# INT234

## [Predictive Analytics]

## CA-3

**SUBMITTED TO : Mr. Sandeep Kumar**

| | |
|---|---|
| **NAME:** | **Vinay Kumar** |
| **REGISTRATION NO.:** | **12104993** |
| **ROLL. NO.:** | **RK-21-BP-A-07** |
| **SECTION:** | **K-21-BP** |

**INTRODUCTION :**

This project focuses on analyzing user behavior patterns using machine learning techniques implemented in R. The goal is to categorize and segment users based on their behavior characteristics to gain insights for applications such as personalized marketing and targeted recommendations

_____

**About the Project :**

The project uses a dataset called "user_behavior_dataset.csv" to explore and classify user behaviors. By leveraging both supervised and unsupervised learning methods, we aim to understand user segments and behavior classes within the data. The project examines the performance of different algorithms, assessing their ability to classify users and form meaningful clusters.

_____

**Data Used and Algorithms Applied :**

The "user_behavior_dataset.csv" contains information on user attributes, including device model, operating system, and demographic details, which form the basis for analysis. Three algorithms are used:
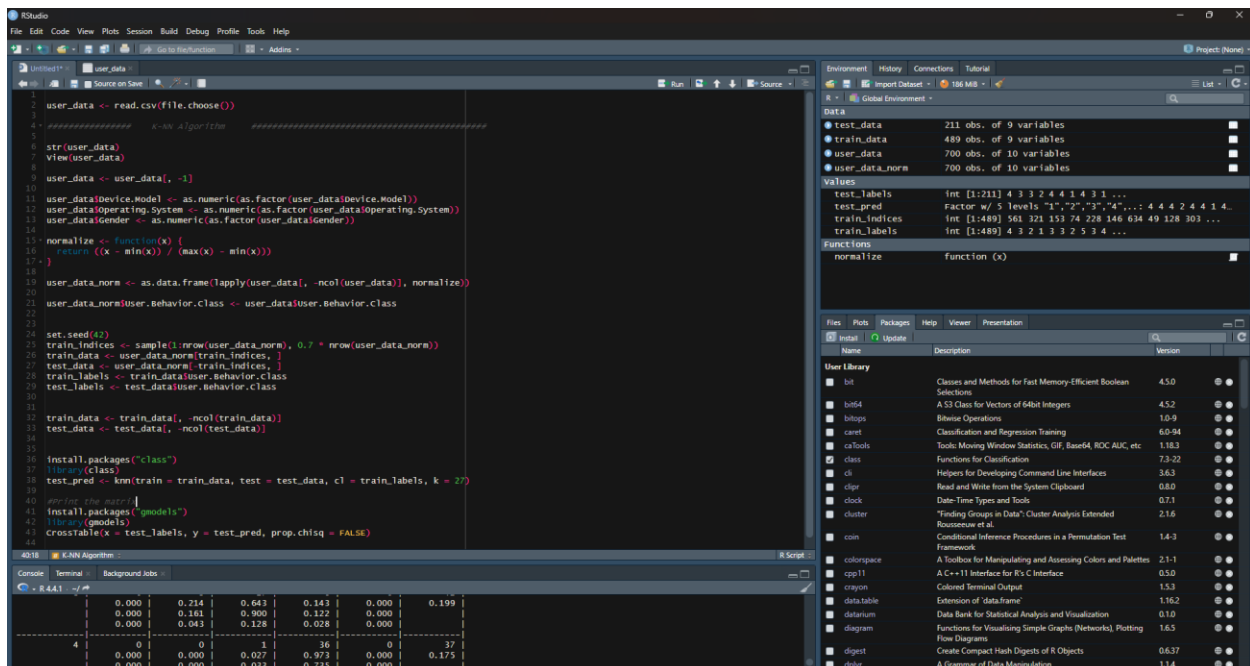
1. **K-Nearest Neighbors (K-NN):** A supervised learning algorithm that classifies users into predefined behavior classes by comparing each user to its nearest neighbors in the dataset.

2. **K-Means Clustering:** An unsupervised learning algorithm that groups users into clusters based on similarity, without any pre-labeled categories, to discover natural segments.

3. **Hierarchical Clustering:** Another unsupervised technique that builds a tree-like structure of user groups based on similarity, providing a visual interpretation of user segments at different levels of detail.

_____

**Performance Comparison :**

The performance of each algorithm is evaluated based on its effectiveness in identifying distinct user behaviors. K-NN is assessed through classification accuracy, measuring how well it predicts user classes. K-Means and Hierarchical Clustering are compared based on cluster cohesion and interpretability, focusing on how consistently each method groups similar users. Through this comparison, we can identify the strengths and best applications of each approach, providing insights into their suitability for behavior analysis.

**Algorithms Code's**

**K-Nearest Neighbors (K-NN):**



This code performs K-Nearest Neighbors (K-NN) classification on a dataset (user_data) to predict user behavior based on features such as device model,

operating system, and gender. The dataset is first loaded and preprocessed by removing any non-informative columns (like an ID column) and converting categorical variables (e.g., Device.Model, Operating.System, and Gender) to numeric form. Then, all features except the target variable (User.Behavior.Class) are normalized to a 0-1 scale to ensure equal weighting in the K-NN algorithm. The data is split into training (70%) and testing (30%) sets, with a seed set for reproducibility. Using the class library's knn function, the model is trained with k=27 and tested on the test set. Finally, the gmodels package is used to generate a confusion matrix, allowing for an assessment of classification accuracy by comparing predicted and actual classes.

```
> CrossTable(x = test_labels, y = test_pred, prop.chisq = FALSE)

   Cell Contents
|-------------------------|
|                       N |
|           N / Row Total |
|           N / Col Total |
|         N / Table Total |
|-------------------------|

Total Observations in Table:  211

             | test_pred
 test_labels |        1 |        2 |        3 |        4 |        5 | Row Total |
-------------|----------|----------|----------|----------|----------|-----------|
           1 |       31 |       11 |        0 |        0 |        0 |        42 |
             |    0.738 |    0.262 |    0.000 |    0.000 |    0.000 |     0.199 |
             |    0.912 |    0.196 |    0.000 |    0.000 |    0.000 |           |
             |    0.147 |    0.052 |    0.000 |    0.000 |    0.000 |           |
-------------|----------|----------|----------|----------|----------|-----------|
           2 |        3 |       36 |        2 |        0 |        0 |        41 |
             |    0.073 |    0.878 |    0.049 |    0.000 |    0.000 |     0.194 |
             |    0.088 |    0.643 |    0.067 |    0.000 |    0.000 |           |
             |    0.014 |    0.171 |    0.009 |    0.000 |    0.000 |           |
-------------|----------|----------|----------|----------|----------|-----------|
           3 |        0 |        9 |       27 |        6 |        0 |        42 |
             |    0.000 |    0.214 |    0.643 |    0.143 |    0.000 |     0.199 |
             |    0.000 |    0.161 |    0.900 |    0.122 |    0.000 |           |
             |    0.000 |    0.043 |    0.128 |    0.028 |    0.000 |           |
-------------|----------|----------|----------|----------|----------|-----------|
           4 |        0 |        0 |        1 |       36 |        0 |        37 |
             |    0.000 |    0.000 |    0.027 |    0.973 |    0.000 |     0.175 |
             |    0.000 |    0.000 |    0.033 |    0.735 |    0.000 |           |
             |    0.000 |    0.000 |    0.005 |    0.171 |    0.000 |           |
-------------|----------|----------|----------|----------|----------|-----------|
           5 |        0 |        0 |        0 |        7 |       42 |        49 |
             |    0.000 |    0.000 |    0.000 |    0.143 |    0.857 |     0.232 |
             |    0.000 |    0.000 |    0.000 |    0.143 |    1.000 |           |
             |    0.000 |    0.000 |    0.000 |    0.033 |    0.199 |           |
-------------|----------|----------|----------|----------|----------|-----------|
Column Total |       34 |       56 |       30 |       49 |       42 |       211 |
             |    0.161 |    0.265 |    0.142 |    0.232 |    0.199 |           |
-------------|----------|----------|----------|----------|----------|-----------|
```
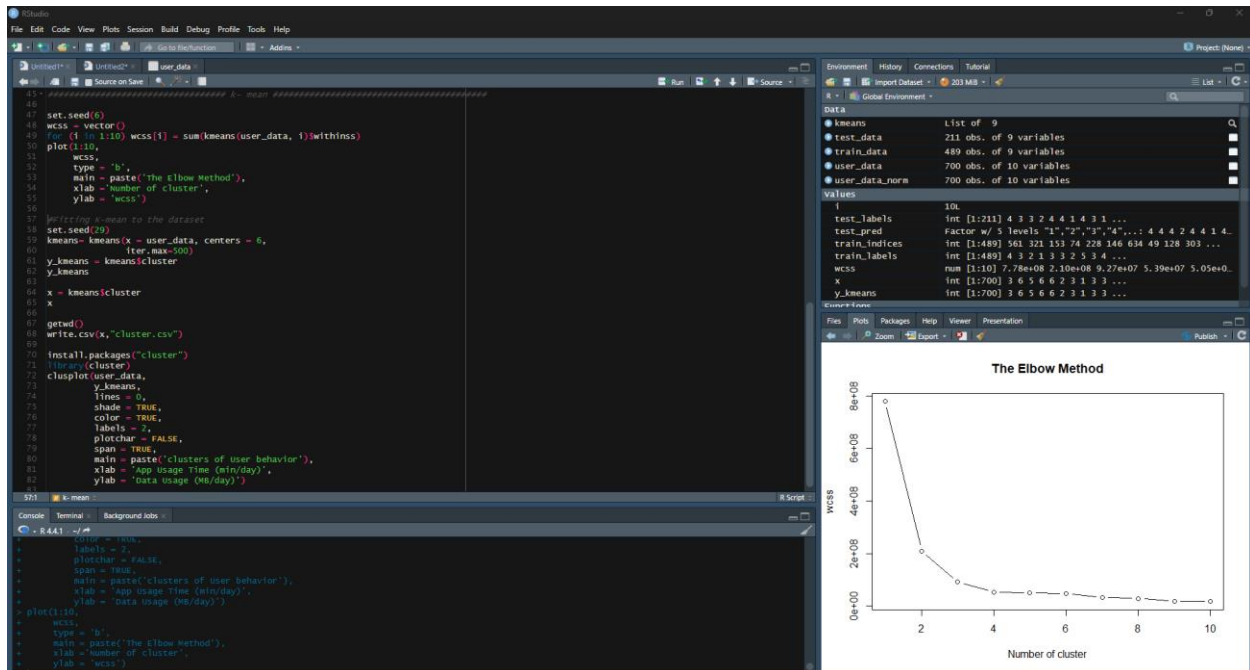
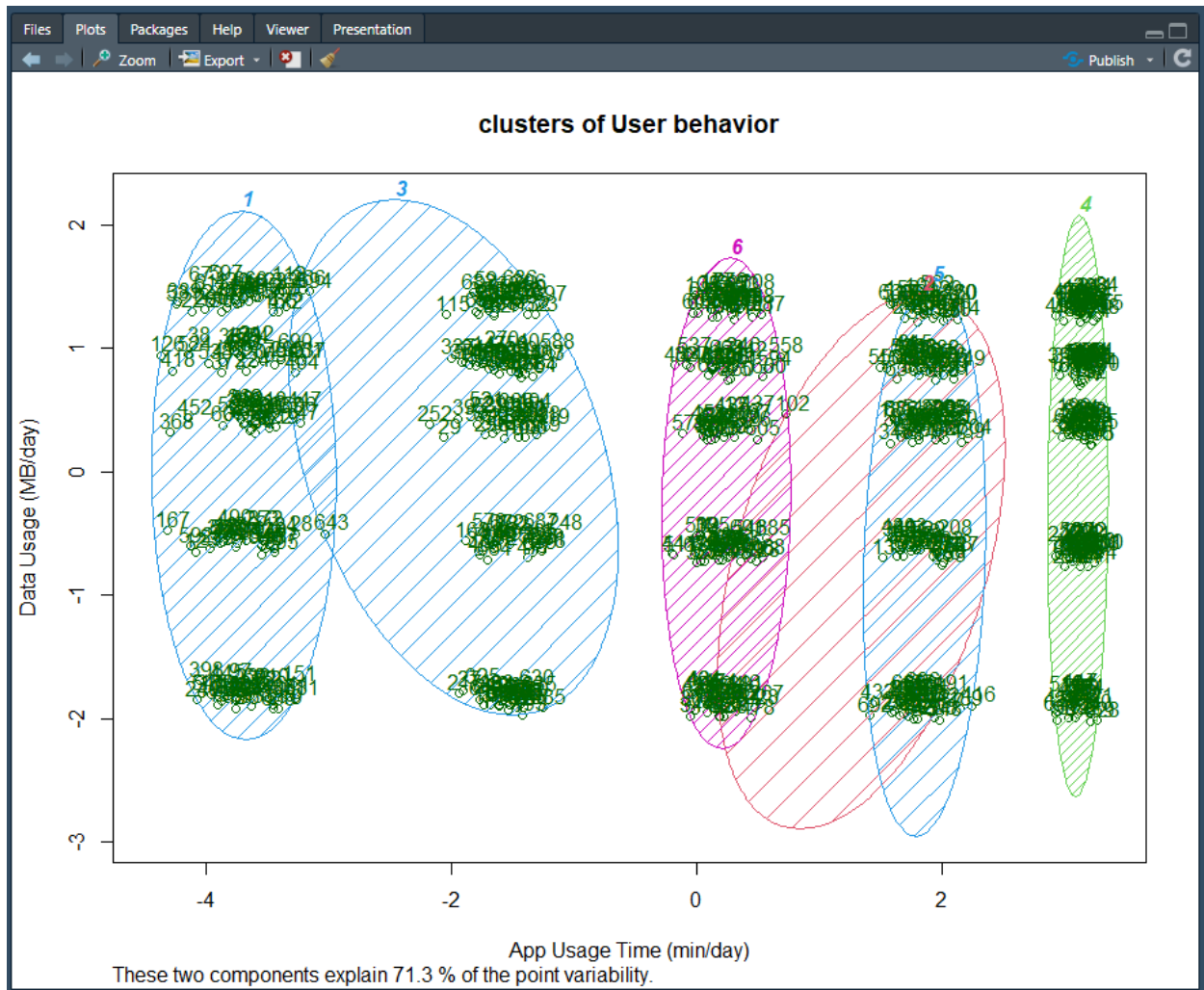Here we have the Confusing Matrix For our data set  with Prediction.

# K-Means Clustering :



This code uses K-Means clustering in R to analyze user behavior patterns in user_data. First, the Elbow Method is applied to find the optimal number of clusters by plotting the within-cluster sum of squares (WCSS) for clusters ranging from 1 to 10. The "elbow" in this plot suggests the optimal cluster count, where adding more clusters provides diminishing returns in WCSS reduction. With the optimal number identified (6 clusters in this case), the K-Means algorithm is run with centers = 6 and a maximum of 500 iterations to generate clusters. The resulting cluster assignments are saved to a CSV file for further analysis.
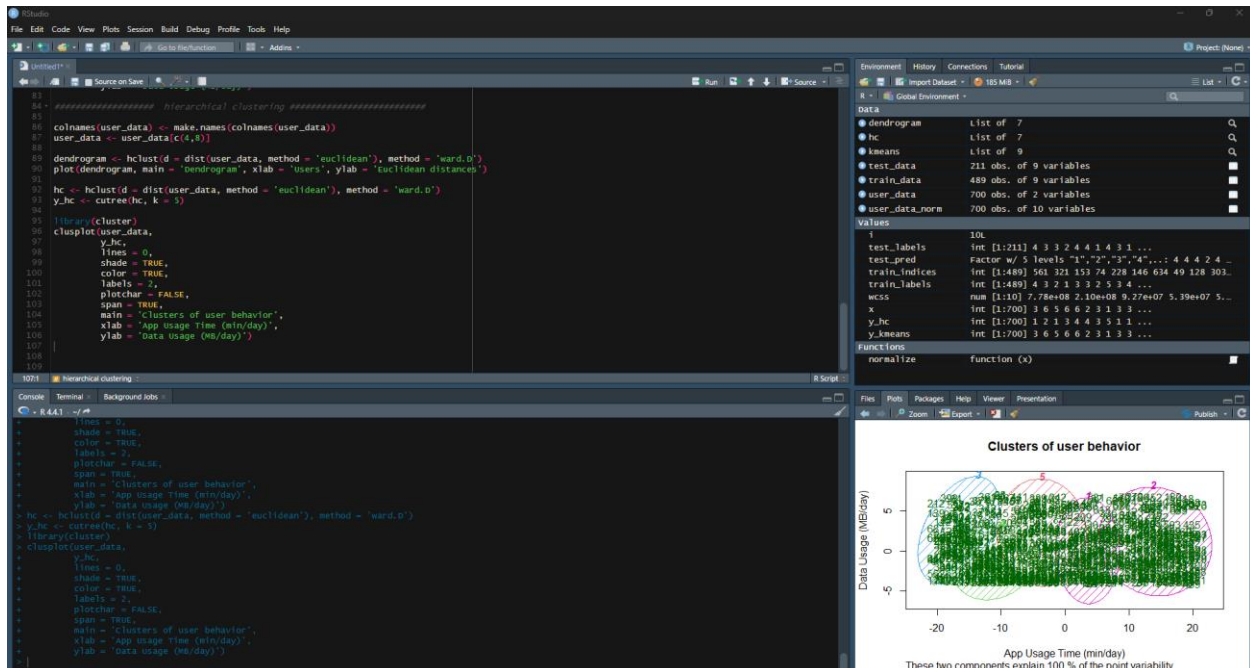
Finally, the clusplot function visualizes the clusters with shading and color differentiation, providing a clear, 2D view of the cluster distribution. This plot helps in understanding the cohesion and separation between clusters based on two attributes, labeled here as "App Usage Time" and "Data Usage." Overall, this workflow identifies distinct user behavior segments, aiding in interpreting and analyzing clustering effectiveness within the dataset.
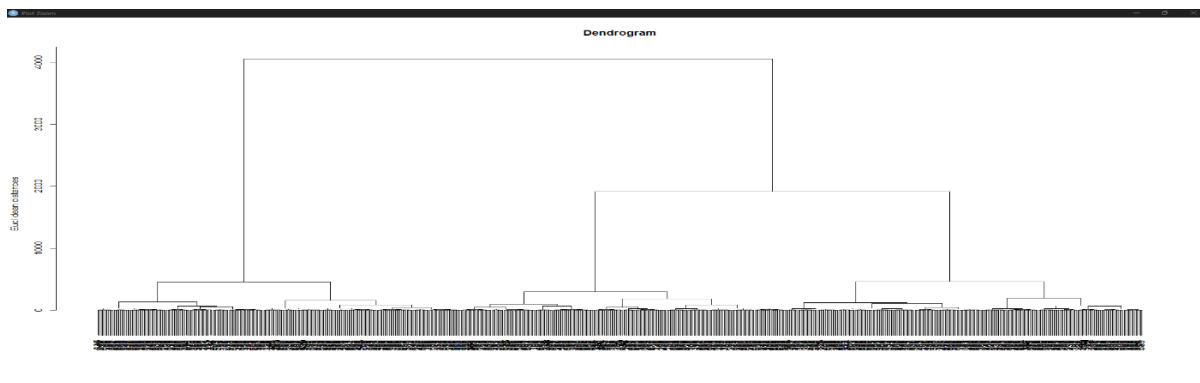


This plot shows the clusters of user behavior, segmented by the K-Means algorithm based on "App Usage Time (min/day)" and "Data Usage (MB/day)." Each color and shaded ellipse represents a distinct cluster, helping to visualize the separation and distribution of user groups.
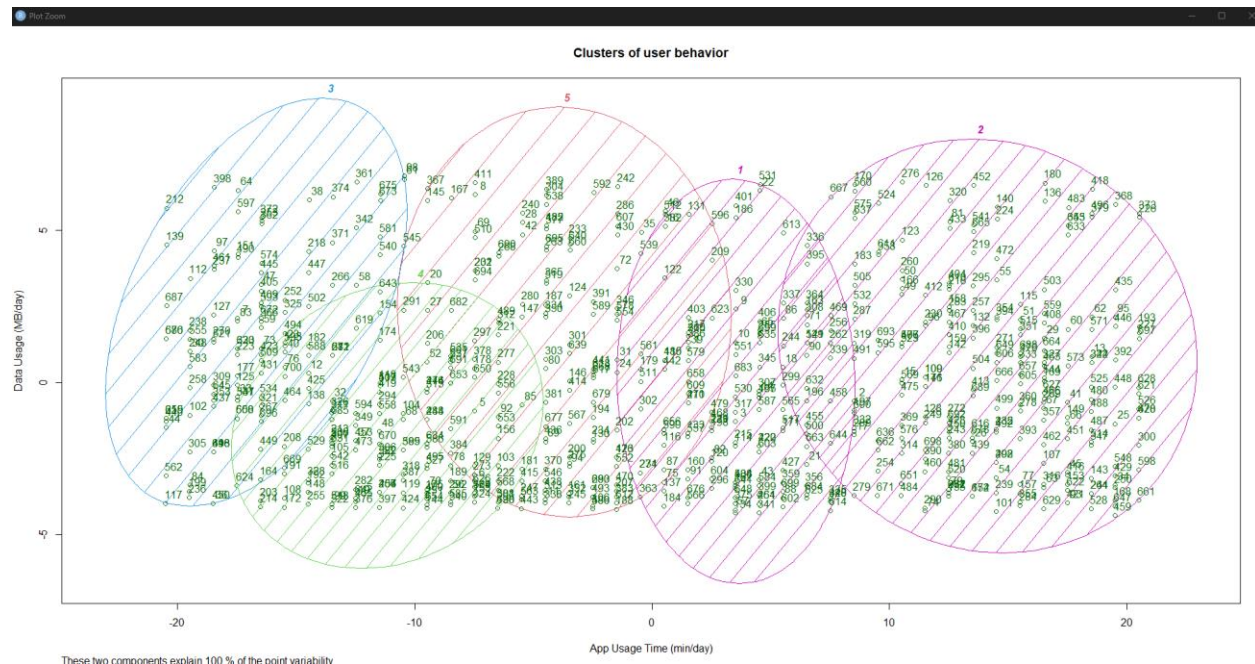
## Hierarchical Clustering:



This code performs **Hierarchical Clustering** on the user_data dataset in R to group users based on selected attributes, followed by visualization with a dendrogram and cluster plot. First, it ensures all column names are syntactically valid and selects two key features (e.g., "App Usage Time" and "Data Usage") for analysis. Using **Euclidean distance** and **Ward's method**, it computes hierarchical clusters that minimize within-cluster variance. The dendrogram displays hierarchical relationships and helps determine the optimal number of clusters, with larger merges indicating greater dissimilarity. After determining the number of clusters (k), the data is segmented accordingly, and the results are visualized in a 2D cluster plot using clusplot(). This plot shows clusters as distinct, color-coded regions, allowing for easy interpretation of how users are grouped based on the selected behavioral features.

This dendrogram represents the hierarchical clustering of users based on Euclidean distances. Two main clusters are visible, with higher distances indicating greater dissimilarity between the groups.



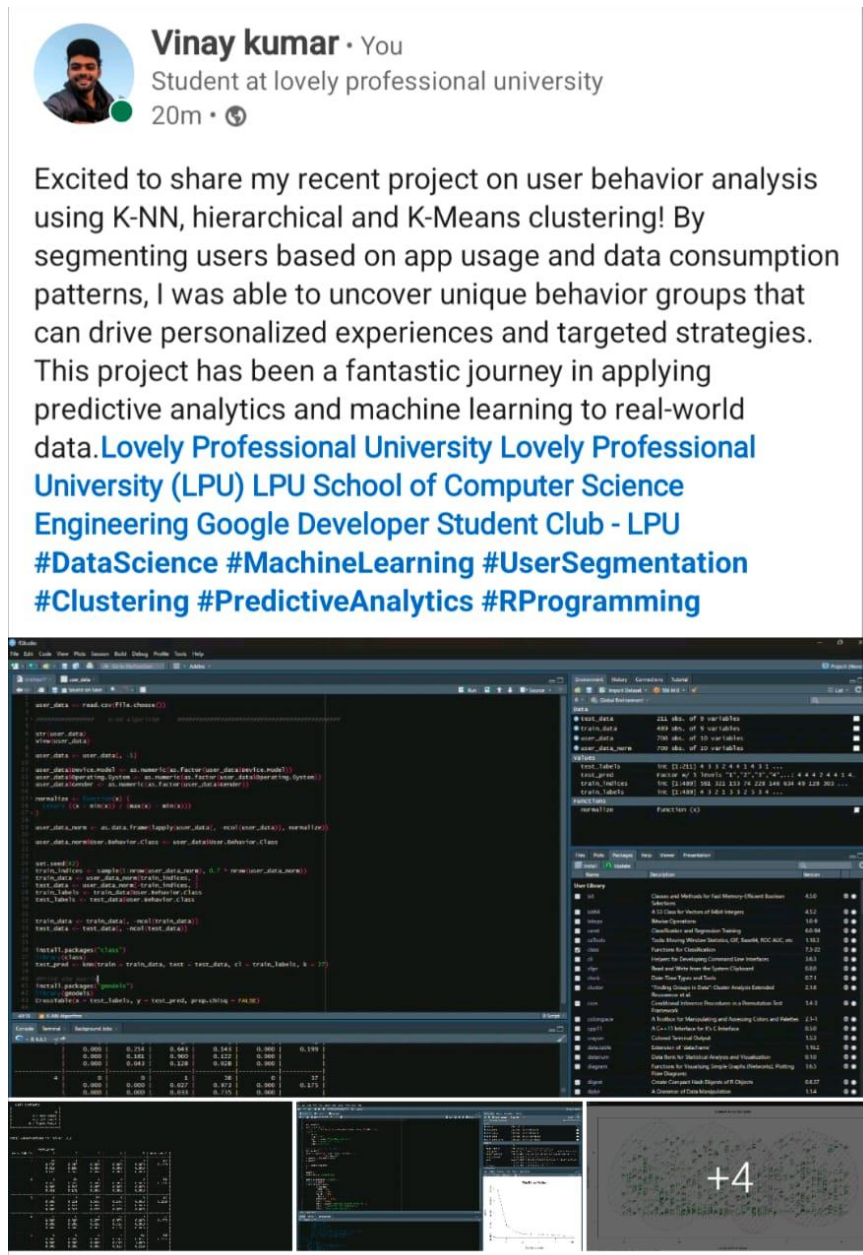These two components explain 100 % of the point variability

This cluster plot visualizes user behavior by grouping users based on "App Usage Time" (x-axis) and "Data Usage" (y-axis). There are five distinct clusters, each represented by a shaded ellipse with a unique color for clear separation. Cluster numbers are labeled at the center of each group. Data points within each cluster are displayed as green numbers, indicating individual users.

The clusters are structured as follows:

1. Cluster 3 (blue) contains users with low app usage time and data usage.

2. Cluster 5 (red) includes users with moderate app usage and data usage levels.

3. Cluster 1 (purple) has users with higher data usage and slightly higher app usage.

4. Cluster 2 (magenta) includes users with the highest app usage and moderately high data usage.

5. Cluster 4 (green) represents a middle range of app usage and data usage.

# LINKEDIN POST



## LINK :

https://www.linkedin.com/posts/vinay-kumar-9b27411b9_datascience-machinelearning-usersegmentation-activity-7259306141975797760-CPgJ?utm_source=screenshot_social_share&utm_medium=android_app&utm_campaign=copy_link