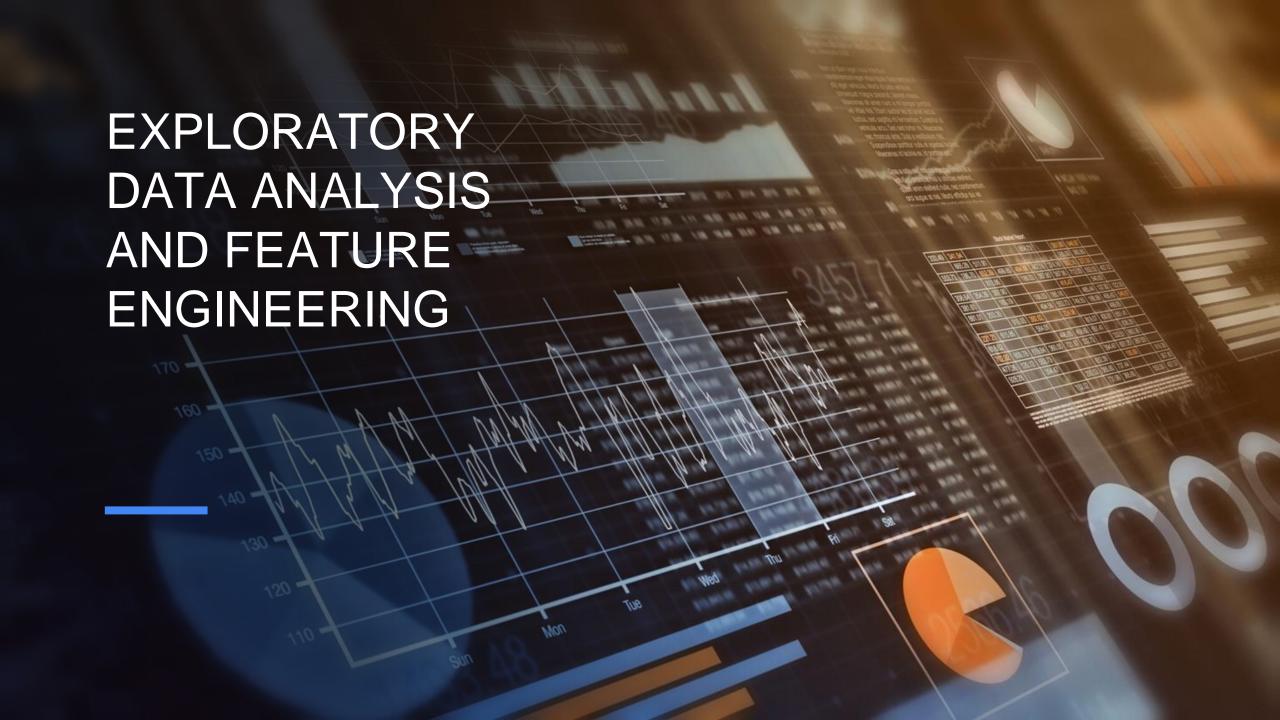
UNSUPERVISED
METHODS TO GROUP
USERS' CONSUMPTION
BEHAVIOUR TO
ENHANCE
PERSONALIZE SERVICE
DEGRADATION
POLICIES

PROJECT MODULE 1

ROBERTO CASALUCE – MACIEJ ZUZIAK



OVERVIEW



ONE DATASET
CONTAINING 1249
STATISTICAL UNITS



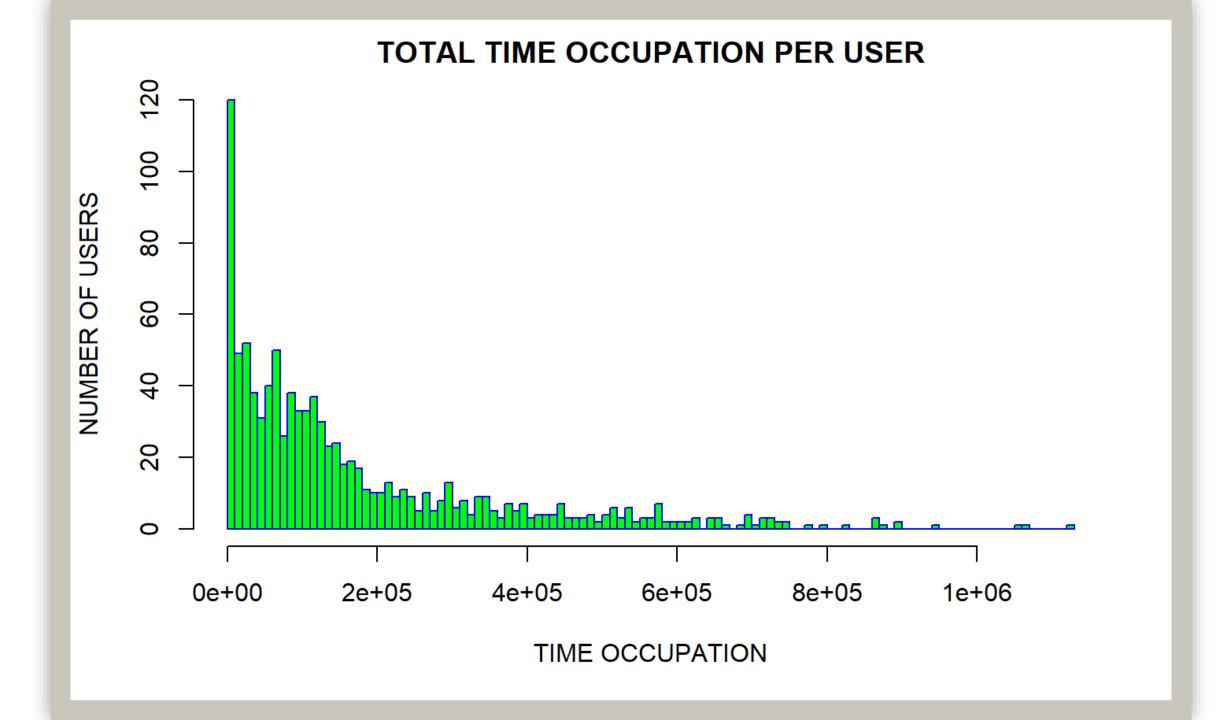
114 DIFFERENT
FEATURES (ONLINE
PLATFORMS AND
MOBILE APPLICATIONS)

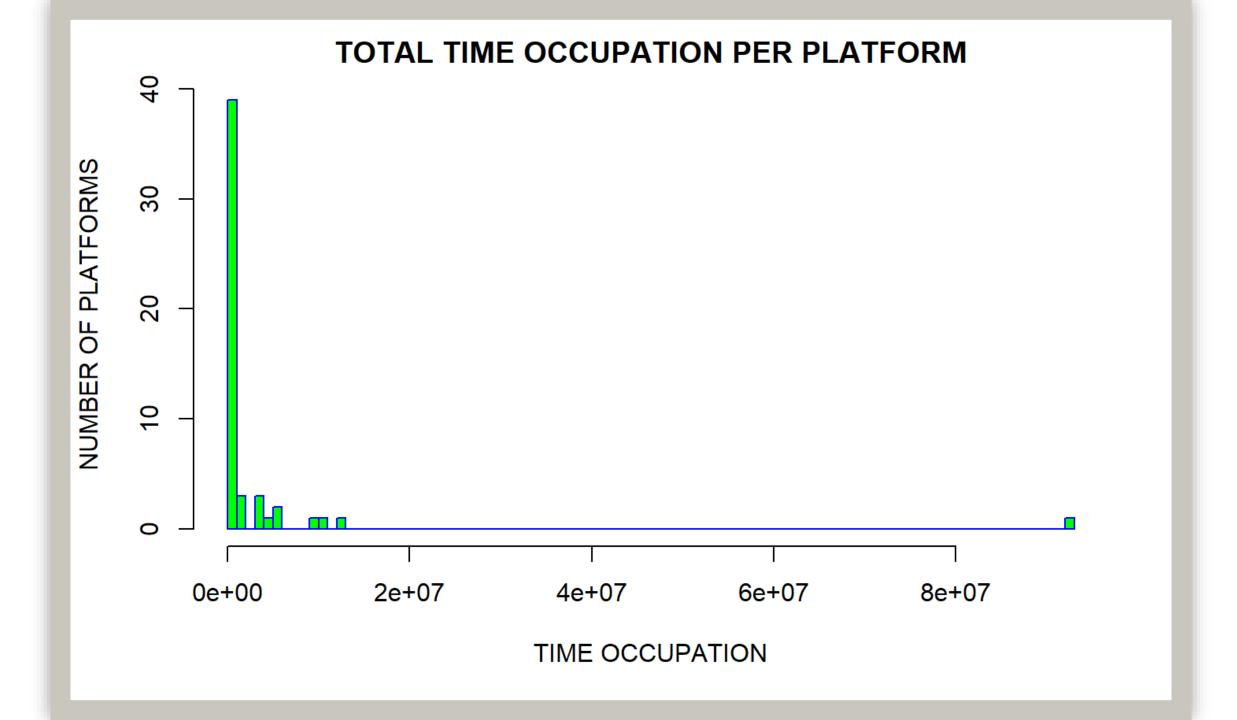


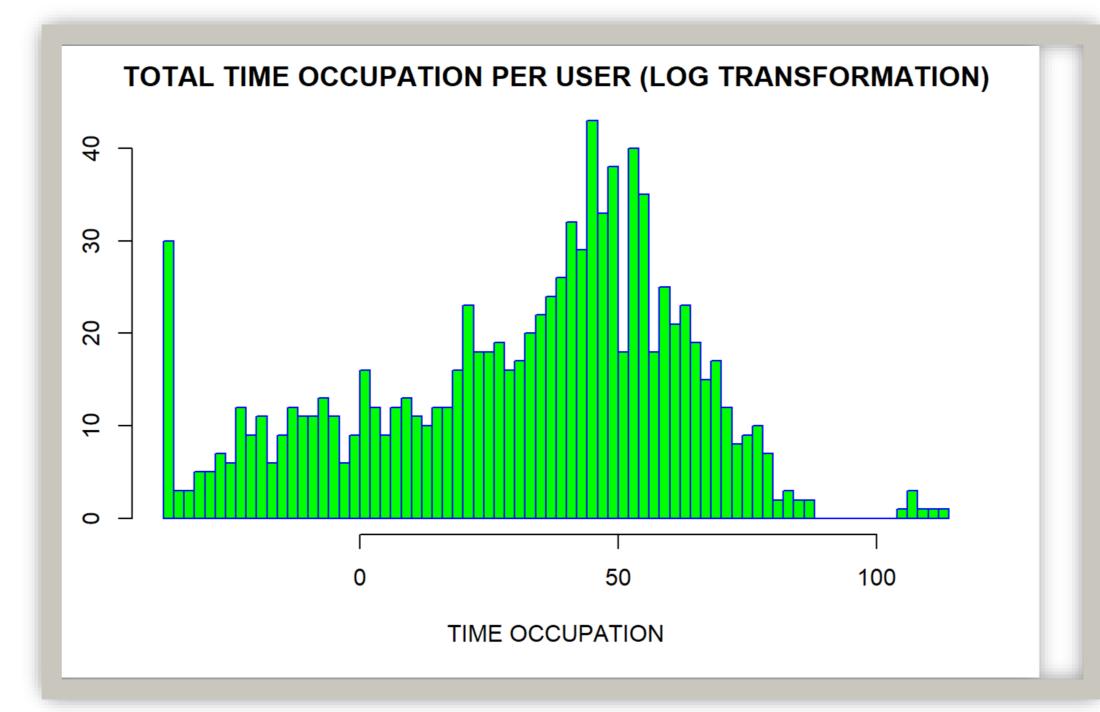
TIME EXPRESSED IN SECOND, DATA CONSUMPTION IN BYTES



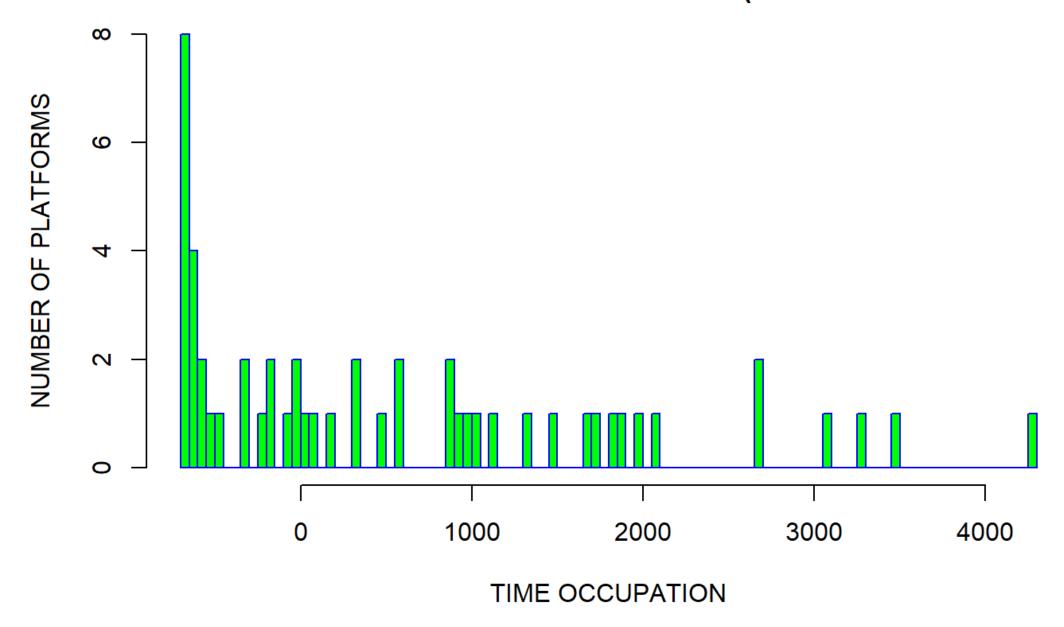
TASK IS TO IDENTIFY THE EXISTENCE OF CONSUMER GROUPS





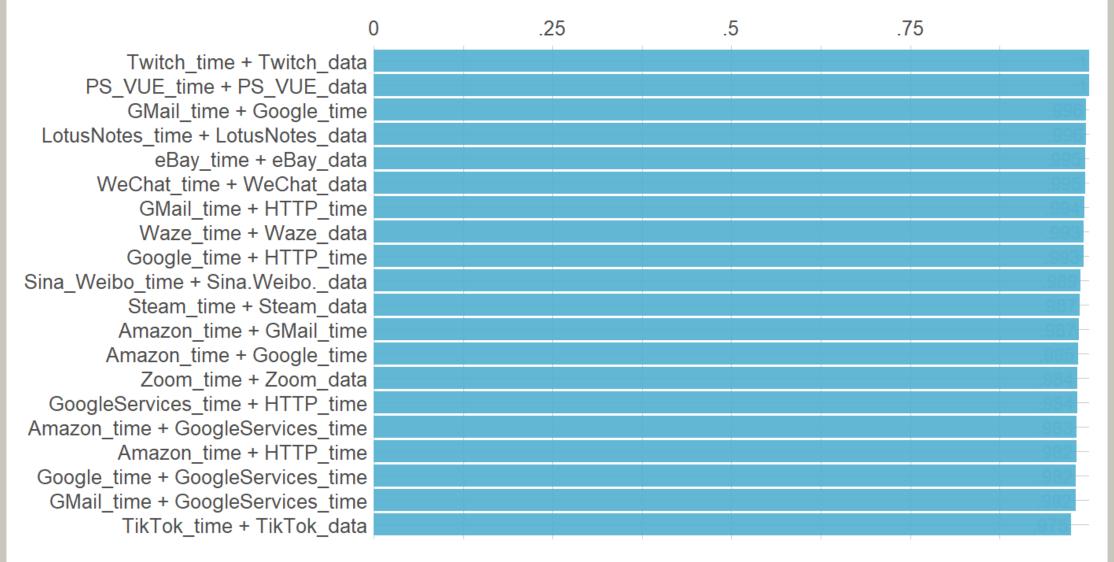


TOTAL TIME OCCUPATION PER PLATFORM (LOG TRANSFORMATION)



Ranked Cross-Correlations

20 most relevant





WE HAVE...

01

DELETED FEATURES RELATED TO DATA USAGE 02

DELETED FEATURES
THAT CONTAINED
LESS THAN TWO
USERS' ENTRIES

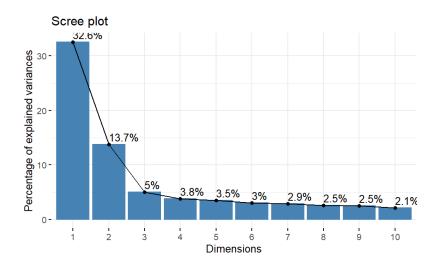
03

DELETED DUPLICATED ENTRIES 04

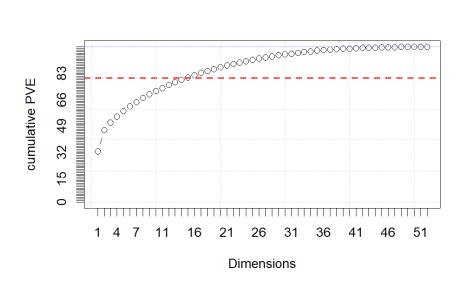
PERFORMED LOG TRANSFORMATION DESCRIBED BEFORE

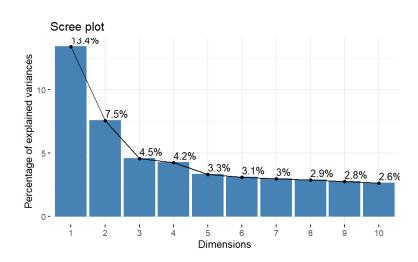


PCA

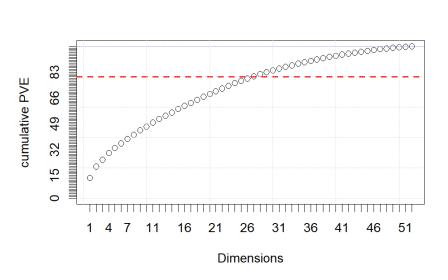


Log transformed data

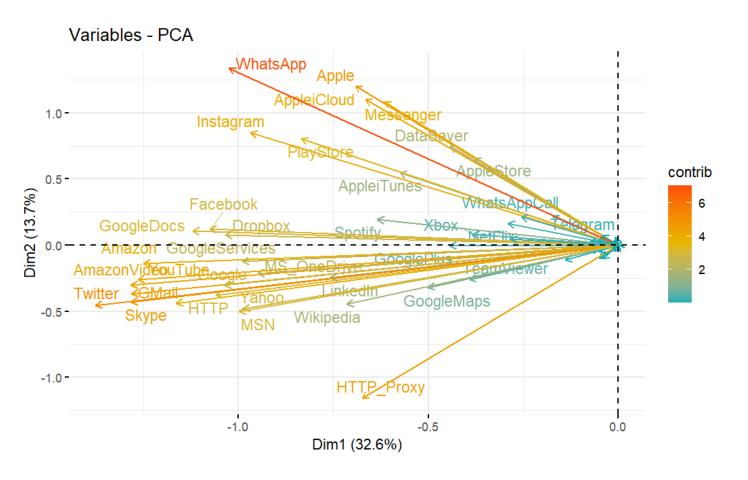


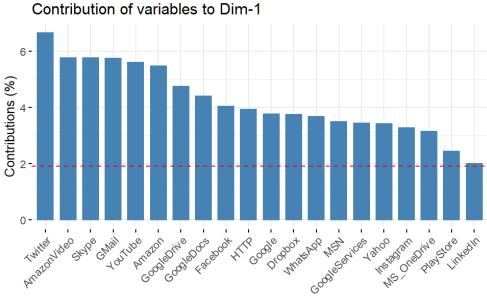


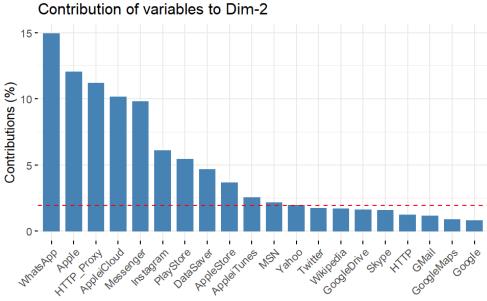
Scaled data



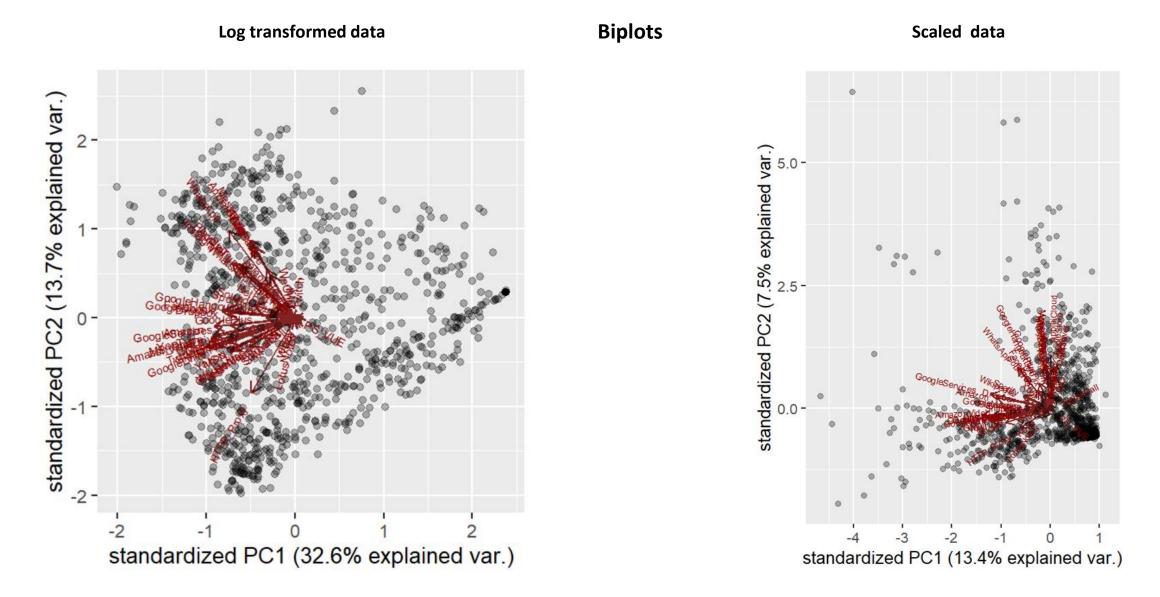
PCA







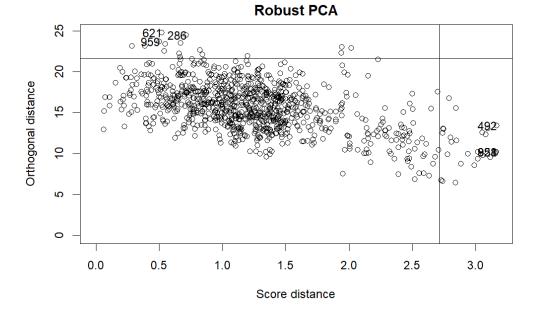
PCA



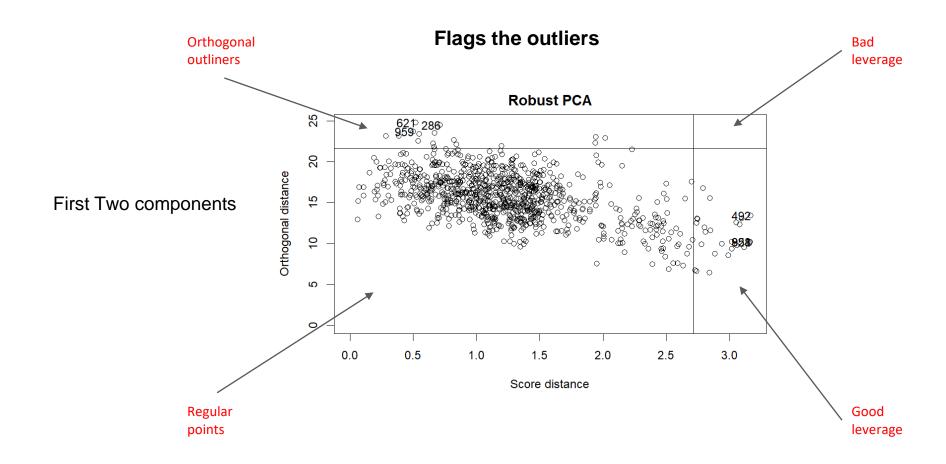
Robust PCA

Flags the outliers

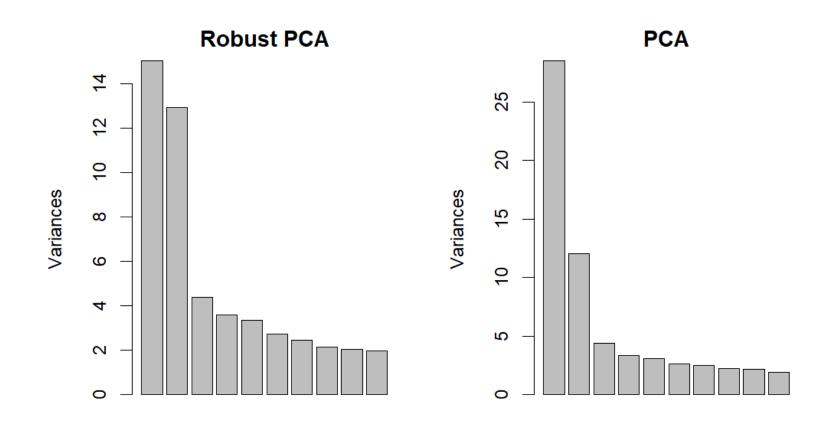
First Two components

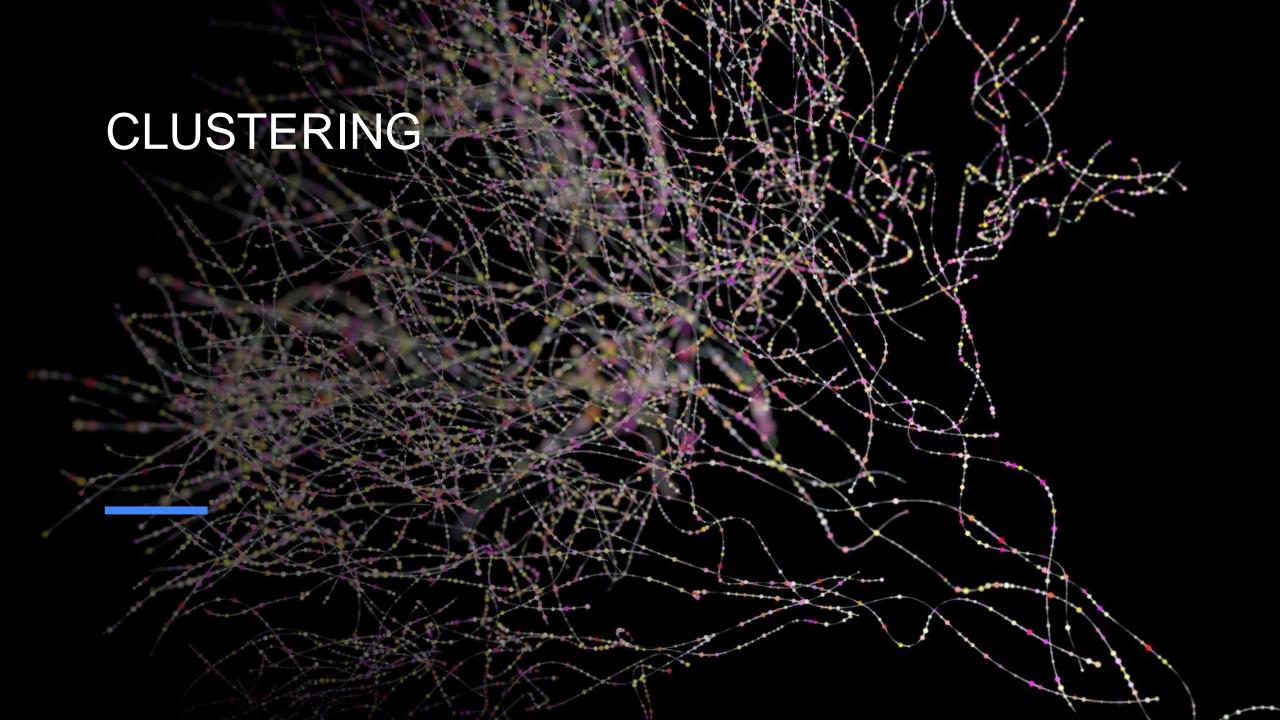


Robust PCA



Variances Robust and Classic PCA





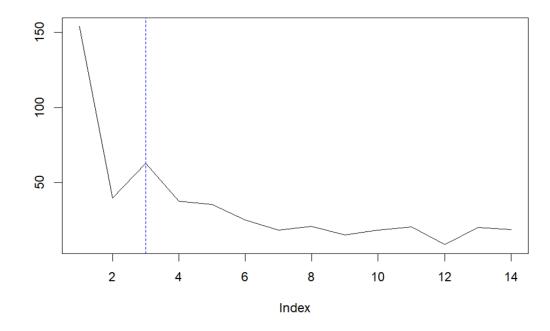
Clustering method selection

From a previous work (Rojas et al., 2020): 3/4 clusters - Low, Medium, High and Very High

Hartigan Index

Agglomerative Hierarchical clustering method

Kmeans clustering method



Evaluating Clustering methods

Cluster **stability** measure:

- The average proportion of non-overlap (APN)
- The average distance (AD)
- The average distance between means (ADM)
- The figure of merit (FOM)

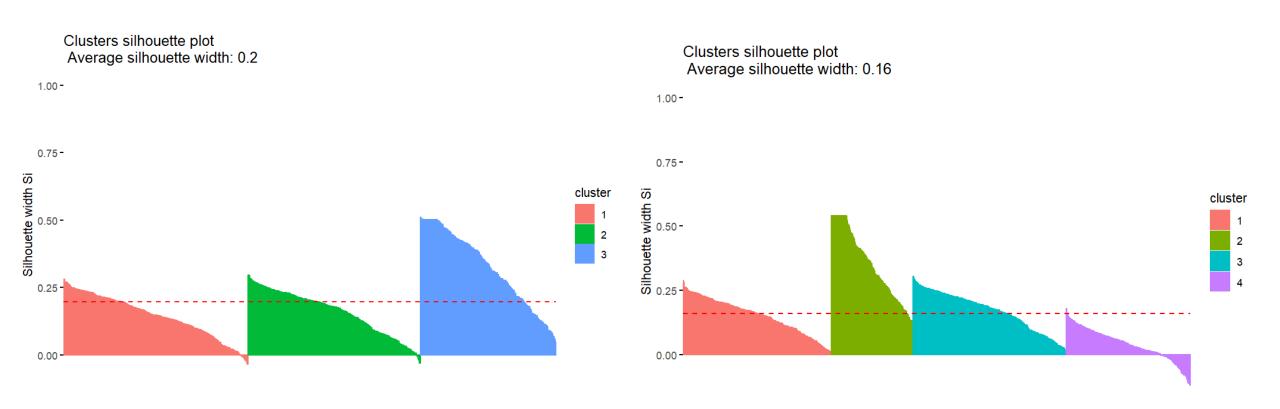
APN	0.009	kmeans	3
AD	10.2	kmeans	4
ADM	0.095	kmeans	3
FOM	0.93	kmeans	4

Connectivity	139.75	hierarchical	3
Dunn	0.26	hierarchical	3
Silhouette	0.20	kmeans	3

Internal measures for cluster validation

Evaluating a Clustering Solution

Silhouette widths kmeans

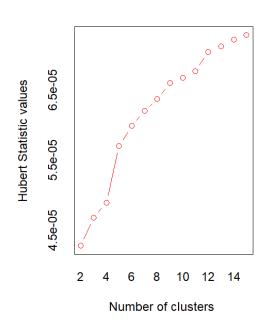


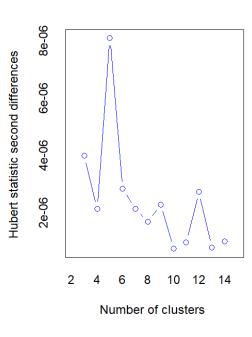
Evaluating a Clustering Solution

Hubert index kmeans

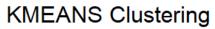
According to the majority rule, the best number of clusters is 3

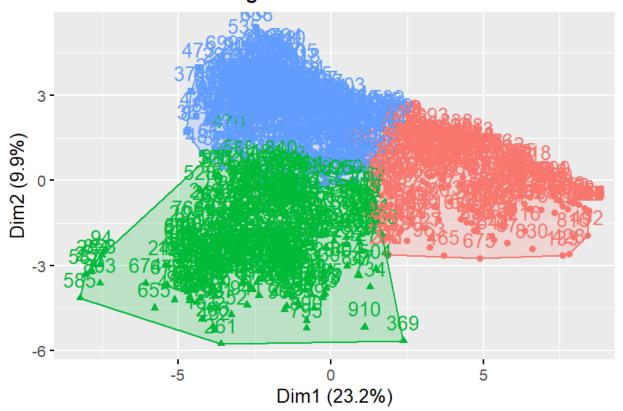
- * 11 proposed 3 as the best number of clusters
- * 1 proposed 4 as the best number of clusters
- * 9 proposed 2 as the best number of clusters





Kmeans 3 clusters





cluster		Size Clusters
1	1	2
2	266	343

Reference

Chen, X., Zhang, B., Wang, T., Bonni, A., & Zhao, G. (2020). Robust principal component analysis for accurate outlier sample detection in RNA-Seq data. *BMC bioinformatics*, *21*(1), 1-20.

Metcalf, L., & Casey, W. (2016). Cybersecurity and applied mathematics. Syngress.

Rojas, J. S., Pekar, A., Rendón, Á., & Corrales, J. C. (2020). Smart user consumption profiling: Incremental learning-based OTT service degradation. *IEEE access*, *8*, 207426-207442.

Forgy, E. W. (1965). Cluster analysis of multivariate data: efficiency vs interpretability of classifications. Biometrics, 21, 768–769.

Hartigan, J. A. and Wong, M. A. (1979). Algorithm AS 136: A K-means clustering algorithm. Applied Statistics, 28, 100–108. doi: 10.2307/2346830.

Lloyd, S. P. (1957, 1982). Least squares quantization in PCM. Technical Note, Bell Laboratories. Published in 1982 in IEEE Transactions on Information Theory, 28, 128–137.

MacQueen, J. (1967). Some methods for classification and analysis of multivariate observations. In Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, eds L. M. Le Cam & J. Neyman, 1, pp. 281–297. Berkeley, CA: University of California Press.

Reference

Charrad M., Ghazzali N., Boiteau V., Niknafs A. (2014). "NbClust: An R Package for Determining the Relevant Number of Clusters in a Data Set.", "Journal of Statistical Software, 61(6), 1-36.", "URL http://www.jstatsoft.org/v61/i06/".

Becker, R. A., Chambers, J. M. and Wilks, A. R. (1988) The New S Language. Wadsworth & Brooks/Cole.

Mardia, K. V., Kent, J. T. and Bibby, J. M. (1979) Multivariate Analysis. Academic Press.

Borg, I. and Groenen, P. (1997) Modern Multidimensional Scaling. Theory and Applications. Springer.