

PART 1

List of Figures

Boxplot

Dendrogram

WSS plot

List of Tables

Dataset Sample

Summary of the data

Data Describe

Data Description

We have a data on The ads24*7, Digital Marketing company which has now got seed funding \$10M. They are expanding wings in Marketing Analytics. There is the data description

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

Where we can see Timestamp, Inventory, Ad Type, Platform, Device Type, Format is a object data type and some of int64 and float64 datatypes are there.

TEST FOR MISSING VALUES AND TREAT THE MISSING VALUES

We can clearly see that in CTR, CPM & CPC has some missing values. Next I use to calculate the missing values using

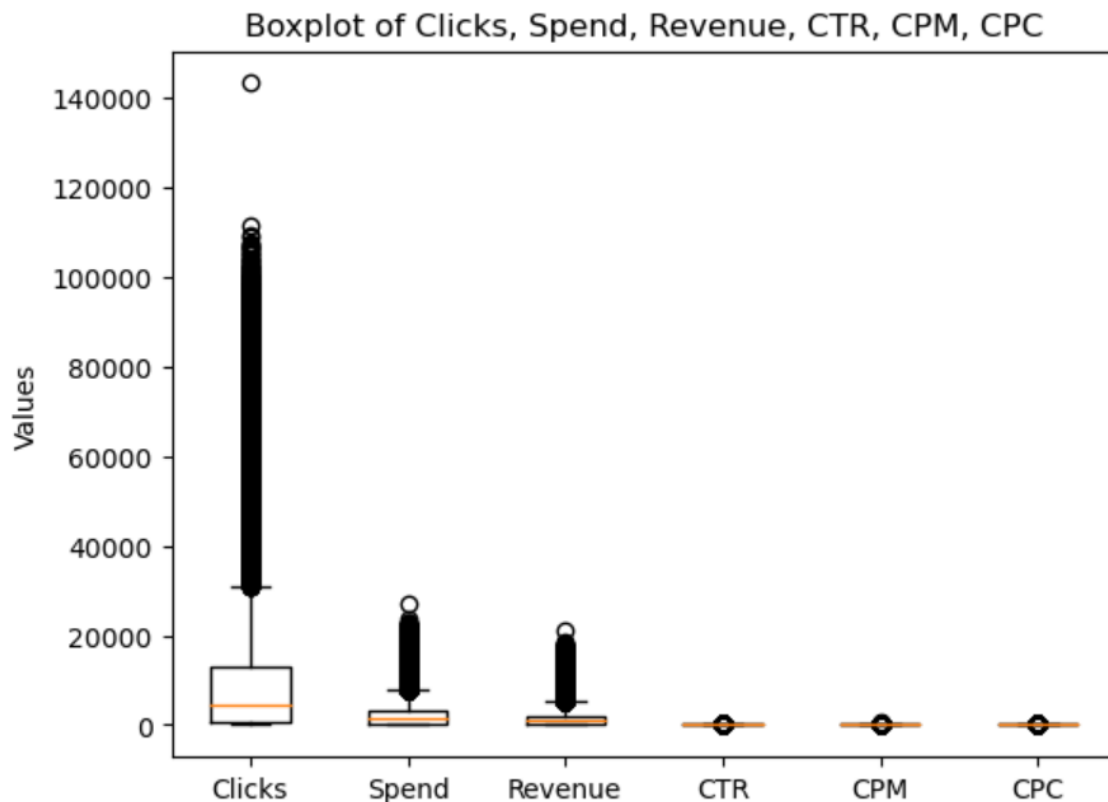
$$\text{CPM} = (\text{Total Campaign Spend} / \text{Number of Impressions}) * 1,000$$

CPC = Total Cost (spend) / Number of Clicks

CTR = Total Measured Clicks / Total Measured Ad Impressions x 100

These formulas.

CHEKING IF THERE ANY OUTLIERS



In this boxplot we can clearly see there is outliers available.

I do not think treating outliers is necessary. Because if we treat the outliers it will impact on the data set. It will change the entire data.

PERFORM Z-SCORE

In the data set there are 23066 rows and 21 columns which are unscaled data. Here is the sample data

	Timestamp	InventoryType	Ad - Length	Ad- Width	Ad Size	Ad Type	Platform	Device Type	Format	Available_Impressions	Matched_Queries	Impressions	Clicks	Sp
0	2020-9-2-17	Format1	300	250	75000	Inter222	Video	Desktop	Display	1806	325	323	1	
1	2020-9-2-10	Format1	300	250	75000	Inter227	App	Mobile	Video	1780	285	285	1	
2	2020-9-1-22	Format1	300	250	75000	Inter222	Video	Desktop	Display	2727	356	355	1	
3	2020-9-3-20	Format1	300	250	75000	Inter228	Video	Mobile	Video	2430	497	495	1	
4	2020-9-4-15	Format1	300	250	75000	Inter217	Web	Desktop	Video	1218	242	242	1	

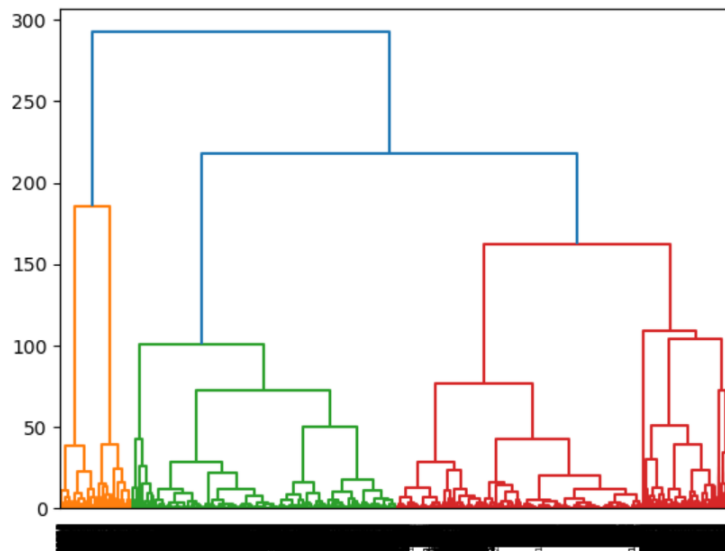
We cannot go through the clustering process with uncalled data so here we apply z-score.

Here is the sample

	Clicks	Spend	Revenue	CTR	CPM	CPC
0	-0.615311	-0.665372	-0.619693	-0.332572	-0.92711	-0.986603
1	-0.615311	-0.665372	-0.619693	-0.332521	-0.92711	-0.986603
2	-0.615311	-0.665372	-0.619693	-0.332610	-0.92711	-0.986603
3	-0.615311	-0.665372	-0.619693	-0.332712	-0.92711	-0.986603
4	-0.615311	-0.665372	-0.619693	-0.332444	-0.92711	-0.986603

HIERARCHICAL CLUSTERING AND DENDROGRAM

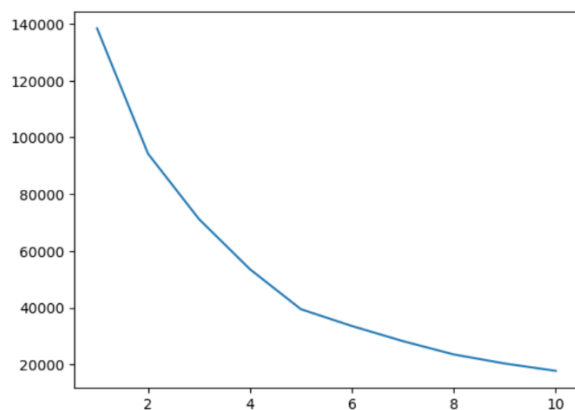
In hierarchical clustering we take all the rows as a cluster and we check the dendrogram.



Dendrogram suggests for 3 cluster.

ELBOW PLOT

For k-means algorithm we take hole data as 1cluster and end with each row as a cluster. Here we observe the WSS plot to take the number of cluster.



In this plot we can see that 5 is the significant optimum number taking as cluster number.

SILHOUETTE SCORES

the silhouette score for a clustering solution is greater than 0.5, it indicates that the clusters are relatively well-separated and the data points within each cluster are more similar to each other than to data points in other clusters.

the clustering results have a moderate to good level of separation and coherence based on the silhouette score

SAMMARY

Clustering technique is used to make some group on different behaviours of customers. It was used to identified them and put them in a group for taking a decision.

In this project I use the clustering technique, make some cluster and find there sum, average and see each clusters different characteristics.

PART 2

LIST OF FIGURES

Bar plot

Boxplot

Heatmap

Scree plot

LIST OF TABLES

Dataset Sample

Summary of the data

Data Describe

Data Description

In this data set we have 640 rows and 61 columns. Here we use PCA for dimensionality reduction . we only take the independent variables for this technique. We consider each variable as a pc and reduce its dimensionality before machine learning technique. Here is the sample data

	State Code	Dist.Code	State	Area Name	No_HH	TOT_M	TOT_F	M_06	F_06	M_SC	...	MARG_CL_0_3_M	MARG_CL_0_3_F	MARG_AL_0_3_M	MARG_...
0	1	1	Jammu & Kashmir	Kupwara	7707	23388	29796	5862	6196	3	...	1150	749	180	
1	1	2	Jammu & Kashmir	Badgam	6218	19585	23102	4482	3733	7	...	525	715	123	
2	1	3	Jammu & Kashmir	Leh(Ladakh)	4452	6546	10964	1082	1018	3	...	114	188	44	
3	1	4	Jammu & Kashmir	Kargil	1320	2784	4206	563	677	0	...	194	247	61	
4	1	5	Jammu & Kashmir	Punch	11654	20591	29981	5157	4587	20	...	874	1928	465	

STATE AND DISTRICT HAS HIGHEST AND LOWEST GENDER RATIO

State with the highest gender ratio: Andhra Pradesh

State with the lowest gender ratio: Lakshadweep

District with the highest gender ratio: 547

District with the lowest gender ratio: 587

OUTLIERS

Here we were asked not to choose outliers. But PCA is a technique which depends on the means. As we know means are highly sensitive of outliers. So treating outliers in PCA is necessary.

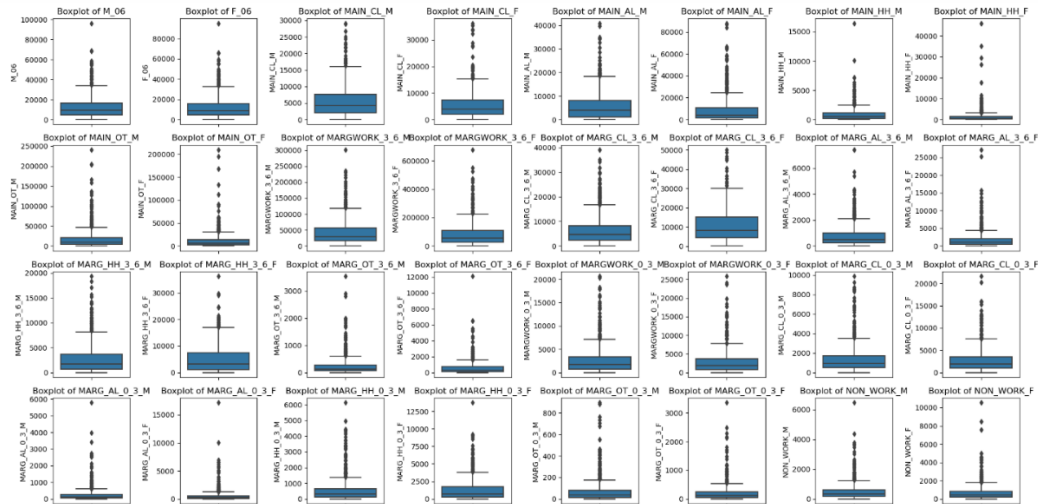
SCALING THE DATA

we use z-score for scaling the data.

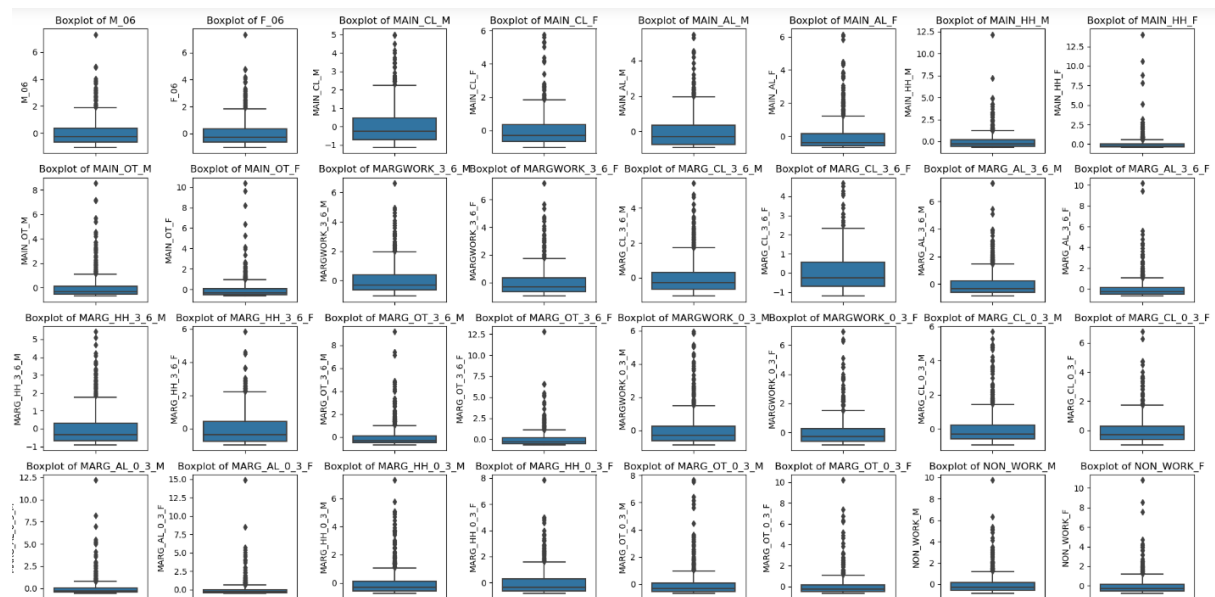
z-score does not impact outliers.

Here is the example

Before z-score



after z-score

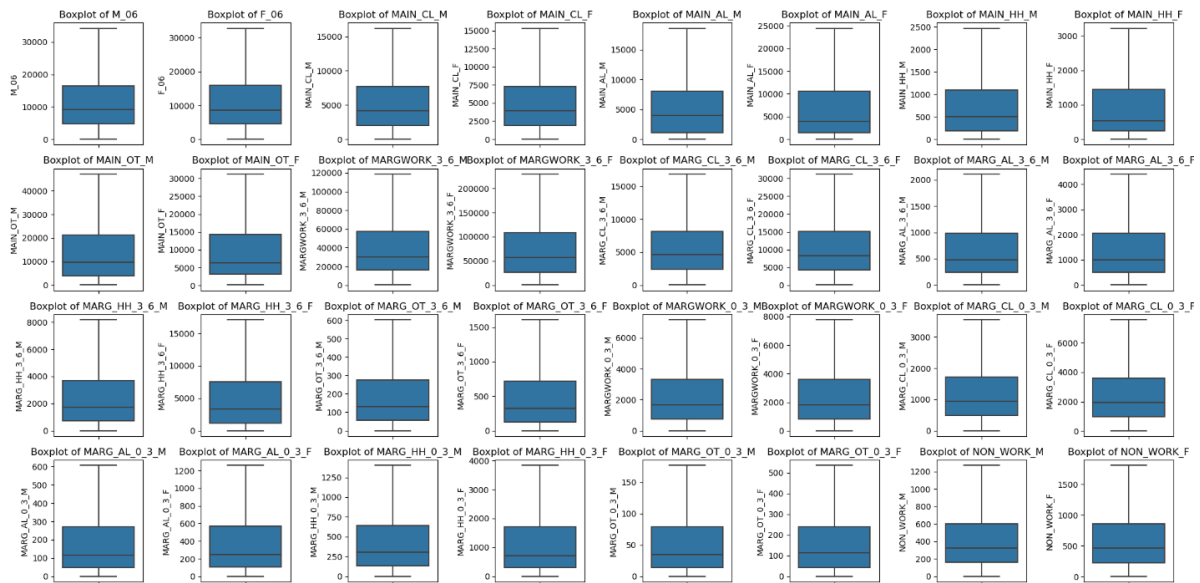


z-score use for scaling the data set, it does not effect the outliers. In PCA treating outliers and z-score both are important.

PERFORMING PCA

PCA cannot go with the outliers and its works on only independent variables.

So fast I use to treat the outliers.



In the boxplot we can see there are only independent variables and containing no outliers.

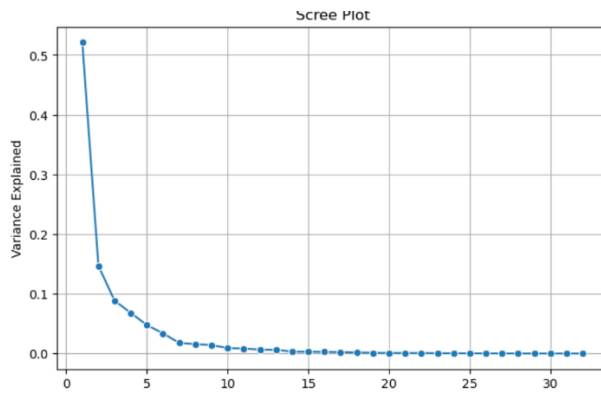
	M_06	-1.00	1.00	0.56	0.38	0.55	0.30	0.66	0.35	0.68	0.56	0.99	0.94	0.86	0.72	0.47	0.25	0.64	0.47	0.69	0.53	0.83	0.74	0.76	0.65	0.27	0.20	0.60	0.51	0.69	0.57	0.78	0.65
	F_06	-1.00	1.00	0.56	0.38	0.55	0.30	0.66	0.36	0.68	0.56	0.99	0.93	0.87	0.72	0.48	0.25	0.65	0.47	0.70	0.54	0.83	0.74	0.76	0.65	0.26	0.19	0.61	0.52	0.70	0.57	0.78	0.65
	MAIN_CL_M	0.56	0.56	1.00	0.67	0.59	0.37	0.41	0.23	0.20	0.16	0.54	0.45	0.54	0.54	0.43	0.30	0.52	0.45	0.47	0.40	0.36	0.31	0.46	0.39	0.24	0.18	0.43	0.36	0.46	0.39	0.34	0.22
	MAIN_CL_F	0.38	0.38	0.67	1.00	0.40	0.49	0.21	0.18	0.11	0.18	0.38	0.27	0.30	0.51	0.42	0.52	0.28	0.42	0.23	0.19	0.16	0.18	0.31	0.37	0.37	0.35	0.23	0.29	0.24	0.18	0.15	0.13
	MAIN_AL_M	0.55	0.55	0.59	0.40	1.00	0.83	0.42	0.39	0.32	0.30	0.55	0.54	0.55	0.61	0.13	-0.03	0.61	0.66	0.40	0.45	0.39	0.38	0.38	0.39	-0.04	-0.07	0.43	0.48	0.37	0.44	0.36	0.31
	MAIN_AL_F	0.30	0.30	0.37	0.49	0.83	1.00	0.21	0.29	0.23	0.31	0.33	0.32	0.25	0.49	-0.03	-0.05	0.31	0.61	0.12	0.19	0.18	0.24	0.12	0.23	-0.11	-0.10	0.16	0.33	0.11	0.17	0.14	0.18
	MAIN_HH_M	0.66	0.66	0.41	0.21	0.42	0.21	1.00	0.52	0.65	0.54	0.69	0.70	0.64	0.51	0.22	0.06	0.40	0.26	0.83	0.69	0.70	0.66	0.53	0.41	0.08	0.03	0.38	0.28	0.73	0.67	0.63	0.48
	MAIN_HH_F	0.35	0.36	0.23	0.18	0.39	0.29	0.52	1.00	0.43	0.41	0.38	0.41	0.33	0.34	0.05	-0.03	0.22	0.19	0.38	0.71	0.37	0.39	0.24	0.23	-0.03	-0.04	0.17	0.15	0.34	0.57	0.34	0.30
	MAIN_OT_M	0.68	0.68	0.20	0.11	0.32	0.23	0.65	0.43	1.00	0.92	0.74	0.86	0.49	0.42	0.02	-0.03	0.12	0.11	0.38	0.40	0.79	0.81	0.37	0.29	-0.04	-0.05	0.09	0.08	0.35	0.38	0.75	0.67
	MAIN_OT_F	0.56	0.56	0.16	0.18	0.30	0.31	0.54	0.41	0.92	1.00	0.63	0.73	0.40	0.42	-0.01	-0.04	0.08	0.14	0.28	0.34	0.68	0.78	0.27	0.24	-0.07	-0.07	0.04	0.08	0.27	0.31	0.62	0.58
	MARGWORK_3_6_M	0.99	0.99	0.54	0.38	0.55	0.33	0.69	0.38	0.74	0.63	1.00	0.96	0.85	0.72	0.44	0.23	0.60	0.45	0.69	0.54	0.86	0.78	0.75	0.64	0.24	0.18	0.57	0.50	0.69	0.57	0.82	0.68
	MARGWORK_3_6_F	0.94	0.93	0.45	0.27	0.54	0.32	0.70	0.41	0.86	0.73	0.96	1.00	0.79	0.66	0.31	0.12	0.50	0.37	0.63	0.54	0.90	0.85	0.68	0.57	0.14	0.09	0.46	0.40	0.61	0.56	0.86	0.76
	MARG_CL_3_6_M	0.86	0.87	0.54	0.30	0.55	0.25	0.64	0.33	0.49	0.40	0.85	0.79	1.00	0.86	0.64	0.31	0.87	0.65	0.83	0.66	0.82	0.73	0.92	0.83	0.34	0.26	0.83	0.75	0.84	0.73	0.80	0.65
	MARG_CL_3_6_F	0.72	0.72	0.54	0.51	0.61	0.49	0.51	0.34	0.42	0.42	0.72	0.66	0.86	1.00	0.63	0.47	0.77	0.85	0.65	0.60	0.66	0.67	0.79	0.88	0.40	0.37	0.70	0.81	0.66	0.64	0.63	0.61
	MARG_AL_3_6_M	0.47	0.48	0.43	0.42	0.13	-0.03	0.22	0.05	0.02	-0.01	0.44	0.31	0.64	0.63	1.00	0.85	0.58	0.41	0.49	0.33	0.32	0.22	0.77	0.77	0.85	0.79	0.64	0.57	0.54	0.42	0.35	0.25
	MARG_AL_3_6_F	0.25	0.25	0.30	0.52	-0.03	-0.05	0.06	-0.03	-0.03	-0.04	0.23	0.12	0.31	0.47	0.85	1.00	0.22	0.21	0.20	0.10	0.11	0.07	0.49	0.59	0.92	0.94	0.27	0.28	0.23	0.15	0.14	0.13
	MARG_HH_3_6_M	0.64	0.65	0.52	0.28	0.61	0.31	0.40	0.22	0.12	0.08	0.60	0.50	0.87	0.77	0.58	0.22	1.00	0.81	0.68	0.55	0.45	0.35	0.80	0.76	0.26	0.17	0.93	0.88	0.72	0.65	0.47	0.35
	MARG_HH_3_6_F	0.47	0.47	0.45	0.42	0.66	0.61	0.26	0.19	0.11	0.14	0.45	0.37	0.65	0.85	0.41	0.21	0.81	1.00	0.44	0.39	0.31	0.31	0.56	0.70	0.15	0.11	0.69	0.88	0.46	0.45	0.32	0.31
	MARG_OT_3_6_M	0.69	0.70	0.47	0.23	0.40	0.12	0.83	0.38	0.38	0.28	0.69	0.63	0.83	0.65	0.49	0.20	0.68	0.44	1.00	0.76	0.69	0.59	0.77	0.63	0.25	0.17	0.70	0.56	0.95	0.81	0.66	0.48
	MARG_OT_3_6_F	0.53	0.54	0.40	0.19	0.45	0.19	0.69	0.71	0.40	0.34	0.54	0.54	0.66	0.60	0.33	0.10	0.55	0.39	0.76	1.00	0.57	0.55	0.59	0.55	0.13	0.07	0.54	0.46	0.73	0.95	0.55	0.47
	MARGWORK_0_3_M	0.83	0.83	0.36	0.16	0.39	0.18	0.70	0.37	0.79	0.68	0.86	0.90	0.82	0.66	0.32	0.11	0.45	0.31	0.69	0.57	1.00	0.95	0.69	0.57	0.13	0.09	0.44	0.36	0.67	0.58	0.93	0.81
	MARGWORK_0_3_F	0.74	0.74	0.31	0.18	0.38	0.24	0.66	0.39	0.81	0.78	0.78	0.85	0.73	0.67	0.22	0.07	0.35	0.31	0.59	0.55	0.95	1.00	0.58	0.52	0.06	0.03	0.33	0.31	0.57	0.55	0.85	0.83
	MARG_CL_0_3_M	0.76	0.76	0.46	0.31	0.38	0.12	0.53	0.24	0.37	0.27	0.75	0.68	0.92	0.79	0.77	0.49	0.80	0.56	0.77	0.59	0.69	0.58	1.00	0.92	0.59	0.49	0.88	0.77	0.83	0.71	0.77	0.62
	MARG_CL_0_3_F	0.65	0.65	0.39	0.37	0.39	0.23	0.41	0.23	0.29	0.24	0.64	0.57	0.83	0.88	0.77	0.59	0.76	0.70	0.63	0.55	0.57	0.52	0.92	1.00	0.62	0.59	0.81	0.85	0.69	0.67	0.64	0.64
	MARG_AL_0_3_M	0.27	0.26	0.24	0.37	-0.04	-0.11	0.08	-0.03	-0.04	0.07	0.24	0.14	0.34	0.40	0.85	0.92	0.26	0.15	0.25	0.13	0.13	0.06	0.59	0.62	1.00	0.96	0.35	0.30	0.29	0.20	0.19	0.14
	MARG_AL_0_3_F	0.20	0.19	0.18	0.35	-0.07	-0.10	0.03	-0.04	-0.05	0.07	0.18	0.09	0.26	0.37	0.79	0.94	0.17	0.11	0.17	0.07	0.09	0.03	0.49	0.59	0.96	1.00	0.25	0.23	0.21	0.13	0.14	0.12
	MARG_HH_0_3_M	0.60	0.61	0.43	0.23	0.43	0.16	0.38	0.17	0.09	0.04	0.57	0.46	0.83	0.70	0.64	0.27	0.93	0.69	0.70	0.54	0.44	0.33	0.88	0.81	0.35	0.25	1.00	0.90	0.77	0.67	0.49	0.36
	MARG_HH_0_3_F	0.51	0.52	0.36	0.29	0.48	0.33	0.28	0.15	0.08	0.08	0.50	0.40	0.75	0.81	0.57	0.28	0.88	0.88	0.56	0.46	0.36	0.31	0.77	0.85	0.30	0.23	0.90	1.00	0.63	0.58	0.41	0.37
	MARG_OT_0_3_M	0.69	0.70	0.46	0.24	0.37	0.11	0.73	0.34	0.35	0.27	0.69	0.61	0.84	0.66	0.54	0.23	0.72	0.46	0.95	0.73	0.67	0.57	0.83	0.69	0.29	0.21	0.77	0.63	1.00	0.82	0.67	0.49
	MARG_OT_0_3_F	0.57	0.57	0.39	0.18	0.44	0.17	0.67	0.57	0.38	0.31	0.57	0.56	0.73	0.64	0.42	0.15	0.65	0.45	0.81	0.95	0.58	0.55	0.71	0.67	0.20	0.13	0.67	0.58	0.82	1.00	0.61	0.52
	NON_WORK_M	0.78	0.78	0.34	0.15	0.36	0.14	0.63	0.34	0.75	0.62	0.82	0.86	0.80	0.63	0.35	0.14	0.47	0.32	0.66	0.55	0.93	0.85	0.77	0.64	0.19	0.14	0.49	0.41	0.67	0.61	1.00	0.88
	NON_WORK_F	0.65	0.65	0.22	0.13	0.31	0.18	0.48	0.30	0.67	0.58	0.68	0.76	0.65	0.61	0.25	0.13	0.35	0.31	0.48	0.47	0.81	0.83	0.62	0.64	0.14	0.12	0.36	0.37	0.49	0.52	0.88	1.00

Here is the covariance matrix.

IDENTIFYING OPTIMUM NUMBER OF PCs

To identify the optimal number of pcs I use the see the scree plot and explained variance ratio

```
[0.52120054, 0.6677232 , 0.75618259, 0.82417856, 0.87170677,
0.90546198, 0.92310679, 0.93879236, 0.952818 , 0.96188109,
0.97010962, 0.97658872, 0.98265229, 0.9857252 , 0.98868869,
0.99133356, 0.9935503 , 0.99511636, 0.99615724, 0.99704083,
0.99780206, 0.99849627, 0.99901355, 0.99936005, 0.99965068,
0.99984142, 0.99996858, 1. , 1. , 1. ,
1. , 1. ])
```

Lets take 90% of explained variance ratio.

PC1 is the most important among all the PCs. Only pc1 contents 52% of the data.

If we take care only pc1 components huge benefits we will get.



LINEAR EQUATION FOR FIRST PC

$$PC1 = w1 * X1 + w2 * X2 + w3 * X3 + ... + wn * Xn$$

Where X1, X2, X3, ..., Xn are the original variables, and w1, w2, w3, ..., wn are the coefficients or loadings of the first PC.