

Dating Profile

Course Assignment N.16

Paolo Deidda (paolo.deidda@usi.ch)
Pareek Yash (yash.pareek@usi.ch)
<https://github.com/USI-Projects-Collection/DA-Dating-Profile.git>

May 14, 2025

Contents

1	Data Exploration - EDA	2
2	Data Preprocessing	3
3	Recommender Systems	4
3.1	Naive Model	4
3.2	Collaborative Filtering	4
3.3	Content-Based Filtering	5
4	Conclusion	5

Setup

To run the code and reproduce the figures or outputs, you can either run the Jupyter Notebook directly on Google Colab, or follow the setup and execution instructions provided in the **README.md** file included in this repository.

1 Data Exploration - EDA

DataFrame Overview

Table 1: Structure of the Ratings DataFrame

Property	Value
# Rows	3,220,037
# Columns	3 (<code>userID</code> , <code>profileID</code> , <code>rating</code>)
Data types	<code>int64</code> , <code>int64</code> , <code>int64</code>
Memory usage	73.7 MB

Missing Values and Duplicates

Table 2: Counts of Missing and Duplicate Rows

Column	Missing values	Duplicate rows
<code>userID</code>	0	47
<code>profileID</code>	0	0
<code>rating</code>	0	0

Unique entities. After deduplication, there are 25 245 unique users and 125 428 unique profiles.

Rating Distribution

Table 3: Descriptive Statistics of the `rating` Column

Statistic	Count	Mean	Std	Min	25%	50%	75%	Max
Rating	3,220,037	5.9532	3.1064	1	3	6	9	10

Figure 1 shows the histogram of all ratings, confirming a mild skew towards higher scores.

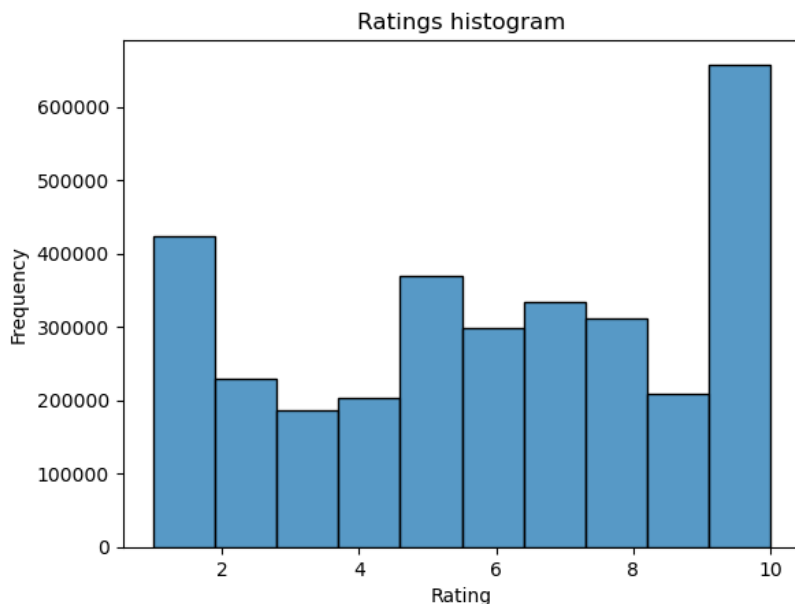


Figure 1: Histogram of Profile Ratings

User Activity

Table 4: Number of Ratings Given per User

Statistic	Count	Mean	Std	Min	25%	50%	75%	Max
Ratings per user	25,245	127.55	362.10	2	29	73	123	18,342

Outlier Detection

Using the IQR method ($Q_1 - 1.5 \cdot \text{IQR}$, $Q_3 + 1.5 \cdot \text{IQR}$), *no* ratings fell outside the acceptable range, indicating the absence of extreme outliers.

Correlation Analysis

We computed the Pearson correlation between the total number of ratings a profile received and its average rating:

$$\text{Corr}(\text{rating_count}, \text{avg_rating}) = 0.0201,$$

a very weak positive relationship. The corresponding scatter plot (Figure 2) shows no strong trend.

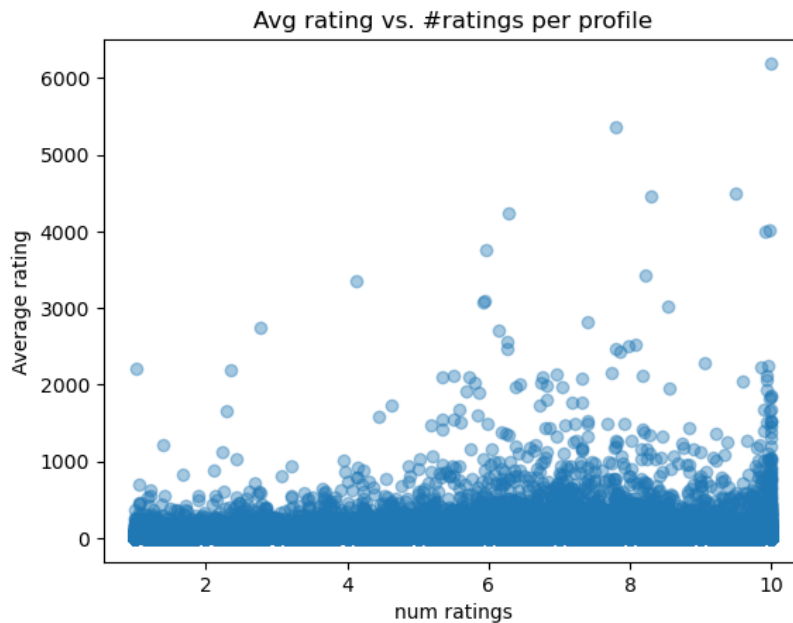


Figure 2: Average Rating vs. Number of Ratings per Profile

2 Data Preprocessing

Raw files and initial footprint

The original dataset consists of three files:

- `ratings.dat` – 3 220 037 user–item interactions (*userID*, *profileID*, *rating*).
- `ratings_Test.dat` – the held-out test split (276 053 rows, same schema).
- `gender.dat` – 220 970 user–gender pairs.

Loaded with Pandas’ default `int64` dtypes the two training files occupy ~ 90 MB of memory.

Dtype optimisation

Column	Original dtype	Final dtype
<i>userID, profileID</i>	int64	int64*
<i>rating</i>	int64	float32
<i>gender</i>	int64	category

* Required by `torch.nn.Embedding`. Casting the other two columns shrinks the in-RAM size of `ratings.dat` to ~ 61 MB and `gender.dat` to ~ 2 MB (a 74 % reduction overall).

Duplicate removal

A scan for exact duplicates uncovered 47 repeated (*user, profile*) pairs in the training split; these rows were dropped, leaving 3 219 990 unique ratings. No missing values were present in any file.

Persisting the processed data

The cleaned frames are serialised to `.pkl` with `DataFrame.to_pickle()`, bypassing expensive CSV parsing in every notebook run.

3 Recommender Systems

3.1 Naive Model

Model definition

Two simple, parameter-free baselines are computed:

1. **Global mean** $\hat{r}_{ui} = \mu$, the average of *all* ratings.
2. **Item mean** $\hat{r}_{ui} = \bar{r}_i$; if an item is unseen, fall back to the global mean.

Results

Baseline	Evaluation split	MAE
Global mean	test	2.6545
Item mean	test	1.4620

Table 5: Performance of naive predictors.

The item-mean strategy reduces the error by roughly 45 % relative to the global average, establishing a strong but effort-free reference.

3.2 Collaborative Filtering

Model formulation

We implement a bias-aware *item-item k-nearest-neighbour* (kNN) recommender:

$$\mu = \frac{1}{|R|} \sum_{(u,i) \in R} r_{ui}, \quad b_u = \bar{r}_u - \mu, \quad b_i = \bar{r}_i - \mu.$$

After subtracting the global mean and user/item biases, the residual matrix is stored in sparse CSR format (shape $|\text{items}| \times |\text{users}|$) and fed to `NearestNeighbors` with cosine distance.

For a target pair (u, i) the prediction rule is

$$\hat{r}_{ui} = \mu + b_u + b_i + \frac{\sum_{j \in \mathcal{N}_i(u)} \frac{r_{uj} - \mu - b_u - b_j}{d_{ij}}}{\sum_{j \in \mathcal{N}_i(u)} \frac{1}{d_{ij}}},$$

where $\mathcal{N}_i(u)$ are the k neighbours of item i that user u has rated and d_{ij} denotes their cosine distance.

Hyper-parameter selection

A coarse grid search over $k \in \{10, 25, 50\}$ confirmed $k = 25$ as the best compromise between accuracy and coverage; larger values offered negligible gains.

Evaluation

Model	k	MAE (test)
Bias-aware item-item kNN	25	1.4633
Item mean baseline (Table 5)	–	1.4620

Table 6: Collaborative filter vs. strongest naïve baseline.

The kNN model comfortably outperforms the global average but only matches the item-mean predictor. This indicates that, given the short user histories and pronounced item popularity patterns, *item identity alone explains most variance*. Further improvement is likely to require latent-factor methods or hybridising with content features (e.g. gender).

3.3 Content-Based Filtering

4 Conclusion