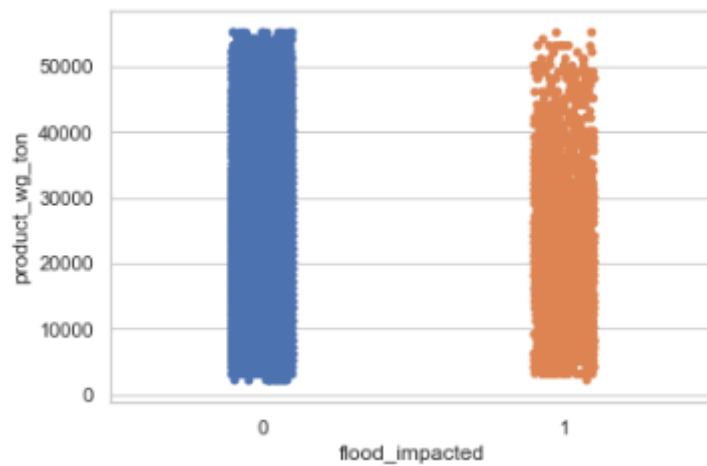


Fig 1.17 Stripplot for Product weight ton vs. flood\_impacted.

The Flood indicator indicates that most of the products in the ware house are affected by the flood.

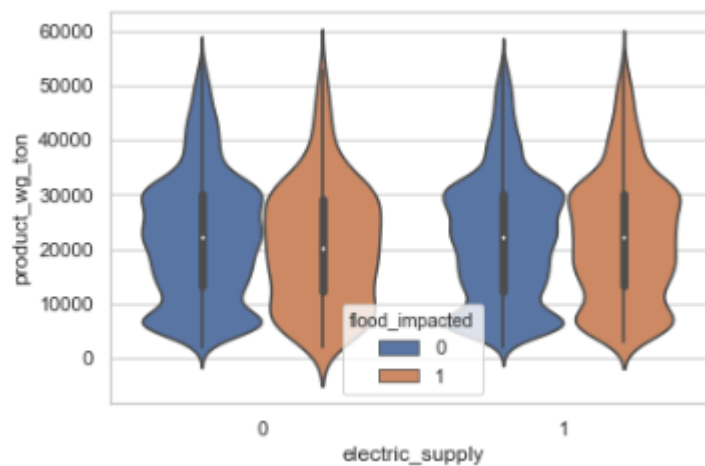


Fig 1.18 Violinplot for Product weight ton vs. electric\_supply with flood\_impacted.

The warehouse is having the backup power supply where the flood impacted area will have the loss of power or when there is power shutdown in the area.

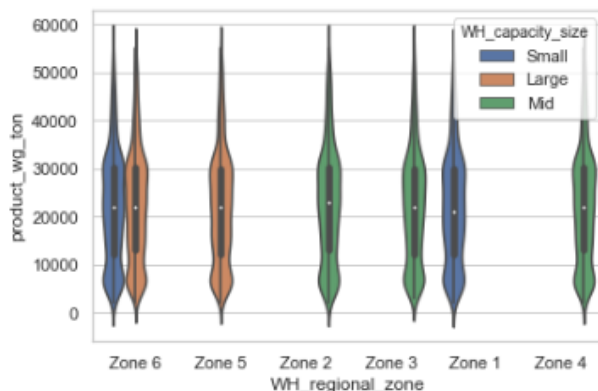


Fig 1.19 Violinplot for Product weight ton vs. wh\_regional\_zone with wh\_capacity\_size.

2 types of ware house capacity size (small and large) are there in zone 6 and other zone has either one type of ware house capacity (small, medium, large).

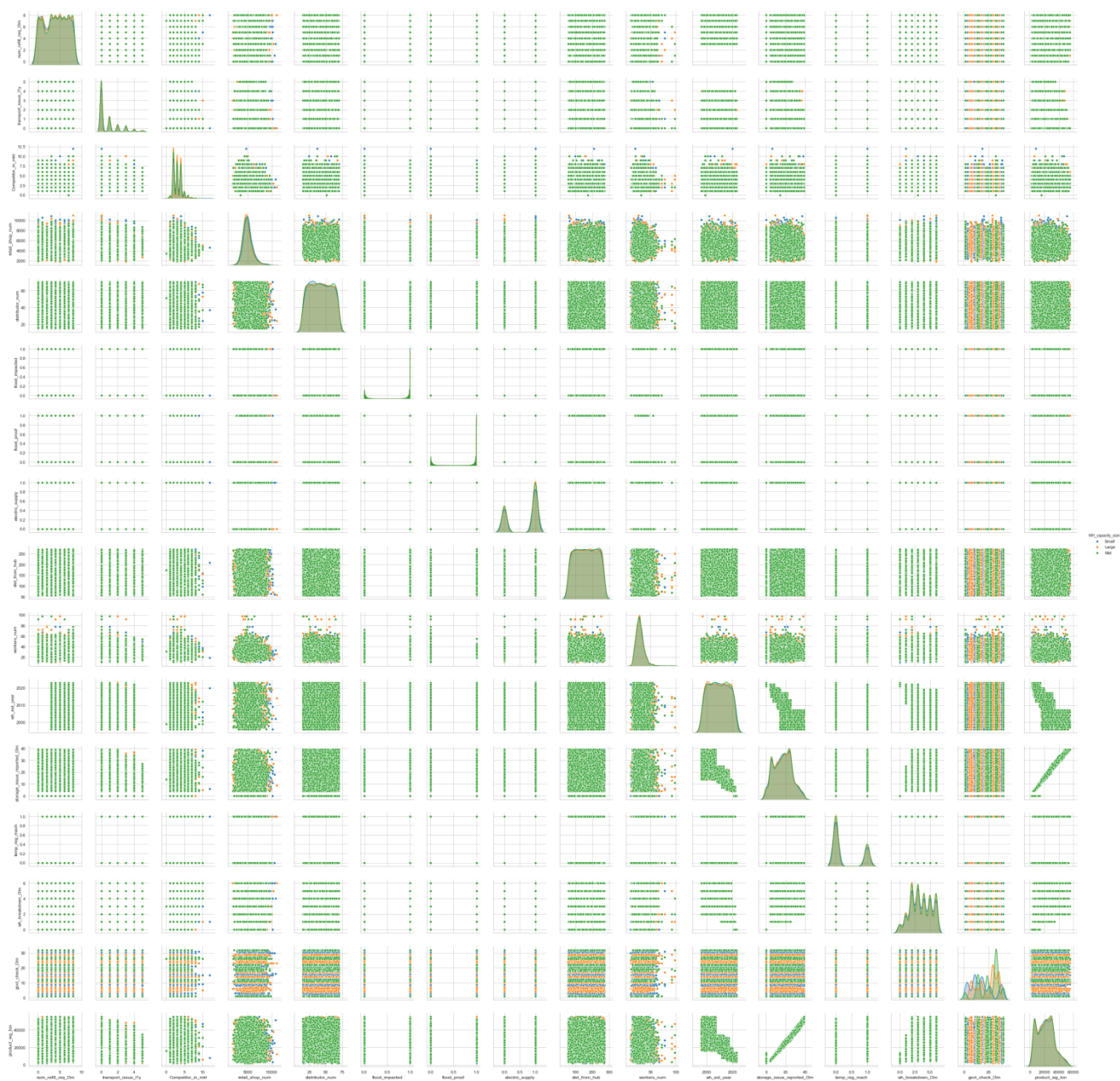


Fig 1.20 Pairplot for the bivariate analysis.

	num_refill_req_13m	transport_issue_11y	Competitor_in_mkt	retail_shop_num	distributor_num	flood_impacted	flood_proof	elect
num_refill_req_13m	1.000000	0.018549	0.002985	-0.001186	0.003995	-0.010548	-0.001123	
transport_issue_11y	0.018549	1.000000	-0.005826	-0.001826	0.008993	-0.009596	0.000022	
Competitor_in_mkt	0.002985	-0.005826	1.000000	-0.156943	-0.001492	0.009338	-0.003444	
retail_shop_num	-0.001186	-0.001826	-0.156943	1.000000	-0.000395	-0.003774	0.007223	
distributor_num	0.003995	0.008993	-0.001492	-0.000395	1.000000	0.004611	-0.003409	
flood_impacted	-0.010548	-0.009596	0.009338	-0.003774	0.004611	1.000000	0.107015	
flood_proof	-0.001123	0.000022	-0.003444	0.007223	-0.003409	0.107015	1.000000	
electric_supply	-0.007959	-0.009299	0.001759	-0.009207	0.000454	0.164815	0.114811	
dist_from_hub	0.000048	0.014336	0.008407	0.000429	-0.011838	0.000749	-0.005315	
workers_num	-0.013764	-0.009004	0.000050	-0.005406	-0.014682	0.168425	0.041228	
wh_est_year	0.015363	-0.012910	-0.011202	0.005721	-0.012295	-0.000668	-0.003329	
storage_issue_reported_13m	-0.006602	-0.144327	0.009543	-0.006632	0.003396	-0.003157	-0.002712	
temp_reg_mach	0.260928	0.018207	0.009524	-0.001273	0.002827	-0.008554	0.005636	
wh_breakdown_13m	0.000608	0.012990	0.012733	-0.008420	0.004286	-0.001744	-0.005151	
govt_check_13m	-0.003302	0.002190	-0.043455	0.045749	-0.007934	0.000587	-0.003600	
product_wg_ton	0.001415	-0.173992	0.008884	-0.006615	0.004999	-0.002299	-0.000441	

Fig 1.21 Sample Correlation data for the bivariate analysis.

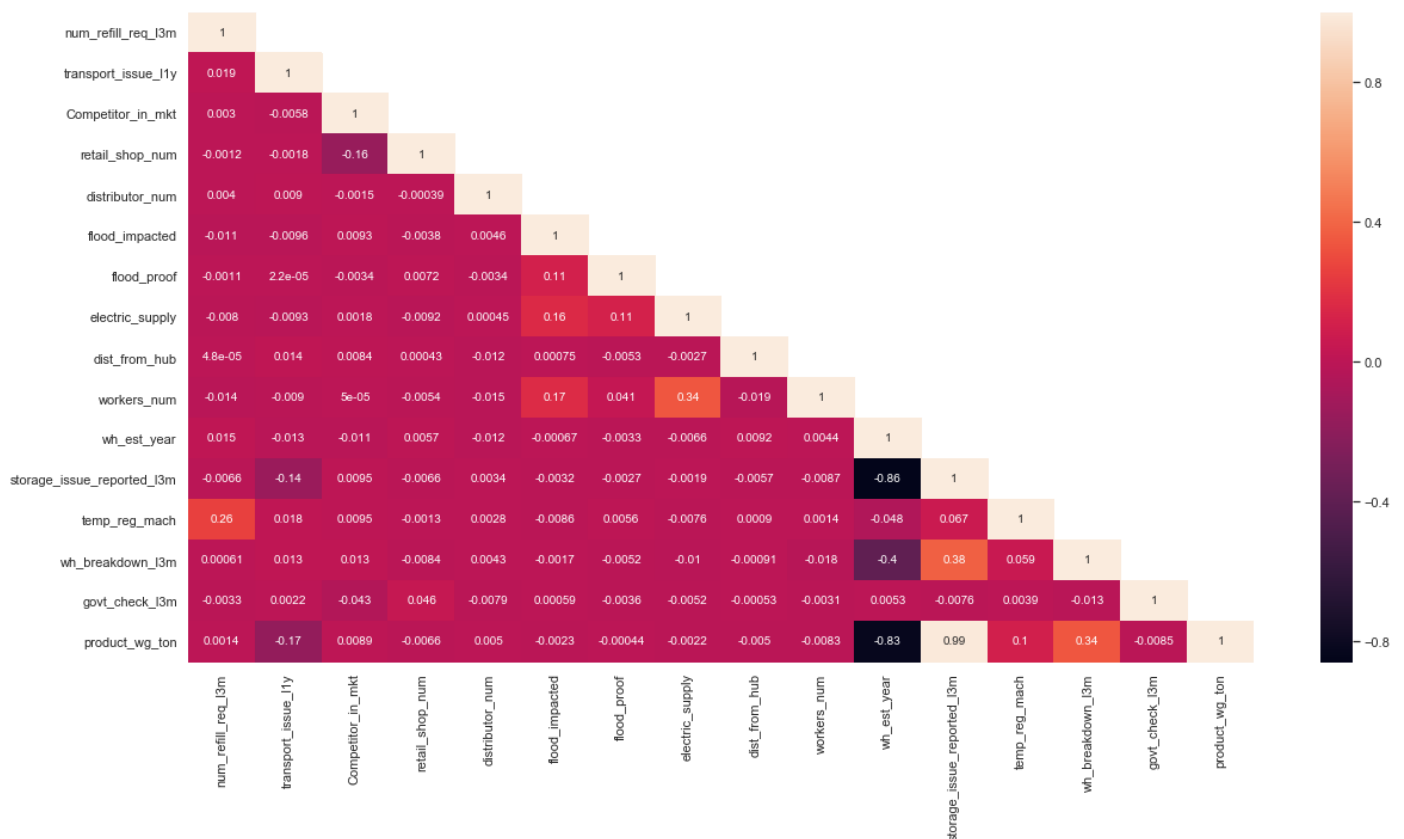


Fig 1.22 Plotting correlation in the heatmap for the bivariate analysis.

From the heat map, we can infer that, '**storage\_issue\_reported\_13m**' (storage issue reported last 3 months) is **highly correlated** with the target variable ('**product\_wg\_ton**').

Wh\_est\_year has **highly negatively correlated** with the target variable ('**product\_wg\_ton**') and **storage\_issue\_reported\_13m** variable

### c) Removal of unwanted variables (if applicable):

As 'WH\_Manager\_ID' and 'Ware\_house\_ID' are unique values, we are dropping 'WH\_Manager\_ID' and setting 'Ware\_house\_ID' as an index value.

As 'wh\_est\_year' is having **48% of null values**, so we are **dropping 'wh\_est\_year'** independent variable from the dataframe.



'storage\_issue\_reported\_l3m' independent variable is highly correlated with the target column. So 'storage\_issue\_reported\_l3m' can also be dropped from the dataframe.

(25000, 20)

Fig 1.23 shape of the dataframe after dropping unwanted variables

```
<class 'pandas.core.frame.DataFrame'>
Index: 25000 entries, WH_100000 to WH_124999
Data columns (total 20 columns):
Location_type                25000 non-null object
WH_capacity_size             25000 non-null object
zone                        25000 non-null object
WH_regional_zone            25000 non-null object
num_refill_req_l3m          25000 non-null int64
transport_issue_l1y         25000 non-null int64
Competitor_in_mkt           25000 non-null int64
retail_shop_num              25000 non-null int64
wh_owner_type                25000 non-null object
distributor_num              25000 non-null int64
flood_impacted               25000 non-null int64
flood_proof                  25000 non-null int64
electric_supply              25000 non-null int64
dist_from_hub                25000 non-null int64
workers_num                  25000 non-null float64
temp_reg_mach                25000 non-null int64
approved_wh_govt_certificate 25000 non-null object
wh_breakdown_l3m             25000 non-null int64
govt_check_l3m               25000 non-null int64
product_wg_ton               25000 non-null int64
dtypes: float64(1), int64(13), object(6)
memory usage: 4.6+ MB
```

Fig 1.24 Info of the dataframe after dropping unwanted variables.

#### d) Missing Value treatment (if applicable):

```
Location_type                0
WH_capacity_size             0
zone                        0
WH_regional_zone            0
num_refill_req_l3m          0
transport_issue_l1y         0
Competitor_in_mkt           0
retail_shop_num              0
wh_owner_type                0
distributor_num              0
flood_impacted               0
flood_proof                  0
electric_supply              0
dist_from_hub                0
workers_num                  990
temp_reg_mach                0
approved_wh_govt_certificate 908
wh_breakdown_l3m             0
govt_check_l3m               0
product_wg_ton               0
dtype: int64
```

Fig 1.25 Missing values before treating.

```

Location_type      0
WH_capacity_size   0
zone               0
WH_regional_zone   0
num_refill_req_13m 0
transport_issue_11y 0
Competitor_in_mkt  0
retail_shop_num    0
wh_owner_type      0
distributor_num    0
flood_impacted     0
flood_proof        0
electric_supply     0
dist_from_hub      0
workers_num        0
temp_reg_mach      0
approved_wh_govt_certificate 0
wh_breakdown_13m   0
govt_check_13m     0
product_wg_ton     0
dtype: int64

```

Fig 1.26 Missing values after treating

Missing values are treated using mean value of integer type and mode value for object type variable and there are no missing values present in the dataset.

#### *e) Outlier treatment (if required):*

From the boxplot we can infer that, outlier treatment is not required for this dataset.

#### *f) Variable transformation (if applicable):*

(25000, 30)

Fig 1.27 Shape of the dataset after variable transformation.

Ware_house_ID	Location_type_Urban	WH_capacity_size_Mid	WH_capacity_size_Small	zone_North	zone_South	zone_West	WH_regional_zone_Zone_2	WH_regional
WH_100000	1	0	1	0	0	1	0	
WH_100001	0	0	0	1	0	0	0	
WH_100002	0	1	0	0	1	0	0	1
WH_100003	0	1	0	1	0	0	0	0
WH_100004	0	0	0	1	0	0	0	0

Fig 1.28 sample dataset after variable transformation.

```

<class 'pandas.core.frame.DataFrame'>
Index: 25000 entries, WH_100000 to WH_124999
Data columns (total 30 columns):
Location_type_Urban                25000 non-null uint8
WH_capacity_size_Mid               25000 non-null uint8
WH_capacity_size_Small             25000 non-null uint8
zone_North                        25000 non-null uint8
zone_South                        25000 non-null uint8
zone_West                         25000 non-null uint8
WH_regional_zone_Zone 2           25000 non-null uint8
WH_regional_zone_Zone 3           25000 non-null uint8
WH_regional_zone_Zone 4           25000 non-null uint8
WH_regional_zone_Zone 5           25000 non-null uint8
WH_regional_zone_Zone 6           25000 non-null uint8
wh_owner_type_Rented              25000 non-null uint8
approved_wh_govt_certificate_A+   25000 non-null uint8
approved_wh_govt_certificate_B    25000 non-null uint8
approved_wh_govt_certificate_B+   25000 non-null uint8
approved_wh_govt_certificate_C    25000 non-null uint8
num_refill_req_13m                25000 non-null int64
transport_issue_11y               25000 non-null int64
Competitor_in_mkt                 25000 non-null int64
retail_shop_num                   25000 non-null int64
distributor_num                   25000 non-null int64
flood_impacted                    25000 non-null int64
flood_proof                       25000 non-null int64
electric_supply                   25000 non-null int64
dist_from_hub                     25000 non-null int64
workers_num                       25000 non-null float64
temp_reg_mach                     25000 non-null int64
wh_breakdown_13m                  25000 non-null int64
govt_check_13m                    25000 non-null int64
product_wg_ton                    25000 non-null int64
dtypes: float64(1), int64(13), uint8(16)
memory usage: 3.9+ MB

```

Fig 1.29 Info of the dataset after variable transformation.

Lambda value used for Transformation: 0.583633892766992

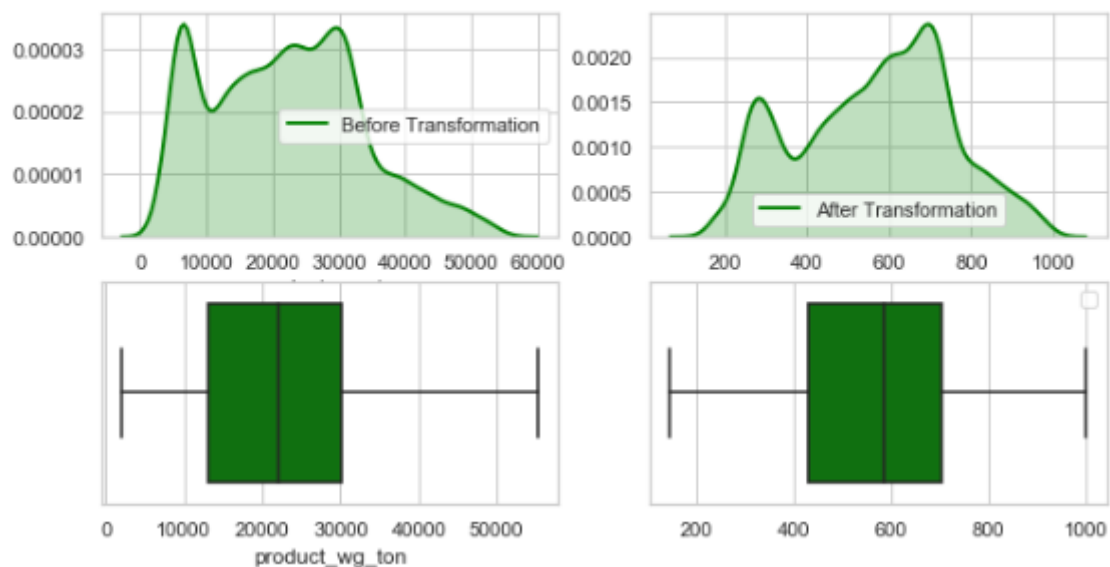


Fig 1.30 Transforming right skewed target variable.

The target variable is right skewed before transformation. After transforming the target variable the values are now changed to a slightly normal distributed value.

Location_type_Urban	3.054017
WH_capacity_size_Mid	0.404872
WH_capacity_size_Small	1.560454
zone_North	0.361295
zone_South	1.127422
zone_West	0.785433
WH_regional_zone_Zone 2	2.360618
WH_regional_zone_Zone 3	2.410081
WH_regional_zone_Zone 4	1.785360
WH_regional_zone_Zone 5	1.635609
WH_regional_zone_Zone 6	0.706068
wh_owner_type_Rented	0.173135
approved_wh_govt_certificate_A+	1.779592
approved_wh_govt_certificate_B	1.560127
approved_wh_govt_certificate_B+	1.526275
approved_wh_govt_certificate_C	1.116087
num_refill_req_13m	-0.075217
transport_issue_1ly	1.610907
Competitor_in_mkt	0.978456
retail_shop_num	0.908302
distributor_num	0.015213
flood_impacted	2.701327
flood_proof	3.919343
electric_supply	-0.660933
dist_from_hub	-0.005999
workers_num	1.081539
temp_reg_mach	0.855960
wh_breakdown_13m	-0.068026
govt_check_13m	-0.363262
product_wg_ton	-0.117054

dtype: float64

Fig 1.31 Finding skewness of the data frame.

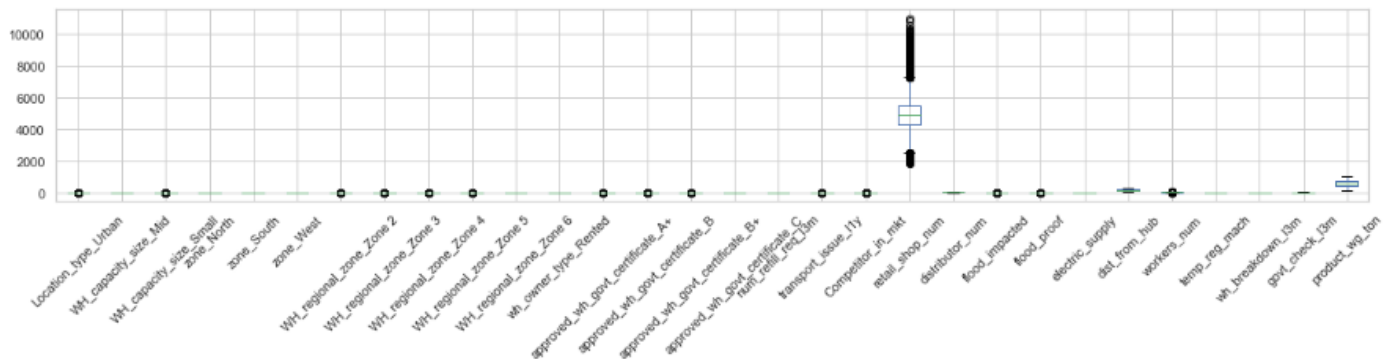


Fig 1.32 Boxplot before scaling

### g) Addition of new variables (if required):

Addition of new variables is not required in this dataset.

## Business insights from EDA

a) Is the data unbalanced? If so, what can be done? Please explain in the context of the business:

If the data is unbalanced, we can use Smote-R and Smote-R Gaussian for treating the unbalanced data.

b) Any business insights using clustering (if applicable)

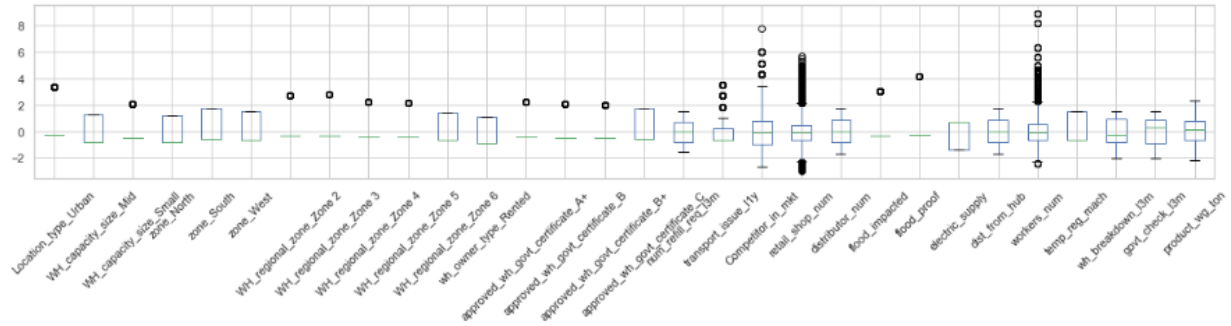


Fig 1.33 Boxplot after scaling

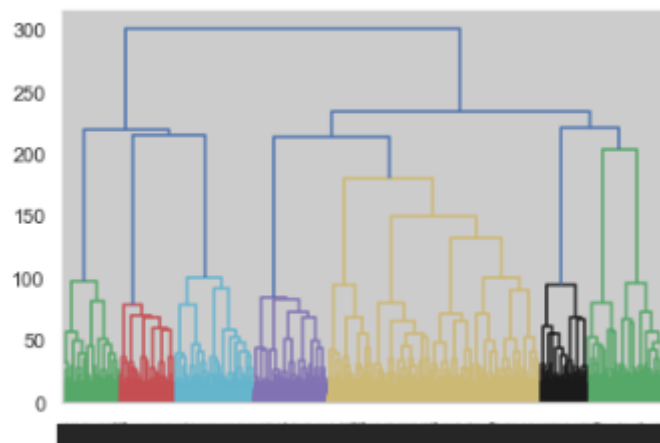


Fig 1.34 Dendrogram of hierarchial clustering

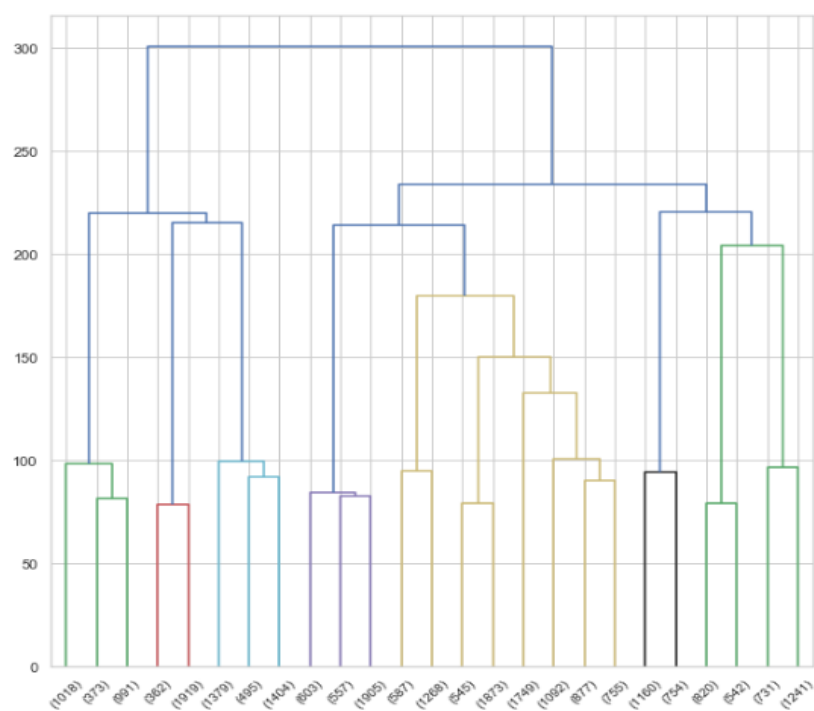


Fig 1.35 Dendrogram of hierarchial clustering after Truncating

The cluster grouping linkage based on the dendrogram, 3 or 4 looks good. The further analysis, and based on the dataset had gone for 3 linkage solution based on the hierarchical clustering. From hierarchial clustering we can find there are 7 clusters identified for the data.

In K-means clustering, the inertia is calculated, from that finding plotting the elbow curve and finding the clusters from the elbow curve

```
[156510180.3596399,
43443968.93680901,
20621288.387078643,
11492506.372049445,
7350408.125392827,
5438510.894673958,
4245456.305741977,
3271363.6982872887,
2672881.0274393656,
2299211.6490169726,
2000273.1450390322,
1776152.5373905546,
1562201.1989981781,
1425965.8490435758,
1308757.9397151927,
1201345.0573615218,
1127460.7549983964,
1052962.5302454599,
1005583.0463379754,
952646.3267051913,
899360.664886525,
853996.2506287948,
837108.1628498367,
793088.2354395089]
```

Fig 1.36 Finding inertia from K-means clustering

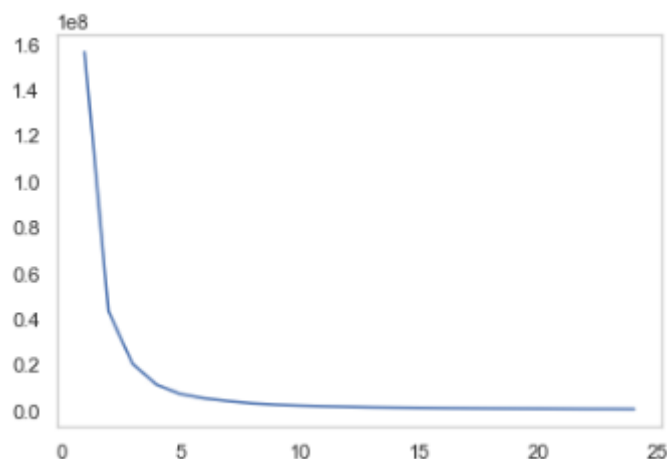


Fig 1.37 Elbow curve for K – means

From the Elbow curve we find that the elbow bends down at 4 cluster linkages are present in the data.

**Silhouette Score and Silhouette width:**

```
0.5454022159943779
```

Fig 1.38 Silhouette score

```
array([0.67117602, 0.67430139, 0.73966249, ..., 0.17528338, 0.51774779,
0.72617875])
```

Fig 1.39 Silhouette width

***Silhouette sample value:***

**0.015491729420163901**

Fig – 1.40 Silhouette Sample value

From this we can clearly identify that we can separate the market zone based on the 4 type of clusters. 4 type of clusters can be

1. High demand & high supply
2. High demand & low supply
3. Low demand & high supply
4. Low demand & low supply

From this we can separate the high performing areas and low performing zones. Based on this we can make promotion, offers, and discounts, make new marketing strategies for making low demand zone to high demand zone and make profits for the company.

**c) Any other business insights:**

**Important business insights are**

- reduce cost,
- Improve the overall organization performance
- customer satisfaction