

Objective: To leverage the previously engineered "Taste DNA" feature to perform an unsupervised segmentation of the user base. The goal was to discover meaningful groups of users with similar taste profiles ("Taste Tribes") to enable a powerful, personalized collaborative filtering strategy for recommendations.

Executive Summary: This final modeling stage successfully transitioned from feature engineering to user segmentation. The core of this procedure was the application of the **HDBSCAN** clustering algorithm to the 32-dimensional "Taste DNA" vectors that were created by our Autoencoder. A critical pre-processing step involved aggregating the interaction-level embeddings to create a single, average taste profile for each unique user.

The clustering algorithm successfully identified **48 distinct "Taste Tribes"** of varying sizes. A key finding was the identification of one massive "mainstream" cluster and a significant number of users (approximately 47%) whose tastes were too unique to be classified, labeled as "noise." The final output is a file, `user_clusters.csv`, which maps each `user_id` to a specific cluster ID. This artifact is the final, actionable output of our modeling pipeline, directly enabling the multi-strategy recommendation engine.

Methodology: A Step-by-Step Breakdown

The clustering procedure was executed in two primary steps, as detailed in the `perform_clustering()` script.

Step 1: Data Aggregation - From Interactions to Profiles

- **What:** The script first loaded the `data_with_embeddings.csv` file, which contained a "Taste DNA" vector for each of the 200,000 user-restaurant interactions. To cluster the *users*, we needed a single, representative taste profile for each of the 100,000 unique users.
- **The Code's Logic (`groupby('user_id').mean()`):** The script grouped the entire dataset by `user_id` and calculated the **mean** of the 32 `taste_dna` columns for each user.
- **Why this is the correct approach:** A user's true taste profile is an aggregate of all their interactions. By averaging their "Taste DNA" vectors across the two restaurants they were paired with, we create a more stable and holistic representation of their overall preferences. This resulted in a clean, 100,000-row DataFrame where each user was represented by a single, powerful 32-dimensional vector.

Step 2: Unsupervised Clustering with HDBSCAN

- **What:** The HDBSCAN algorithm was applied to the 100,000 aggregated user taste profiles.
- **The Choice of Algorithm (HDBSCAN):** As your sparring partner, it's important to state why this choice was strategically superior to a simpler algorithm like K-Means.
 1. **No Need to Guess k:** HDBSCAN does not require us to specify the number of clusters beforehand. It discovers the "natural" number of clusters present in the data, which is more intellectually honest.

2. **Ability to Identify "Noise":** This is the most critical feature. HDBSCAN is a density-based algorithm. It identifies dense regions of users with similar tastes as clusters. Crucially, it labels users in sparse, low-density regions as "noise" (Cluster ID -1). This is a realistic acknowledgment that not every user fits neatly into a well-defined group. K-Means, by contrast, would have incorrectly forced every single user into a cluster, even if they didn't belong.
- **The Code's Logic (`hdbscan.HDBSCAN(...).fit_predict(...)`):** The script initialized HDBSCAN with a `min_cluster_size` of 100 (meaning it would not recognize any "tribe" with fewer than 100 users) and then fit it to the user profile data. The output was a cluster label for each of the 100,000 users.

Results and Interpretation

The output of the clustering process was highly informative and revealed a classic real-world distribution of user preferences.

- **Quantitative Results:**
 - **Number of "Taste Tribes" Discovered:** 48
 - **Number of "Noise" Users (Cluster -1):** 47,459
- **Qualitative Analysis:**
 - **The "Noise" Points are a Feature, Not a Bug:** The fact that nearly 47% of users were unclassifiable is a sign of the algorithm's strength. It correctly identified that a large portion of the user base has eclectic, inconsistent, or highly unique tastes that do not conform to any specific tribe. This prevents us from making inaccurate, stereotyped recommendations to them.
 - **The "Mainstream" Cluster:** The results showed one enormous cluster (#46 with 43,570 users). This group represents the "average" or "mainstream" user, whose tastes are not strongly defined in any one direction.
 - **The "Niche Taste Tribes":** The other 47 smaller clusters represent the true, high-value user segments. These are groups of users with very specific and consistent preferences (e.g., the ramen lovers, the biryani fanatics). These are the users for whom personalized, collaborative filtering will be most effective.

File Output

A	B
user_id	cluster_id
1	-1
2	46
3	-1
4	46
5	46
6	44
7	-1
8	-1
9	-1
10	-1
11	-1
12	46
13	46
14	46
15	-1
16	46
17	-1
18	-1
19	-1
20	46
21	46
22	-1
23	-1
24	-1
25	46

Conclusion: The clustering stage was a success. It effectively translated the abstract "Taste DNA" embeddings into a concrete, actionable user segmentation scheme. The final output, `user_clusters.csv`, is the key that unlocks our multi-strategy recommendation logic, allowing us to provide hyper-personalized recommendations to users in niche taste tribes while serving robust, popular recommendations to mainstream and unclassified users.