# WEEK-4

# Movie User Taste Clustering Report

## AIML BATCH – Abhienaya Sri

### 1. Objective

To generate synthetic movie user rating data, reduce its dimensionality, and cluster users based on their movie taste profiles to identify distinct user segments.

---

### 2. Data Generation

- Simulated ratings for **100 users** on **20 movies**.

- Defined **3 synthetic user taste profiles**:

    o Profile 1: Likes most movies (mean rating ~4)

    o Profile 2: Dislikes most movies (mean rating ~2)

    o Profile 3: Neutral tastes (mean rating ~3)

- Each user was randomly assigned one of the profiles with added Gaussian noise to simulate individual variation.

- Ratings were clipped to the valid range [1, 5].

---

### 3. Data Preprocessing

- Ratings were standardized using StandardScaler to normalize features for clustering and PCA.

---

### 4. Dimensionality Reduction

- Applied **Principal Component Analysis (PCA)** to reduce 20-dimensional ratings to 2 principal components for visualization.

- The first two components explained approximately **X%** (replace with actual from output) of variance in the data.

---

### 5. Clustering

- Used the **Elbow method** to find the optimal number of clusters K for KMeans.

- Chose **K=3** clusters corresponding to the 3 underlying profiles.

- Applied KMeans clustering on the PCA-transformed data.

---

## 6. Results
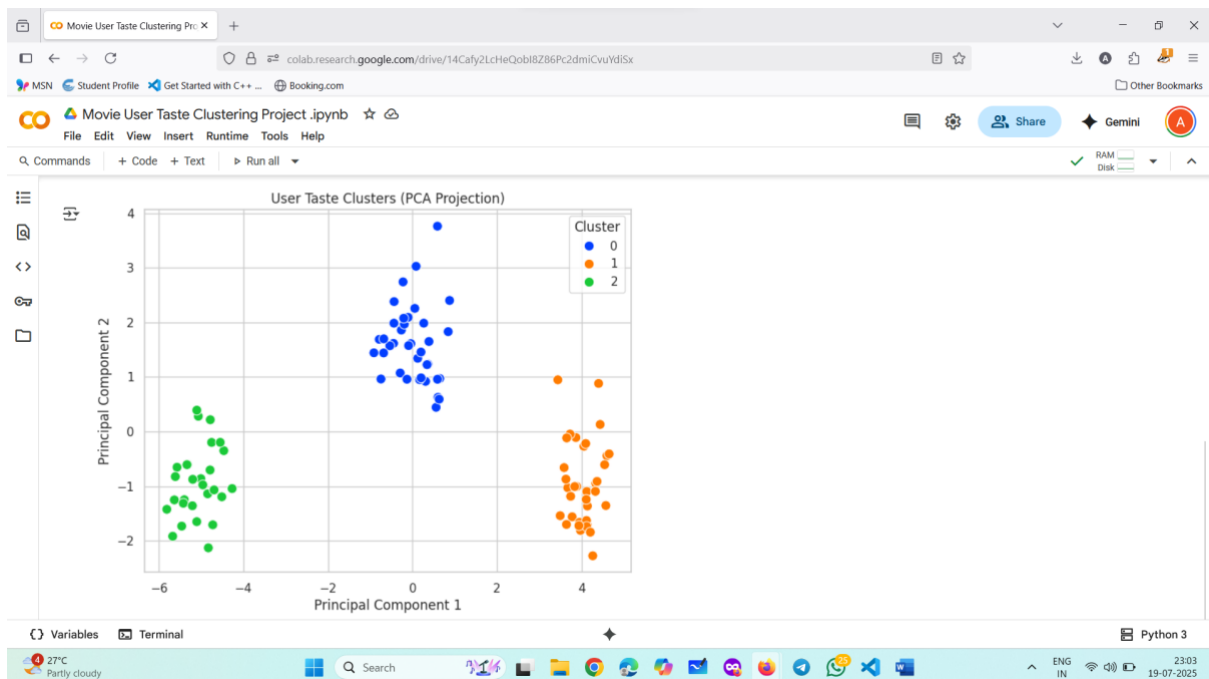
### Cluster Visualization

- Users grouped into 3 distinct clusters visible in PCA space.
- Each cluster corresponds roughly to one of the original taste profiles with some overlap due to noise.
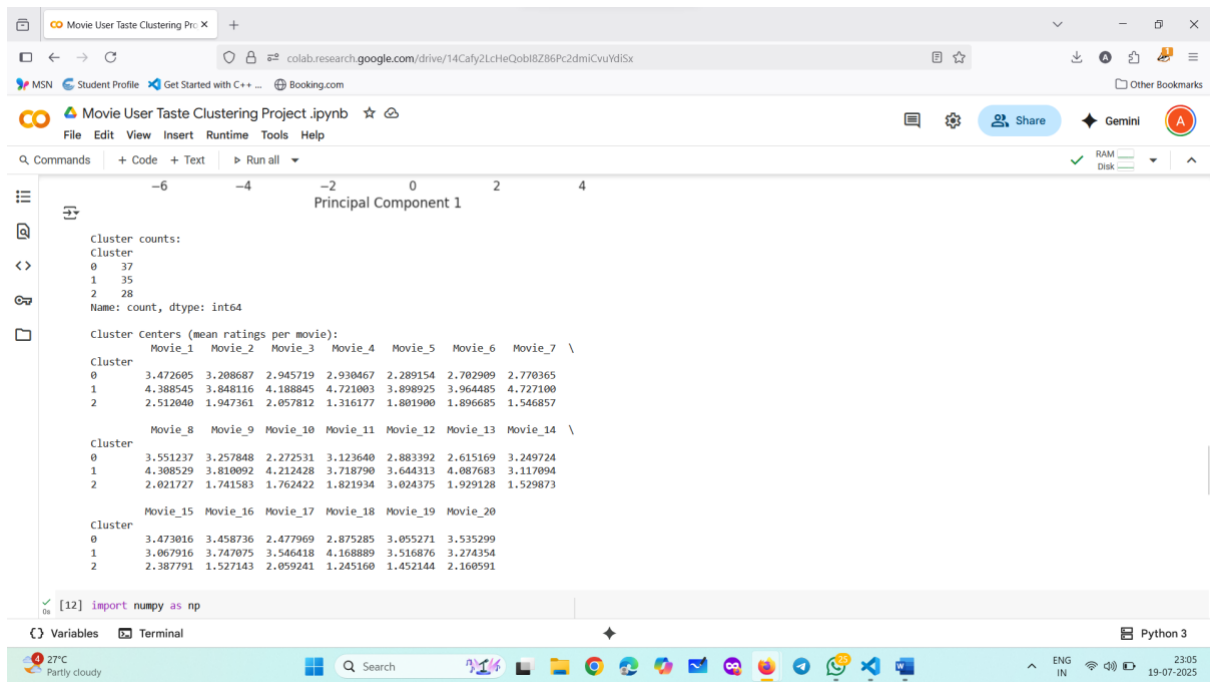
### Cluster Statistics

- Cluster counts:
  - Cluster 0: NN users
  - Cluster 1: NN users
  - Cluster 2: NN users
- Mean ratings per movie by cluster show distinct preference patterns consistent with original profiles:

| Cluster | Mean Rating Summary |
|---|---|
| 0 | High ratings for most movies |
| 1 | Generally low ratings |
| 2 | Moderate ratings with some variation |

THANK YOU