

1 Recommendation Steps

1. List set W_u of movies watched by user U_i , extract tag-genome matrix for W_u .
2. Calculate Silhouette scores for different possible cluster sizes, different sizes are from 2 to $|W_u|$.
3. Select $n_clusters$ such that it yields the highest Silhouette score.
4. Rank clusters by their sizes, i.e. number of movies in each cluster. Rank in such a fashion that highest size cluster gets the low rank. These ranks are used by the 'Novel Ranking' algorithm.
5. **N_Similar_movies** is a configurable parameter that controls the number of movies in Recommendation List **RL** that are most similar to the user profile. Select $N_Similar_movies$ from the dense clusters, using **Algorithm 1**. Dense clusters are the clusters \geq mean cluster size.
6. In case, above step returns less number of movies, the difference should be covered by the next movie selection **Algorithm 2**. Hence, recalculate the value of **N_Novel_Movies** as:

$$N_Novel_Movies = K - N_Similar_Movies \quad (1)$$

7. Find all possible **N_Novel_Movies** using the **Algorithm 2**.
8. Append **N_Similar_Movies** to the Recommendation List **RL**.
9. Choose **N_Novel_Movies** from the Novel Recommendation List **NRL**, and append them to the final Recommendation List **RL**.
10. Recommend the Recommendation List **RL**.

Algorithm 1 Most Similar Movies Selection Algorithm

- 1: $D \leftarrow dense_clusters$ ▷ D is a set of Dense Clusters, input argument
 - 2: $N \leftarrow N_Movies_Per_Dense_Cluster$ ▷ Input argument N
 - 3: **while** D has more clusters **do**
 - 4: $C \leftarrow next_cluster$ from D
 - 5: $W \leftarrow list$ of watched movies from C
 - 6: $similar_movies \leftarrow top$ similar movies ▷ Cosine similarity
 - 7: $RL1 \leftarrow$ append top N movies most similar to user's profile
 - 8: **return** **RL1**
-

Algorithm 2 Novel Re-Ranking Algorithm

```
1:  $W_u \leftarrow$  watched movies from all sparse clusters
2:  $R_{df} \leftarrow$  A dataframe object used for Ranking
3: while  $W_u$  has more movies do
4:    $W_i \leftarrow$  next movie
5:    $R_{df} \leftarrow$  append similar N movies
6:    $R_{df}[C_i] \leftarrow$  cluster score for  $W_i$ 
7:    $R_{df}[S_{wi}] \leftarrow$  similarity with  $W_i$ 
8:    $R_{df}[S_u] \leftarrow$  similarity of all movies to user profile
9:    $R_{df}[R_{wi}] \leftarrow$  rating of user i for watched movie  $W_i$ 
10:  $R_{df}[diversity] \leftarrow 1 - R_{df}[S_u]$ 
     $\triangleright$  Diversity = 1 - sim(user_profile)
11:  $R_{df}[rank(R_{cu})] \leftarrow$  Rank for column values  $R_{cu}$ 
     $\triangleright$  Dense rank as available in Pandas
12:  $R_{df}[rank(diversity)] \leftarrow$  Rank for column values  $diversity$ 
13:  $R_{df}[rank(S_u)] \leftarrow$  Rank for column values  $S_u$ 
14:  $R_{df}[rank(S_c)] \leftarrow$  Rank for column values  $S_c$ 
15:  $R_{df}[RN_c] \leftarrow$  Composite Rank using equation 2
16:  $NRL \leftarrow$  Sort movies in descending order based on  $R_{df}[RN_c]$  - composite rank.
17: return NRL  $\triangleright$  Sorted Novel Recommendation List
```

Components of the algorithm 2,

1. R_{wi} , is Rating given by the user to the watched movie.
2. S_u , is the Similarity to the user's profile, calculated using cosine similarity with the user's tag-genome based term vector.
3. $diversity$, is the diversity to the user profile, calculated as '1 - S_u '
4. C_i , is the score for the cluster which has a watched movie W_i .

This is the ranking equation:

$$\begin{aligned} R_{df}[RN_{comp}] = & (R_{wi_weight} * R_{df}[rank(R_{wi})]) + \\ & (diversity_weight * R_{df}[rank(diversity)]) + \\ & (S_u_weight * R_{df}[rank(S_u)]) + \\ & (C_i_weight * R_{df}[rank(C_i)]) \end{aligned} \quad (2)$$

The $R_{df}[RN_{comp}]$ is a composite rank for a movie.