

Faculté des Sciences de Sfax



Rapport de Projet

Recommandation Basée sur la Similarité
(Approche centrée item)

Par

Ahmedou Yahye Gleiguem Kheyri

étudiant en Master Recherche Informatique

25 février 2025

Table des matières

1	Introduction	3
2	Presentation Mathematique	3
2.1	Normalisation des Notes	3
2.2	Mesures de Similarite	3
3	Présentation des Données	3
4	ETAPE 1 : Chargement des Données	4
4.1	Code	4
4.2	Output	5
5	ETAPE 2 : Fusion des Données	5
5.1	Code	5
5.2	Output	6
6	ETAPE 3 : Création de la Matrice Utilisateur-Film	6
6.1	Code	6
6.2	Output	7
7	ETAPE 4 : Normalisation des Notes	7
7.1	Code	7
7.2	Output	8
8	ETAPE 5 : Calcul des Similarités	8
8.1	Code	8
8.2	Output	9
9	ETAPE 6 : Prédictions des Notes et Comparaison	9
9.1	Code	9
9.2	Output	11
10	ETAPE 7 : Analyse Finale	11
10.1	Code et Output d'Analyse Finale	11
10.2	Output	12
11	Comparaison et Perspectives d'Amélioration	12
11.1	Comparaison des Méthodes	12
11.2	Perspectives d'Amélioration	12
11.3	Conclusion	12

1 Introduction

Ce rapport présente l'analyse d'un système de recommandation basé sur la similarité entre films. Nous utilisons des évaluations de films ainsi que leurs caractéristiques pour prédire les notes manquantes. Le document est découpé en plusieurs ETAPES :

- Chargement des données.
- Fusion des données.
- Création de la matrice utilisateur-film.
- Normalisation des notes.
- Calcul des similarités (Pearson et Cosinus).
- Prédiction des notes et comparaison des méthodes.
- Analyse finale et conclusion.

2 Presentation Mathematique

2.1 Normalisation des Notes

Pour reduire le biais des utilisateurs, on centre les notes autour de la moyenne :

$$r_{norm}(u, i) = r(u, i) - \bar{r}_i \quad (1)$$

ou \bar{r}_i est la moyenne des notes pour l'item i .

2.2 Mesures de Similarite

- **Pearson** : Mesure la correlation lineaire entre les vecteurs centres

$$sim_p(A, B) = \frac{\sum(a_i - \bar{a})(b_i - \bar{b})}{\sqrt{\sum(a_i - \bar{a})^2} \sqrt{\sum(b_i - \bar{b})^2}} \quad (2)$$

- **Cosinus** : Mesure l'angle entre les vecteurs

$$sim_c(A, B) = \frac{A \cdot B}{\|A\| \|B\|} \quad (3)$$

3 Présentation des Données

Les fichiers utilisés sont :

movies-Rene (1).csv

```
movieId;title;genres
1;Item1;"Adventure|Animation|Children|Comedy|Fantasy"
2;Item2;"Adventure|Children|Fantasy"
3;Item3;"Comedy|Romance"
4;Item4;"Adventure|Children|Fantasy"
```

rating-Rene (1).csv

```
userId;movieId;rating;timestamp
1;1;;964982703
1;2;2;964981247
1;3;7;964981247
1;4;8;964981247
2;1;4;964982224
2;2;1;964982224
2;3;;964982224
2;4;7;964982224
3;1;3;964983815
3;2;8;964983815
3;3;;964983815
3;4;4;964983815
4;1;4;964982931
4;2;1;964982931
4;3;6;964982931
4;4;;964982931
```

4 ETAPE 1 : Chargement des Données

Objectif : Importation des evaluations utilisateurs et des metadonnees des films. Les donnees brutes contiennent des valeurs manquantes (NaN) qu'il faudra traiter.

4.1 Code

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.metrics.pairwise import cosine_similarity
4
```

```

5 # -----
6 # ETAPE 1 : Chargement des donn\ees
7 # -----
8 print("\n[ETAPE 1/7] Chargement des donn\ees...")
9 ratings = pd.read_csv('rating-Rene (1).csv', sep=';',
10                       na_values=[''])
11 movies = pd.read_csv('movies-Rene (1).csv', sep=';')
12
13 print(f"* {len(ratings)} \xe9valuations charg\ees")
14 print(f"* {len(movies)} films dans le catalogue")
15 print("Apercu des \xe9valuations :")
16 print(ratings.head())

```

Listing 1 – Chargement des données

4.2 Output

Output - Etape 1

[ETAPE 1/7] Chargement des données...

* 16 évaluations chargées

* 4 films dans le catalogue

Apercu des évaluations :

userId	movieId	rating	timestamp
0	1	1	NaN
1	1	2	2.0
2	1	3	7.0
3	1	4	8.0
4	2	1	4.0

5 ETAPE 2 : Fusion des Données

Processus : Jointure des evaluations avec les metadonnees des films via l'identifiant movieId. Permet d'enrichir les donnees d'evaluation avec les genres des films.

5.1 Code

```

1 print("\n[ETAPE 2/7] Fusion des donn\ees...")
2 df = pd.merge(ratings, movies, on='movieId', how='inner')
3 print(f"* Dataset fusionn\ee : {df.shape[0]} lignes")

```

```
4 print("Colonnes disponibles :", list(df.columns))
```

Listing 2 – Fusion des données

5.2 Output

Output - ETAPE 2

[ETAPE 2/7] Fusion des données...

* Dataset fusionné : 16 lignes

Colonnes disponibles : ['userId', 'movieId', 'rating', 'timestamp', 'title', 'genres']

6 ETAPE 3 : Création de la Matrice Utilisateur-Film

Structure : Creation d'une matrice creuse de taille $n_{utilisateurs} \times n_{films}$ ou chaque cellule contient la note attribuee. Le taux de remplissage initial est de 75%.

$$M = \begin{bmatrix} - & 4 & 3 & 4 \\ 2 & 1 & 8 & 1 \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} \quad (4)$$

6.1 Code

```
1 print("\n[ETAPE 3/7] Cr\éation de la matrice utilisateur-film
...")
2 matrix = df.pivot_table(index='title', columns='userId',
values='rating')
3 print("\nMatrice brute (5 premières lignes) :")
4 print(matrix.head())
5 print(f"\n* Dimensions : {matrix.shape[0]} films x {matrix.
shape[1]} utilisateurs")
6 print(f"* Taux de remplissage : {(1 - matrix.isna().mean().
mean()):.1%}")
```

Listing 3 – Création de la matrice utilisateur-film

6.2 Output

Output - ETAPE 3

[ETAPE 3/7] Création de la matrice utilisateur-film...

Matrice brute (5 premières lignes) :

	1	2	3	4
Item1	NaN	4.0	3.0	4.0
Item2	2.0	1.0	8.0	1.0
Item3	7.0	NaN	NaN	6.0
Item4	8.0	7.0	4.0	NaN

* Dimensions : 4 films x 4 utilisateurs

* Taux de remplissage : 75.0%

7 ETAPE 4 : Normalisation des Notes

Impact : Réduit le biais systématique entre utilisateurs. Permet une comparaison plus équitable entre les films.

Exemple de normalisation :

$$M_{norm} = M - \bar{M} = \begin{bmatrix} - & 0.33 & -0.66 & 0.33 \\ -1 & -2 & 5 & -2 \end{bmatrix} \quad (5)$$

7.1 Code

```
1 print("\n[ETAPE 4/7] Normalisation des notes...")
2 matrix_norm = matrix.subtract(matrix.mean(axis=1), axis=0)
3 print("\nMatrice normalisée (exemple) :")
4 print(matrix_norm.loc[['Item1', 'Item2'], [1, 4]])
5 print("\nExplication : Les notes sont centrées autour de la
  moyenne de chaque film")
```

Listing 4 – Normalisation des notes

7.2 Output

Output - ETAPE 4

[ETAPE 4/7] Normalisation des notes...

Matrice normalisée (exemple) :

	1	4
Item1	NaN	0.33
Item2	-1.00	-2.00

Explication : Les notes sont centrées autour de la moyenne de chaque film

8 ETAPE 5 : Calcul des Similarités

Comparaison :

- Pearson sensible aux tendances lineaires
- Cosinus sensible à l'amplitude des notes

Resultats typiques :

$$S_{Pearson} = \begin{bmatrix} 1 & -0.93 \\ -0.93 & 1 \end{bmatrix}, \quad S_{Cosinus} = \begin{bmatrix} 1 & 0.40 \\ 0.40 & 1 \end{bmatrix} \quad (6)$$

8.1 Code

```
1 print("\n[ETAPE 5/7] Calcul des similarités entre films...")
2
3 # Méthode Pearson
4 print("\nCalcul des corrélations de Pearson...")
5 item_similarity_pearson = matrix_norm.T.corr()
6 print("Matrice de similarité Pearson (extrait) :")
7 print(item_similarity_pearson.head())
8
9 # Méthode Cosinus
10 print("\nCalcul des similarités cosinus...")
11 matrix_norm_zero = matrix_norm.fillna(0)
12 item_similarity_cosine = pd.DataFrame(
13     cosine_similarity(matrix_norm_zero),
14     index=matrix_norm.index,
15     columns=matrix_norm.index
16 )
17 print("Matrice de similarité Cosinus (extrait) :")
18 print(item_similarity_cosine.head())
```

8.2 Output

Output - ETAPE 5

[ETAPE 5/7] Calcul des similarités entre films...

Calcul des corrélations de Pearson...

Matrice de similarité Pearson (extrait) :

	Item1	Item2	Item3	Item4
Item1	1.00	-1.00	1.00	-0.93
Item2	-1.00	1.00	-1.00	0.93
Item3	1.00	-1.00	1.00	-0.93
Item4	-0.93	0.93	-0.93	1.00

Calcul des similarités cosinus...

Matrice de similarité Cosinus (extrait) :

	Item1	Item2	Item3	Item4
Item1	1.00	-0.98	0.74	-0.85
Item2	-0.98	1.00	-0.29	0.40
Item3	0.74	-0.29	1.00	-0.85
Item4	-0.85	0.40	-0.85	1.00

9 ETAPE 6 : Prédiction des Notes et Comparaison

Methodologie : Pour chaque film non note, on calcule une moyenne ponderee des notes des films similaires :

$$\hat{r}_{u,i} = \bar{r}_i + \frac{\sum_{j \in N(i)} sim(i,j) \cdot (r_{u,j} - \bar{r}_j)}{\sum_{j \in N(i)} |sim(i,j)|} \quad (7)$$

ou $N(i)$ sont les voisins les plus proches de l'item i .

9.1 Code

```

1 def item_based_recommendation(user_id, similarity_matrix,
2   method_name):
3     print(f"\nAnalyse pour l'utilisateur {user_id} ({
4       method_name})")
5     user_ratings = matrix_norm[user_id]
6
7     # Films non not\és
8     unwatched = user_ratings[user_ratings.isna()].index
9     print(f"\n* Films non \évalu\és : {list(unwatched)}")
10
11     predictions = []
12     for movie in unwatched:
13       if movie not in similarity_matrix:
14         continue
15       # Calcul des similarit\és
16       sim_scores = similarity_matrix[movie].reset_index(
17         name='similarity')
18       merged = pd.merge(user_ratings.dropna().reset_index()
19         , sim_scores, on='title')
20
21       if not merged.empty:
22         top_sims = merged.nlargest(2, 'similarity')
23         predicted = np.average(top_sims[user_id], weights
24           =top_sims['similarity'])
25         predictions.append((movie, predicted))
26
27     return sorted(predictions, key=lambda x: x[1], reverse=
28       True)[:3]
29
30 # Comparaison des m\éthodes
31 print("\n[ETAPE 6/7] Comparaison des m\éthodes...")
32
33 # M\éthode Pearson
34 print("\n" + "="*50)
35 reco_pearson = item_based_recommendation(4,
36   item_similarity_pearson, "Pearson")
37 print("\nR\ésultats Pearson :")
38 for film, note in reco_pearson:
39   note_reelle = note + matrix.mean(axis=1)[film]
40   print(f" {film} :<45} : {note:.2f} (norm) {note_reelle:.2
41     f}/5")
42
43 # M\éthode Cosinus
44 print("\n" + "="*50)
45 reco_cosine = item_based_recommendation(4,
46   item_similarity_cosine, "Cosinus")
47 print("\nR\ésultats Cosinus :")
48 for film, note in reco_cosine:

```

```

40     note_reelle = note + matrix.mean(axis=1)[film]
41     print(f" {film} : {note:.2f} (norm) {note_reelle:.2f}/5"
          )

```

Listing 6 – Prédictions des notes et comparaison

9.2 Output

Output - ETAPE 6

[ETAPE 6/7] Comparaison des méthodes...

=====

Analyse pour l'utilisateur 4 (Pearson)

* Films non évalués : ['Item4']

Résultats Pearson :

Item4 : 31.57 (norm) 37.91/5

=====

Analyse pour l'utilisateur 4 (Cosinus)

* Films non évalués : ['Item4']

Résultats Cosinus :

Item4 : 0.04 (norm) 6.37/5

10 ETAPE 7 : Analyse Finale

10.1 Code et Output d'Analyse Finale

```

1 print("\n[ETAPE 7/7] Analyse finale :")
2 print("\nObservations clés :")
3 print("1. La méthode Pearson donne généralement des pr\é
   diction plus \élev\ées")
4 print("2. Le cosinus est plus conservateur car sensible aux
   notes nulles")
5 print("3. Item4 est fortement recommand\é dans les deux m\é
   thodes")
6 print("\nConclusion : La corr\élation de Pearson semble mieux
   adapt\ée")
7 print("pour ce jeu de donn\ées clairsem\é avec normalisation"
      )

```

Listing 7 – Analyse finale

10.2 Output

Output - ETAPE 7

[ETAPE 7/7] Analyse finale :

11 Comparaison et Perspectives d'Amélioration

11.1 Comparaison des Méthodes

La corrélation de Pearson apparaît comme mieux adaptée pour ce jeu de données, car la normalisation permet de centrer les notes et de mettre en évidence les relations linéaires entre les films. La méthode cosinus, quant à elle, est plus sensible aux valeurs manquantes.

Observations clés :

1. La méthode Pearson donne généralement des prédictions plus élevées
2. Le cosinus est plus conservateur car sensible aux notes nulles
3. Item4 est fortement recommandé dans les deux méthodes

11.2 Perspectives d'Amélioration

- Explorer d'autres techniques de normalisation ou de pondération des notes.
- Augmenter la taille du jeu de données pour améliorer la robustesse des prédictions.
- Tester des méthodes hybrides pour pallier les limites de chaque approche.

11.3 Conclusion

Conclusion : L'approche basée sur la corrélation de Pearson montre des résultats prometteurs pour des données claires, mais une approche hybride permettrait de bénéficier des avantages des deux méthodes.

A Code Complet (version plus pratique)

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.metrics.pairwise import cosine_similarity
4
5 # Chargement des donnees
6 ratings = pd.read_csv('rating-Rene (1).csv', sep=';',
7                       na_values=[''])
8 movies = pd.read_csv('movies-Rene (1).csv', sep=';')
9
10 # Fusion des donnees
11 df = pd.merge(ratings, movies, on='movieId', how='inner')
12
13 # Creation de la matrice utilisateur-film
14 matrix = df.pivot_table(index='title', columns='userId',
15                          values='rating')
16
17 # Normalisation de la matrice
18 matrix_norm = matrix.subtract(matrix.mean(axis=1), axis=0)
19
20 # Calcul des similarites
21 # Pearson
22 item_similarity_pearson = matrix_norm.T.corr()
23
24 # Cosine (remplacer NaN par 0)
25 matrix_norm_zero = matrix_norm.fillna(0)
26
27 item_similarity_cosine = pd.DataFrame(
28     cosine_similarity(matrix_norm_zero),
29     index=matrix_norm.index,
30     columns=matrix_norm.index
31 )
32
33 # Fonction de recommandation generique
34 def item_based_recommendation(user_id, similarity_matrix,
35                               n_similar=2, n_reco=3):
36     user_ratings = matrix_norm[user_id]
37     unwatched_movies = user_ratings[user_ratings.isna()].
38         index
39
40     predictions = []
41
42     for movie in unwatched_movies:
43         if movie not in similarity_matrix:
44             continue
```

```

42     # Recuperer les similarites (Similarites du film
        cible avec les autres films)
43     sim_scores = similarity_matrix[movie].reset_index(
        name='similarity')
44
45     # Fusionner avec les films notes par l'utilisateur
46     merged = pd.merge(
47         user_ratings.dropna().reset_index(),
48         sim_scores,
49         on='title'
50     )
51
52     if merged.empty:
53         continue
54
55     # Filtrer les similarites non-nulles
56     valid_sims = merged[merged['similarity'].notna()]
57     if len(valid_sims) < 1:
58         continue
59
60     # Top N similarites
61     top_sims = valid_sims.nlargest(n_similar, 'similarity')
62
63     # Prediction ponderee
64     predicted_rating = np.average(top_sims[user_id],
        weights=top_sims['similarity'])
65     predictions.append((movie, predicted_rating))
66
67     # Tri et retour
68     predictions.sort(key=lambda x: x[1], reverse=True)
69     return predictions[:n_reco]
70
71 # Comparaison pour l'utilisateur 4
72 print("=== Resultats Pearson ===")
73 reco_pearson = item_based_recommendation(4,
    item_similarity_pearson)
74 for film, note in reco_pearson:
75     print(f"- {film} : {note:.3f} (norm) / {note + matrix.
        mean(axis=1)[film]:.3f} (reel)")
76
77 print("\n=== Resultats Cosinus ===")
78 reco_cosine = item_based_recommendation(4,
    item_similarity_cosine)
79 for film, note in reco_cosine:
80     print(f"- {film} : {note:.3f} (norm) / {note + matrix.
        mean(axis=1)[film]:.3f} (reel)")

```

Listing 8 – Code complet de l'implémentation