

MACHINE LEARNING-1- REPORT-CODED

FASNA-PP



TABLE OF CONTENTS: Problem 1:

- 1:1)-Clustering: Define the problem and perform Exploratory Data Analysis
- 1:2) -Clustering: Data Preprocessing
- 1: 3)-Clustering: Hierarchical Clustering
- 1: 4)-Clustering: K-means Clustering
- 1: 5)-Clustering: Actionable Insights & Recommendations

Problem-2:

- 2:1)-PCA: Define the problem and perform Exploratory Data Analysis
- 2:2)-PCA: Data Preprocessing
- 2:3)-PCA: PCA

TABLE OF FIGURES

Fig-1 histogram of Ad_length-----	9
Fig-2-histogram and boxplot of Ad-length-----	9
Fig-3-histogram and boxplot of Ad-Width-----	10
Fig-4-histogram and boxplot of Ad-size-----	11
Fig-5-histogram and boxplot of Available_impression-----	12
Fig-6-histogram and boxplot of matched_queries-----	13
Fig-7-histogram and boxplot of impression-----	14
Fig-8-histogram and boxplot of clicks-----	15
Fig-9-histogram and boxplot of spend-----	16
Fig-10-histogram and boxplot of Fee-----	17
Fig-11-histogram and boxplot of Revenue-----	17
Fig-12-histogram and boxplot of CTR-----	18
Fig-13-histogram and boxplot of CPM-----	19
Fig-14-histogram and boxplot of CPC-----	20
Fig-15- UNIVARIATE ANALYSIS FOR CATEGORICAL COLUMNS -----	21
Fig-16-correlation map-----	24

Fig-17-barplot of impression with ad type-----	25
Fig-18-barplot of CTR with Device Type-----	25
Fig-19-barplot CPM with Device Type-----	26
Fig-20-barplot of clicks with device Type-----	26
Fig-21-check outliers-----	28
Fig-22-dendrogram-----	30
Fig-23-Elbow curve-----	31
Fig-24- plot silhouette score-----	32
Fig-25-barplot for <i>Group the data by clusters and take sum or mean to identify trends in clicks, spend, r venue, CPM, CTR, & CPC based on Device Type</i> -----	34
Fig-26-barplot for state based by gender ratio-----	43
Fig-27-checking outliers-----	45
Fig-28-treated outliers-----	47
Fig-29-before scaling boxplot for checking outliers-----	51
Fig-30-after scaling boxplot for checking outliers-----	52
Fig-31-scree plot-----	60
Fig-32-compare pcs with actual columns-----	67

Tables:

<i>Table -1-first 5 rows</i>	6
<i>Table-2-describe the data set</i>	7
Table-3-scaled data	29
Table-4-calculating mean and median of the original data for each label	33
Table-5- statical summary	41
Table-6-scaled data	50
Table-7-before scaling	52
Table-8-after scaling	55
Table-9-compare pcs with actual columns and identify which is explaining most variance	64
Table-10-compare pcs with actual columns	72

Problem 1:**Clustering:****Digital Ads Data:**

The ads24x7 is a Digital Marketing company which has now got seed funding of \$10 Million. They are expanding their wings in Marketing Analytics. They collected data from their Marketing Intelligence team and now wants you (their newly appointed data analyst) to segment type of ads based on the features provided. Use Clustering procedure to segment ads into homogeneous groups.

The following three features are commonly used in digital marketing:

CPM = (Total Campaign Spend / Number of Impressions) * 1,000. Note that the Total Campaign Spend refers to the 'Spend' Column in the dataset and the Number of Impressions refers to the 'Impressions' Column in the dataset.

CPC = Total Cost (spend) / Number of Clicks. Note that the Total Cost (spend) refers to the 'Spend' Column in the dataset and the Number of Clicks refers to the 'Clicks' Column in the dataset.

CTR = Total Measured Clicks / Total Measured Ad Impressions x 100. Note that the Total Measured Clicks refers to the 'Clicks' Column in the dataset and the Total Measured Ad Impressions refers to the 'Impressions' Column in the dataset.

DATA DICTIONARY

1-**Timestamp**-The Timestamp of the particular Advertisement.

2-**InventoryType**-The Inventory Type of the particular Advertisement. Format 1 to 7. This is a Categorical Variable.

3-**Ad - Length**-The Length Dimension of the particular Advertisement.

4-**Ad- Width**-The Width Dimension of the particular Advertisement.

5-**Ad Size**-The Overall Size of the particular Advertisement. Length*Width.

6-**Ad Type**-The type of the particular Advertisement. This is a Categorical Variable.

7-**Platform**-The platform in which the particular Advertisement is displayed. Web, Video or App. This is a Categorical Variable.

8-**Device Type**-The type of the device which supports the particular Advertisement. This is a Categorical Variable.

9-**Format**-The Format in which the Advertisement is displayed. This is a Categorical Variable.

10-Available Impressions-How often the particular Advertisement is shown. An impression is counted each time an Advertisement is shown on a search result page or other site on a Network.

11-Matched Queries-Matched search queries data is pulled from Advertising Platform and consists of the exact searches typed into the search Engine that generated clicks for the particular Advertisement.

12- Impressions-The impression count of the particular Advertisement out of the total available impressions.

13-Clicks-It is a marketing metric that counts the number of times users have clicked on the particular advertisement to reach an online property.

14-Spend-It is the amount of money spent on specific ad variations within a specific campaign or ad set. This metric helps regulate ad performance.

15-Fee-The percentage of the Advertising Fees payable by Franchise Entities.

16-Revenue-It is the income that has been earned from the particular advertisement.

17-CTR-CTR stands for "Click through rate". CTR is the number of clicks that your ad receives divided by the number of times your ad is shown. Formula used here is $CTR = \text{Total Measured Clicks} / \text{Total Measured Ad Impressions} \times 100$. Note that the Total Measured Clicks refers to the 'Clicks' Column and the Total Measured Ad Impressions refers to the 'Impressions' Column.

18-CPM-CPM stands for "cost per 1000 impressions." Formula used here is $CPM = (\text{Total Campaign Spend} / \text{Number of Impressions}) \times 1,000$. Note that the Total Campaign Spend refers to the 'Spend' Column and the Number of Impressions refers to the 'Impressions' Column.

19-CPC-CPC stands for "Cost-per-click". Cost-per-click (CPC) bidding means that you pay for each click on your ads. The Formula used here is $CPC = \text{Total Cost (spend)} / \text{Number of Clicks}$. Note that the Total Cost (spend) refers to the 'Spend' Column and the Number of Clicks refers to the 'Clicks' Column.

Part 1:--1)Clustering: Define the problem and perform Exploratory Data Analysis

- Problem definition - Check shape, Data types, statistical summary - Univariate analysis - Bivariate analysis - Key meaningful observations on individual variables and the relationship between variables

ANSWRS:Import all the necessary libraries

read the dataset, print first 5 rows

	Time stamp	Inventory Type	Ad- Length	Ad- Width	Ad Size	Ad Type	Pla tform	De vice Type	For mat	Availabl e_Impre ssions	Match ed_Qu eries	Imp ress ions	Cl ic ks	S p e n d	F e e	Re ve nu e	CT R	C P M	C P C
0	2020-9-2-17	Format1	300	250	75000	Inter222	Video	Desktop	Display	1806	325	323	1	0.00	0.35	0.0	0.0031	0.0	0.0
1	2020-9-2-10	Format1	300	250	75000	Inter227	App	Mobile	Video	1780	285	285	1	0.00	0.35	0.0	0.0035	0.0	0.0
2	2020-9-1-22	Format1	300	250	75000	Inter222	Video	Desktop	Display	2727	356	355	1	0.00	0.35	0.0	0.0028	0.0	0.0
3	2020-9-3-20	Format1	300	250	75000	Inter228	Video	Mobile	Video	2430	497	495	1	0.00	0.35	0.0	0.0020	0.0	0.0
4	2020-9-4-15	Format1	300	250	75000	Inter217	Web	Desktop	Video	1218	242	242	1	0.00	0.35	0.0	0.0041	0.0	0.0

Table -1-first 5 rows**1---1)-a)- Check shape**

#shape of the data set: (23066, 19)

1---1)-b)-Data types

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 23066 entries, 0 to 23065
Data columns (total 19 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Timestamp                             23066 non-null  object
1   InventoryType                         23066 non-null  object
2   Ad - Length                           23066 non-null  int64
3   Ad- Width                             23066 non-null  int64
4   Ad Size                               23066 non-null  int64
5   Ad Type                               23066 non-null  object
6   Platform                              23066 non-null  object
7   Device Type                           23066 non-null  object
8   Format                                 23066 non-null  object
9   Available_Impressions                 23066 non-null  int64
10  Matched_Queries                       23066 non-null  int64
11  Impressions                           23066 non-null  int64
12  Clicks                                23066 non-null  int64
13  Spend                                 23066 non-null  float64
14  Fee                                    23066 non-null  float64
15  Revenue                               23066 non-null  float64
16  CTR                                   18330 non-null  float64
17  CPM                                   18330 non-null  float64
18  CPC                                   18330 non-null  float64
dtypes: float64(6), int64(7), object(6)
memory usage: 3.3+ MB
```

1---1)-c)-statistical summary

	count	mean	std	min	25%	50%	75%	max
Ad - Length	23066.0	385.16	233.65	120.00	120.00	300.00	720.00	728.00
Ad- Width	23066.0	337.90	203.09	70.00	250.00	300.00	600.00	600.00
Ad Size	23066.0	96674.47	61538.33	33600.00	72000.00	72000.00	84000.00	216000.00
Available_Impressions	23066.0	2432043.67	4742887.76	1.00	33672.25	483771.00	2527711.75	27592861.00
Matched_Queries	23066.0	1295099.14	2512969.86	1.00	18282.50	258087.50	1180700.00	14702025.00
Impressions	23066.0	124151	2429399.9	1.00	7990.50	225290.0	1112428.50	14194774.00

	0	9.52	6			0		
Clicks	23066.0	10678.52	17353.41	1.00	710.00	4425.00	12793.75	143049.00
Spend	23066.0	2706.63	4067.93	0.00	85.18	1425.12	3121.40	26931.87
Fee	23066.0	0.34	0.03	0.21	0.33	0.35	0.35	0.35
Revenue	23066.0	1924.25	3105.24	0.00	55.37	926.34	2091.34	21276.18
CTR	18330.0	0.07	0.08	0.00	0.00	0.08	0.13	1.00
CPM	18330.0	7.67	6.48	0.00	1.71	7.66	12.51	81.56
CPC	18330.0	0.35	0.34	0.00	0.09	0.16	0.57	7.26

Table-2-describe the data set

#check the null values
 null values in CTR,CPM,CPC
#check the total duplicated values in the data set
 no duplicated values

#create data frame with categorical variables

#create data frame with numerical variables

#check the counts of each columns for checking data irregularities

1---1-d)Univariate analysis

for numerical columns

```
#plot histogram of Ad_length
count      23066.000000
mean       385.163097
std        233.651434
min         120.000000
25%         120.000000
50%         300.000000
75%         720.000000
max         728.000000
Name: Ad - Length, dtype: float64
```

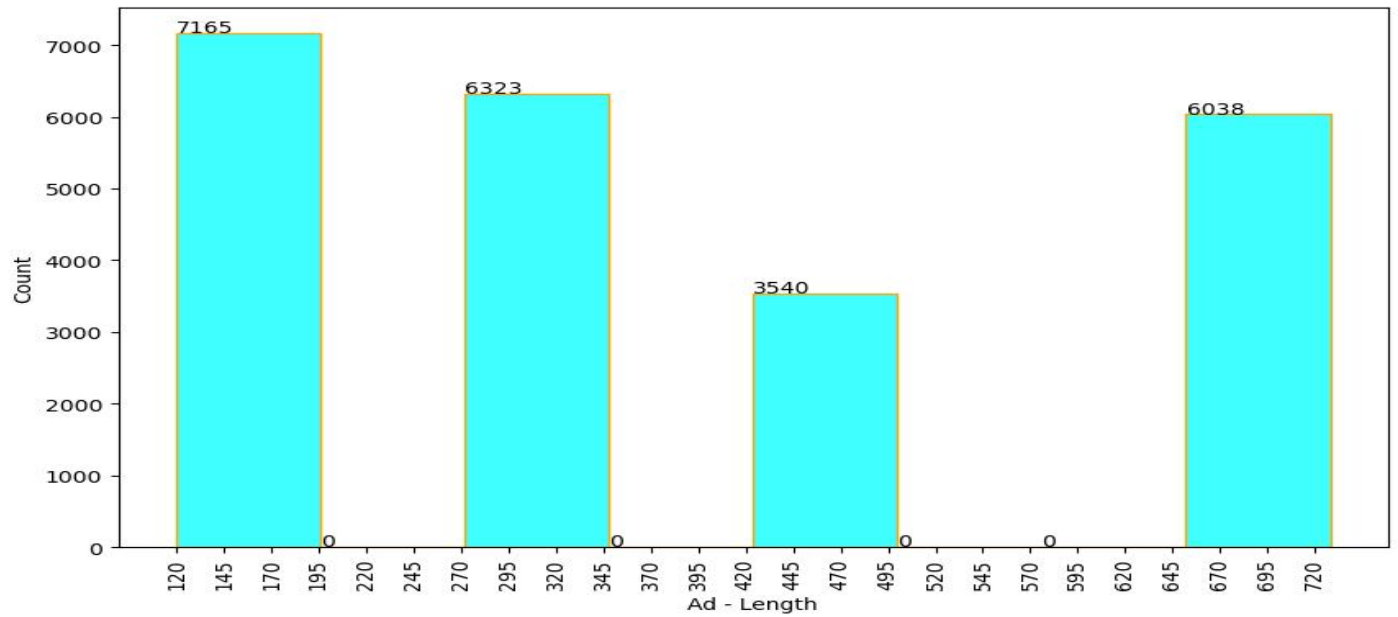


Fig-1 histogram of Ad_length

#plot histogram and boxplot of Ad-length

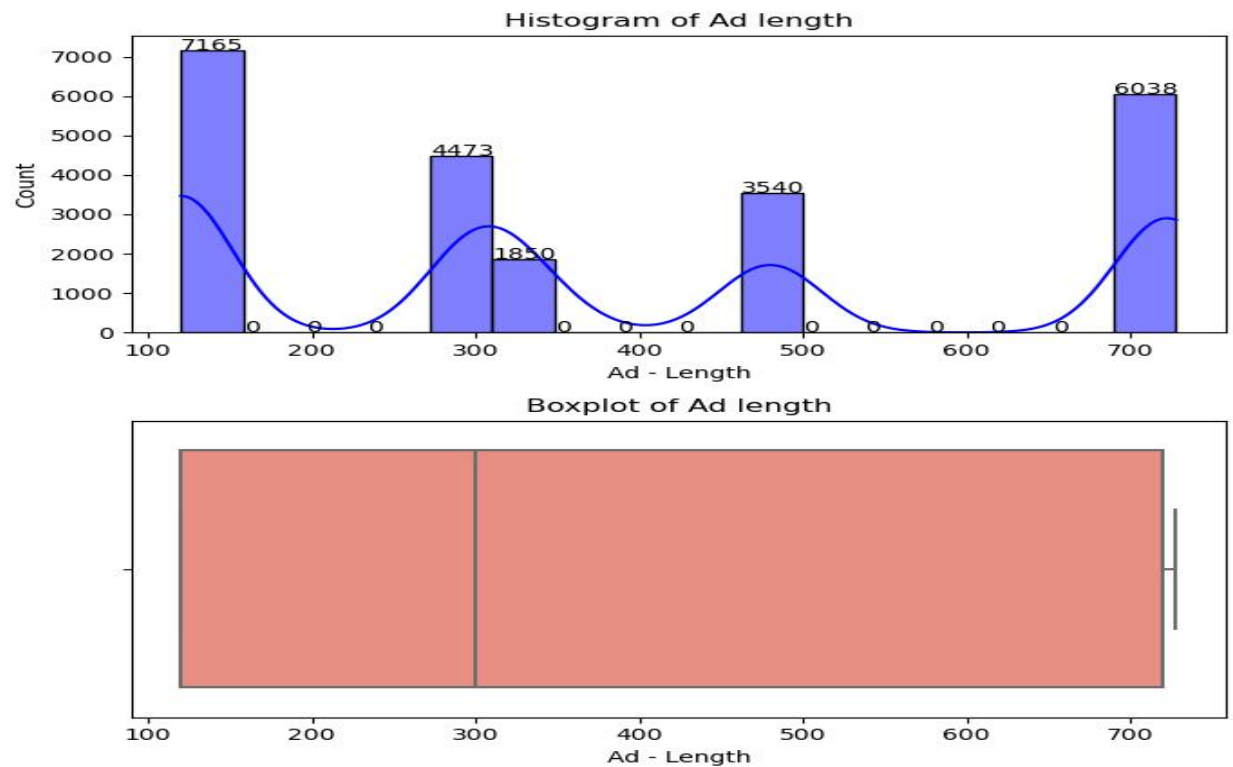


Fig-2-histogram and boxplot of Ad-length

plot histogram and boxplot of Ad-width

count 23066.000000

```

mean      337.896037
std       203.092885
min        70.000000
25%       250.000000
50%       300.000000
75%       600.000000
max       600.000000
Name: Ad- Width, dtype: float64

```

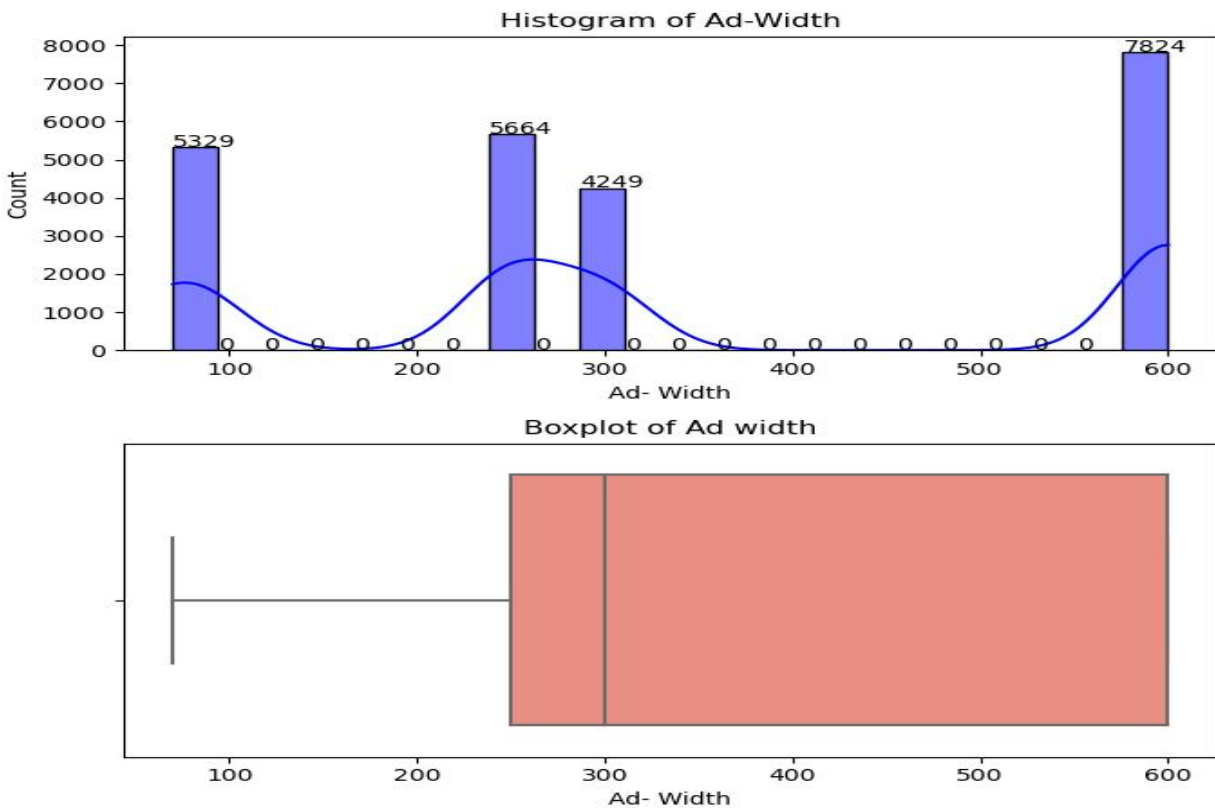


Fig-3-histogram and boxplot of Ad-width

#plot histogram and boxplot of Ad-size

```

count      23066.000000
mean      96674.468048
std       61538.329557
min       33600.000000
25%       72000.000000
50%       72000.000000
75%       84000.000000
max      216000.000000

```

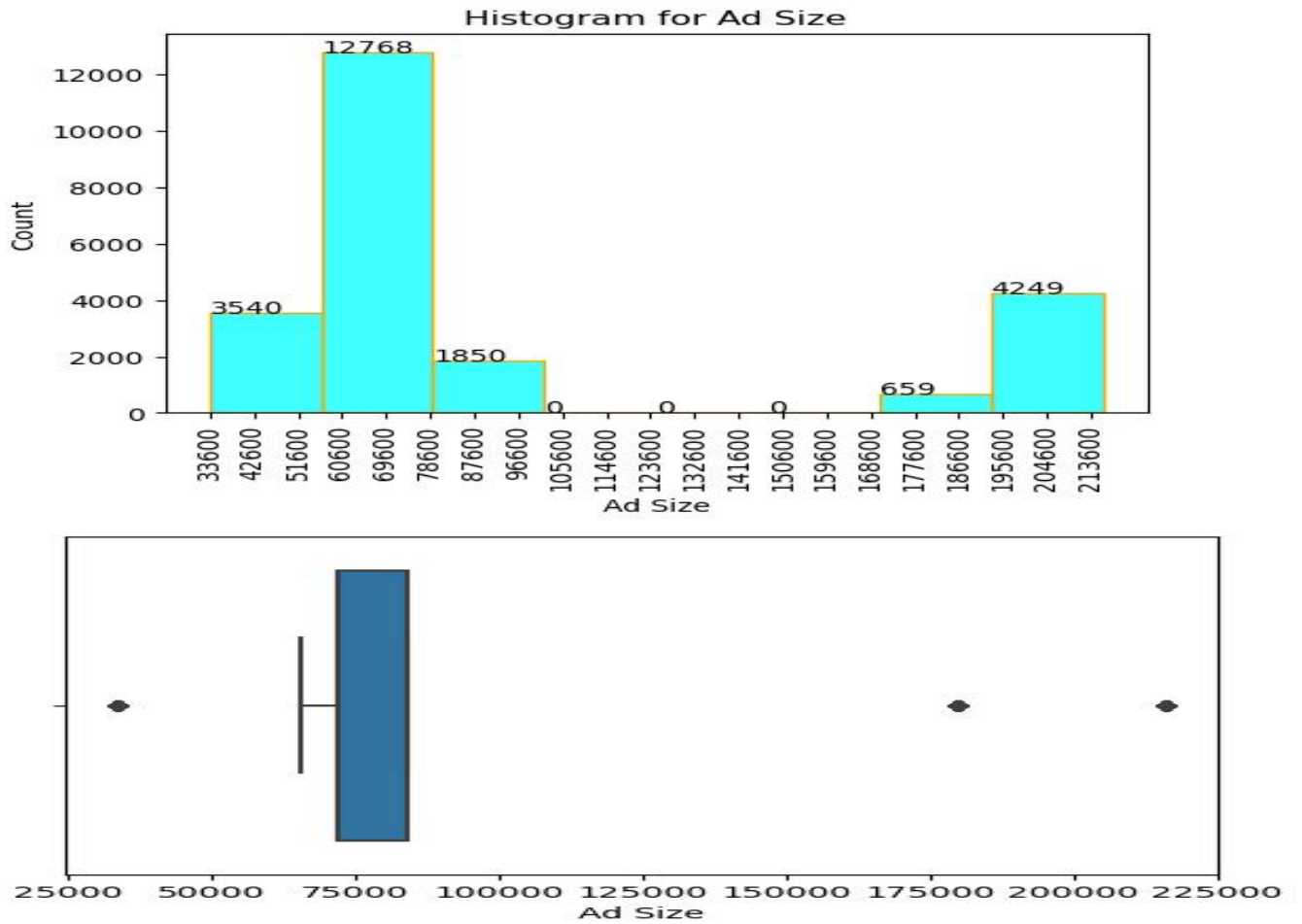


Fig-4-histogram and boxplot of Ad-size

#plot histogram and boxplot of available impression

```
count    2.306600e+04
mean     2.432044e+06
std      4.742888e+06
min      1.000000e+00
25%      3.367225e+04
50%      4.837710e+05
75%      2.527712e+06
max      2.759286e+07
```

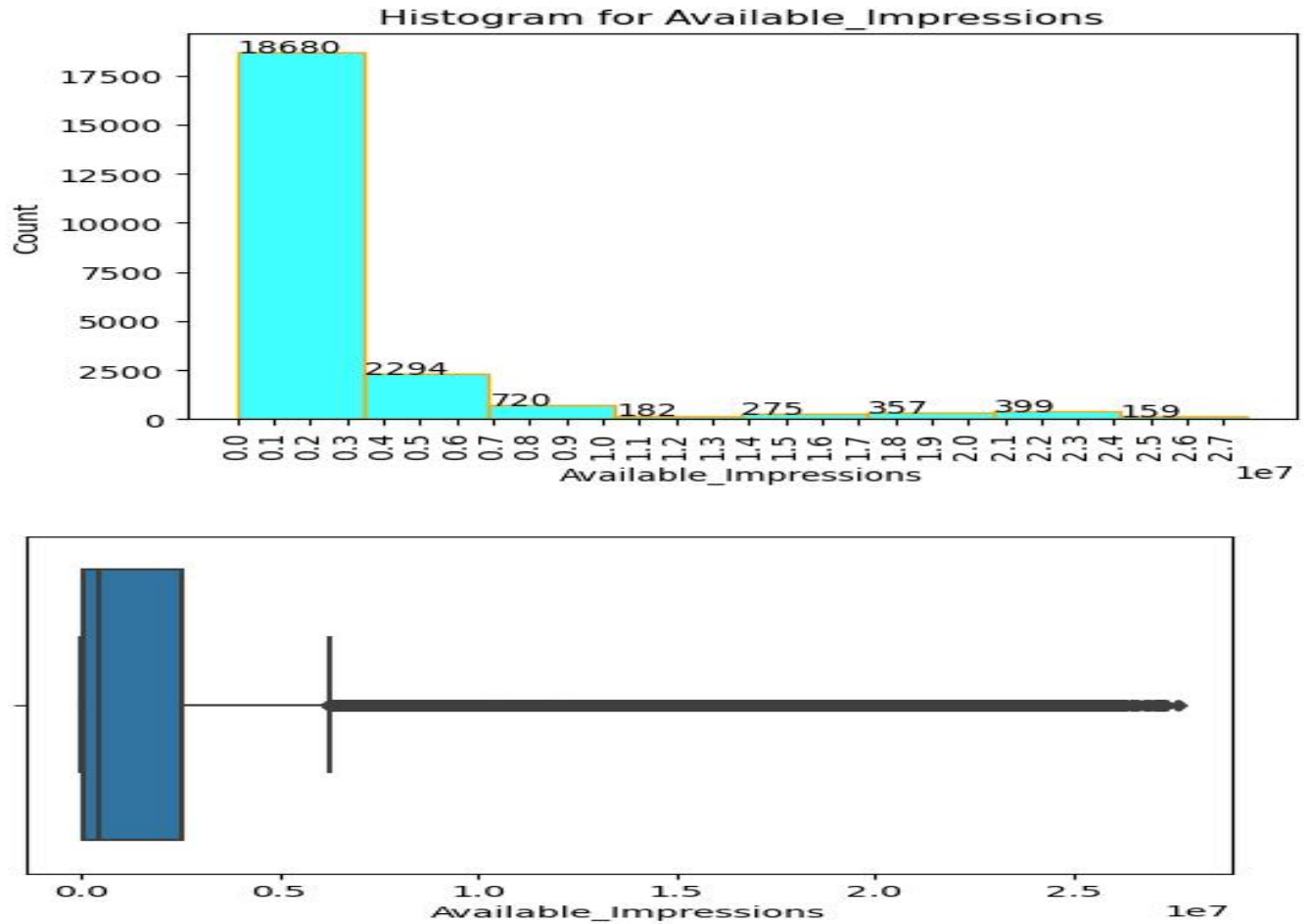


Fig-5-histogram and boxplot of available impression

#plot histogram and boxplot of matched queries

```

count      2.306600e+04
mean       1.295099e+06
std        2.512970e+06
min        1.000000e+00
25%        1.828250e+04
50%        2.580875e+05
75%        1.180700e+06
max        1.470202e+07

```

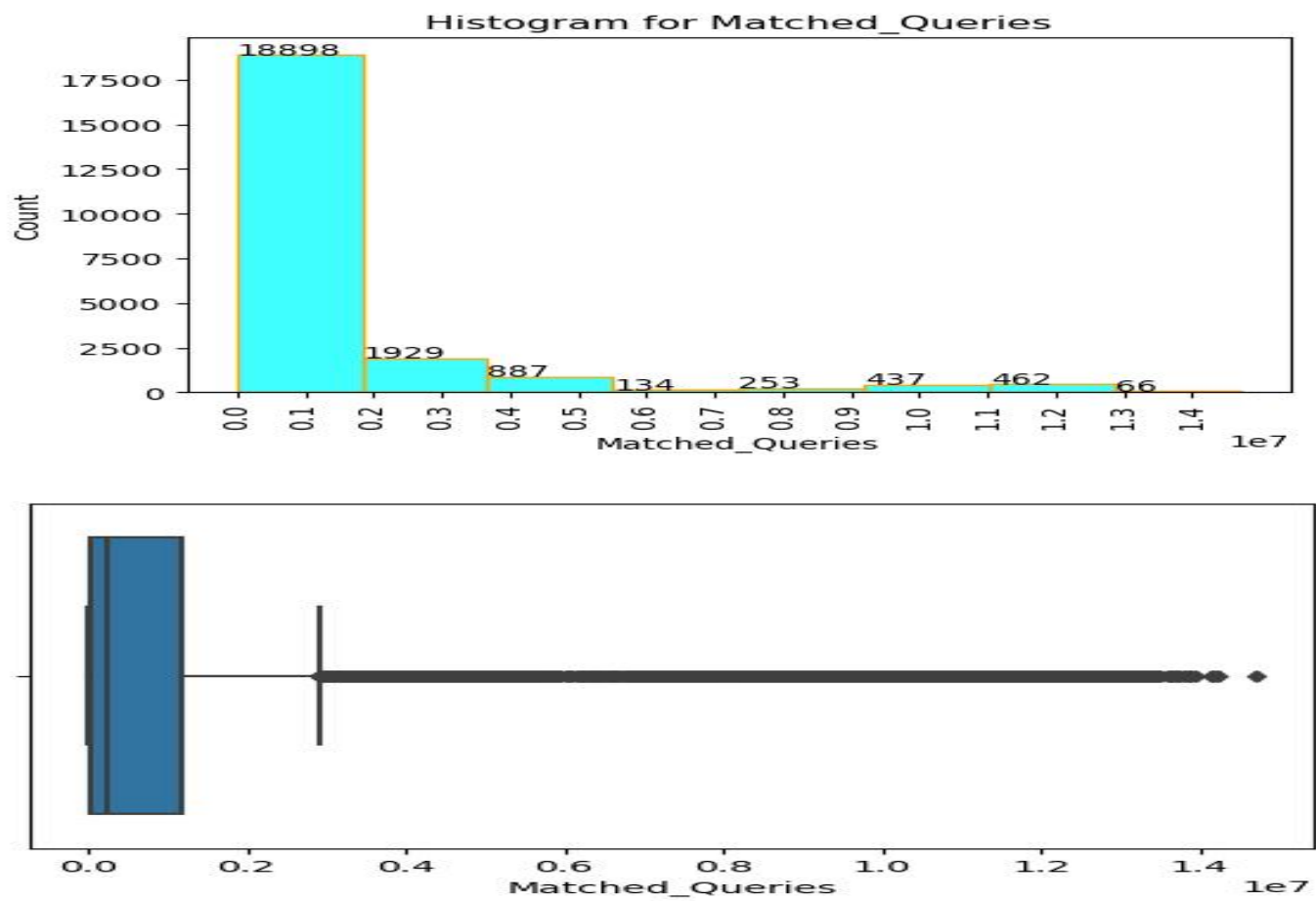


Fig-6-histogram and boxplot of matched queries

#plot histogram and boxplot of impression

```
count    2.306600e+04
mean     1.241520e+06
std      2.429400e+06
min      1.000000e+00
25%      7.990500e+03
50%      2.252900e+05
75%      1.112428e+06
max      1.419477e+07
```

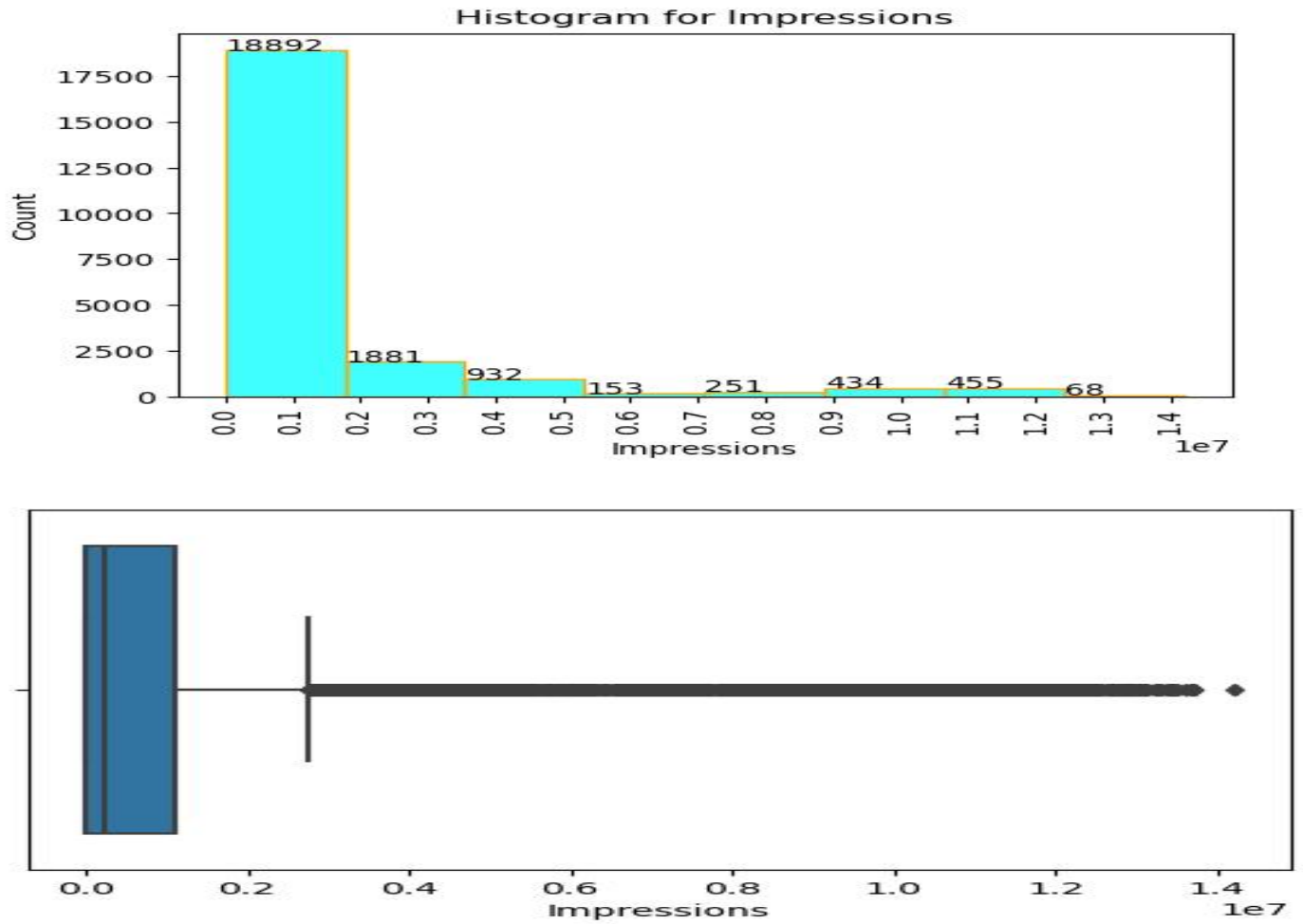


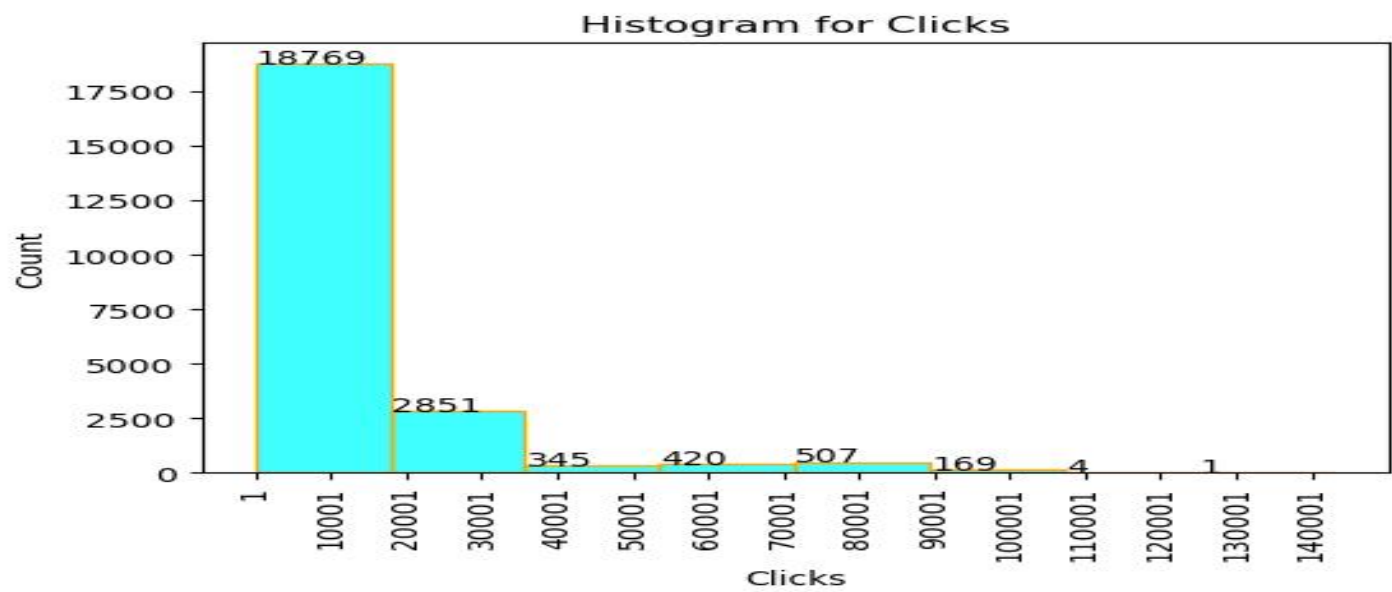
Fig-7-histogram and boxplot of impression

#plot histogram and boxplot of clicks

```

count      23066.000000
mean       10678.518816
std        17353.409363
min         1.000000
25%         710.000000
50%        4425.000000
75%       12793.750000
max       143049.000000
Name: Clicks, dtype: float64 Distribution of Clicks

```



BoxPlot of clicks
<Figure size 640x480 with 0 Axes>

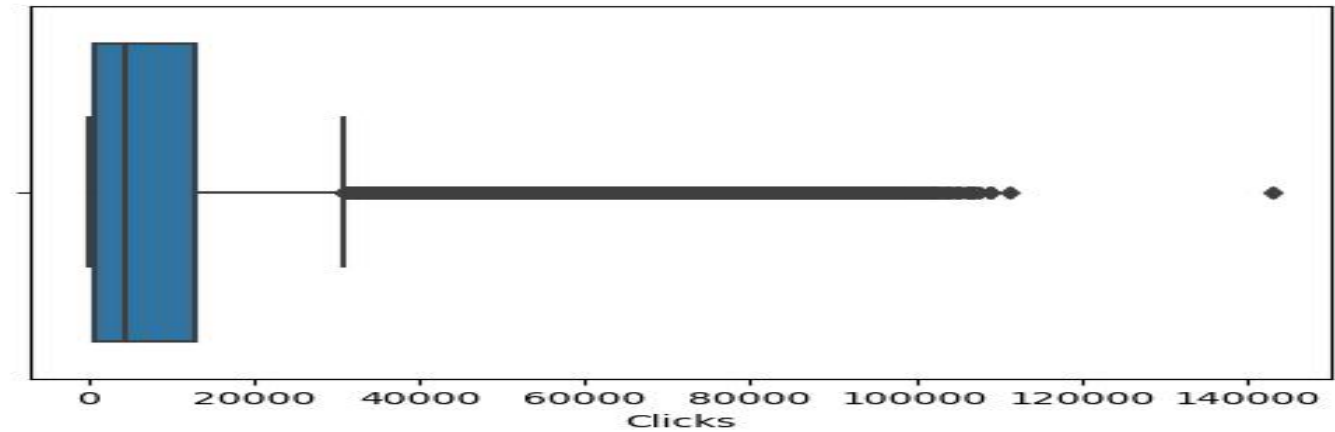
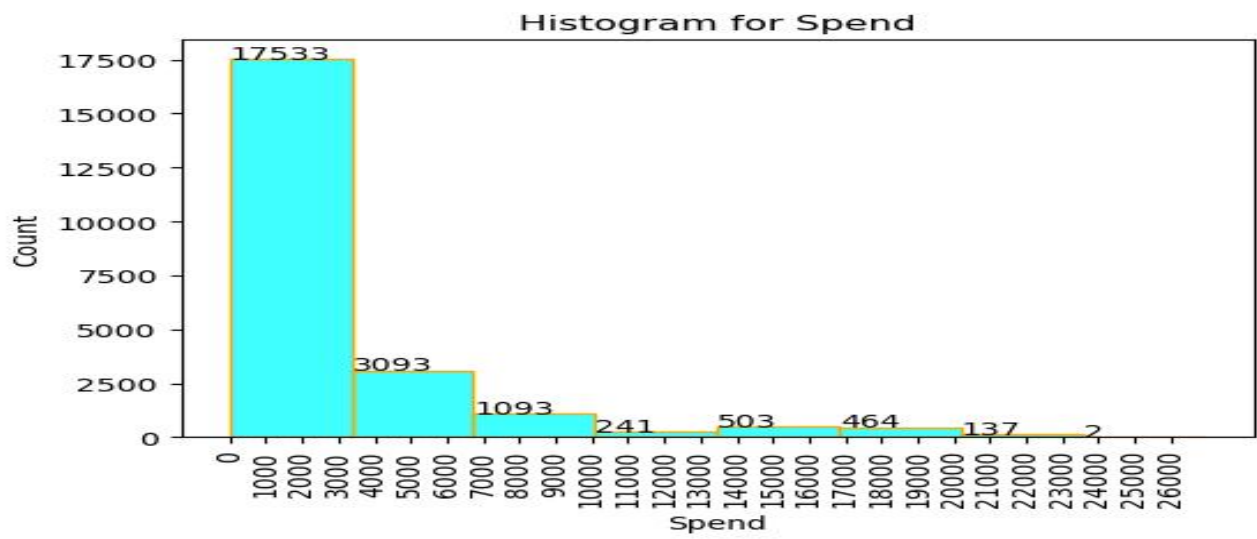


Fig-8-histogram and boxplot of clicks

#plot histogram and boxplot of spend

```
count    23066.000000
mean      2706.625689
std       4067.927273
min        0.000000
25%       85.180000
50%      1425.125000
75%      3121.400000
max      26931.870000
```

Name: Spend, dtype: float64 Distribution of Spend



BoxPlot of Spend
<Figure size 640x480 with 0 Axes>

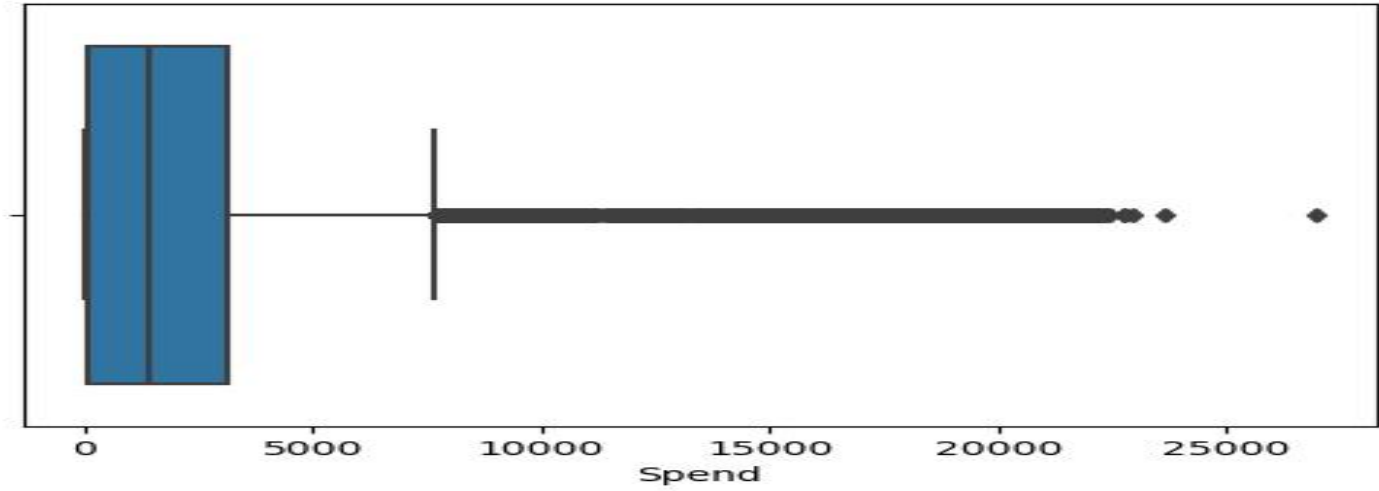


Fig-9-histogram and boxplot of spend
plot histogram and boxplot of Fee

```
count    23066.000000
mean      0.335123
std       0.031963
min       0.210000
25%       0.330000
50%       0.350000
75%       0.350000
max       0.350000
Name: Fee, dtype: float64
```

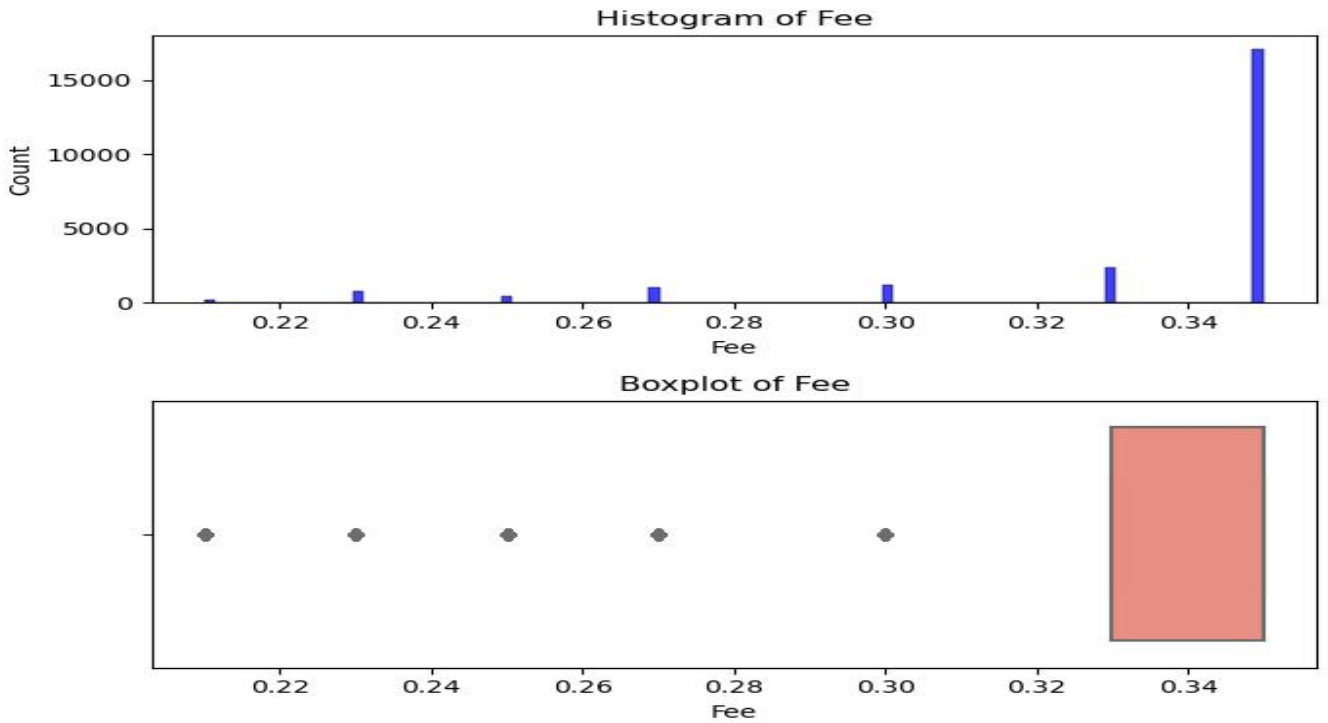
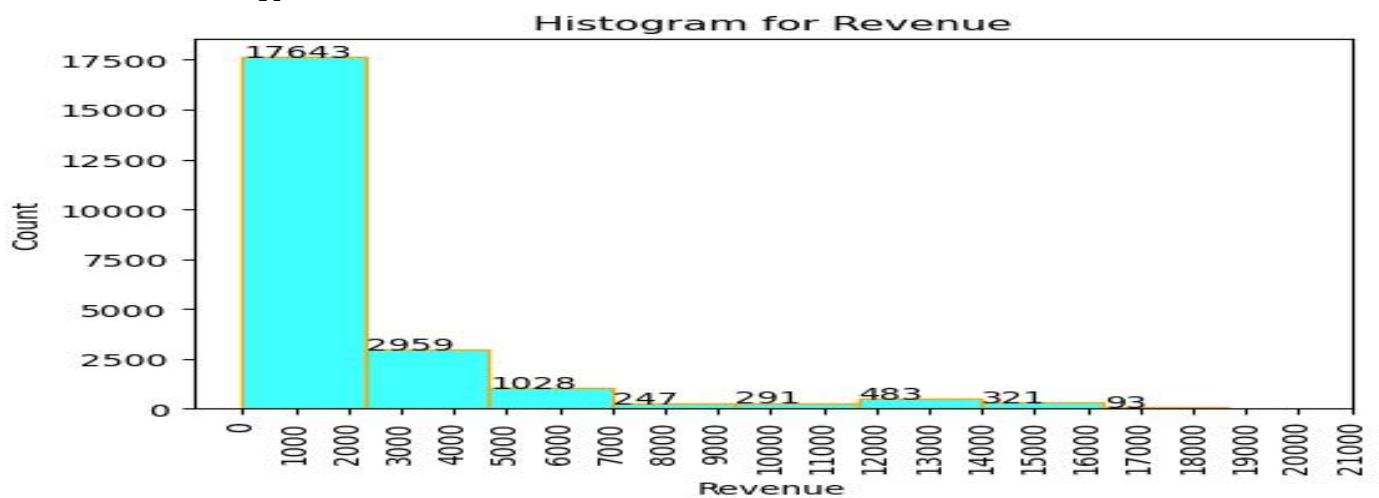


Fig-10- histogram and boxplot of Fee

#plot histogram and boxplot of revenue

```
count    23066.000000
mean      1924.252331
std       3105.238410
min        0.000000
25%       55.365375
50%      926.335000
75%     2091.338150
max     21276.180000
Name: Revenue, dtype: float64 Distribution of Revenue
```



BoxPlot of Revenue

<Figure size 640x480 with 0 Axes>

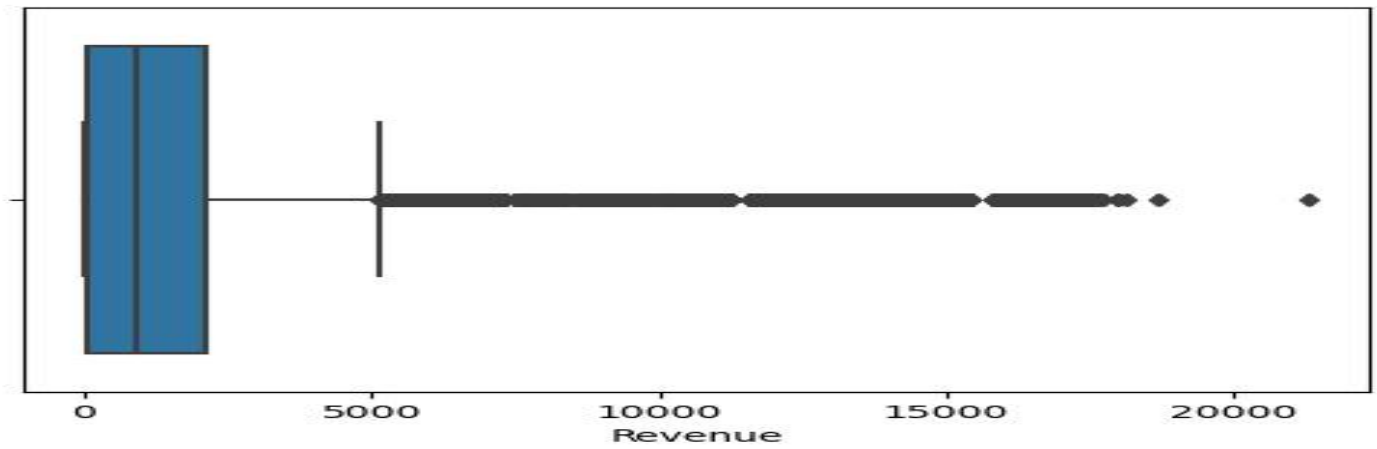


Fig-11-histogram and boxplot of revenue
plot histogram and boxplot of CTR

```
count    18330.000000
mean      0.073661
std       0.075160
min       0.000100
25%       0.002600
50%       0.082550
75%       0.130000
max       1.000000
Name: CTR, dtype: float64
```

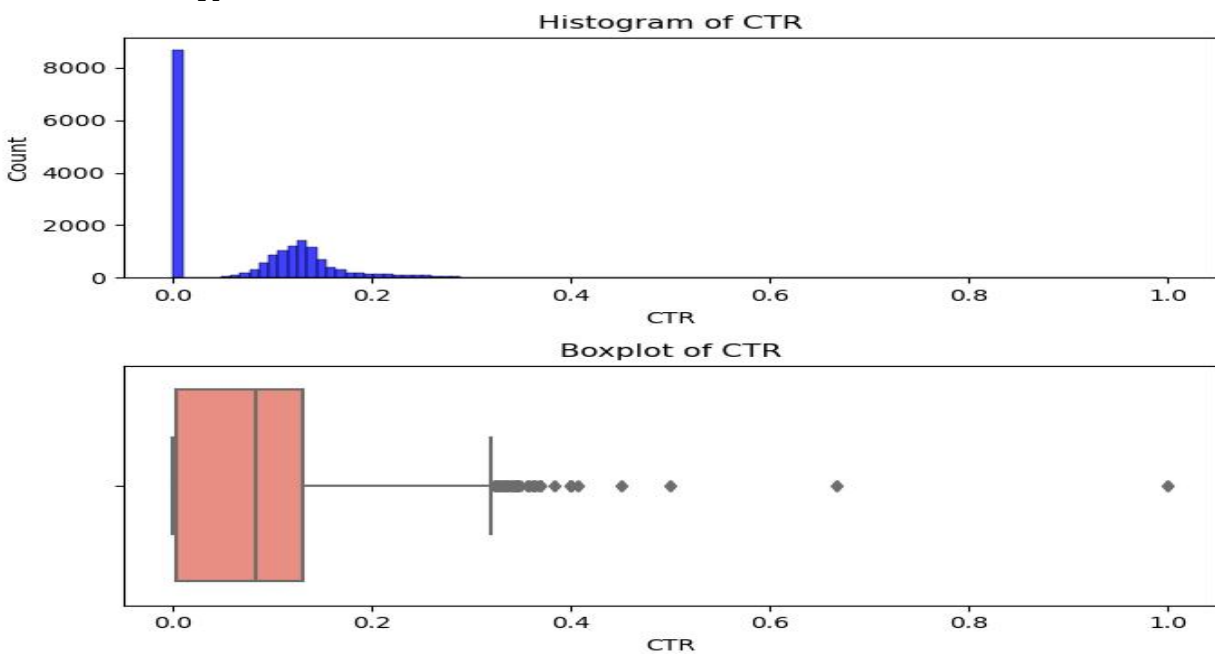


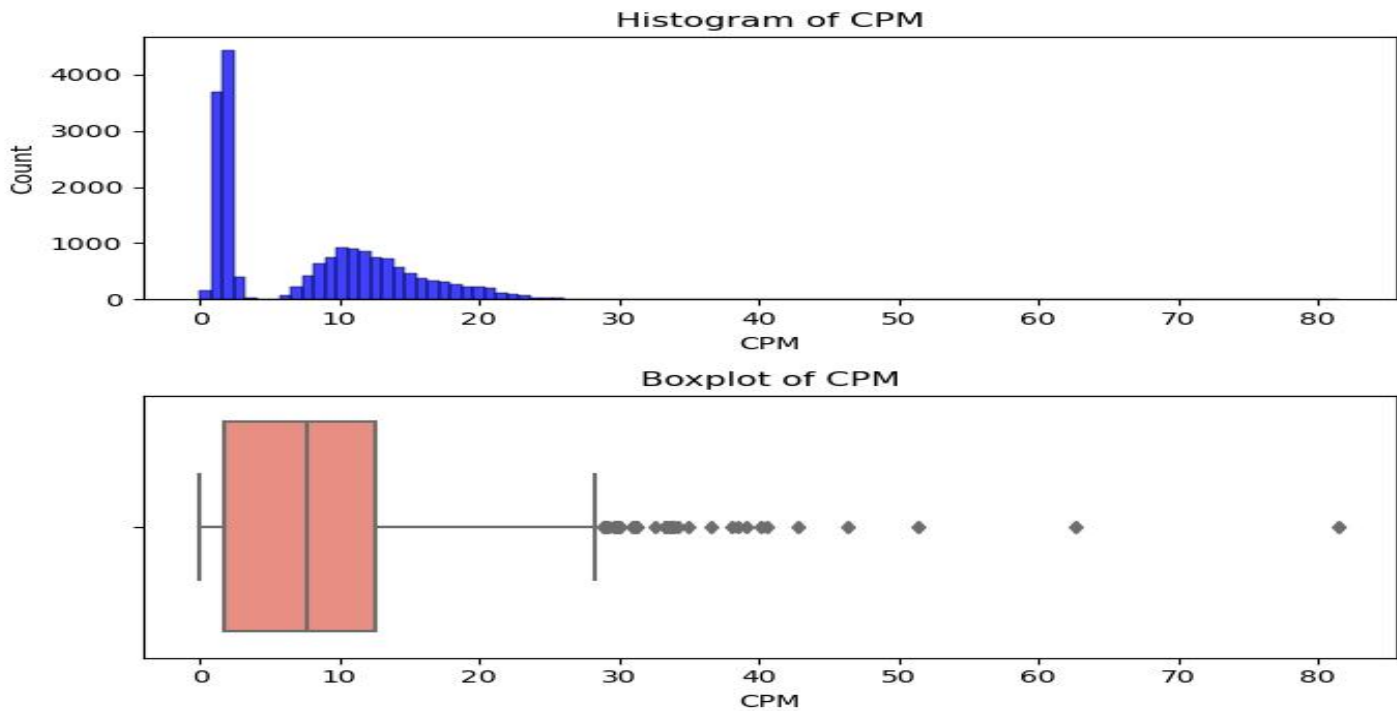
Fig-12- histogram and boxplot of CTR

#plot histogram and boxplot of CPM

```

count    18330.000000
mean      7.672045
std       6.481391
min       0.000000
25%      1.710000
50%      7.660000
75%     12.510000
max      81.560000
Name: CPM, dtype: float64

```

**Fig-13- histogram and boxplot of CPM****#plot histogram and boxplot of CPC**

```

count    18330.000000
mean      0.351061
std       0.343334
min       0.000000
25%      0.090000
50%      0.160000
75%      0.570000
max       7.260000
Name: CPC, dtype: float64

```

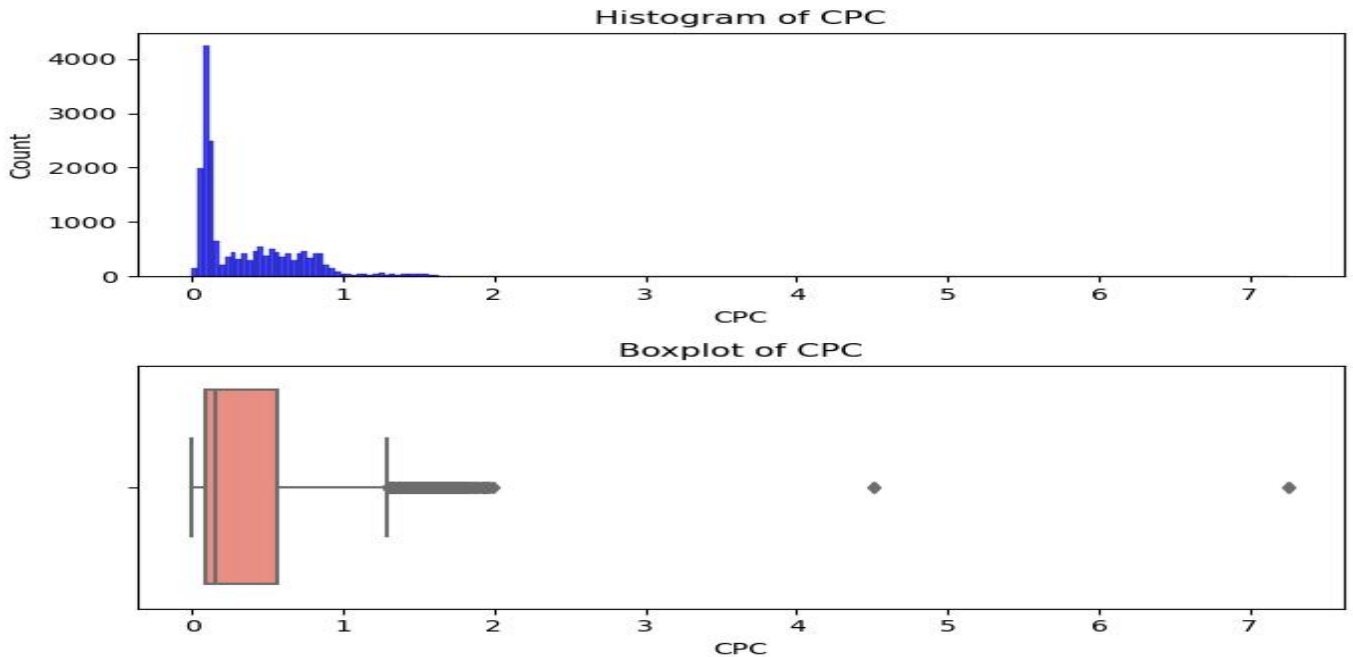


Fig-14- histogram and boxplot of CPC

#DESCRIBE THE FEE COLUMN

```
count    23066.000000
mean      0.335123
std       0.031963
min       0.210000
25%       0.330000
50%       0.350000
75%       0.350000
max       0.350000
```

Observations

- There are 13 numeric fields in the data
- Customer ad length ranges from 120 to 728
- maximum of Ad_width is 600
- Ad size ranges from 33600 to 216000 with an average 7200
- Available_Impressions ranges from 1.000000e+00 to 2.759286e+07
- only few matched queries above 1.4
- Average impression is around 2.252900e+05
- more Ads are zero clicks very less ads are more than 80000 Clicks

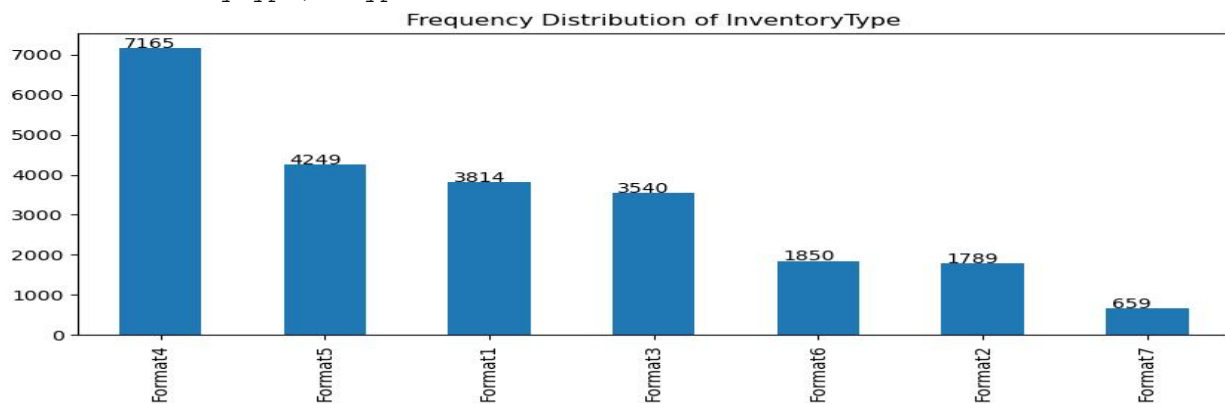
- Spend range is 0.00 to 26931.870000
- 50%,75%,maximum of fee of Ads is 350000
- maximum of Revenue of Ad is 21276.180000
- CTR range in between 0.00 to 1.00
- maximum of CPM is 81.560000
- Range of CPC 0.000 to 7.260000 with average 0.160000
- Outliers to be treated

#UNIVARIATE ANALYSIS FOR CATEGORICAL COLUMNS

Distribution of InventoryType

```
Format4      7165
Format5      4249
Format1      3814
Format3      3540
Format6      1850
Format2      1789
Format7       659
```

Name: InventoryType, dtype: int64



Distribution of Ad Type

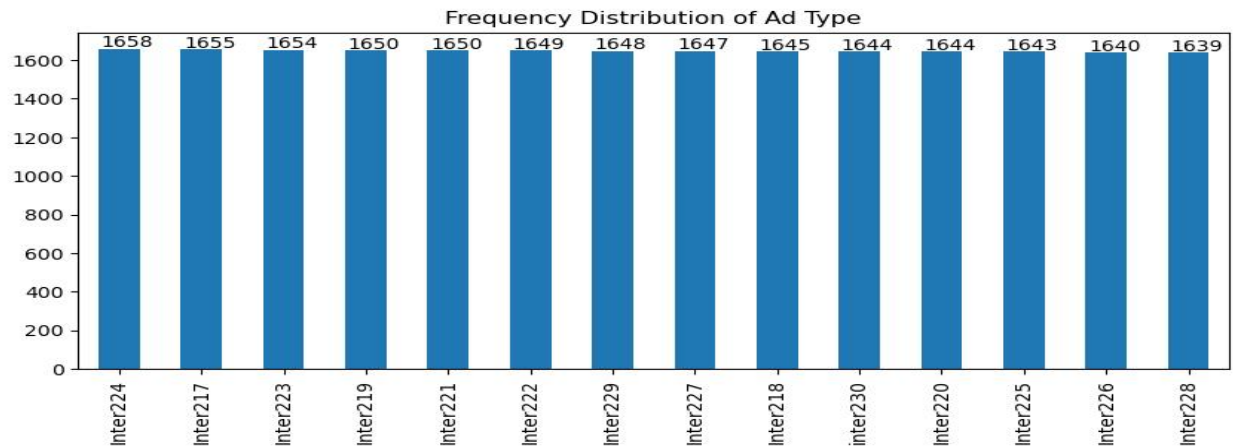
```
Inter224      1658
Inter217      1655
Inter223      1654
Inter219      1650
Inter221      1650
```

```

Inter222    1649
Inter229    1648
Inter227    1647
Inter218    1645
inter230    1644
Inter220    1644
Inter225    1643
Inter226    1640
Inter228    1639

```

Name: Ad Type, dtype: int64



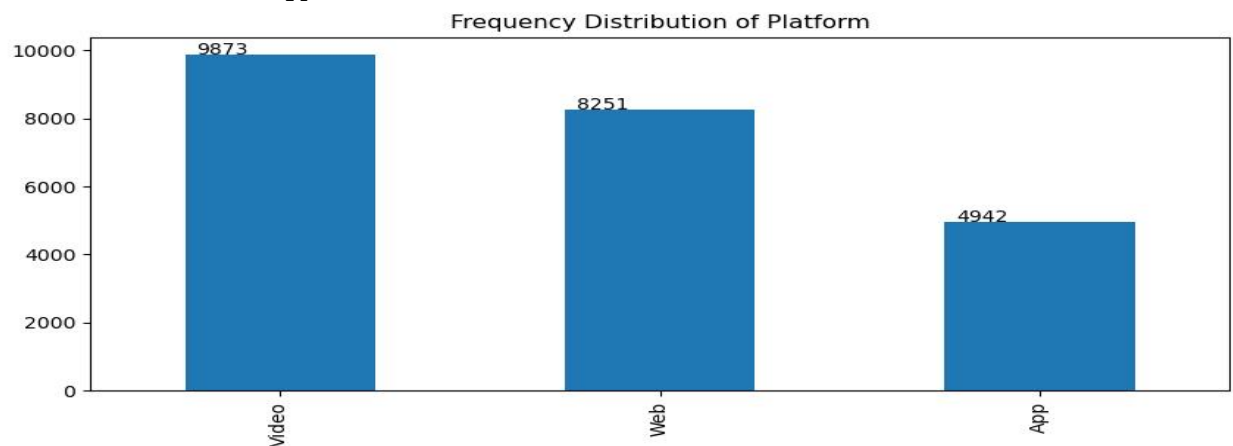
Distribution of Platform

```

Video      9873
Web        8251
App        4942

```

Name: Platform, dtype: int64



Distribution of Device Type

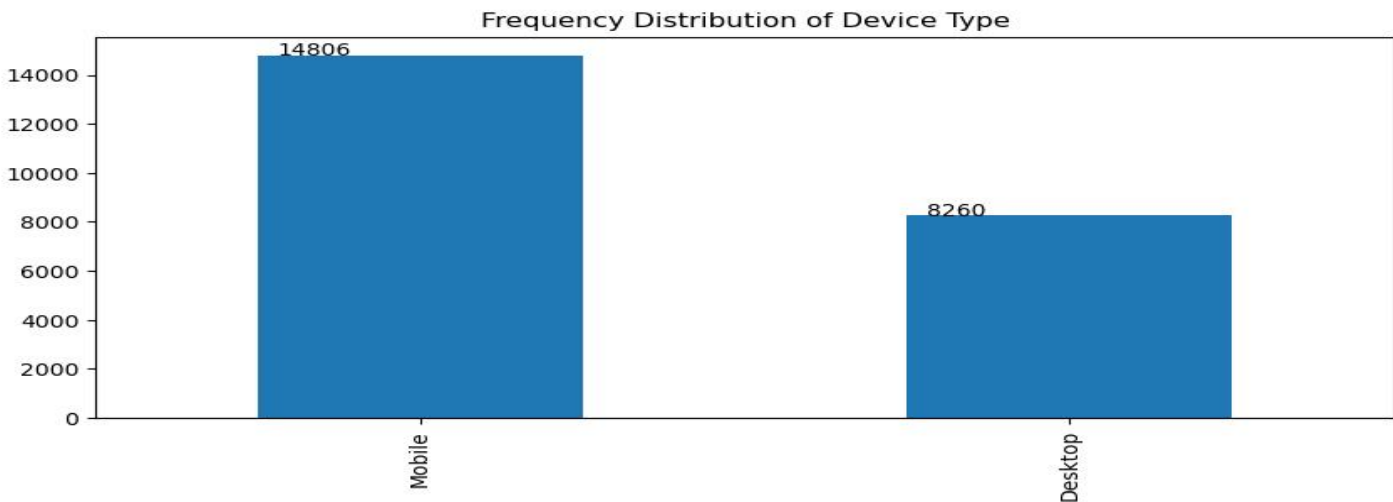
```

Mobile      14806

```

Desktop 8260

Name: Device Type, dtype: int64



Distribution of Format

Video 11552

Display 11514

Name: Format, dtype: int64

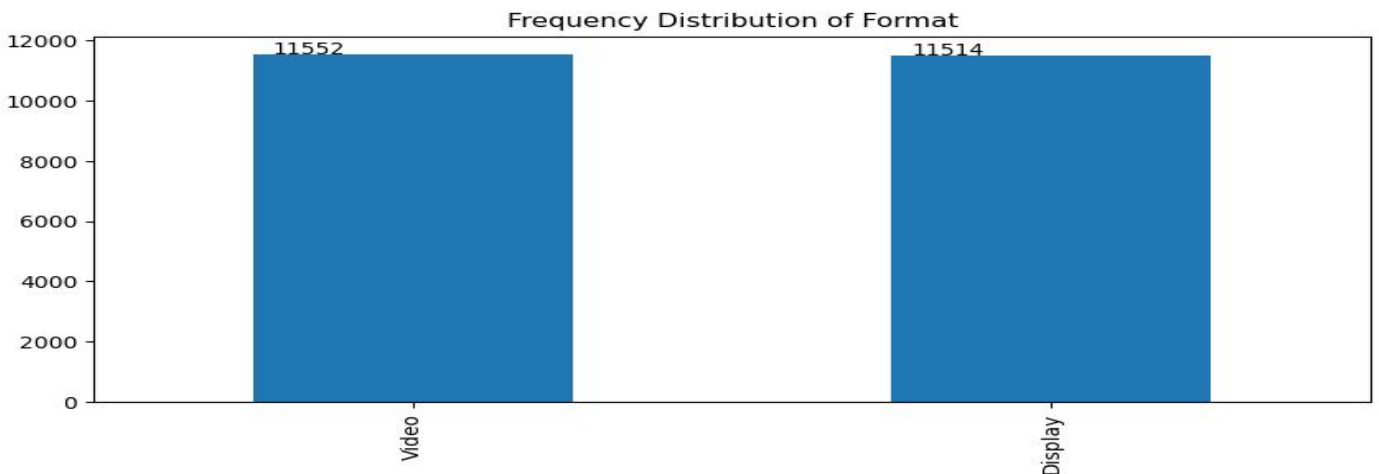


Fig-15- UNIVARIATE ANALYSIS FOR CATEGORICAL COLUMNS

- more inventory type is Format1
- counts of All Ad Type is approximately same
- most count Ad are in video platform, less counts in App platform
- mobile device have high Ad than desktop
- Format of Ads video type and display type are approximately same

1---1--e)-Bivariate analysis

#plot the correlation map

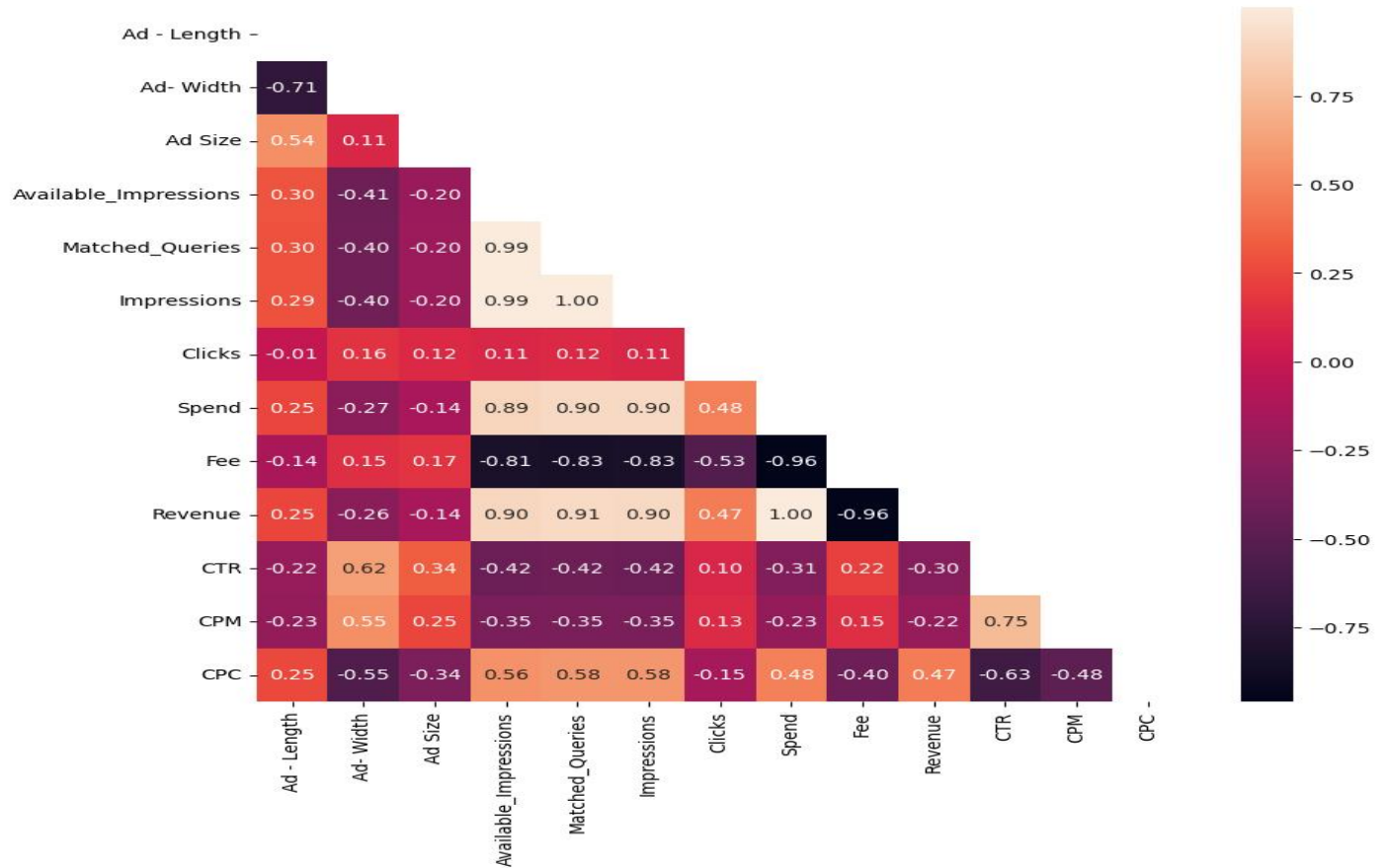


Fig-16-correlation map

Observation

- 'Spend' shows high correlation with 'Available impression','Matched Queries','impression'
- "Revenue"shows high correlation with 'Available impression','Matched Queries','impression','Spend'
- "Fee"negatively correlated with 'Available impression','Matched Queries','impression','Spend','Revenue'
- 'CPC'negatively correlated with 'CTR','CPM'

#barplot of impression with Ad Type

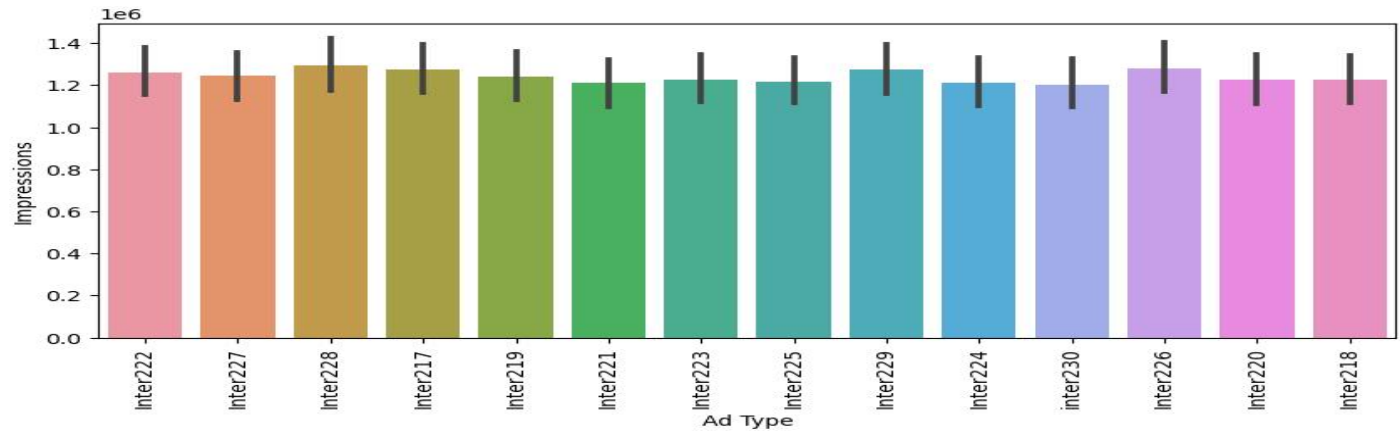


Fig-17-barplot of impression with Ad Type

impression for every Ad Types are approximately same

#barplot of CTR with Device Type

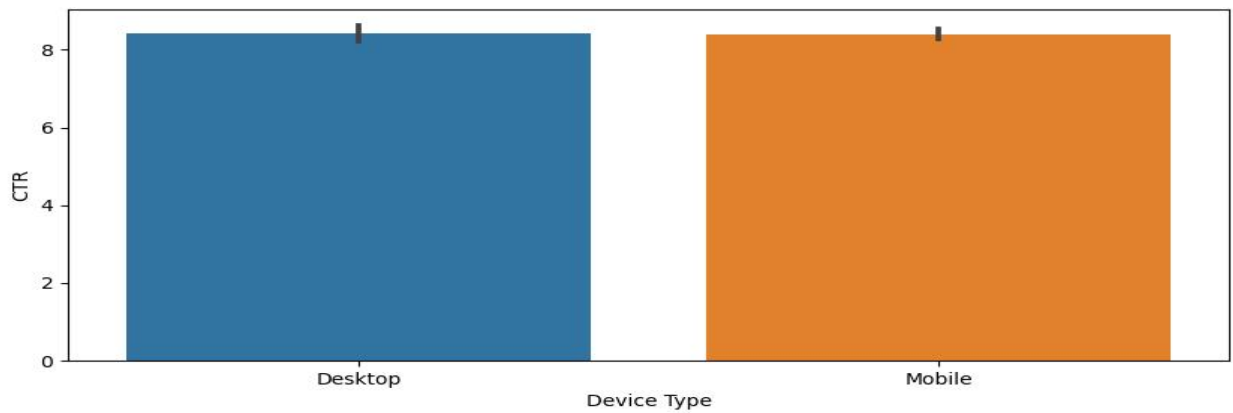


Fig-18-barplot of CTR with Device Type

CTR for both device type are same

#barplot with CPM with Device Type

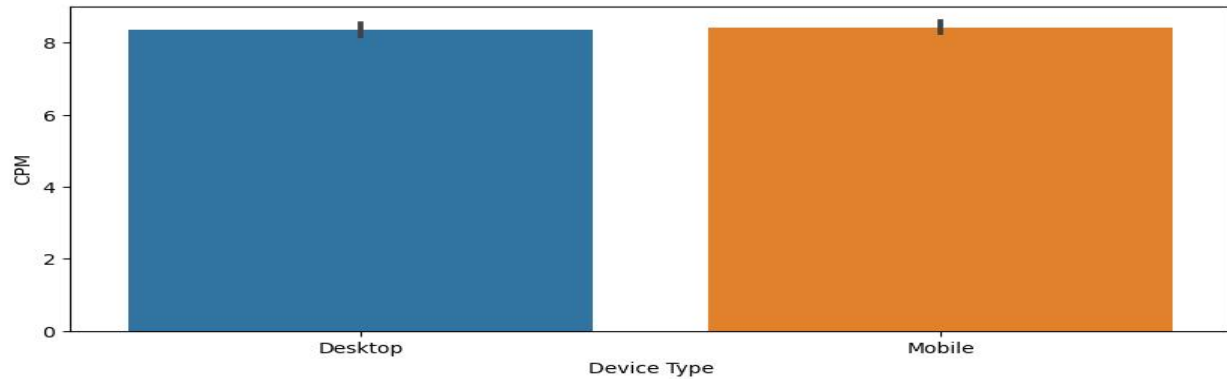


Fig-19-barplot with CPM with Device Type

CPM for both device type are same

#barplot of clicks with device type

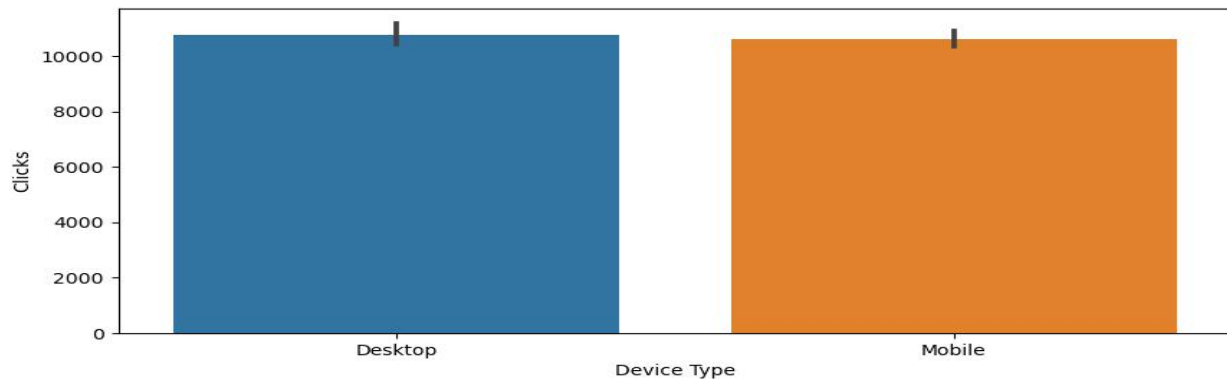


Fig-20-barplot of clicks with device type

Clicks for both device type are same.

Key meaningful observations on individual variables and the relationship between variables

- There are 13 numeric fields in the data
- Customer ad length ranges from 120 to 728
- maximum of Ad_width is 600
- Ad size ranges from 33600 to 216000 with an average 7200
- Available_Impressions ranges from 1.000000e+00 to 2.759286e+07
- only few matched queries above 1.4
- Average impression is around 2.252900e+05
- more Ads are zero clicks very less ads are more than 80000 Clicks
- Spend range is 0.00 to 26931.870000

- 50%,75%,maximum of fee of Ads is 350000
- maximum of Revenue of Ad is 21276.180000
- CTR range in between 0.00 to 1.00
- maximum of CPM is 81.560000
- Range of CPC 0.000 to 7.260000 with average 0.160000
- Outliers to be treated
- more inventory type is Format1
- counts of All Ad Type is approximately same
- most count Ad are in video platform,less counts in App platform
- mobile device have high Ad than desktop
- Format of Ads video type and display type are approximately same
- *relationship between variables*
- 'Spend' shows high correlation with 'Available impression','Matched Queries','impression'
- 'Revenue' shows high correlation with 'Available impression','Matched Queries','impression','Spend'
- 'Fee' negatively correlated with 'Available impression','Matched Queries','impression','Spend','Revenue'
- 'CPC' negatively correlated with 'CTR','CPM'
- CTR for both device type are same
- CPM for both device type are same
- Clicks for both device type are same

Part 1:--2) Clustering: Data Preprocessing

- Missing value check and treatment - Outlier Treatment - z-score scaling Note: Treat missing values in CPC, CTR and CPM using the formula given.

1---2)-a)Missing value check and treatment

#check missing values

CPM,CTR,CPC columns have 4736 null values,this null values treat using these equations.

$CPM = (\text{Total Campaign Spend} / \text{Number of Impressions}) * 1000$

$CPC = \text{Total Cost (spend)} / \text{Number of Clicks}$.

$CTR = \text{Total Measured Clicks} / \text{Total Measured Ad Impressions} \times 100$

1--2)-b) Check if there are any outliers

#Check for presence of outliers in each feature

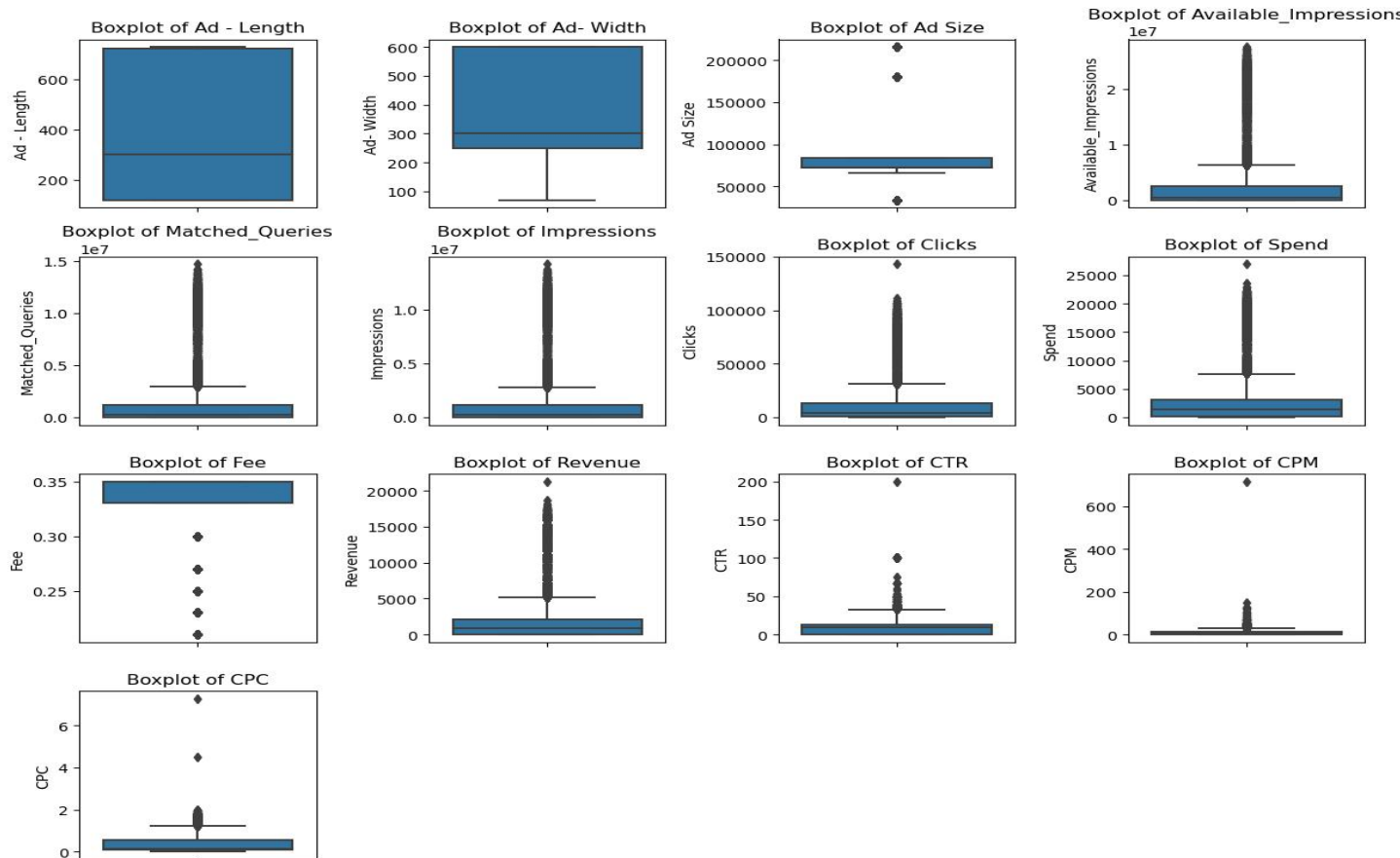


Fig-21-check outliers

treat outliers using IQR method

1--2-c)-z-score scaling

#scaling data frame using z-score method, Scaled data

Ad - Length	Ad-Width	Ad Size	Available_Impressions	Matched_Queries	Impression_s	Clicks	Spend	Fee	Revenue	CTR	CPM	CPC	
0	-0.36449	-0.43279	-0.352218	0.512407	-0.515248	-0.51091	-0.61531	-0.66537	0.465447	-0.61969	-0.87459	-0.92705	-0.98661

Ad - Length	Ad-Width	Ad Size	Available_Impressions	Matched_Queries	Impressions	Clicks	Spent	Fee	Revenue	CTR	CPM	CPC	
	6	7				8	1	2		3	3	4	5
1	0.364496	0.432797	-0.352218	0.512413	0.515264	0.510933	0.615311	0.665372	0.465447	0.619693	0.870136	0.927054	0.986615
2	0.364496	0.432797	-0.352218	0.512213	0.515235	0.510905	0.615311	0.665372	0.465447	0.619693	0.877606	0.927054	0.986615
3	0.364496	0.432797	-0.352218	0.512276	0.515179	0.510847	0.615311	0.665372	0.465447	0.619693	0.886208	0.927054	0.986615
4	0.364496	0.432797	-0.352218	0.512531	0.515281	0.510951	0.615311	0.665372	0.465447	0.619693	0.863404	0.927054	0.986615
...
23061	1.433093	0.186599	1.939086	0.512788	0.515377	0.511050	0.615311	0.665355	0.465447	0.619678	9.888962	6.801294	0.781484
23062	1.433093	0.186599	1.939086	0.512787	0.515376	0.511050	0.615311	0.665362	0.465447	0.619684	4.490471	1.281046	0.869397
23063	1.433093	0.186599	1.939086	0.512788	0.515377	0.511050	0.615311	0.665360	0.465447	0.619682	9.888962	4.593195	0.840092
23064	1.134891	1.290590	-0.400970	0.512787	0.515377	0.511050	0.615311	0.665355	0.465447	0.619678	9.888962	6.801294	0.781484
23065	1.433093	0.186599	1.939086	0.512788	0.515376	0.511050	0.615311	0.665350	0.465447	0.619674	4.490471	4.041170	0.722875

23066 rows × 13 columns

Table-3- Scaled data

Part 1:--3) Clustering: Hierarchical Clustering

- Construct a dendrogram using Ward linkage and Euclidean distance - Identify the optimum number of Clusters

1---3-a)-Construct a dendrogram using Ward linkage and Euclidean distance

#plot dendrogram

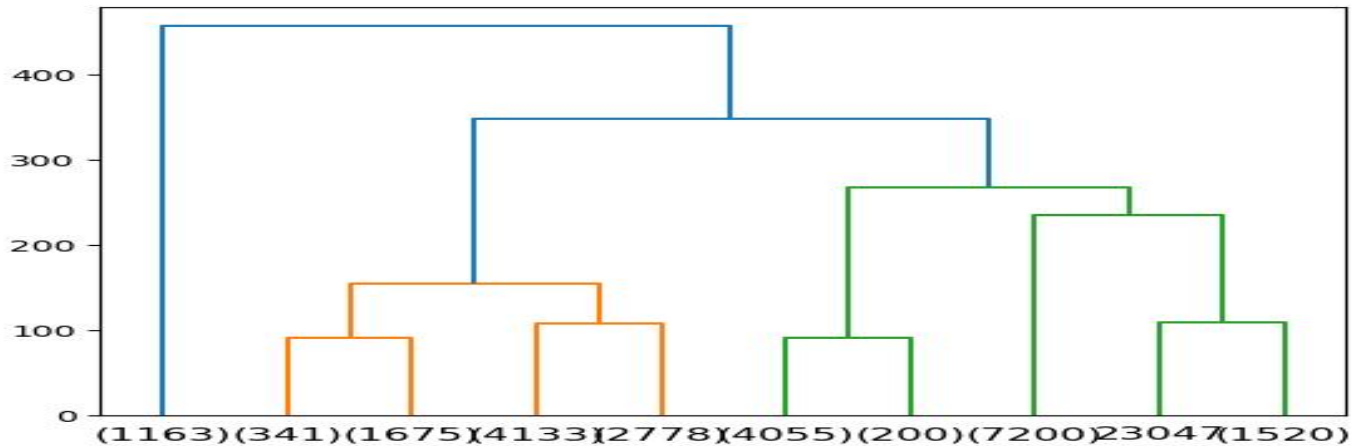


Fig-22-dendrogram

1-3-b) Identify the optimum number of Clusters

In this example, you might decide to cut the dendrogram at a certain height (horizontal line) where the vertical lines are longest, and you would consider the number of clusters as the number of branches you've crossed. we can extract the optimal number of clusters by looking dendrogram. By looking at the dendrogram ,we can say that the optimmal number of clusters is 2.because the vertical distance is high for 2 clusters

Optimum number of clusters=2

Part 1--4)-: Clustering: K-means Clustering

- Apply K-means Clustering - Plot the Elbow curve - Check Silhouette Scores - Figure out the appropriate number of clusters - Cluster Profiling

1--4--a)Apply K-means Clustering

Elbow plot(up to n=10)

Wss=

```
[299858.0,
188281.49,
130712.22,
94685.77,
66289.08,
55262.54,
49134.2,
44094.3,
40161.91,
34878.56]
```

1-4-b)-Plot the Elbow curve

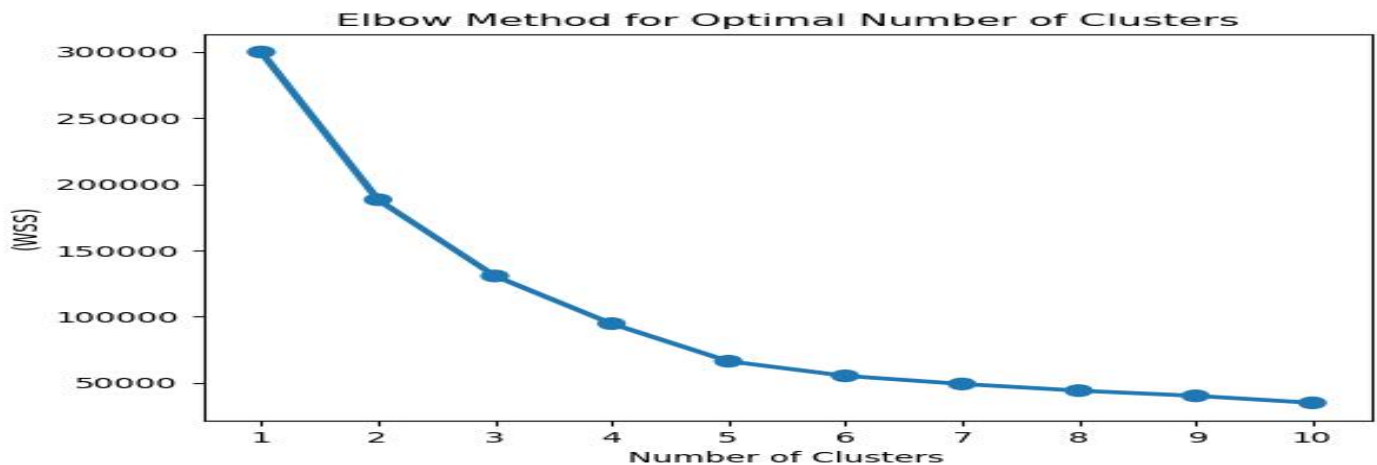


Fig-23-Elbow curve

We can see from the plot that there is a consistent dip from 2 to 10 and there doesn't seem to be a clear 'elbow' here. We may choose any from 2 to 8 as our # of clusters. So, let's look at another method to get a 'second opinion from maths'. Let's create a plot with Silhouette scores to see how it varies with k.

1-4-c)-Check Silhouette Scores

Let us now find the Silhouette Score for the values of K from 2 to 10

Silhouette Analysis

```
For n_clusters=2, the silhouette score is 0.6122159547090302
For n_clusters=3, the silhouette score is 0.38939835242959714
For n_clusters=4, the silhouette score is 0.5033463240941357
For n_clusters=5, the silhouette score is 0.557117025455947
For n_clusters=6, the silhouette score is 0.5291811252973383
For n_clusters=7, the silhouette score is 0.5367029944777137
For n_clusters=8, the silhouette score is 0.5381765938866728
For n_clusters=9, the silhouette score is 0.5531121077559789
For n_clusters=10, the silhouette score is 0.5534060301063022
```

1-4-d)Figure out the appropriate number of clusters

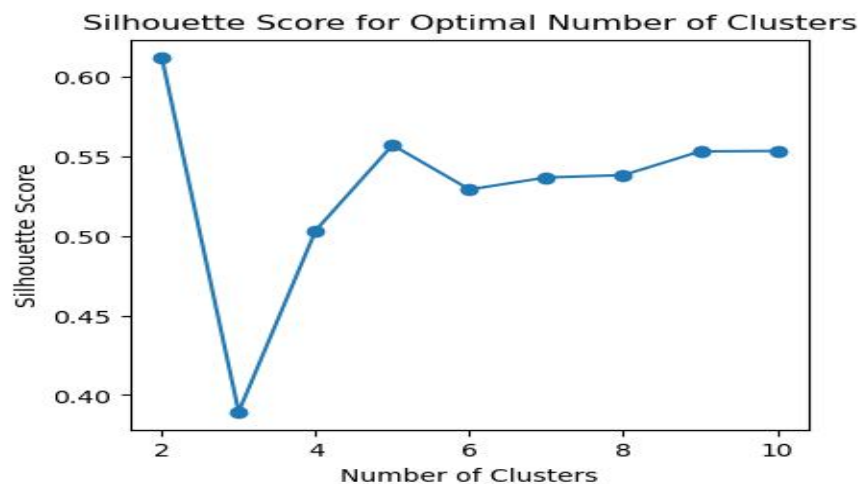


Fig-24-plot of silhouette score

1--4--e)Cluster Profiling

Profile the ads based on optimum number of clusters using silhouette score and your domain understanding

[Hint: Group the data by clusters and take sum or mean to identify trends in clicks, spend, revenue, CPM, CTR, & CPC based on Device Type. Make bar plots.]

We can see from the plot that silhouette score is highest for k=2. Well that makes it slightly easy for us and we can start with first understanding these 2 clusters. So let's take the number of clusters as 2.

#let's take the number of clusters as 2.

#Adding predicted labels to the original data and scaled data

#counts for clusters

Observations:

This looks like a very skewed clustering . Let's check out the profiles of these clusters

#Calculating mean and median of the original data for each label

	group_0 Mean	group_1 Mean	group_0 Median	group_1 Median
Ad - Length	3.645779e+02	6.841050e+02	300.000000	7.280000e+02
Ad- Width	3.531738e+02	1.160296e+02	300.000000	9.000000e+01
Ad Size	9.852045e+04	6.986673e+04	72000.000000	6.552000e+04
Available_Impressions	1.349530e+06	1.815254e+07	408718.000000	1.894390e+07
Matched_Queries	7.173359e+05	9.685496e+06	221739.000000	1.024037e+07
Impressions	6.848326e+05	9.325842e+06	188945.500000	9.857084e+06
Clicks	1.019593e+04	1.768678e+04	3912.000000	1.896250e+04

	group_0 Mean	group_1 Mean	group_0 Median	group_1 Median
Spend	1.818878e+03	1.559868e+04	1299.980000	1.607917e+04
Fee	3.417966e-01	2.382100e-01	0.350000	2.300000e-01
Revenue	1.234364e+03	1.194295e+04	844.985000	1.238096e+04
CTR	8.976118e+00	1.878036e-01	10.178568	1.904888e-01
CPM	8.857767e+00	1.703308e+00	9.106158	1.667699e+00
CPC	2.968454e-01	9.151287e-01	0.127471	8.573520e-01

Table-4- Calculating mean and median of the original data for each label

Observations:

It looks like Cluster 1 with high impression and CPC, Cluster 0 is of less impression,CPC than cluster 1 .

Hint: Group the data by clusters and take sum or mean to identify trends in clicks, spend, revenue, CPM, CTR, & CPC based on Device Type. Make bar plots.]

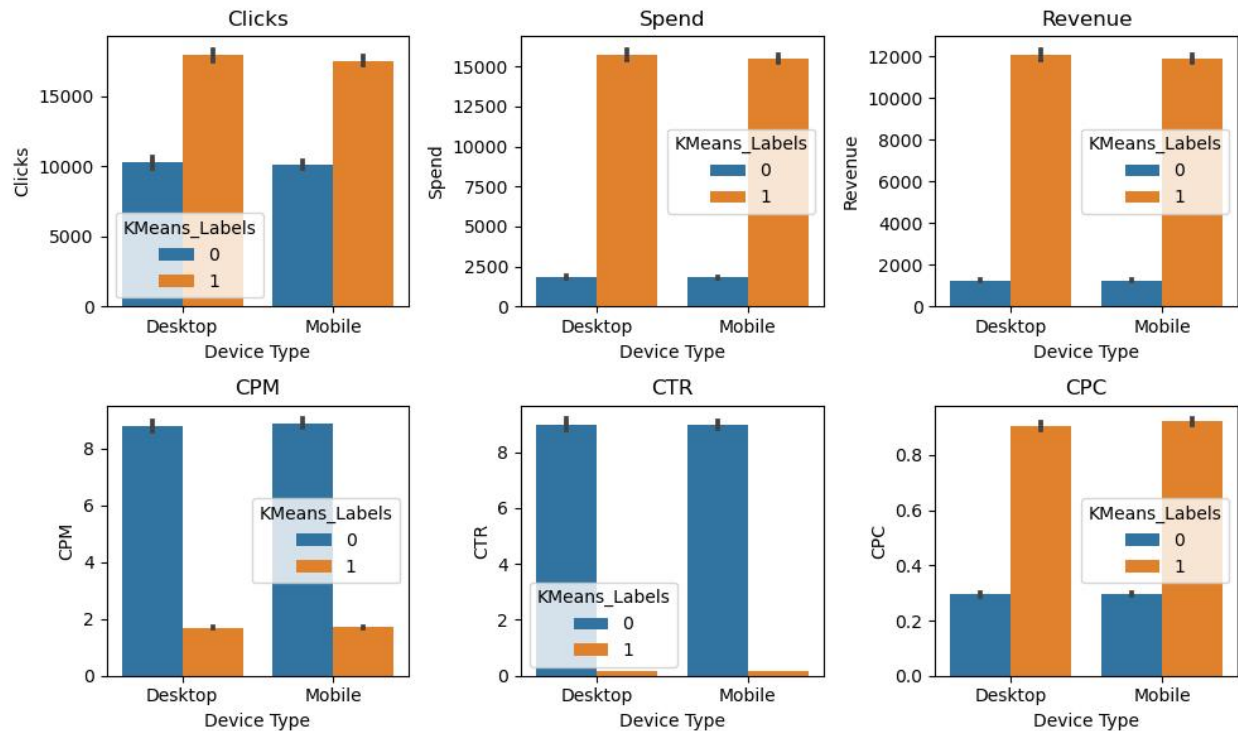


Fig-25-barplot for Group the data by clusters and take sum or mean to identify trends in clicks, spend, revenue, CPM, CTR, & CPC based on Device Type

Conclude the project by providing summary of your learnings.

- There are 23066 rows and 19 columns into the dataset.
- There are no duplicated values in data frame.
- There are 4736 null values in CTR,CPM,CPC columns.
- I have treated missing values in CPC,CTR,CPM columns using the given formula
- It seems that there are outliers into the dataset
- We treated outliers using IQR method
- I have applied z-score method on data frame for scaling .
- I have plotted Dendrogram for value of P=10
- plotted elbow plot and got the optimum value of cluster is 2.
- As per Elbow plot/scree plot ,we conclude that the optimal number of clusters should be 5.
- I have create 2 clusters for the data set.

conclusion after clustering ``

- When clicks on Ads gets increases then revenue is also increases.
- When amount of money spent on specific Ad variation within a specific campaign or ad set is increases then revenue is also increases.
- when impression count of the particular Advertisement increasa then revenue is also increases.
- the clicks for both device type have approximately same in both clusters.
- all variables are same in both devices type.

- click, spend, revenue, cpc are high in cluster 1 type.
- CPM and CTR are high in cluster 0 type .

Problem 2:

PCA FH (FT): Primary census abstract for female headed households excluding institutional households (India & States/UTs - District Level), Scheduled tribes - 2011 PCA for Female Headed Household Excluding Institutional Household. The Indian Census has the reputation of being one of the best in the world. The first Census in India was conducted in the year 1872. This was conducted at different points of time in different parts of the country. In 1881 a Census was taken for the entire country simultaneously. Since then, Census has been conducted every ten years, without a break. Thus, the Census of India 2011 was the fifteenth in this unbroken series since 1872, the seventh after independence and the second census of the third millennium and twenty first century. The census has been uninterruptedly continued despite of several adversities like wars, epidemics, natural calamities, political unrest, etc. The Census of India is conducted under the provisions of the Census Act 1948 and the Census Rules, 1990. The Primary Census Abstract which is important publication of 2011 Census gives basic information on Area, Total Number of Households, Total Population, Scheduled Castes, Scheduled Tribes Population, Population in the age group 0-6, Literates, Main Workers and Marginal Workers classified by the four broad industrial categories, namely, (i) Cultivators, (ii) Agricultural Laborers, (iii) Household Industry Workers, and (iv) Other Workers and also Non-Workers. The characteristics of the Total Population include Scheduled Castes, Scheduled Tribes, Institutional and Houseless Population and are presented by sex and rural-urban residence. Census 2011 covered 35 States/Union Territories, 640 districts, 5,924 sub-districts, 7,935 Towns and 6,40,867 Villages. The data collected has so many variables thus making it difficult to find useful details without using Data Science Techniques. You are tasked to perform detailed EDA and identify Optimum Principal Components that explains the most variance in data. Use Sklearn only.

Note: The 24 variables given in the Rubric is just for performing EDA. You will have to consider the entire dataset, including all the variables for performing PCA. Data file - PCA India Data Census.xlsx

DATA DICTIONARY

State Code: Code representing the state.

Dist. Code: Code representing the district within the state.

State, Area Name: Name of the state AND district.

No_HH: Number of households.

TOT_M: Total male population.

TOT_F: Total female population.

M_06, F_06: Male and female population in the age group 0-6.

M_SC, F_SC: Male and female population belonging to the Scheduled Caste.

M_ST, F_ST: Male and female population belonging to the Scheduled Tribe.

M_LIT, F_LIT: Male and female literate population.

M_ILL, F_ILL: Male and female illiterate population.

TOT_WORK_M, TOT_WORK_F: Total working male and female population.

MAINWORK_M, MAINWORK_F: Main working male and female population.

MAIN_CL_M, MAIN_CL_F: Main working male and female population engaged in cultivation.

MAIN_AL_M, MAIN_AL_F: Main working male and female population engaged in agriculture and allied activities.

MAIN_HH_M, MAIN_HH_F: Main working male and female population engaged in household industry.

MAIN_OT_M, MAIN_OT_F: Main working male and female population engaged in other occupations.

MARGWORK_M, MARGWORK_F: Marginal working male and female population.

MARG_CL_M, MARG_CL_F: Marginal working male and female population engaged in cultivation.

MARG_AL_M, MARG_AL_F: Marginal working male and female population engaged in agriculture and allied activities.

MARG_HH_M, MARG_HH_F: Marginal working male and female population engaged in household industry.

MARG_OT_M, MARG_OT_F: Marginal working male and female population engaged in other occupations.

MARGWORK_3_6_M, MARGWORK_3_6_F: Marginal working male and female population in the age group 3-6.

MARG_CL_3_6_M, MARG_CL_3_6_F: Marginal working male and female population in the age group 3-6 engaged in cultivation.

MARG_AL_3_6_M, MARG_AL_3_6_F: Marginal working male and female population in the age group 3-6 engaged in agriculture and allied activities.

MARG_HH_3_6_M, MARG_HH_3_6_F: Marginal working male and female population in the age group 3-6 engaged in household industry.

MARG_OT_3_6_M, MARG_OT_3_6_F: Marginal working male and female population in the age group 3-6 engaged in other occupations.

MARGWORK_0_3_M, MARGWORK_0_3_F: Marginal working male and female population in the age group 0-3.

MARG_CL_0_3_M, MARG_CL_0_3_F: Marginal working male and female population in the age group 0-3 engaged in cultivation.

MARG_AL_0_3_M, MARG_AL_0_3_F: Marginal working male and female population in the age group 0-3 engaged in agriculture and allied activities.

MARG_HH_0_3_M, MARG_HH_0_3_F: Marginal working male and female population in the age group 0-3 engaged in household industry.

MARG_OT_0_3_M, MARG_OT_0_3_F: Marginal working male and female population in the age group 0-3 engaged in other occupations.

NON_WORK_M, NON_WORK_F: Non-working male and female population.

2-1: PCA: Define the problem and perform Exploratory Data Analysis

- Problem Definition - Check shape, Data types, statistical summary - Perform an EDA on the data to extract useful insights Note: 1. Pick 5 variables out of the given 24 variables below for EDA: No_HH, TOT_M, TOT_F, M_06, F_06, M_SC, F_SC, M_ST, F_ST, M_LIT, F_LIT, M_ILL, F_ILL, TOT_WORK_M, TOT_WORK_F, MAINWORK_M, MAINWORK_F, MAIN_CL_M, MAIN_CL_F, MAIN_AL_M, MAIN_AL_F, MAIN_HH_M, MAIN_HH_F, MAIN_OT_M, MAIN_OT_F 2. Example questions to answer from EDA - (i) Which state has highest gender ratio and which has the lowest? (ii) Which district has the highest & lowest gender ratio?

ANSWER

#read the data set

2-1-b) Check shape

(640, 61)-640 rows, 61 columns

2-1-c) Data types

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 640 entries, 0 to 639
```

```
Data columns (total 61 columns):
```

#	Column	Non-Null Count	Dtype
0	State Code	640 non-null	int64
1	Dist.Code	640 non-null	int64
2	State	640 non-null	object
3	Area Name	640 non-null	object
4	No_HH	640 non-null	int64
5	TOT_M	640 non-null	int64
6	TOT_F	640 non-null	int64
7	M_06	640 non-null	int64
8	F_06	640 non-null	int64
9	M_SC	640 non-null	int64
10	F_SC	640 non-null	int64

11	M_ST	640 non-null	int64
12	F_ST	640 non-null	int64
13	M_LIT	640 non-null	int64
14	F_LIT	640 non-null	int64
15	M_ILL	640 non-null	int64
16	F_ILL	640 non-null	int64
17	TOT_WORK_M	640 non-null	int64
18	TOT_WORK_F	640 non-null	int64
19	MAINWORK_M	640 non-null	int64
20	MAINWORK_F	640 non-null	int64
21	MAIN_CL_M	640 non-null	int64
22	MAIN_CL_F	640 non-null	int64
23	MAIN_AL_M	640 non-null	int64
24	MAIN_AL_F	640 non-null	int64
25	MAIN_HH_M	640 non-null	int64
26	MAIN_HH_F	640 non-null	int64
27	MAIN_OT_M	640 non-null	int64
28	MAIN_OT_F	640 non-null	int64
29	MARGWORK_M	640 non-null	int64
30	MARGWORK_F	640 non-null	int64
31	MARG_CL_M	640 non-null	int64
32	MARG_CL_F	640 non-null	int64
33	MARG_AL_M	640 non-null	int64
34	MARG_AL_F	640 non-null	int64
35	MARG_HH_M	640 non-null	int64
36	MARG_HH_F	640 non-null	int64
37	MARG_OT_M	640 non-null	int64
38	MARG_OT_F	640 non-null	int64
39	MARGWORK_3_6_M	640 non-null	int64
40	MARGWORK_3_6_F	640 non-null	int64
41	MARG_CL_3_6_M	640 non-null	int64
42	MARG_CL_3_6_F	640 non-null	int64
43	MARG_AL_3_6_M	640 non-null	int64
44	MARG_AL_3_6_F	640 non-null	int64
45	MARG_HH_3_6_M	640 non-null	int64
46	MARG_HH_3_6_F	640 non-null	int64
47	MARG_OT_3_6_M	640 non-null	int64
48	MARG_OT_3_6_F	640 non-null	int64
49	MARGWORK_0_3_M	640 non-null	int64
50	MARGWORK_0_3_F	640 non-null	int64
51	MARG_CL_0_3_M	640 non-null	int64
52	MARG_CL_0_3_F	640 non-null	int64
53	MARG_AL_0_3_M	640 non-null	int64
54	MARG_AL_0_3_F	640 non-null	int64

```
55 MARG_HH_0_3_M      640 non-null      int64
56 MARG_HH_0_3_F      640 non-null      int64
57 MARG_OT_0_3_M      640 non-null      int64
58 MARG_OT_0_3_F      640 non-null      int64
59 NON_WORK_M          640 non-null      int64
60 NON_WORK_F          640 non-null      int64
```

dtypes: int64(59), object(2)

memory usage: 305.1+ KB

2-1-d)statistical summary

St	D																			
at	is																			
C	C	N	T	T	M	F	M	F	M	.	M	M	M	M	M	M	M	N	N	
o	o	o	O	O	_	_	_	_	_	.	AR	AR	AR	AR	AR	AR	AR	O	O	
d	d	H	T	T	6	06	C	SC	T	.	G_	G_	G_	G_	G_	G_	G_	N_	_	
e	e	H	M	F						.	CL	CL	AL	AL	HH	HH	OT	W	W	
										.	_0	_0	_0	_0	_0	_0	_0	ORK	O	
										.	_3	_3	_3	_3	_3	_3	_3		R	
										.	-	-	-	-	-	-	-	-	K_	
										.	M	F	M	F	M	F	M	M	F	

	6	6																		
	4	4																		
c	0.	0.	64	64	64	64	64	64	64	64	64	64	64	64	64	64	64	64	64	64
o	0	0	0.	0.	0.	0.	0.	0.	0.	.	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.	0.	0.
u	0	0	00	00	00	00	00	00	00	.	00	00	00	00	00	00	00	00	00	00
n	0	0	00	00	00	00	00	00	00	.	00	00	00	00	00	00	00	00	00	00
t	0	0	00	00	00	00	00	00	00	.	0	0	0	0	0	0	0	00	00	00
	0	0																		
	0	0																		

	1	3																		
	7.	2			12															
m	1	0.	51	79	23	12	11	13	20	61	13	27	25	55	56	12	71.	20	51	70
e	1	5	22	94	72	30	94	82	77	91	92.	57.	0.8	8.0	0.6	93.	37	0.7	0.	4.
a	1	0	2.	0.	.0	9.	2.	0.	8.	.8	97	05	89	98	90	43	96	42	01	77
n	4	0	87	57	84	09	30	94	39	07	34	00	06	43	62	12	88	18	40	81
	0	0	18	65	37	84	00	68	21	81	38	00	2	8	5	50		8	63	25
	6	0	75	63	5	38	00	75	88	3										
	2	0																		

	9.	1	48	73	11	11	14	21	99											
	4	8	13	38	36	50	32	42	72	12	14	27	45	11	76	15	10	30	61	91
s	2	4.	5.	4.	00	0.	6.	6.	7.	.6	89.	88.	3.3	17.	2.5	85.	7.8	9.7	0.	0.
t	6	8	40	51	.7	90	29	37	88	68	70	77	36	64	78	37	97	40	60	20
d	4	9	54	11	17	68	45	31	77	94	70	66	59	27	99	79	62	85	31	92
	8	6	75	14	28	81	67	30	13	8	52	76	4	48	1	36	7	4	87	25

	St at e C o d e	D is t. C o d e	N o_ H H	T O T_ M	T O T_ F	M _0 6	F _06	M _S C	F _SC	M _S T	.	M AR G_ CL _0 _3 _M	M AR G_ CL _0 _3 _F	M AR G_ AL _0 _3 _M	M AR G_ AL _0 _3 _F	M AR G_ HH _0 _3 _M	M AR G_ HH _0 _3 _F	M AR G_ OT _0 _3 _M	M AR G_ OT _0 _3 _F	N O N_ W O RK _M	N O N_ _W O R K_ F
	6	3 6 7			2																
m i n	1. 0 0 0 0 0 0	1. 0 0 0 0 0 0	35 0. 00 00 00 00	39 1. 00 00 00 00	69 8. 00 00 00 00	56 .0 00 00 00 0	56 .0 00 00 00 0	0. 00 00 00 00	0. 00 00 00 00	0. 00 00 00 00	.	4.0 00 00 00 0	30. 00 00 00 00	0.0 00 00 00 0	0.0 00 00 00 0	0.0 00 00 00 0	0.0 00 00 00 0	0.0 00 00 00 0	0.0 00 00 00 0	0. 00 00 00 00	5. 00 00 00 00
2 5 %	9. 0 0 0 0 0 0	1 6 0. 7 5 0 0 0 0 0	19 48 4. 00 00 00 00	30 22 8. 00 00 00 00	46 51 7. 75 50 00 00	47 33 .7 50 00 00 0	46 72 .2 50 00 00 0	34 66 .2 50 00 00 0	56 03 .2 50 00 00 0	29 3. 75 00 00 00	.	48 9.5 00 00 0	95 7.2 50 00 0	47. 00 00 00	10 9.0 00 00 0	13 6.5 00 00 0	29 8.0 00 00 0	14. 00 00 00 00	43. 00 00 00 00	16 1. 00 00 00 00	22 0. 50 00 00 00
5 0 %	1 8. 0 0 0 0 0 0 0	3 2 0. 5 0 0 0 0 0 0	35 83 7. 00 00 00 00	58 33 9. 00 00 00 00	87 72 4. 50 00 00 00	91 59 .0 00 00 0	86 63 .0 00 00 0	95 91 .5 00 00 0	13 70 9. 00 00 00 00	23 33 .5 00 00 00 0	.	94 9.0 00 00 0	19 28. 00 00 00	11 4.5 00 00 0	24 7.5 00 00 0	30 8.0 00 00 0	71 7.0 00 00 0	35. 00 00 00 00	11 3.0 00 00 0	32 6. 00 00 00	46 4. 50 00 00
7 5 %	2 4. 0 0	4 8 0. 2	68 89 2. 00	10 79 18 .5	16 42 51 .7	16 52 0. 25	15 90 2. 25	19 42 9. 75	29 18 0. 00	76 58 .0 00	.	17 14. 00 00	35 99. 75 00	27 0.7 50 00	56 8.7 50 00	64 2.0 00 00	17 10. 75 00	79. 00 00	24 0.0 00 00	60 4. 50 00	85 3. 50 00

State	Dist.code	No_HH	TOT_M	TOT_F	M_06	F_06	M_SC	F_SC	M_ST	F_ST	M_LIT	F_LIT	M_ILL	F_ILL	TOT_WORK_M	TOT_WORK_F	MAINWORK_M	MAINWORK_F	MAIN_CL_M	MAIN_CL_F	MAIN_AL_M	MAIN_AL_F	MAIN_HH_M	MAIN_HH_F	MAIN_OT_M	MAIN_OT_F	NONWORK	NONWORK

0	5	00	00	50	00	00	00	00	00	00	00	00	0	0	0	00	00	0	00	00
0	0	00	00	00	00	00	00	00	00	0										
0	0		0	0																
0	0																			
	0																			

	3	6																										
	5.	4	31	48	75	96	95	10	15	96																		
	0	0.	04	54	03	22	12	33	64	78																		
	0	0	50	17	92	3.	9.	07	29	5.																		
	0	0	.0	.0	.0	3.	9.	.0	.0	5.																		
	0	0	00	00	00	00	00	00	00	00																		
	0	0	00	00	00	00	00	00	00	00																		
	0	0	00	00	00	00	00	00	00	00																		
	0	0	0	0	0	00	00	0	0	00																		
	0	0																										

8 rows × 59 columns

Table-5-static summary

2-1-e)-Perform an EDA on the data to extract useful insights Note:

2-1-e)-1. Pick 5 variables out of the given 24 variables below for EDA: No_HH, TOT_M, TOT_F, M_06, F_06, M_SC, F_SC, M_ST, F_ST, M_LIT, F_LIT, M_ILL, F_ILL, TOT_WORK_M, TOT_WORK_F, MAINWORK_M, MAINWORK_F, MAIN_CL_M, MAIN_CL_F, MAIN_AL_M, MAIN_AL_F, MAIN_HH_M, MAIN_HH_F, MAIN_OT_M, MAIN_OT_F

Answers

I have picked 5 variables such as TOT_M, TOT_F, M_LIT, F_LIT and TOT_WORK_M and comparing those 5 variables against "state", "Dist.code"

TOT_M-total population male

TOT_F-total population female

M_LIT-literates population male

F_LIT-literates population female

TOT_WORK_M-total work population male

State-state code

District-district code

2-1-e)-2. Example questions to answer from EDA -

(i) Which state has highest gender ratio and which has the lowest?

(ii) Which district has the highest & lowest gender ratio?

Answers:

2-1-e)2-(i) Which state has highest gender ratio and which has the lowest?

using the bar plot we can find which state has highest gender ratio and which has the lowest

ans:

gender ratio=(total population of female/total population of male)*1000

create a column in data for store the gender ratio

#barplot for state based by gender ratio

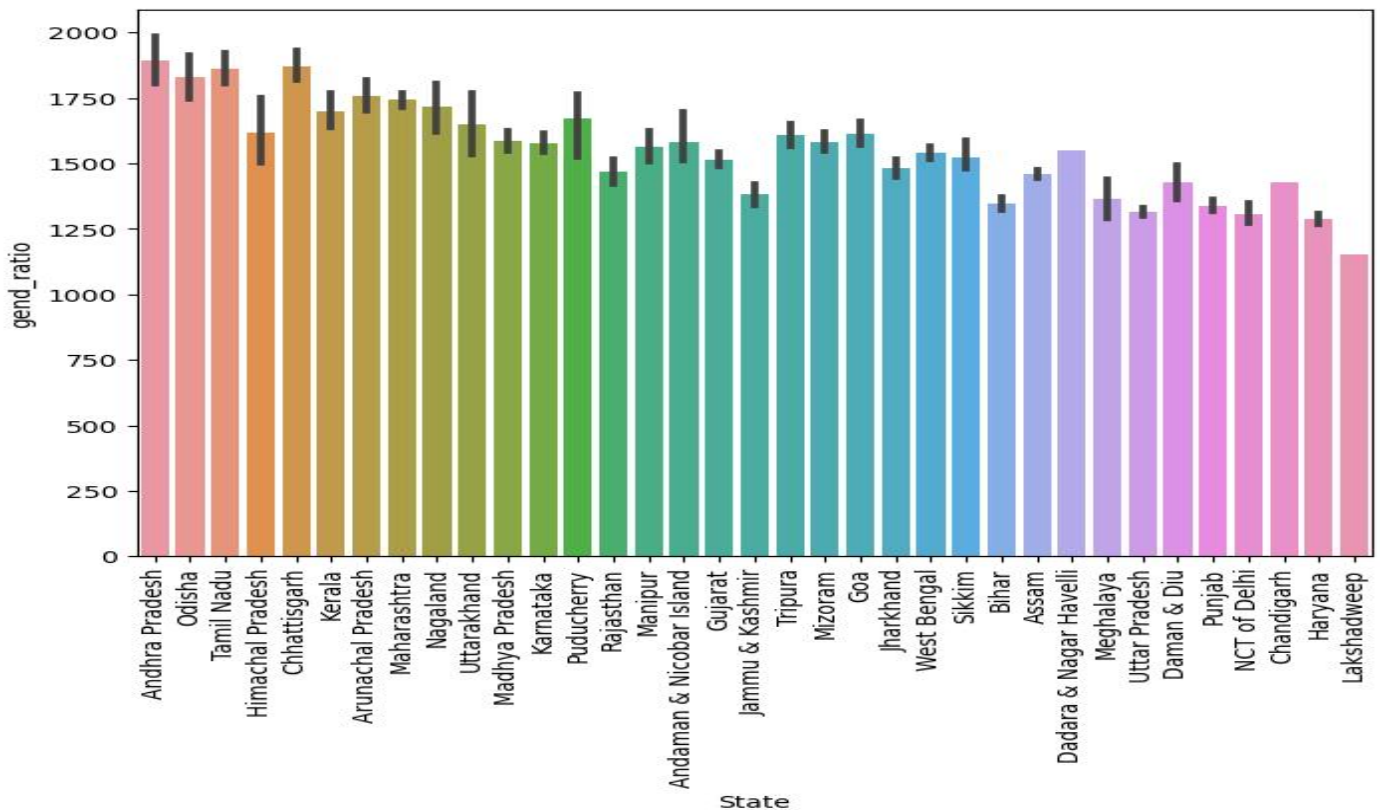


Fig-26-barplot for state based by gender ratio

```

546      Andhra Pradesh
397              Odisha
624          Tamil Nadu
545      Andhra Pradesh
390              Odisha

```

```

...
138      Uttar Pradesh
105      Rajasthan
143      Uttar Pradesh
1        Jammu & Kashmir
586      Lakshadweep
Name: State, Length: 640, dtype: object
Andhra pradesh state have highest gender ratio

```

Lakshadweep state have lowest gender ratio

2-1-e)2-(ii) Which district has the highest & lowest gender ratio?

Ans: sort values Area name based on gender ratio

```

546      Krishna
397      Koraput
624      Virudhunagar
545      West Godavari
390      Baudh
...
138      Baghpat
105      Dhaulpur
143      Mahamaya Nagar
1        Badgam
586      Lakshadweep
Name: Area Name, Length: 640, dtype: object
krishna have high gender ratio

```

lakshadweep have low gender ratio

Part 2-2: PCA: Data Preprocessing

- Check for and treat (if needed) missing values - Check for and treat (if needed) data irregularities - Scale the Data using the z-score method - Visualize the data before and after scaling and comment on the impact on outliers

2-2-a)Check for and treat (if needed) missing values

there is no null values in the dataset

2-2-b)Check for and treat (if needed) data irregularities

there is no duplicated values in the data set

check count the counts and unique values of columns for checking irregularities

#create ne data frame with drop the State Code","Dist.Code","State","Area Name","gend_ratio"]

Check the outliers

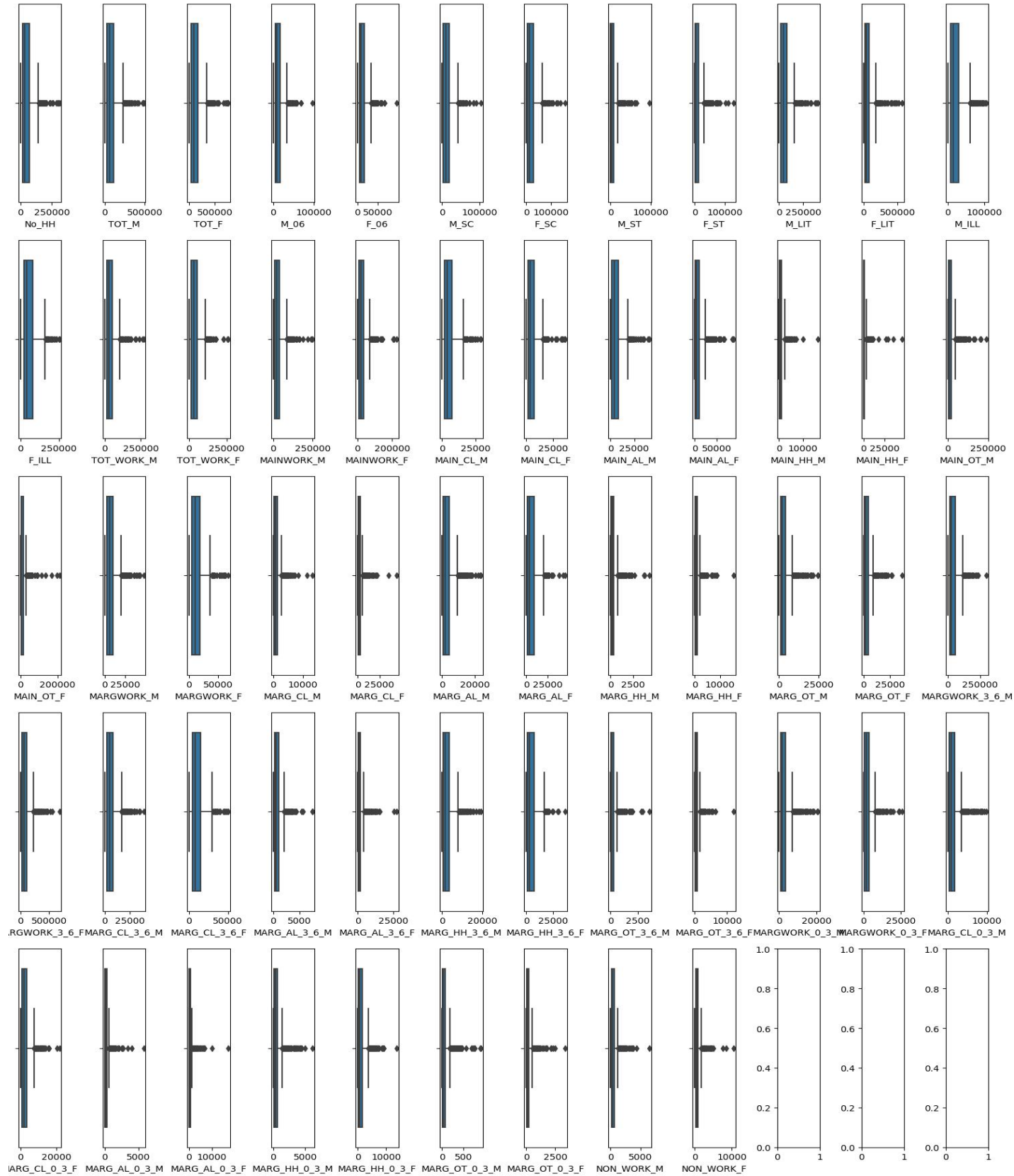
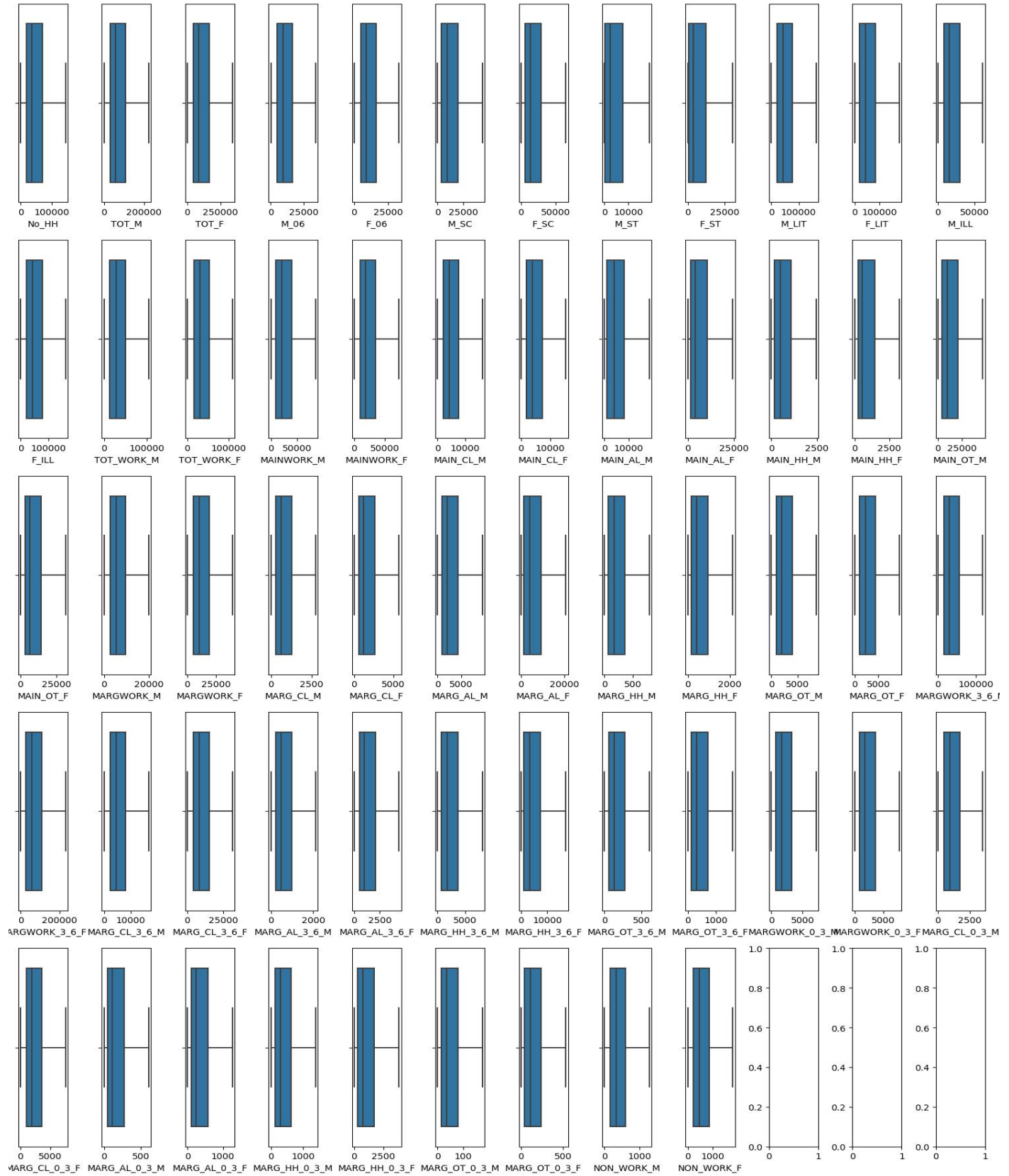


Fig-27-checking outliers

#treat the outliers



Scaling the data using z-score method

N o H	T O T M	T O T F	M 0 6	F 0 6	M S C	F S C	M S T	F S T	M L I T	..	M AR G CL 0 3 M	M AR G CL 0 3 F	M AR G AL 0 3 M	M AR G AL 0 3 F	M AR G HH 0 3 M	M AR G HH 0 3 F	M AR G OT 0 3 M	M AR G OT 0 3 F	N O N W O R K M	N O N W O R K F
4	0 9 3 8 4 9 5	0 9 2 1 3 0 9	0 9 3 5 0 1 8	0 7 0 0 9 3 1	0 7 4 0 5 2 3	1 0 7 8 8 0 7	1 0 7 8 1 6 0	0 4 2 5 3 0 4	0 2 4 4 8 9 7	0.3 698 44	0.2 986 17	1.4 843 98	1.6 331 30	0.5 899 42	0.7 498 82	0.5 892 34	0.3 79 13	0.7 06 20
..
6 3 5	1 1 5 0 3 4 8	1 1 2 7 9 4 9	1 1 3 2 6 6 7	1 1 3 4 7 2 0	1 1 1 7 4 9 5	1 0 7 8 7 2 5	1 0 7 8 3 2 4	0 8 4 2 3 4 4	0 8 3 3 4 1 9	1.2 126 28	1.1 956 66	1.0 057 14	1.0 280 83	1.0 704 66	1.0 144 72	1.0 267 79	1.0 51 59	1.1 02 22
6 3 6	0 9 6 5 0 2 4	1 0 5 8 2 9 9	1 0 5 1 0 4 6	1 0 8 1 6 7 6	1 0 8 0 7 3 4	0 8 9 7 2 1 8	0 8 5 2 9 6	0 8 4 2 3 4 4	0 8 3 2 3 1 7	1.0 895 13	1.0 573 65	0.9 896 49	0.9 923 62	0.9 813 93	0.9 005 75	0.9 538 55	0.9 04 29	0.8 93 50
6 3 7	1 2 0 2 7 4 5	1 2 3 7 6 9 0	1 2 2 2 0 4 5	1 2 3 4 1 0 0	1 2 2 6 4 3 7	1 0 0 0 4 4 3	1 0 7 9 6 3 9	0 6 4 5 5 9 9	0 6 5 8 1 4 0	1.1 405 61	1.1 541 75	0.9 575 19	1.0 178 77	1.0 657 78	1.0 092 15	0.7 168 51	0.7 50 58	0.9 84 48

N o H	T O T H	T O T H	M 6	F 6	M C	F C	M T	F T	M T	..	M AR G CL 0 3 M	M AR G CL 0 3 F	M AR G AL 0 3 M	M AR G AL 0 3 F	M AR G HH 0 3 M	M AR G HH 0 3 F	M AR G OT 0 3 M	M AR G OT 0 3 F	N O N W O R K M	N O N W O R K F
6	1	1	1	1	1	1	1	0	0	1										
3	1	1	1	1	1	0	0	8	8	1	...	1.1	1.1	0.8	0.9	1.0	0.9	1.0	1.0	0.9
8	3	7	7	8	7	8	7	3	2	4		085	360	771	158	446	960	085	25	20
	9	7	3	0	7	0	9	8	8	9		31	53	93	17	82	73	48	97	26
	4	0	5	4	5	4	9	1	7	0								6	7	3
	2	2	9	6	3	4	6	7	4	4										3
	5	9	7	8	9	7	3	8	6	8										6

6	0	0	0	0	1	0	0	8	8	0	...	1.0	1.1	0.9	1.0	1.0	0.9	0.9	1.0	0.7
3	3	6	6	9	0	8	7	1	0	1		714	598	735	229	306	995	903	25	91
9	2	4	4	9	0	0	9	6	7	4		96	98	84	80	17	78	17	97	82
	1	4	5	6	2	4	9	0	3	6								6	9	1
	6	3	9	8	2	4	6	6	6	2										7
	2	0	7	0	9	7	3	3	7	4										3

640 rows × 57 columns

Table-6-scaled data

2-2-4) Visualize the data before and after scaling and comment on the impact on outliers

i have already treated outliers in above question 2-2-b) ,but still i applied the z-score for the scaling of the data set .please find below outputs by boxplot and describe function for before and after

before scaling

Assuming dx_new is your DataFrame, plot boxplot for checking outliers

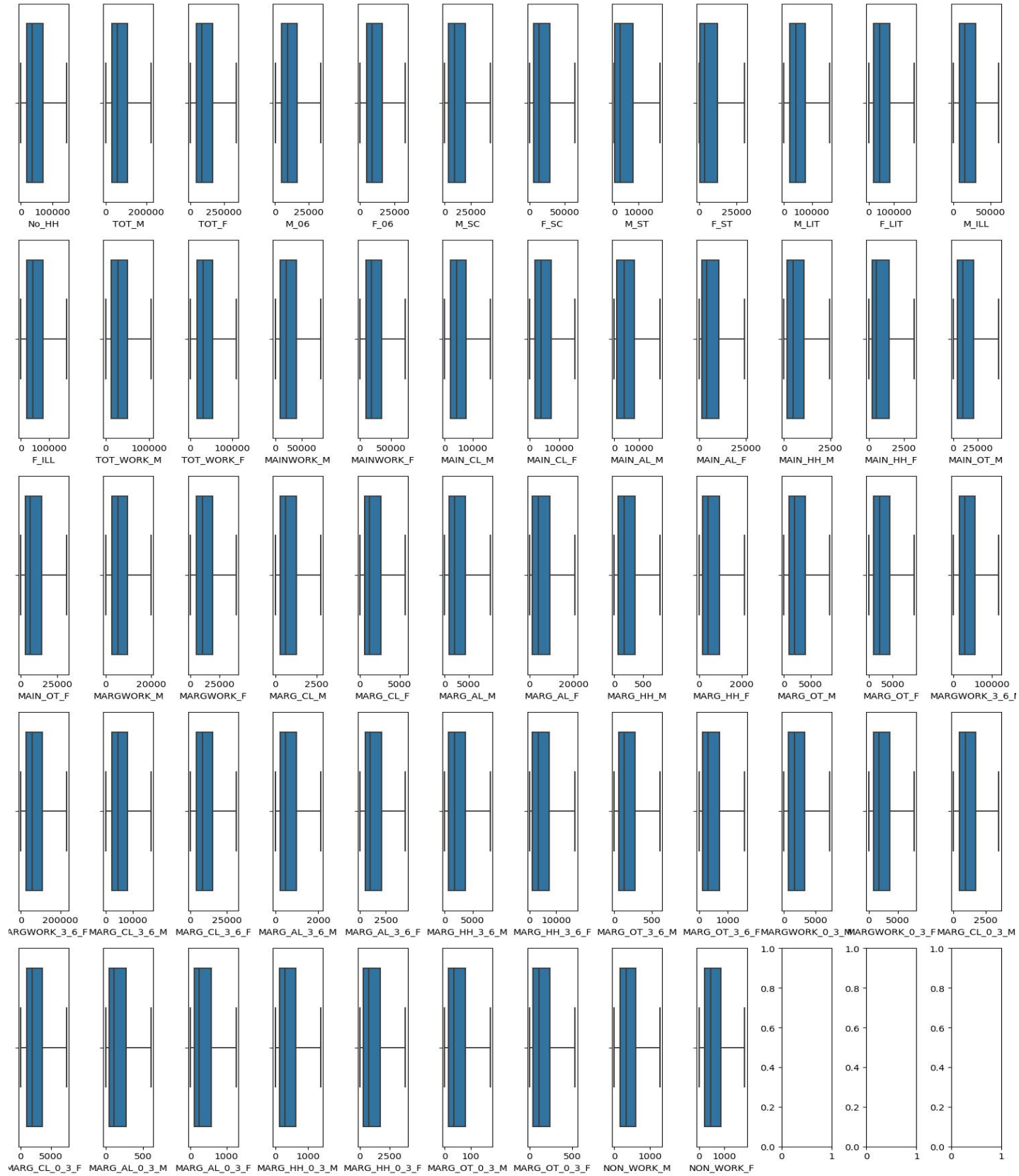


Fig-29-before scalingboxplot for checking outliers

After scaling: *# Assuming dx_new is your DataFrame, after scaling*
boxplot for checking outliers

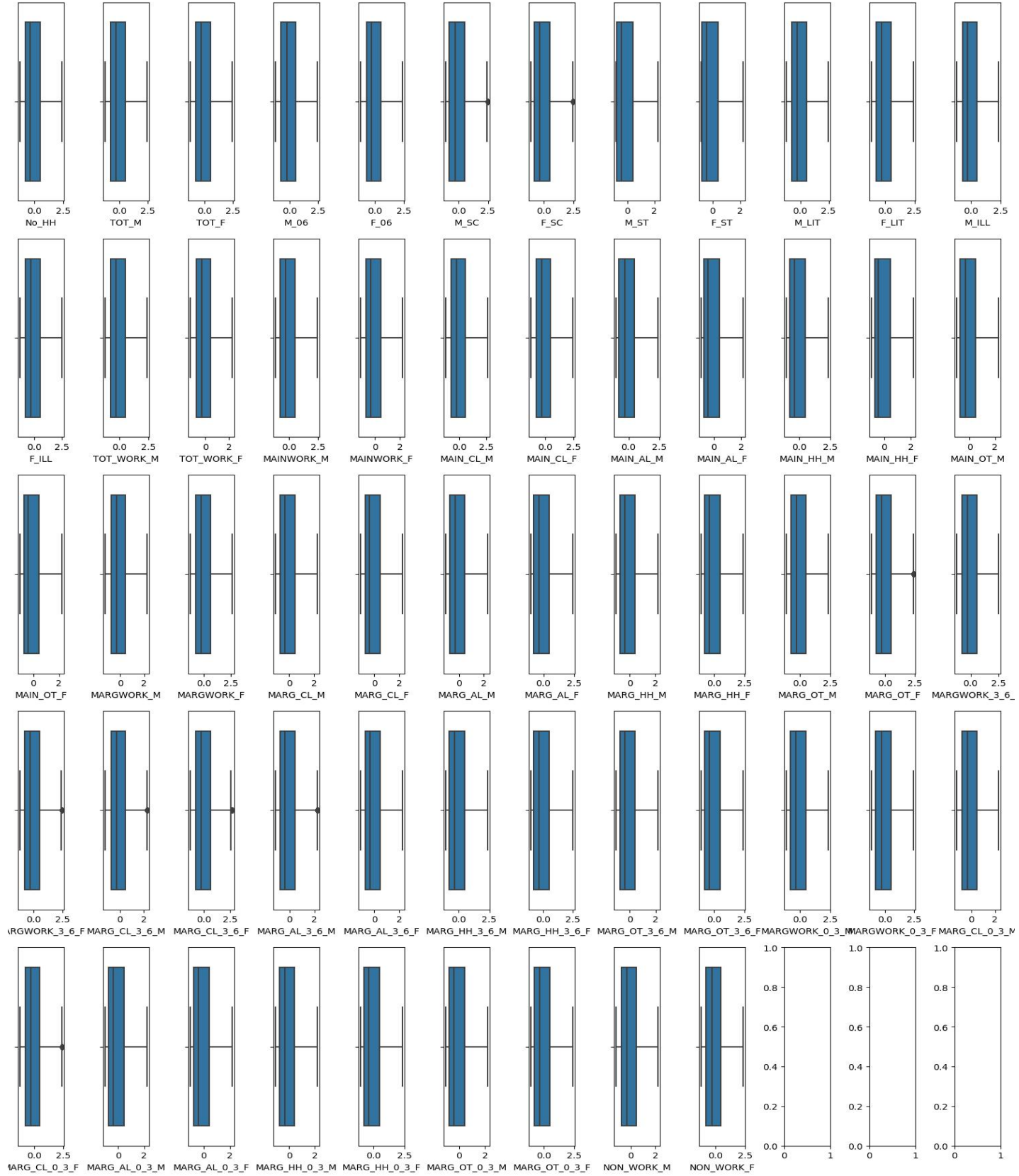


Fig-30-after scaling boxplot for checking outliers

The range of the outliers changes after the scaling.
before scaling

No_HH	TOT_M	TOT_F	M_06	F_06	M_SC	F_SC	M_ST	F_ST	M_LIT	...	MARG_C_L03_M	MARG_C_L03_F	MARG_A_L03_M	MARG_A_L03_F	MARG_H_03_M	MARG_H_03_F	MARG_O_T03_M	MARG_O_T03_F	NON_WORK_M	NON_WORK_F		
										
count	640000000	640000000	640000000	640000000	640000000	640000000	640000000	640000000	640000000	640000000	...	640000000	640000000	640000000	640000000	640000000	640000000	640000000	640000000	640000000	640000000	
	48515.542188	76041.601953	116079.808594	11638.096875	11234.508203	13173.19685	19764.365039	5068.761133	8345.648047	5444.874219	...	1243.50000	2554.161719	187.805664	402.933008	456.679297	1157.905078	563.20312	164.198438	43.924219	609.501562	
	39308.008223	60233.862206	92154.544396	9253.649941	8983.799265	12201.89225	18315.52608	6018.652465	10001.77451	43843.6970	...	999.851461	2098.515606	186.884611	392.233200	426.951049	1142.279691	548.94326	156.264559	374.014651	510.812596	
	35000	39100	69800	56000	56000	00000	00000	00000	00000	28600	...	40000	30000	00000	00000	00000	00000	00000	00000	00000	50000	
min	35000	39100	69800	56000	56000	00000	00000	00000	00000	28600	...	40000	30000	00000	00000	00000	00000	00000	00000	00000	50000	

											MARG_C L_0_3_M	MARG_C L_0_3_F	MARG_A L_0_3_M	MARG_A L_0_3_F	MARG_H H_0_3_M	MARG_H H_0_3_F	MARG_O T_0_3_M	MARG_O T_0_3_F	NON_WORK_M	NON_WORK_F	
NO_HH	TOT_M	TOT_F	M_06	F_06	M_SC	F_SC	M_ST	F_ST	M_LIT	...											
	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000									0000	0000	
25%	194840000000	302280000000	465177000000	473375000000	467225000000	346625000000	560325000000	293750000000	429500000000	212980000000		48950000	95725000	47000000	10900000	13650000	29800000	14000000	43000000		
																				16100000	22050000
																				00000000	00000000
																				00000000	00000000
																				00000000	00000000
																				00000000	00000000
																				00000000	00000000
																				00000000	00000000
																				00000000	00000000
																				00000000	00000000
50%	358370000000	583390000000	877245000000	915900000000	866300000000	959150000000	137090000000	233350000000	383450000000	426935000000		94900000	192800000	11450000	24750000	30800000	71700000	35000000	11300000	36000000	46450000
75%	688920000000	1079185000000	1642517500000	1652250000000	1592250000000	1942750000000	2918000000000	7658000000000	1248250000000	7798500000000		171400000	3599750000	2707500000	5687500000	6420000000	17107500000	7900000000	2400000000	6045000000	8535000000
m	1	2	3	3	3	4	6	1	3	1	...	35	75	60	12	14	38	17	53	1	1

N o _ H H	T O T _ M	T O T _ F	M _ 0 6	F _ 0 6	M _ S C	F _ S C	M _ S T	F _ S T	M _ L I T	...	M A R G _ C _ L 0 _ 3 _ M	M A R G _ C _ L 0 _ 3 _ F	M A R G _ A _ L 0 _ 3 _ M	M A R G _ A _ L 0 _ 3 _ F	M A R G _ H _ 0 _ 3 _ M	M A R G _ H _ 0 _ 3 _ F	M A R G _ O _ T 0 _ 3 _ M	M A R G _ O _ T 0 _ 3 _ F	N O N _ W O R K _ M	N O N _ W O R K _ F	
											a	x									
	4	2	4	4	2	3	4	8	0	6		50	63	6.	58	00	29	6.	5.	2	8
	3	4	0	2	7	3	5	7	5	3		.7	.5	37	.3	.2	.8	50	50	6	0
	0	4	8	0	4	7	4	0	5	0		50	00	50	75	50	75	00	00	9.	3.
	0	5	5	0.	7.	5.	5.	4.	6.	2		00	00	00	00	00	00	00	00	7	0
	4.	4.	2.	0	2	0	1	3	3	6.		0	0		0	0	0			5	0
	0	2	7	0	5	0	2	7	7	7										0	0
	0	5	5	0	0	0	5	5	5	5										0	0
	0	0	0	0	0	0	0	0	0	0										0	0
	0	0	0	0	0	0	0	0	0	0										0	0
	0	0	0	0	0	0	0	0	0	0										0	0
	0	0	0	0	0	0	0	0	0	0											
	0	0	0							0											

8 rows × 57 columns

Table-7-before scaling

After scaling

N O _ H H	T O T _ M	T O T _ F	M _ 0 6	F _ 0 6	M _ S C	F _ S C	M _ S T	F _ S T	M _ L I T	...	M A R G _ C _ L 0 _ 3 _ M	M A R G _ C _ L 0 _ 3 _ F	M A R G _ A _ L 0 _ 3 _ M	M A R G _ A _ L 0 _ 3 _ F	M A R G _ H _ 0 _ 3 _ M	M A R G _ H _ 0 _ 3 _ F	M A R G _ O _ T 0 _ 3 _ M	M A R G _ O _ T 0 _ 3 _ F	N O N _ W O R K _ M	N O N _ W O R K _ F									
c o u n t	6. 4 0 0 0 0 0 e + 0 2	6. 4 0 0 0 0 0 e + 0 2	6. 4 0 0 0 0 0 e + 0 2	6. 4 0 0 0 0 e + 0 2	6. 4 0 0 0 0 e + 0 2	6. 4 0 0 0 0 e + 0 2	6. 4 0 0 0 0 e + 0 2	6. 4 0 0 0 0 e + 0 2	6. 4 0 0 0 0 e + 0 2	6. 4 0 0 0 0 e + 0 2		6. 40 00 00 e+ 02	6. 40 00 00 e+ 02	6. 40 00 00 e+ 02	6. 40 00 00 e+ 02	6. 40 00 00 e+ 02	64 0. 00 00 00	6. 40 00 00 e+ 02	6. 40 00 00 e + 02	6. 4 0 0 0 0 e + 0 2	6. 4 0 0 0 0 e + 0 2								
	- 6. 6 6	- 1. 3 3	- 2. 2 2	5. 5 1	- 3. 3 3	2. 2 0	- 2. 2 2	- 4. 4 4	0 . 0 0	- 7. 7 7	...	- 1. 38 77	- 3. 33 06	- 2. 22 04	- 6. 66 13	2. 22 46	0. 00 00	- 2. 22 04	- 7. 21 64	- 6. 6 6	4. 4 0								

N o _H H	T O T _M	T O T _F	M _0_6	F _0_6	M _S_C	F _S_C	M _S_T	F _S_T	M _L_I_T	...	M A R G_C _L_0_3_M	M A R G_C _L_0_3_F	M A R G_A _L_0_3_M	M A R G_A _L_0_3_F	M A R G_H_0_3_M	M A R G_H_0_3_F	M A R G_O_0_3_M	M A R G_O_0_3_F	N O N _W_O_R_K_M	N O N _W_O_R_K_F		
	1 3 3 8 e- 1 7	2 2 6 8 e- 1 6	0 4 4 6 e- 1 7	1 1 5 e- 1 7	0 6 6 9 e- 1 7	4 4 6 e- 1 7	0 4 4 6 e- 1 7	0 8 9 2 e- 1 7	0 0 0 0	1 5 6 1 e- 1 7		79 e- 17	69 e- 17	46 e- 17	38 e- 17	e- 17		46 e- 17	50 e- 17	1 3 3 8 e- 1 7	8 9 2 e- 1 7	
s t d	1. 0 0 0 7 8 2 e + 0 0	1. 0 0 0 7 8 2 e + 0 0	1. 0 0 0 7 8 2 e + 0 0	1. 0 0 0 7 8 2 e + 0 0	1. 0 0 0 7 8 2 e + 0 0	1. 0 0 0 7 8 2 e + 0 0	1. 0 0 0 7 8 2 e + 0 0	1. 0 0 0 7 8 2 e + 0 0	1 . 0 0 0 7 8 2	1. 0 0 0 7 8 2 e + 0 0	...	1. 00 07 82 e+ 00	1. 00 07 82 e+ 00	1. 00 07 82 e+ 00	1. 00 07 82 e+ 00	1. 00 07 82 e+ 00	1. 00 07 82 e+ 00	1. 00 07 82 e+ 00	1. 00 07 82 e+ 00	1. 0 0 0 7 8 2 e + 0 0	1. 0 0 0 7 8 2 e + 0 0	
	- 1. 2 2 6 2 9 5 e + 0 0	- 1. 2 5 6 9 3 0 e + 0 0	- 1. 2 5 3 0 6 2 e + 0 0	- 1. 2 5 2 6 0 4 e + 0 0	- 1. 2 4 5 2 7 0 e + 0 0	- 1. 0 8 0 4 7 e + 0 0	- 1. 0 7 9 4 6 3 e + 0 0	- 8. 4 2 8 3 4 1 e- 0 1	- 0 8 3 3 7 4 1	1. 0 3 8 5 2 7 e + 0 0	...	- 1. 24 06 54 e+ 00	- 1. 20 37 73 e+ 00	- 1. 00 57 14 e+ 00	- 1. 02 80 83 e+ 00	- 1. 07 04 66 e+ 00	- 1. 01 67 79 e+ 00	- 1. 05 15 94 e+ 00	- 1. 08 7 8 4 5 e+ 0 0	- 1. 1 8 3 7 e+ 0 0	- 1. 1 8 4 3 7 e+ 0 0	
	- 7. 3 9 1 4 3 3 e- 0 1	- 7. 6 1 9 0 4 4 e- 0 1	- 7. 5 5 4 3 7 e- 0 1	- 7. 4 6 7 0 1 e- 0 1	- 7. 3 1 0 2 0 e- 0 1	- 7. 9 6 1 5 0 e- 0 1	- 7. 7 3 9 0 8 e- 0 1	- 7. 9 3 8 4 e- 0 1	- 0 7 9 0 8 3 4	7. 5 8 9 0 1 6 e- 0 1	...	- 7. 54 70 19 e- 01	- 7. 61 56 73 e- 01	- 7. 54 02 57 e- 01	- 7. 49 96 94 e- 01	- 7. 50 50 70 e- 01	- 0. 75 33 86	- 7. 71 54 45 e- 01	- 7. 76 20 44 e- 01	- 7. 5 7 0 4 3 9 e- 0 1	- 7. 6 2 1 3 0 4 e- 0 1	

N o _ H H	T O T _ M	T O T _ F	M _ 0 6	F _ 0 6	M _ S C	F _ S C	M _ S T	F _ S T	M _ L I T	...	M A R G C _ L 0 _ 3 _ M	M A R G C _ L 0 _ 3 _ F	M A R G A _ L 0 _ 3 _ M	M A R G A _ L 0 _ 3 _ F	M A R G H _ 0 _ 3 _ M	M A R G H _ 0 _ 3 _ F	M A R G _ O T 0 _ 3 _ M	M A R G _ O T 0 _ 3 _ F	N O N _ W O R K _ M	N O N _ W O R K _ F	
5 0 %	-	-	-	-	-	-	-	-	-	-										-	-
	3.	2.	3.	2.	2.	2.	3.	4.	-	2.									3.	2.	
	2	9	0	6	8	9	3	5	0	7									1	8	
	2	4	7	8	6	3	0	4	.	0									5	4	
	7	1	9	1	4	7	8	8	4	5									5	0	
	9	2	3	1	6	6	7	1	5	2	...								3	8	
	5	7	3	4	2	5	6	9	0	2									9	6	
	8	7	7	3	3	8	9	5	6	5									7	5	
e-	e-	e-	e-	e-	e-	e-	e-	7	e-									e-	e-		
0	0	0	0	0	0	0	0	0	0									0	0		
1	1	1	1	1	1	1	1	1	1									1	1		
7 5 %	5.	5.	5.	5.	5.	5.	5.	4.		5.										4.	4.
	1	2	2	2	1	1	1	3	0	3									4.	2	7
	8	9	3	8	9	3	4	0	.	5									9	8	
	7	6	1	0	9	1	4	5	4	1									6	0	
	8	3	3	0	7	5	8	3	1	5	...								6	4	
	4	2	8	4	9	3	8	8	3	3									6	0	
	8	8	8	8	6	7	5	9	0	0									0	8	
	e-	e-	e-	e-	e-	e-	e-	e-	5	e-									e-	e-	
0	0	0	0	0	0	0	0	2	0									0	0		
1	1	1	1	1	1	1	1	1	1									1	1		
m a x	2.	2.	2.	2.	2.	2.	2.	2.		2.										2.	2.
	4	4	4	4	3	4	4	2	2	4									2	3	
	0	6	4	4	9	7	4	6	.	7									0	3	
	5	5	0	0	6	7	6	7	2	6									9	8	
	6	8	9	0	4	1	9	3	1	2									7	2	
	7	6	9	7	8	1	0	3	8	3	...								3	9	
	7	8	5	0	8	0	7	1	8	5									1	8	
	e	e	e	e	e	e	e	e	8	e									e	e	
+	+	+	+	+	+	+	+	8	+									+	+		
0	0	0	0	0	0	0	0	1	0									0	0		
0	0	0	0	0	0	0	0	0	0									0	0		

8 rows × 57 columns

table-8-after scaling

Part 2-3); PCA: PCA

- Create the covariance matrix - Get eigen values and eigen vectors - Identify the optimum number of PCs - Show Scree plot - Compare PCs with Actual Columns and identify which is explaining most variance - Write inferences about all the PCs in terms of actual variables - Write linear equation for first PC Note: For the scope of this project, take at least 90% explained variance.

ANSWER:

#Check for presence of correlations

Bartlett's Test of Sphericity

Bartlett's test of sphericity tests the hypothesis that the variables are uncorrelated in the population.

H0: All variables in the data are uncorrelated

Ha: At least one pair of variables in the data are correlated If the null hypothesis cannot be rejected, then PCA is not advisable.

If the p-value is small(Reject H0 if p-value < 0.05), then we can reject the null hypothesis and agree that there is atleast one pair of variables in the data which are correlated hence PCA is recommended

p-value:0.0

the p-value is small(Reject H0 if p-value < 0.05), then we can reject the null hypothesis and agree that there is atleast one pair of variables in the data which are correlated hence PCA is recommended

KMO Test

The Kaiser-Meyer-Olkin (KMO) - measure of sampling adequacy (MSA) is an index used to examine how appropriate PCA is.

Generally, if MSA is less than 0.5, PCA is not recommended, since no reduction is expected. On the other hand, MSA > 0.7 is expected to provide a considerable reduction in the dimension and extraction of meaningful components.

MSA= 0.936189616665265

MSA > 0.7 is expected to provide a considerable reduction in the dimension and extraction of meaningful components.

2-3-a)-Create the covariance matrix

#Apply PCA taking all features

#co variance matrix

```
array([[ -5.52816148e+00,  4.30377559e-01, -1.47382695e+00, ...,
         6.51060294e-03,  2.38391484e-03,  4.36606130e-04],
       [ -5.49201646e+00, -1.06110331e-01, -2.01564100e+00, ...,
        -2.82839348e-04,  8.13609312e-03, -6.60428796e-03],
       [ -7.47464297e+00, -2.17193764e-01, -2.47428211e-01, ...,
        -1.35201718e-03, -1.10109765e-03,  8.64566393e-05],
       ...,
       [ -7.88626804e+00, -1.00353656e+00, -9.09284569e-01, ...,
        -2.15313673e-03,  1.45549207e-03, -4.60053251e-04],
       [ -7.86425952e+00, -9.99337996e-01, -8.51569237e-01, ...,
        -2.06885382e-03, -1.22502335e-03,  1.81303381e-03],
       [ -7.41622568e+00, -1.41214300e+00, -8.65921210e-01, ...,
        -1.06417476e-03, -1.66377584e-03,  1.78275792e-03]])
```

2-3-b) Get eigen values and eigen vectors**#Check the eigen values**

```
array([3.56488638e+01, 7.64357559e+00, 3.76919551e+00, 2.77722349e+00,
       1.90694892e+00, 1.15490310e+00, 9.87726707e-01, 4.64629906e-01,
       3.96708513e-01, 3.22346888e-01, 2.73207369e-01, 2.35647574e-01,
       1.81401107e-01, 1.69243770e-01, 1.38592325e-01, 1.31505852e-01,
       1.03809666e-01, 9.55333831e-02, 8.58580407e-02, 8.09138742e-02,
       6.60179067e-02, 6.30797999e-02, 4.82756124e-02, 4.59506197e-02,
       4.37747566e-02, 3.19339710e-02, 2.86194563e-02, 2.75481445e-02,
       2.34340044e-02, 2.20296816e-02, 1.87487040e-02, 1.59004895e-02,
       1.39957919e-02, 1.18916465e-02, 1.11133495e-02, 9.07842645e-03,
       7.25127869e-03, 6.27213692e-03, 4.95541908e-03, 4.60667097e-03,
       3.45902033e-03, 2.18408510e-03, 2.13514664e-03, 1.92111328e-03,
       1.43840980e-03, 1.09968912e-03, 9.65752052e-04, 8.62630267e-04,
       6.51634478e-04, 5.76658846e-04, 4.35790607e-04, 3.70037468e-04,
       3.06660171e-04, 2.07854170e-04, 1.38286484e-04, 8.97034441e-05,
       4.61745385e-05])
```

#Extract eigen vectors

```
array([[ 0.14922158,  0.15916917,  0.15820921, ...,  0.14136961,
         0.14762899,  0.14210263],
       [-0.11548673, -0.08023879, -0.09371751, ...,  0.03510934,
        -0.04912234, -0.03984815],
       [ 0.1015276 , -0.03866173,  0.0289595 , ..., -0.10217491,
        -0.12667281, -0.02854464],
       ...,
       [ 0.00112879, -0.00673066,  0.02298648, ..., -0.01159627,
        0.05608352, -0.00610478],
       [ 0.00070908,  0.04637872,  0.00402434, ...,  0.01406358,
        -0.07729171, -0.00056173],
       [-0.00461221, -0.00370327,  0.00963954, ...,  0.00227908,
        0.00539901,  0.00130606]])
```

2-3-c) Identify the optimum number of PCs

```
array([6.24441446e-01, 1.33888289e-01, 6.60229147e-02, 4.86470891e-02,
       3.34029704e-02, 2.02297994e-02, 1.73014629e-02, 8.13866529e-03,
       6.94892379e-03, 5.64637229e-03, 4.78562250e-03, 4.12770833e-03,
       3.17750294e-03, 2.96454958e-03, 2.42764517e-03, 2.30351534e-03,
       1.81837655e-03, 1.67340548e-03, 1.50392785e-03, 1.41732362e-03,
       1.15639919e-03, 1.10493400e-03, 8.45617224e-04, 8.04891611e-04,
       7.66778221e-04, 5.59369722e-04, 5.01311201e-04, 4.82545623e-04,
       4.10480504e-04, 3.85881758e-04, 3.28410688e-04, 2.78520087e-04,
       2.45156553e-04, 2.08299401e-04, 1.94666401e-04, 1.59021779e-04,
       1.27016642e-04, 1.09865556e-04, 8.68013375e-05, 8.06925096e-05,
       6.05897475e-05, 3.82574118e-05, 3.74001838e-05, 3.36510796e-05,
       2.51958296e-05, 1.92626466e-05, 1.69165450e-05, 1.51102177e-05,
       1.14143210e-05, 1.01010143e-05, 7.63350323e-06, 6.48174183e-06,
       5.37159674e-06, 3.64086663e-06, 2.42228792e-06, 1.57128566e-06,
       8.08813873e-07])
```

#The percentage of variance explained by each principal component

#Obtaining the Cumulative Sum of the Explained Variance

```
Cumulative Variance Explained in Percentage: [ 62.44  75.83  82.44  87.3   90.
64  92.66  94.39  95.21  95.9   96.47
    96.95  97.36  97.68  97.97  98.22  98.45  98.63  98.79  98.95  99.09
    99.2   99.31  99.4   99.48  99.55  99.61  99.66  99.71  99.75  99.79
    99.82  99.85  99.87  99.89  99.91  99.93  99.94  99.95  99.96  99.97
    99.98  99.98  99.98  99.99  99.99  99.99  99.99 100.   100.   100.
100.   100.   100.   100.   100.   100.   100. ]
```

We can see above that more than 90% of the variance is explained by 5 Principal Components.

optimum number of PCs is 5

2-3-d)-Show Scree plot

Plot to identify the number of components to be built

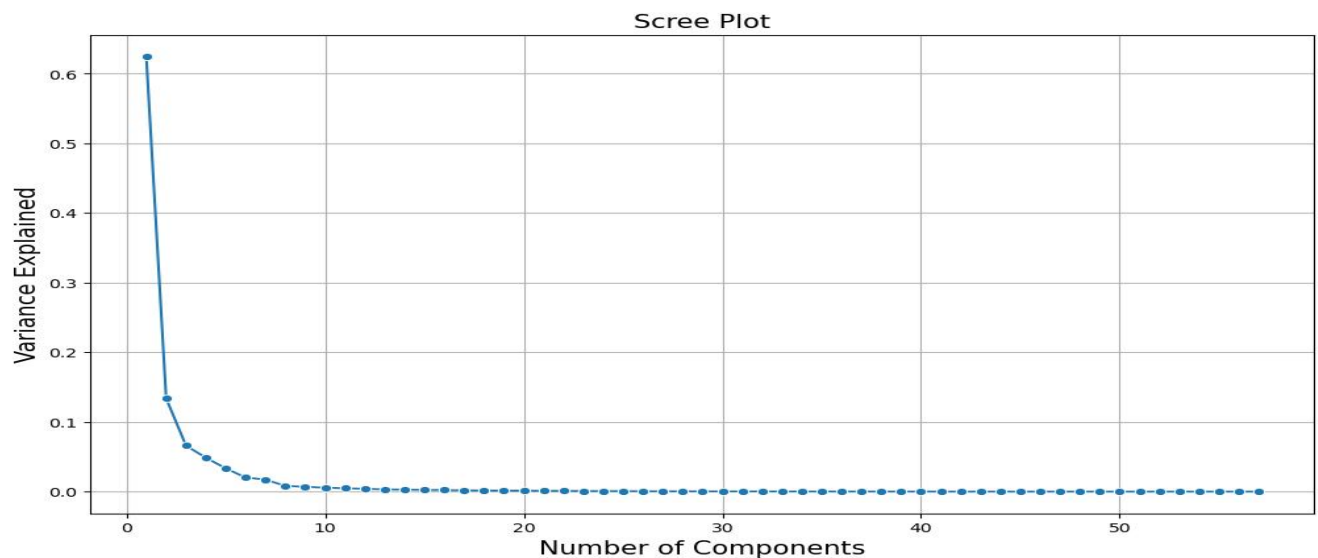


Fig-31- Scree plot

The number of components can be decided based upon the explained variance. Here, it is decided to keep the number of components as 5 as the cumulative explained variance is above 90%

find the least number of components that can explain more than 90% variance

Number of PCs that explain at least 90% variance: 5

Observations:

We can see that out of the 59 original features, we reduced the number of features through principal components to 5, these components explain more than 90% of the original variance. Let us now look at these principal components as a linear combination of original features.

2-3-e)-Compare PCs with Actual Columns and identify which is explaining most variance,Write inferences about all the PCs in terms of actual variables

The number of components can be decided based upon the explained variance. Here, it is decided to keep the number of components as 5 as the cumulative explained variance is above 90%

	PC1	PC2	PC3	PC4	PC5
No_HH	0.15	-0.12	0.10	0.08	-0.01
TOT_M	0.16	-0.08	-0.04	0.05	-0.04
TOT_F	0.16	-0.09	0.03	0.07	-0.02
M_06	0.16	-0.02	-0.07	0.03	-0.08
F_06	0.16	-0.01	-0.07	0.02	-0.08
M_SC	0.14	-0.08	-0.04	0.01	-0.17
F_SC	0.14	-0.09	0.02	0.02	-0.16
M_ST	0.02	0.07	0.32	0.09	0.42
F_ST	0.02	0.07	0.34	0.08	0.42
M_LIT	0.16	-0.11	-0.03	0.09	-0.01

	PC1	PC2	PC3	PC4	PC5
F_LIT	0.15	-0.13	-0.01	0.13	0.03
M_ILL	0.15	-0.01	-0.05	-0.03	-0.10
F_ILL	0.16	-0.02	0.08	-0.01	-0.11
TOT_WORK_M	0.15	-0.12	-0.00	0.07	-0.02
TOT_WORK_F	0.14	-0.08	0.19	0.11	-0.02
MAINWORK_M	0.14	-0.17	0.02	0.10	-0.04
MAINWORK_F	0.13	-0.14	0.21	0.13	-0.05
MAIN_CL_M	0.11	0.04	0.03	0.08	-0.30
MAIN_CL_F	0.08	0.10	0.19	0.27	-0.26
MAIN_AL_M	0.12	-0.05	0.23	-0.12	-0.25
MAIN_AL_F	0.09	-0.07	0.36	-0.02	-0.20
MAIN_HH_M	0.14	-0.10	-0.10	-0.02	-0.06

	PC1	PC2	PC3	PC4	PC5
MAIN_HH_F	0.13	-0.11	0.02	-0.05	-0.02
MAIN_OT_M	0.12	-0.20	-0.03	0.15	0.07
MAIN_OT_F	0.12	-0.21	0.07	0.16	0.11
MARGWORK_M	0.16	0.08	-0.07	-0.08	0.07
MARGWORK_F	0.15	0.11	0.10	0.02	0.08
MARG_CL_M	0.09	0.27	-0.10	0.16	-0.02
MARG_CL_F	0.07	0.28	-0.04	0.29	-0.06
MARG_AL_M	0.13	0.16	0.07	-0.25	-0.05
MARG_AL_F	0.12	0.14	0.26	-0.15	-0.01
MARG_HH_M	0.15	0.04	-0.14	-0.17	0.01
MARG_HH_F	0.14	0.01	-0.09	-0.15	0.04
MARG_OT_M	0.15	-0.07	-0.13	0.02	0.15

	PC1	PC2	PC3	PC4	PC5
MARG_OT_F	0.15	-0.09	-0.05	0.06	0.19
MARGWORK_3_6_M	0.16	-0.04	-0.07	0.04	-0.06
MARGWORK_3_6_F	0.16	-0.09	-0.06	0.05	-0.02
MARG_CL_3_6_M	0.16	0.07	-0.06	-0.09	0.06
MARG_CL_3_6_F	0.15	0.09	0.13	0.02	0.06
MARG_AL_3_6_M	0.09	0.26	-0.10	0.13	-0.01
MARG_AL_3_6_F	0.07	0.27	-0.02	0.29	-0.06
MARG_HH_3_6_M	0.13	0.15	0.08	-0.25	-0.06
MARG_HH_3_6_F	0.11	0.12	0.28	-0.14	-0.03
MARG_OT_3_6_M	0.15	0.04	-0.14	-0.17	0.00
MARG_OT_3_6_F	0.14	-0.00	-0.09	-0.14	0.04
MARGWORK_0_3_M	0.15	-0.08	-0.13	0.02	0.13

	PC1	PC2	PC3	PC4	PC5
MARGWORK_0_3_F	0.15	-0.10	-0.06	0.06	0.17
MARG_CL_0_3_M	0.14	0.14	-0.10	-0.02	0.09
MARG_CL_0_3_F	0.13	0.17	0.03	0.01	0.11
MARG_AL_0_3_M	0.06	0.28	-0.12	0.21	-0.02
MARG_AL_0_3_F	0.06	0.29	-0.09	0.24	-0.04
MARG_HH_0_3_M	0.12	0.18	0.03	-0.24	0.02
MARG_HH_0_3_F	0.11	0.18	0.16	-0.19	0.05
MARG_OT_0_3_M	0.14	0.05	-0.14	-0.17	0.01
MARG_OT_0_3_F	0.14	0.04	-0.10	-0.17	0.05
NON_WORK_M	0.15	-0.05	-0.13	0.02	0.19
NON_WORK_F	0.14	-0.04	-0.03	0.06	0.25

Table-9-)-Compare PCs with Actual Columns and identify which is explaining most variance

#Check as to how the original features matter to each PC

#Note: Here we are only considering the absolute values

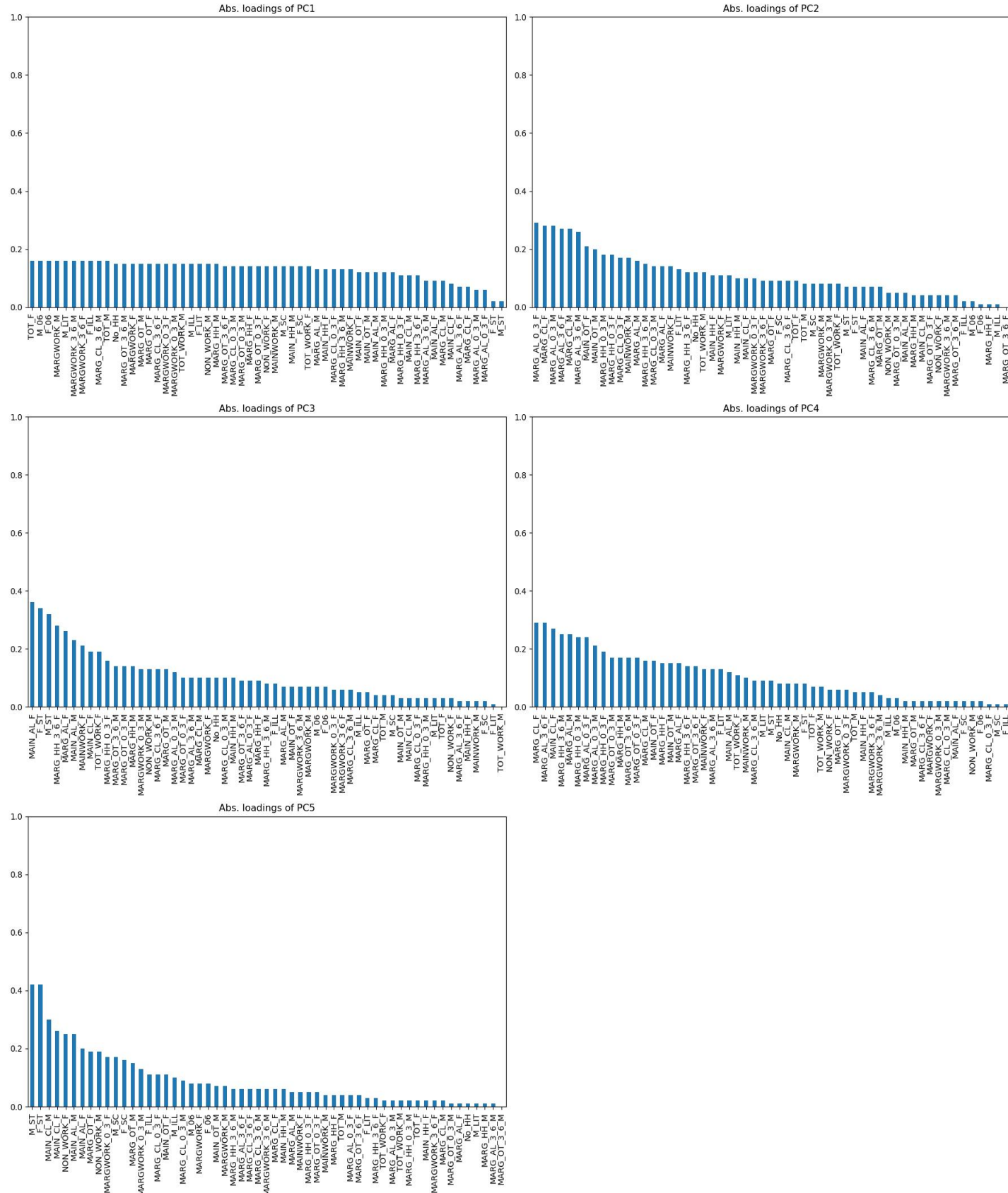




Fig-32-Compare PCs with Actual Columns

	PC1	PC2	PC3	PC4	PC5
No_HH	0.150000	-0.120000	0.100000	0.080000	-0.010000
TOT_M	0.160000	-0.080000	-0.040000	0.050000	-0.040000
TOT_F	0.160000	-0.090000	0.030000	0.070000	-0.020000
M_06	0.160000	-0.020000	-0.070000	0.030000	-0.080000
F_06	0.160000	-0.010000	-0.070000	0.020000	-0.080000
M_SC	0.140000	-0.080000	-0.040000	0.010000	-0.170000
F_SC	0.140000	-0.090000	0.020000	0.020000	-0.160000
M_ST	0.020000	0.070000	0.320000	0.090000	0.420000
F_ST	0.020000	0.070000	0.340000	0.080000	0.420000
M_LIT	0.160000	-0.110000	-0.030000	0.090000	-0.010000
F_LIT	0.150000	-0.130000	-0.010000	0.130000	0.030000
M_ILL	0.150000	-0.010000	-0.050000	-0.030000	-0.100000

	PC1	PC2	PC3	PC4	PC5
F_ILL	0.160000	-0.020000	0.080000	-0.010000	-0.110000
TOT_WORK_M	0.150000	-0.120000	-0.000000	0.070000	-0.020000
TOT_WORK_F	0.140000	-0.080000	0.190000	0.110000	-0.020000
MAINWORK_M	0.140000	-0.170000	0.020000	0.100000	-0.040000
MAINWORK_F	0.130000	-0.140000	0.210000	0.130000	-0.050000
MAIN_CL_M	0.110000	0.040000	0.030000	0.080000	-0.300000
MAIN_CL_F	0.080000	0.100000	0.190000	0.270000	-0.260000
MAIN_AL_M	0.120000	-0.050000	0.230000	-0.120000	-0.250000
MAIN_AL_F	0.090000	-0.070000	0.360000	-0.020000	-0.200000
MAIN_HH_M	0.140000	-0.100000	-0.100000	-0.020000	-0.060000
MAIN_HH_F	0.130000	-0.110000	0.020000	-0.050000	-0.020000
MAIN_OT_M	0.120000	-0.200000	-0.030000	0.150000	0.070000

	PC1	PC2	PC3	PC4	PC5
MAIN_OT_F	0.120000	-0.210000	0.070000	0.160000	0.110000
MARGWORK_M	0.160000	0.080000	-0.070000	-0.080000	0.070000
MARGWORK_F	0.150000	0.110000	0.100000	0.020000	0.080000
MARG_CL_M	0.090000	0.270000	-0.100000	0.160000	-0.020000
MARG_CL_F	0.070000	0.280000	-0.040000	0.290000	-0.060000
MARG_AL_M	0.130000	0.160000	0.070000	-0.250000	-0.050000
MARG_AL_F	0.120000	0.140000	0.260000	-0.150000	-0.010000
MARG_HH_M	0.150000	0.040000	-0.140000	-0.170000	0.010000
MARG_HH_F	0.140000	0.010000	-0.090000	-0.150000	0.040000
MARG_OT_M	0.150000	-0.070000	-0.130000	0.020000	0.150000
MARG_OT_F	0.150000	-0.090000	-0.050000	0.060000	0.190000
MARGWORK_3_6_M	0.160000	-0.040000	-0.070000	0.040000	-0.060000

	PC1	PC2	PC3	PC4	PC5
MARGWORK_3_6_F	0.160000	-0.090000	-0.060000	0.050000	-0.020000
MARG_CL_3_6_M	0.160000	0.070000	-0.060000	-0.090000	0.060000
MARG_CL_3_6_F	0.150000	0.090000	0.130000	0.020000	0.060000
MARG_AL_3_6_M	0.090000	0.260000	-0.100000	0.130000	-0.010000
MARG_AL_3_6_F	0.070000	0.270000	-0.020000	0.290000	-0.060000
MARG_HH_3_6_M	0.130000	0.150000	0.080000	-0.250000	-0.060000
MARG_HH_3_6_F	0.110000	0.120000	0.280000	-0.140000	-0.030000
MARG_OT_3_6_M	0.150000	0.040000	-0.140000	-0.170000	0.000000
MARG_OT_3_6_F	0.140000	-0.000000	-0.090000	-0.140000	0.040000
MARGWORK_0_3_M	0.150000	-0.080000	-0.130000	0.020000	0.130000
MARGWORK_0_3_F	0.150000	-0.100000	-0.060000	0.060000	0.170000
MARG_CL_0_3_M	0.140000	0.140000	-0.100000	-0.020000	0.090000

	PC1	PC2	PC3	PC4	PC5
MARG_CL_0_3_F	0.130000	0.170000	0.030000	0.010000	0.110000
MARG_AL_0_3_M	0.060000	0.280000	-0.120000	0.210000	-0.020000
MARG_AL_0_3_F	0.060000	0.290000	-0.090000	0.240000	-0.040000
MARG_HH_0_3_M	0.120000	0.180000	0.030000	-0.240000	0.020000
MARG_HH_0_3_F	0.110000	0.180000	0.160000	-0.190000	0.050000
MARG_OT_0_3_M	0.140000	0.050000	-0.140000	-0.170000	0.010000
MARG_OT_0_3_F	0.140000	0.040000	-0.100000	-0.170000	0.050000
NON_WORK_M	0.150000	-0.050000	-0.130000	0.020000	0.190000
NON_WORK_F	0.140000	-0.040000	-0.030000	0.060000	0.250000

Table-10--Compare PCs with Actual Columns

OBSERVATION

PC1 is explaining most variance.

- The first principal component, PC1, is a measure of mpg, cylinders, displacement, horsepower, and weight. PC1 is associated with high scores of almost all variables('No_HH', 'TOT_M', 'TOT_F', 'M_06', 'F_06', 'M_SC', 'F_SC', 'M_LIT', 'F_LIT', 'M_ILL', 'F_ILL', 'TOT_WORK_M', 'TOT_WORK_F', 'MAINWORK_M', 'MAINWORK_F', 'MAIN_CL_M', 'MAIN_AL_M', 'MAIN_HH_M', 'MAIN_OT_M', 'MAIN_OT_F', 'MARGWORK_M', 'MARGWORK_F', 'MARG_AL_M', 'MARG_AL_F', 'MARG_HH_M', 'MARG_HH_F', 'MARG_OT_M', 'MARG_OT_F', 'MARGWORK_3_6_M', 'MARGWORK_3_6_F',

'MARG_CL_3_6_M', 'MARG_CL_3_6_F', 'MARG_HH_3_6_M', 'MARG_HH_3_6_F',
 'MARG_OT_3_6_M', 'MARG_OT_3_6_F', 'MARGWORK_0_3_M', 'MARGWORK_0_3_F',
 'MARG_CL_0_3_M', 'MARG_CL_0_3_F', 'MARG_HH_0_3_M', 'MARG_HH_0_3_F',
 'MARG_OT_0_3_M', 'MARG_OT_0_3_F', 'NON_WORK_M', 'NON_WORK_F' M_ST, F_ST
 no negative scores in PC1.

- PC2 have less explain variance than PC2, The second principal component, PC2, is a measure of model year. PC2 is associated with low values of "No_HH", 'TOT_F', 'M_LIT', 'F_LIT', 'TOT_WORK_M', 'MAINWORK_M', 'MAINWORK_F', 'MAIN_OT_M', 'MAIN_OT_F', 'MARG_OT_F', 'MARGWORK_3_6_F', 'MARGWORK_0_3_F'

optimum number of PCs is 5

2-3-f) Write linear equation for first PC Note

PC1 have explain all variance and all are high

To write the linear equation for PC1 (Principal Component 1), you can express PC1 as a linear combination of the original variables. Let's denote the original variables as ($X_1, X_2, X_3, \dots, X_n$) and their corresponding coefficients in PC1 as ($a_1, a_2, a_3, \dots, a_n$). The linear equation for PC1 can be written as:

$$PC1 = a_1x_1 + a_2x_2 + \dots + a_nx_n.$$

Using the coefficients provided in your data, the linear equation for PC1 would be

This equation represents the linear combination of the original variables that make up PC1. The coefficients indicate the contribution of each variable to the overall value of PC1.

$$PC1 = a_1x_1 + a_2x_2 + \dots + a_nx_n.$$

$$\begin{aligned} &0.16No_HH + 0.17TOT_M + 0.17TOT_F + 0.16M_06 + 0.16F_06 + 0.15M_SC + 0.15F_SC + \\ &0.03M_ST + 0.03F_ST + 0.16M_LIT + 0.15F_LIT + 0.16M_ILL + 0.17F_ILL + 0.16TOT_WORK_M + \\ &0.15TOT_WORK_F + \\ &0.15MAINWORK_M + \\ &0.12MAINWORK_F + \\ &0.10MAIN_CL_M + \\ &0.07MAIN_CL_F + \\ &0.11MAIN_AL_M + \\ &0.07MAIN_AL_F + \\ &0.13MAIN_HH_M + \\ &0.08MAIN_HH_F + \\ &0.12MAIN_OT_M + \\ &0.11MAIN_OT_F + \\ &0.16MARGWORK_M + \\ &0.16MARGWORK_F + \\ &0.16MARG_CL_M + \\ &0.05MARG_CL_F + \\ &0.13MARG_AL_M + \\ &0.11MARG_AL_F + \\ &0.14MARG_HH_M + \\ &0.13MARG_HH_F + \\ &0.16MARG_OT_M + \\ &0.15MARG_OT_F + \end{aligned}$$

$0.16MARGWORK_3_6_M + 0.16MARGWORK_3_6_F +$
 $0.17MARG_CL_3_6_M + 0.16MARG_CL_3_6_F +$
 $0.09MARG_AL_3_6_M +$
 $0.05MARG_AL_3_6_F +$
 $0.13MARG_HH_3_6_M +$
 $0.13MARG_HH_3_6_F +$
 $0.14MARG_OT_3_6_M + 0.12MARG_OT_3_6_F + 0.15MARGWORK_0_3_M +$
 $0.15MARGWORK_0_3_F + 0.15MARG_CL_0_3_M + 0.14MARG_CL_0_3_F +$
 $0.05MARG_AL_0_3_M + 0.04MARG_AL_0_3_F + 0.12MARG_HH_0_3_M +$
 $0.12MARG_HH_0_3_F + 0.14MARG_OT_0_3_M + 0.13MARG_OT_0_3_F + 0.15NON_WORK_M$
 $+ 0.13*NON_WORK_F$