

Index

| | | |
|-----|--|---|
| 1 | Description | 1 |
| 2 | Clustering..... | 1 |
| 2.1 | Dimensionality reduction..... | 1 |
| 2.2 | Elbow Method | 4 |
| 2.3 | Kmeans clustering..... | 6 |
| 2.4 | DBscan clustering | 7 |
| 2.5 | Difference between DBscan and Kmeans | 7 |
| 3 | Descriptive Analysis | 8 |

1 Description

Analysis of users of various social media.

In the first part (using Python, sklearn and Google Colab), classes of users with similar habits and characteristics were identified using clustering algorithms.

In the second part with tableau a descriptive analysis was carried out.

2 Clustering

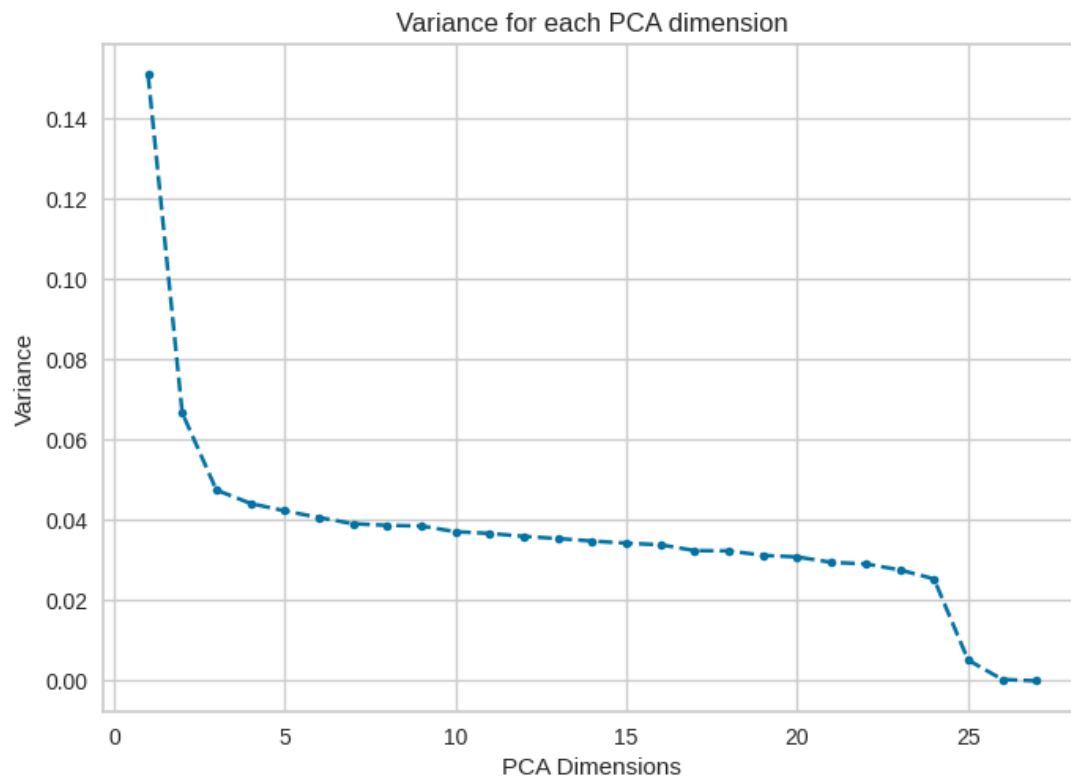
After studying the dataset with the `head()`, `info()`, and `describe ()` commands, we moved on to preparing the data for the application of clustering algorithms.

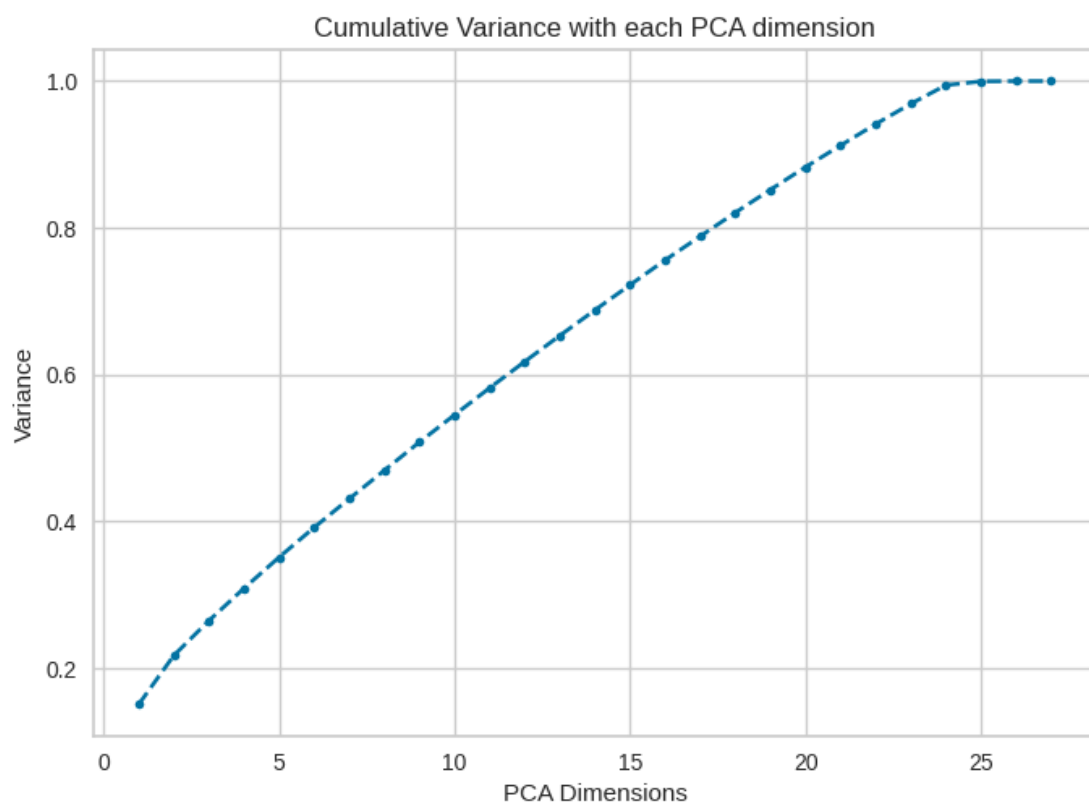
After briefly seeing the correlation between the features, the encoding of the categorical features and the scaling of the data were carried out. This was necessary to increase the performance of the algorithms applied subsequently.

2.1 Dimensionality reduction

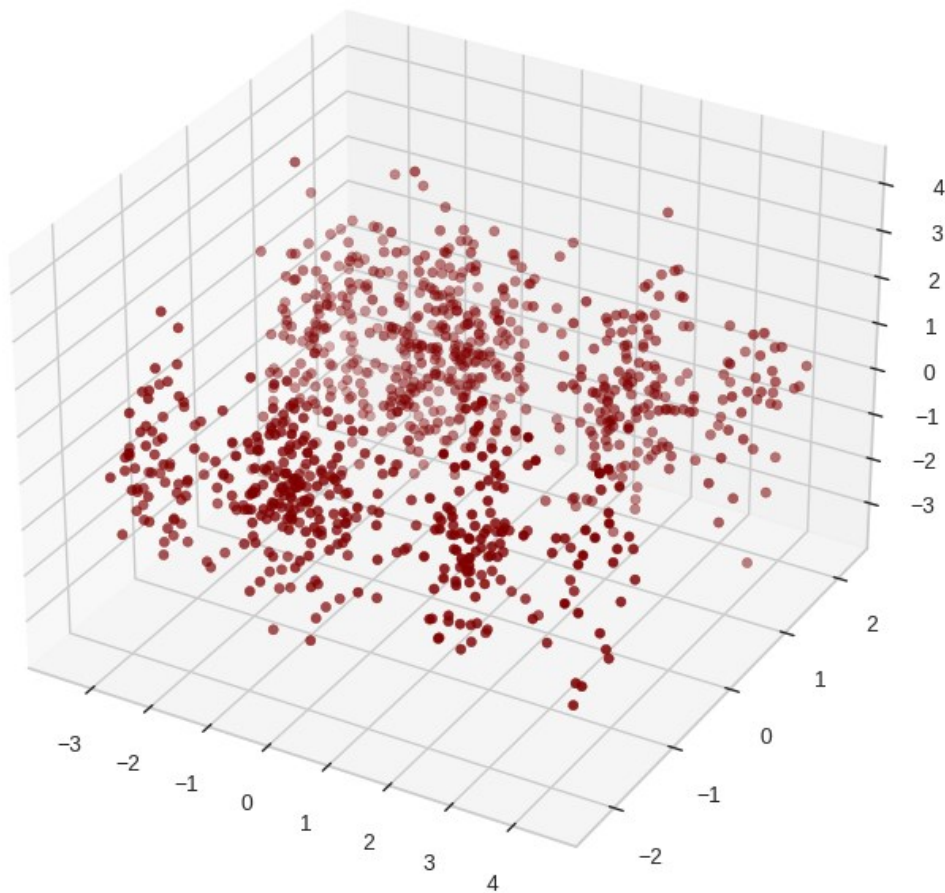
The dataset has many features, therefore to improve the training and performance of the algorithms it was necessary to reduce the dimensionality with the PCA method.

After seeing the results by evaluating the variance based on the number of features it was decided to proceed with the parameter `n_components=3`. Thus the model is simplified in the training and evaluation phase but the loss of information is minimal





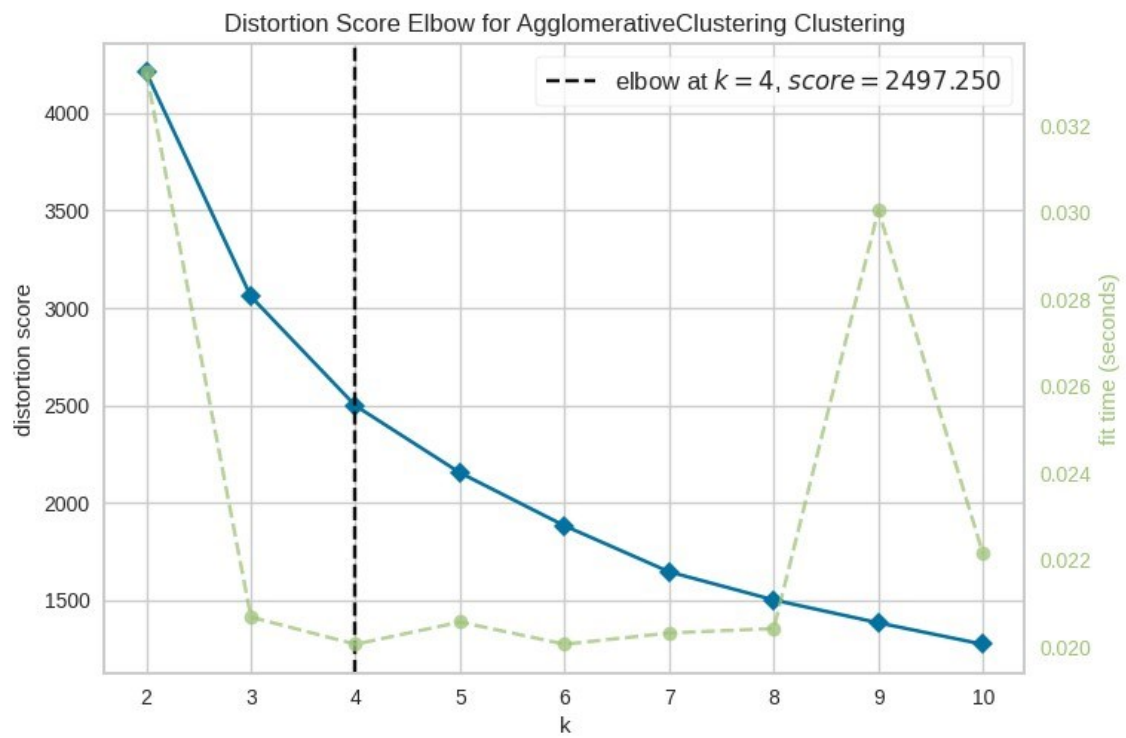
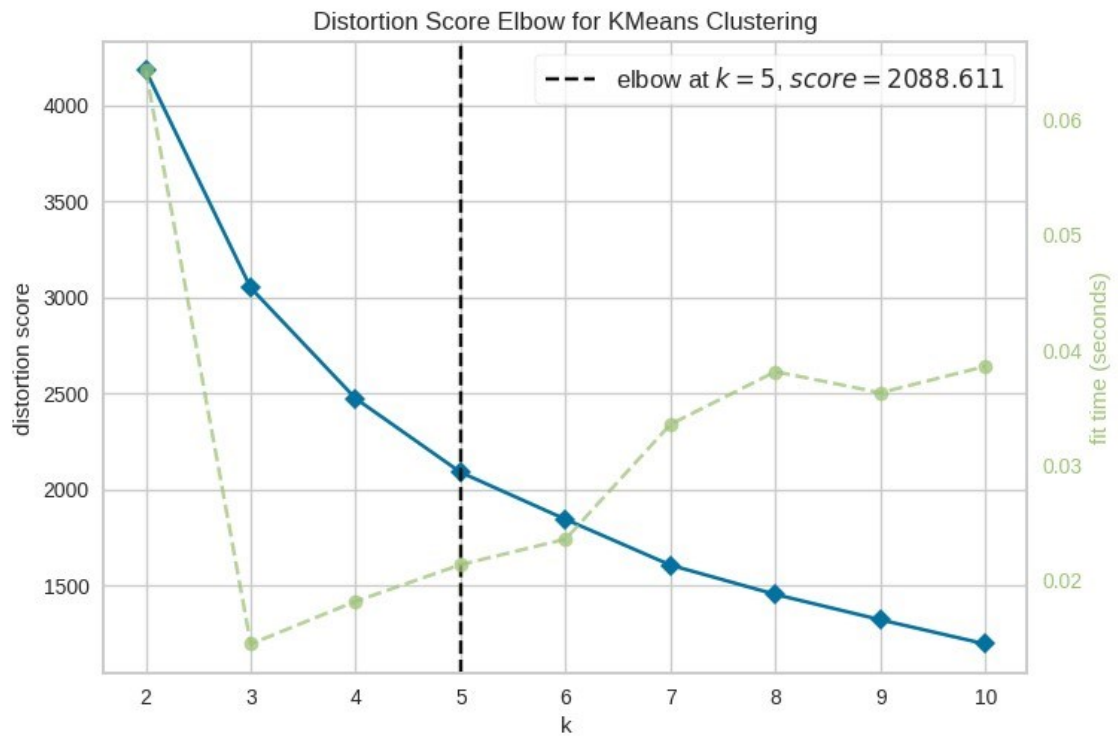
3D visualization after pca



2.2 Elbow Method

For algorithms where the parameter k indicating the number of clusters must be provided, this method provides the value to obtain optimal results.

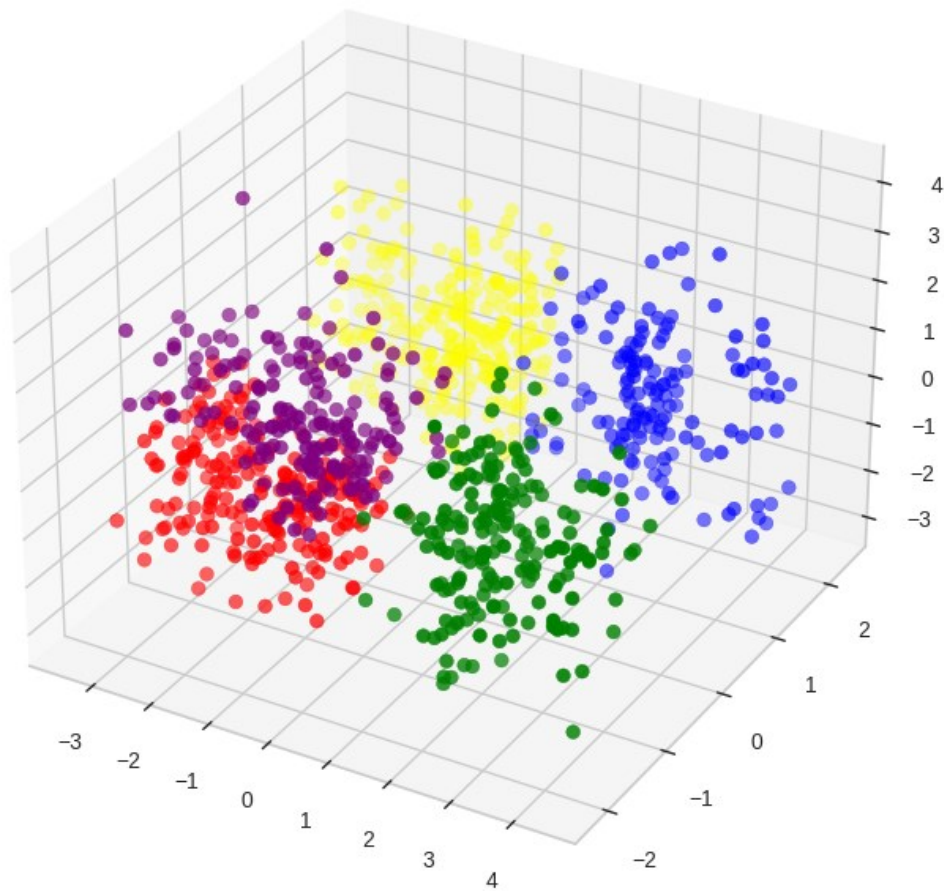
It has been evaluated for the following algorithms:



2.3 Kmeans clustering

We chose to proceed with this algorithm by specifying the parameter $k=5$.

3D Kmeans cluster visualization

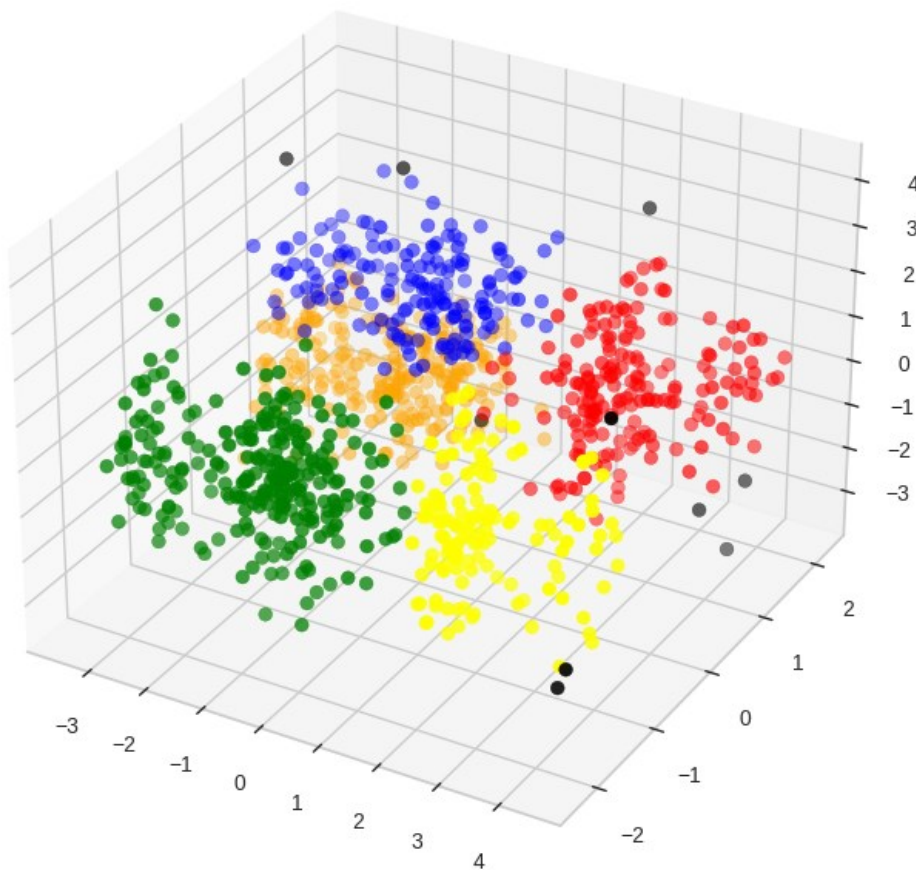


2.4 DBscan clustering

We also chose to see the results of an algorithm with a different logic. In black are the noise points,

which are points not assigned to any cluster as they are too distant from the clusters.

3D DBscan cluster visualization



2.5 Difference between DBscan and Kmeans

Different algorithms offer different results as highlighted below

```
[72] km=dwc['Clusters']
      db=dwc['DB_Clusters']

      cont=0

      for i in km:
          if km[i]!=db[i]:
              cont=cont+1

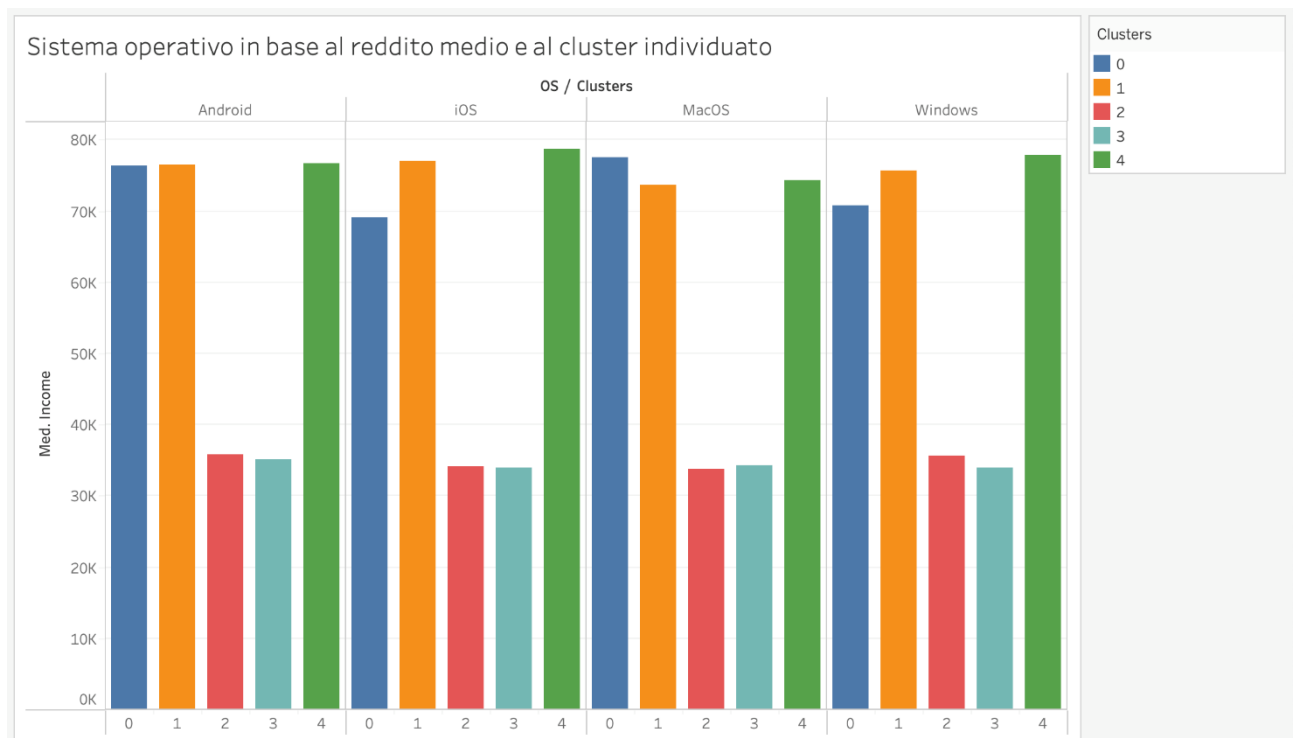
      cont=str(cont)
      print('records with different cluster between dbcan and km are : ' +cont)
```

→ records with different cluster between dbcan and km are : 410

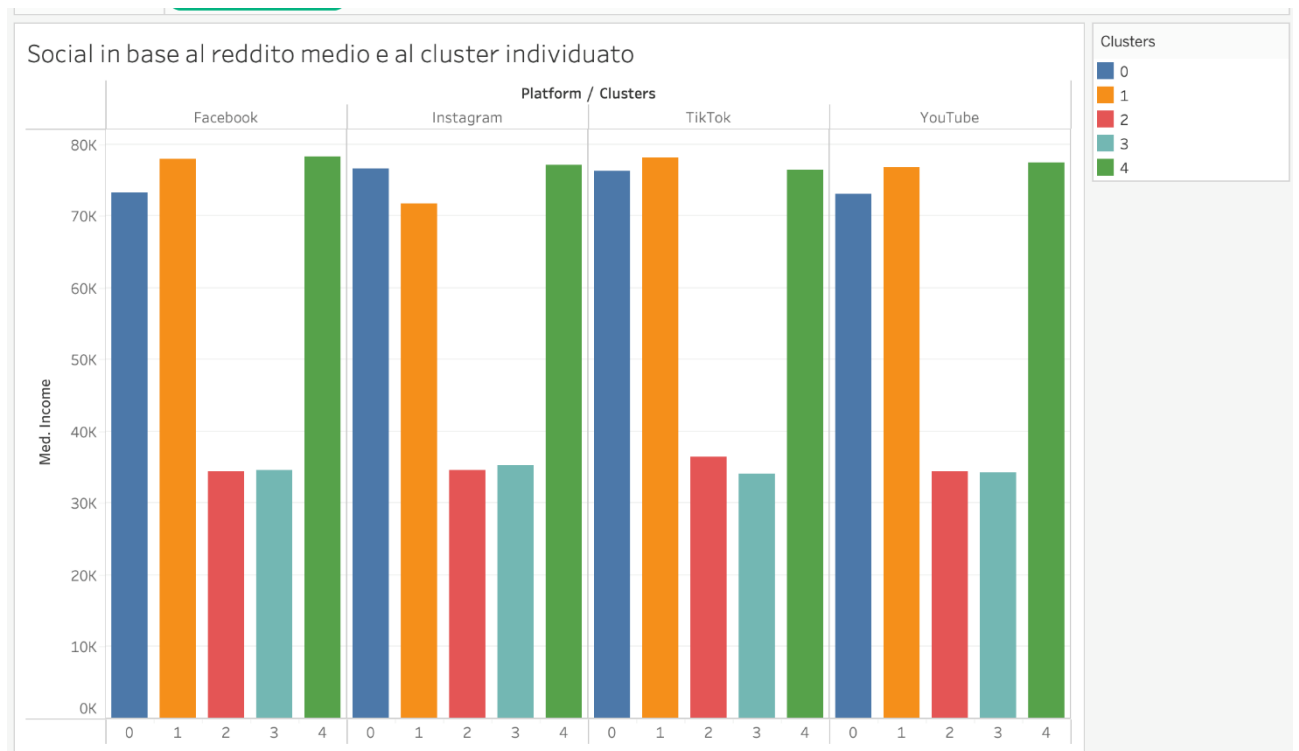
3 Descriptive Analysis

After obtaining 5 clusters, a descriptive analysis is carried out based on the features present in the database.

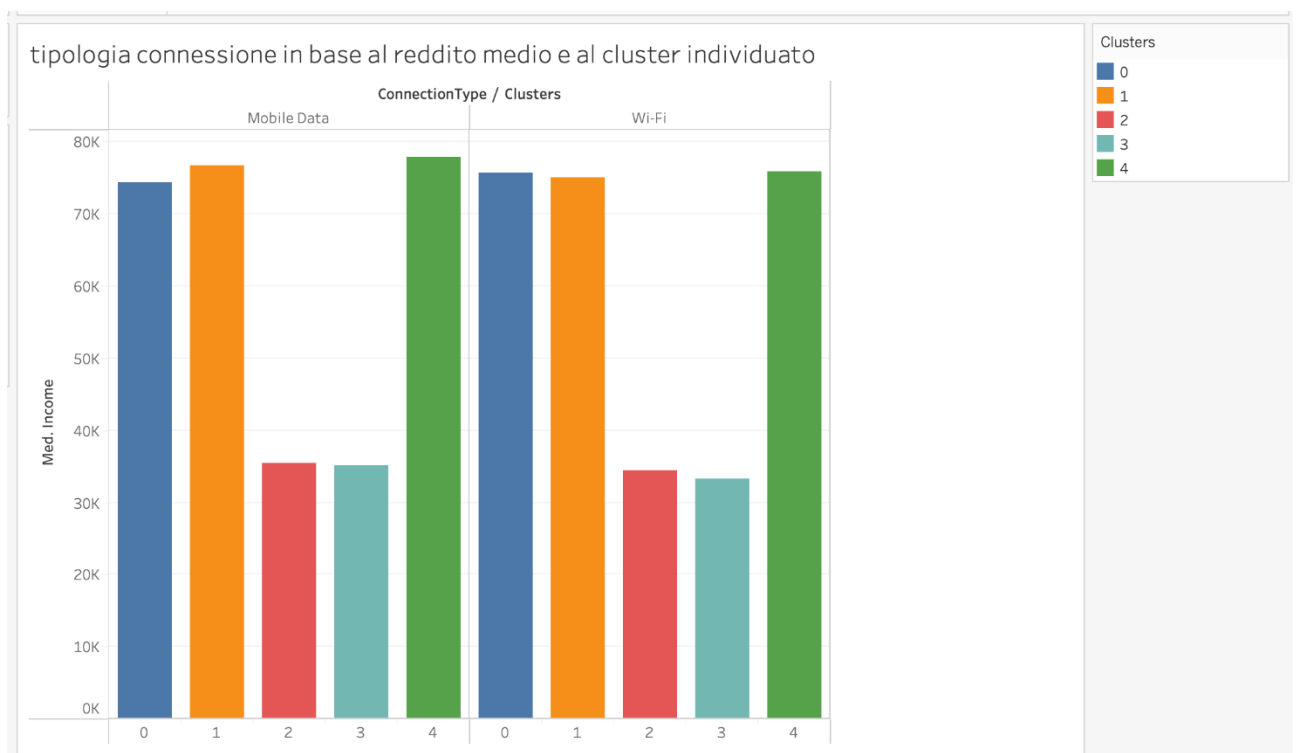
- Operating system



- Social Platform

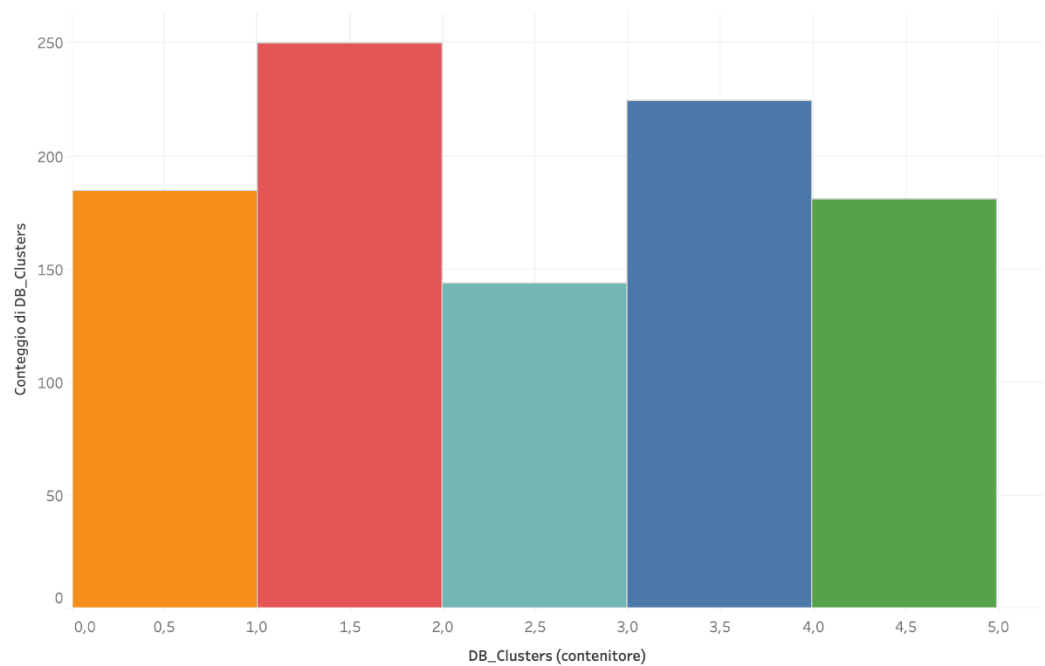


- Type of internet connection

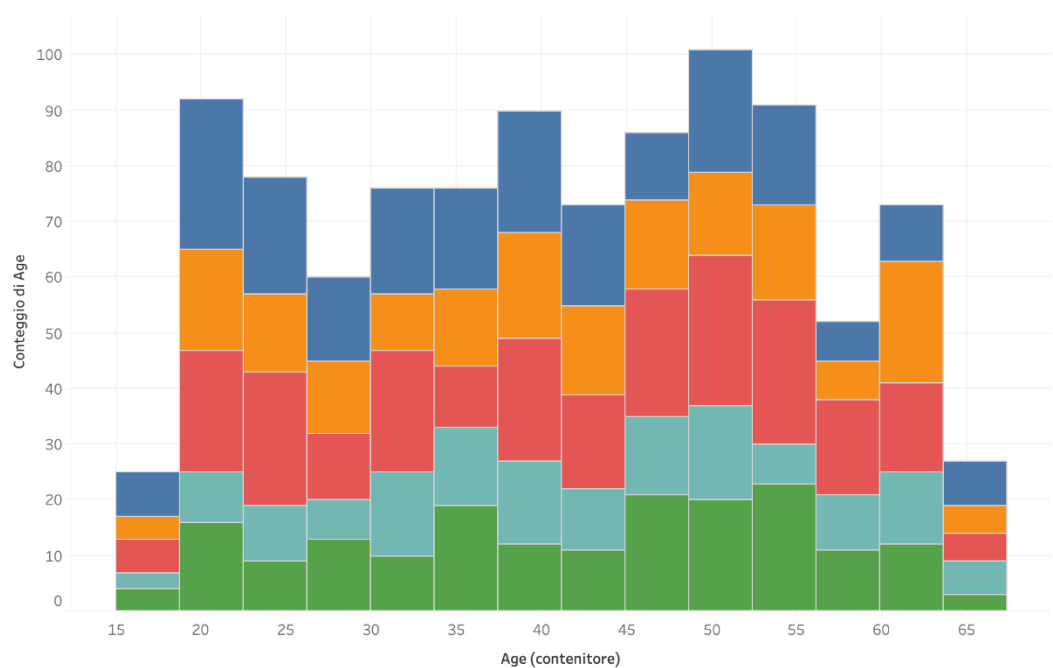


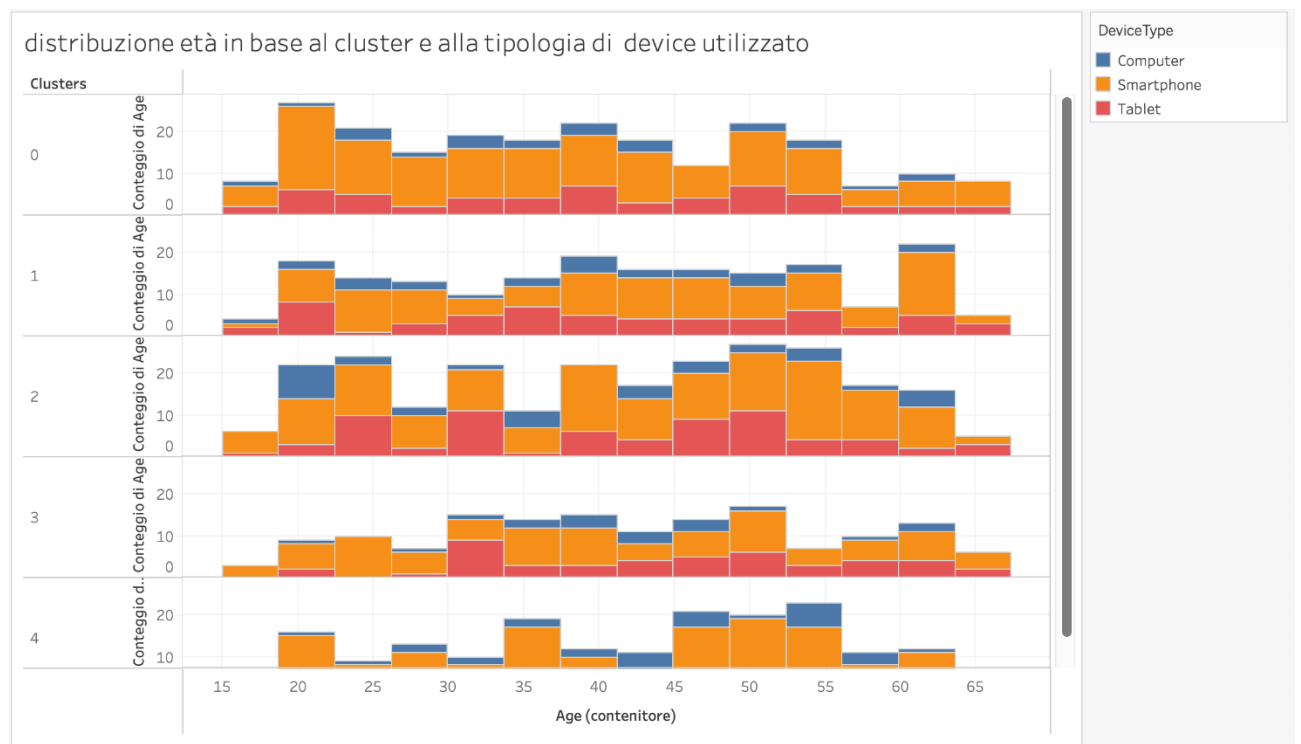
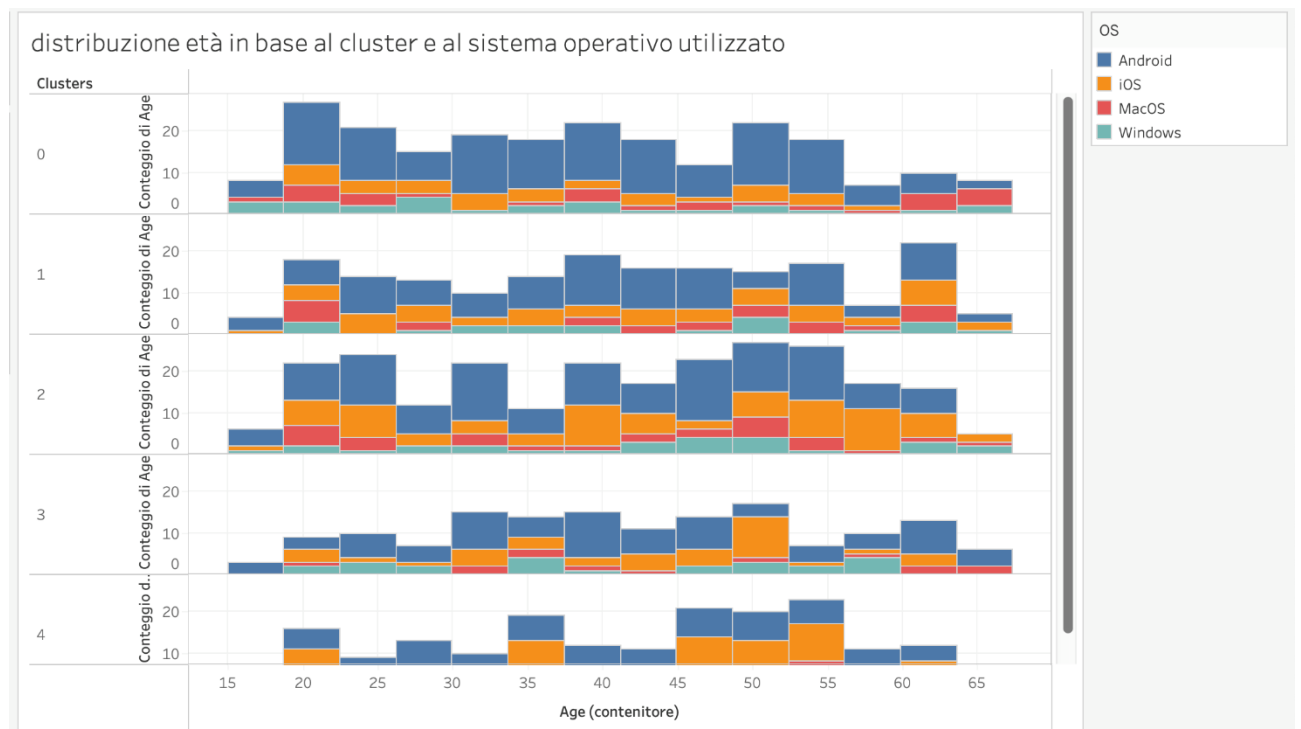
- Cluster distribution

distribuzione cluster individuato

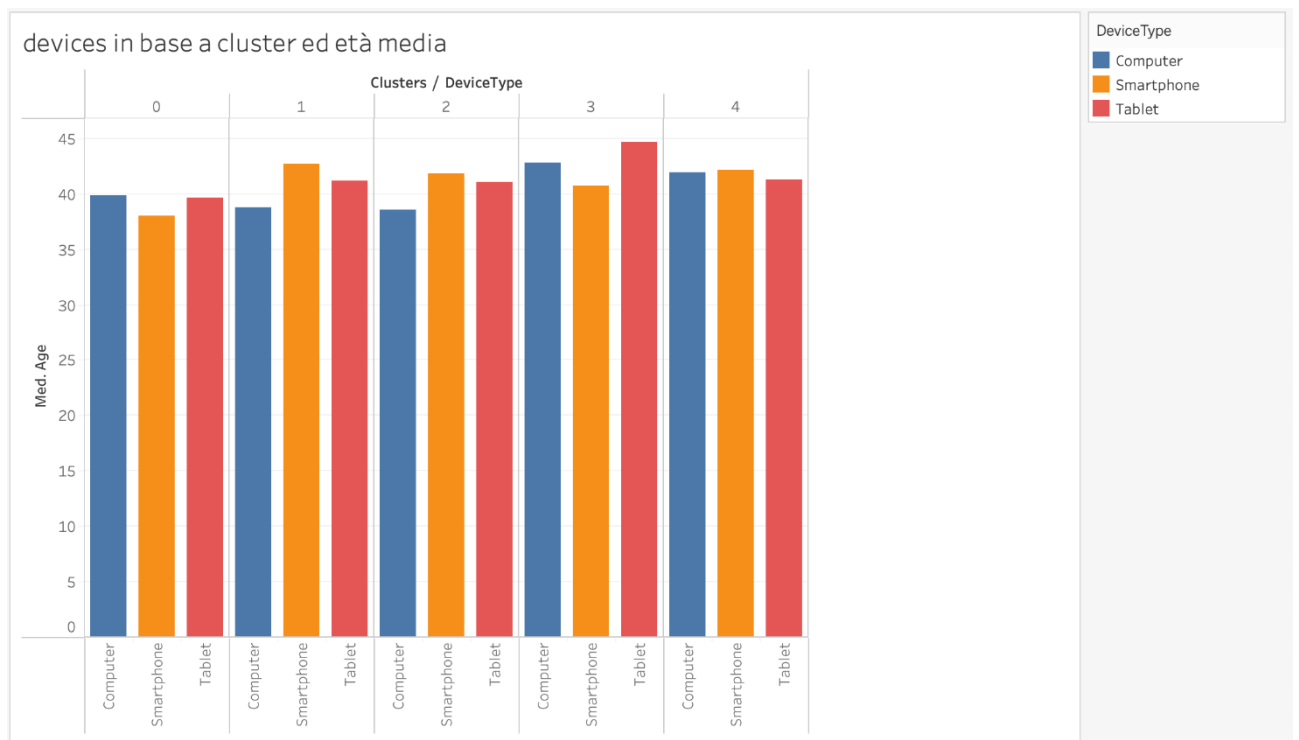


distribuzione età in base al cluster

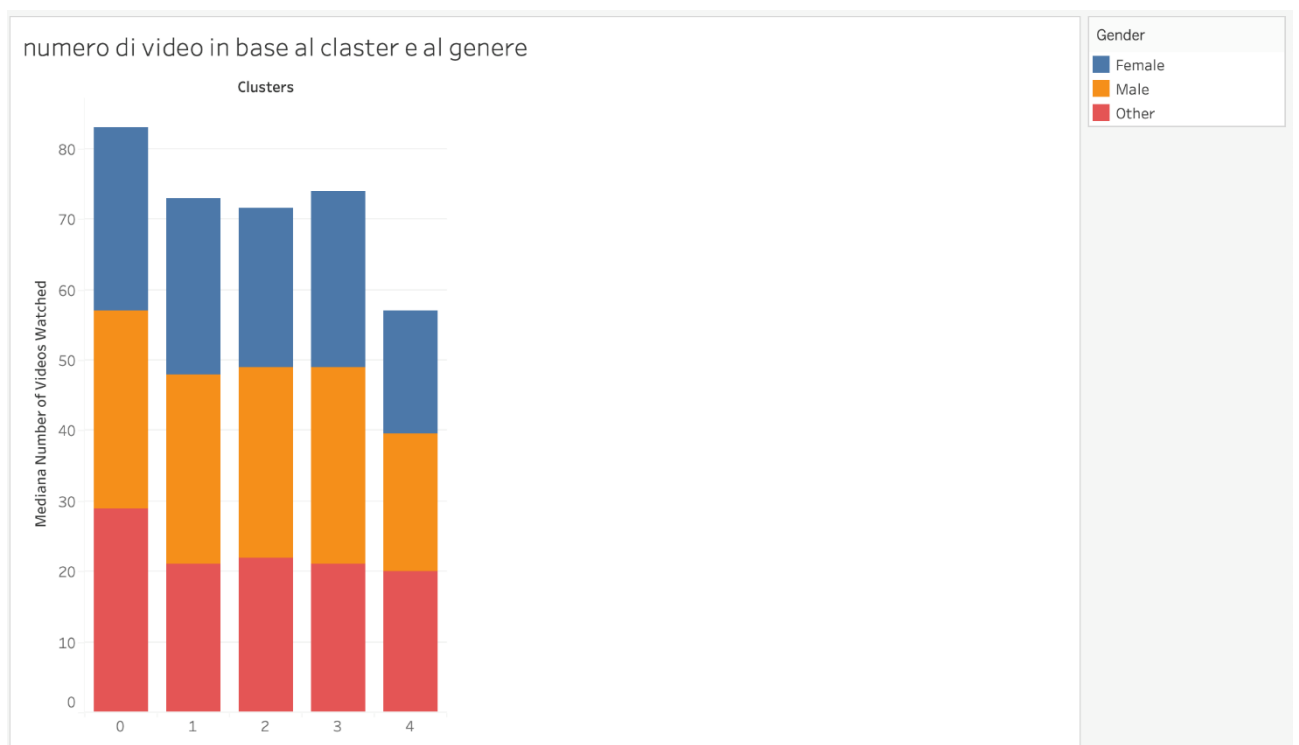




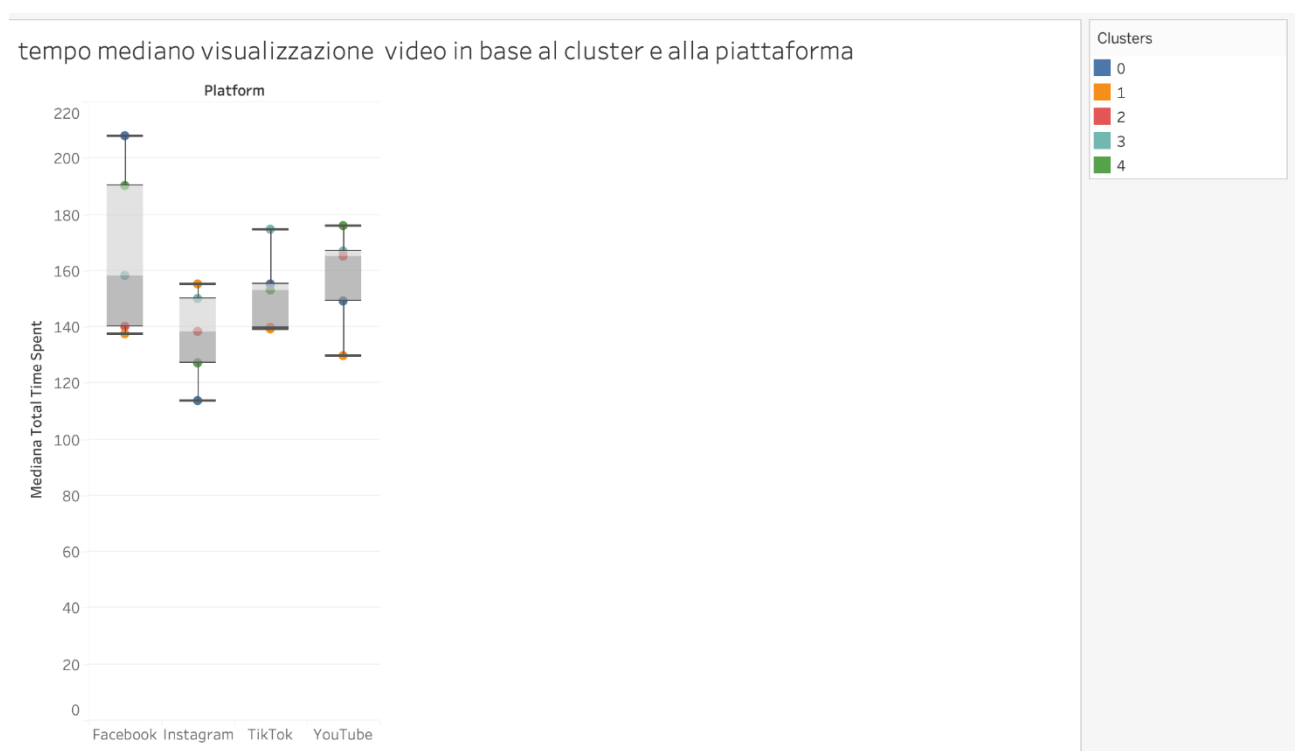
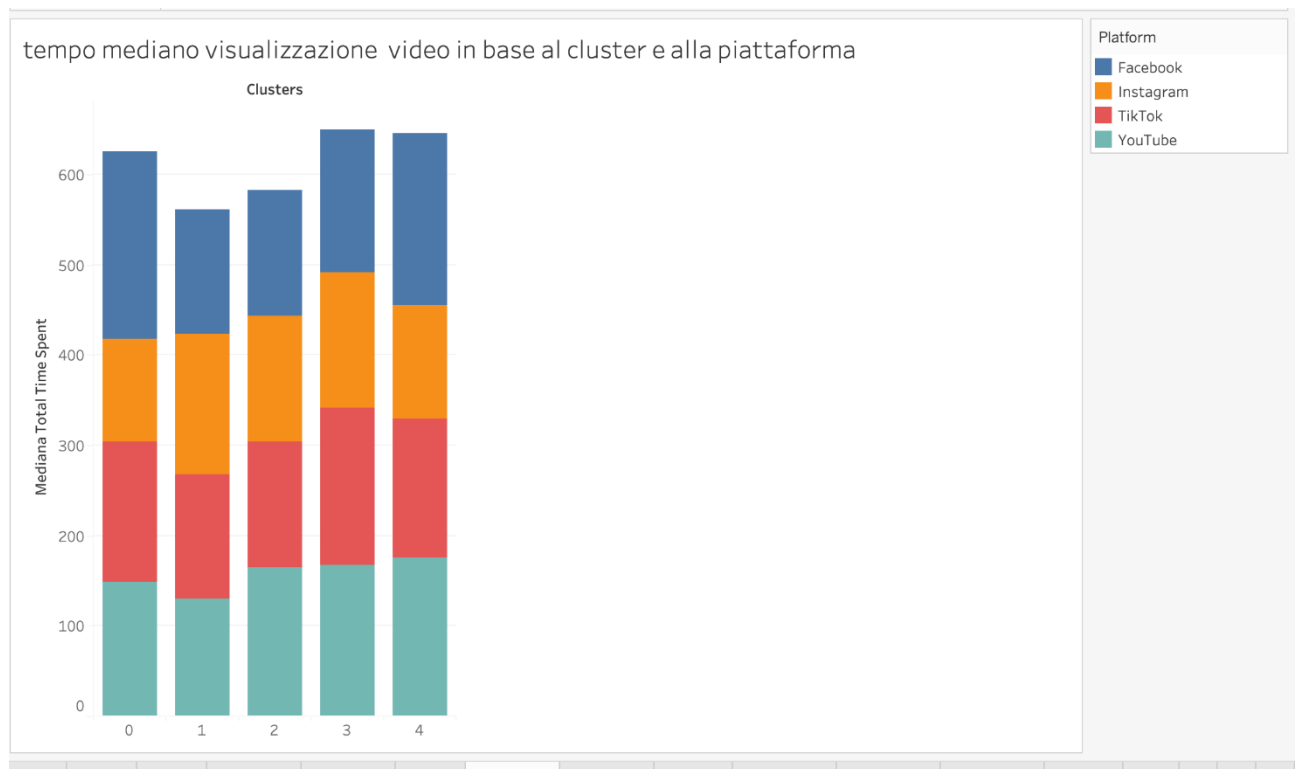
- Type of device based on user age



- Number of videos based on sex and cluster

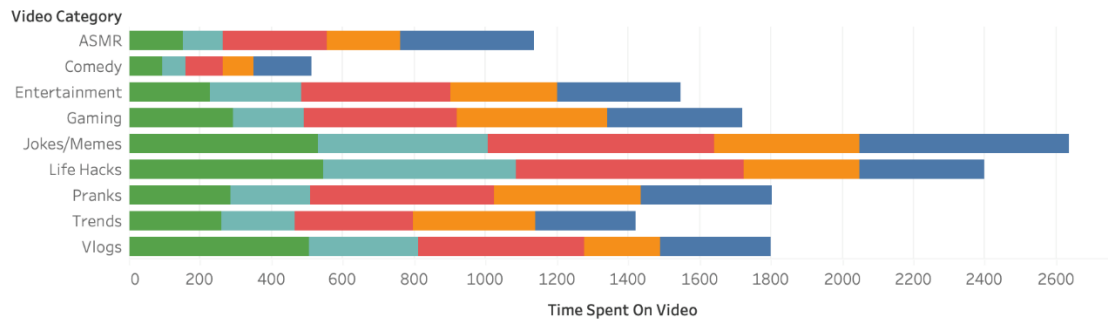


- Average time on platform

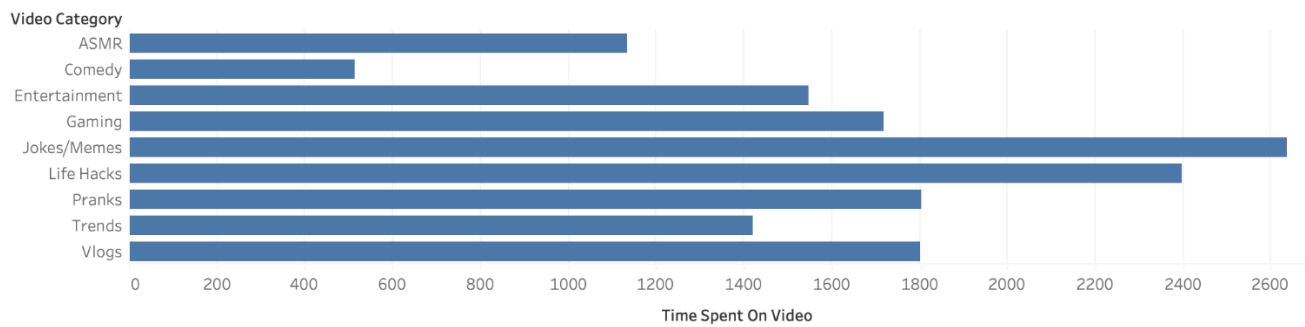


- Total time based on video category

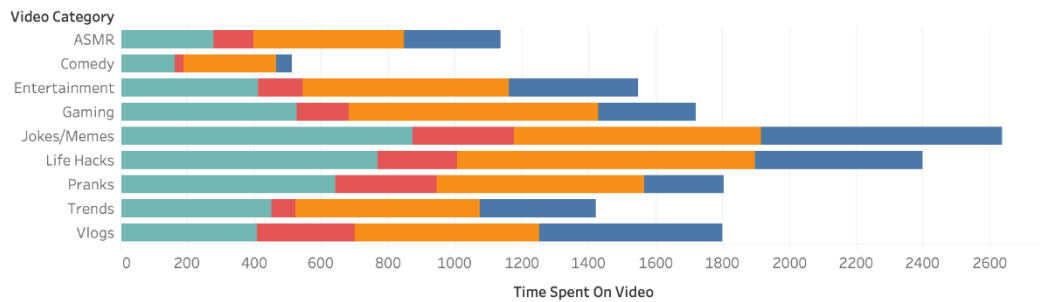
totale tempo in base alla categoria del video (individuo le categorie con più tendenza)



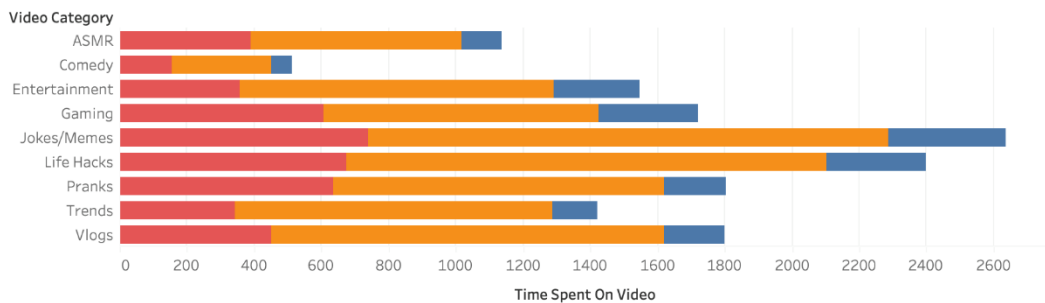
totale tempo in base alla categoria del video (individuo le categorie con più tendenza)



totale tempo in base alla categoria del video (individuo le categorie con più tendenza)



totale tempo in base alla categoria del video (individuo le categorie con più tendenza)



- Reason for use of social media

