

République Algérienne Démocratique et Populaire  
Ministère de l'Enseignement Supérieur et de la Recherche Scientifique  
Université des Sciences et de la Technologie Houari Boumediene



## Faculté de Mathématique

Département de Probabilités et de Statistiques

Domaine Mathématiques et Informatique

## Mémoire

Pour l'obtention du Diplôme de Master

## Probabilités et Statistiques Appliqués

Thème

**Recommendation Systems Based On Covariance Distance  
Similarity Case Study YASSIR EXPRESS**

Présente par : Nihad Senhadji

Soutenu le 29 juin 2022, devant le jury composé de :

Président : Mrs. CHOUEH Meriem , Usthb

Rapporteur : Mr. Medkour Tarek , Usthb

Examineur : Mr. MEZIANI Aymen , Usthb

Invité : Mr. Hocine Abdelouaheb , Yassir Express





---

## ACKNOWLEDGEMENTS

First and foremost, I owe my gratitude to my advisor, Professor Medkour Tarek, for believing in me by providing me with this opportunity to make this project work and for guiding me through my master thesis journey. Above everything else, I appreciate that he patiently allowed me to follow my interests and pursue my own goals. I am also incredibly thankful for his continuous support, constructive feedback, and great ideas.

Whenever I felt stuck, he was there to help me and my next step and make progress.

I want to thank and express my deep appreciation to Dr.Hocine Abdelouaheb from the YASSIR EXPRESS company for guiding my project with unlimited patience, purely selfless help, and compelling support and encouragement.

I am thankful to all my YASSIR EXPRESS data scientists who made this work possible, working with many bright youths throughout my master thesis journey was a privilege.

Also, I want to thank my family members for their support and encouragement. My parents for being there for me no matter what and always believing in me and my brother and sister, without whom this journey would not have been possible.

Last but not least, I would like to thank every person who has contributed or helped me to complete this work. There are not enough words to describe how thankful I am for all the opportunities I have been given and the people I have met while working on this thesis.

---

## Abstract

Recommendation Systems (RSs) is an information filtering system that aims to provide users with relevant items, Their successful application in online retail is reflected in increased user satisfaction and sales revenue. Collaborative filtering is the most used model in recommendation systems one of the main crucial components of these models is the similarity measurement between users or items, Several similarity measures have been proposed, with different performances in terms of accuracy and quality of recommendations.

This thesis provides the results of a building recommendation systems project dealing with new similarity measure distance covariance in terms of their effectiveness. We implement this systems using the YASSIR EXPRESS data set and build a restaurant recommendation systems as a case study.

---

## Notation

*RSs* Recommendation Systems.

*CF* Collaborative Filtering.

$r_{ui}$  Rating of user  $u$  for item  $i$ .

$\bar{r}_u$  Average rating of the user  $u$ .

$U$  Set of users in the system.

$I$  Set of items in the system.

$U_i$  the subset of users that have rated an item  $i$ .

$I_u$  the subset of items that have been rated by a user  $u$ .

$I_u \cap I_v$  the items that have been rated by two users  $u$  and  $v$ .

$U_{ij}$  the set of users that have rated both items  $i$  and  $j$ .

$U = u_1, u_2, u_3, \dots, u_m$  the set of  $m$  users.

$I = i_1, i_2, i_3, \dots, i_n$  the set of  $n$  users.

$m$  Total number of all users.

$n$  Total number of all items.

$R$  the matrix of the logs of ratings  $r_{ui}$ .

$i, j$  any item in the system.

$u, v$  any user in the system.

# Contents

<b>1</b>	<b>YASSIR EXPRESS</b>	<b>13</b>
<b>2</b>	<b>An Introduction To Recommendation Systems</b>	<b>18</b>
2.1	Introduction . . . . .	18
2.2	Recommendation Systems . . . . .	19
2.3	Data and Knowledge Sources . . . . .	20
2.4	Functions of Recommendation Systems . . . . .	21
2.5	Recommendation Systems Applications in The industry . . . . .	22
2.6	Each Domain Has Its Unique Challenges . . . . .	24
2.7	Implementations Of Recommendation Systems In the Industry . . . . .	25
2.7.1	Netflix Recommendation System . . . . .	25
2.7.2	Amazon Recommendation System . . . . .	25
2.7.3	Google News Personalization Recommendation System . . . . .	26
2.7.4	Facebook Friend Recommendations System . . . . .	27
2.8	Models of Recommendation Systems . . . . .	28
2.8.1	Non-Personalized Recommendation System . . . . .	28
2.8.2	Personalized Recommendation System . . . . .	28
2.9	Evaluating Recommendation Systems . . . . .	34
2.9.1	Accuracy . . . . .	34
2.9.2	Non Accuracy Metrics . . . . .	38
2.10	Challenges in Recommendation Systems . . . . .	40
2.10.1	Cold Start . . . . .	40
2.10.2	Sparsity . . . . .	40
2.10.3	Scalability . . . . .	40
<b>3</b>	<b>Neighborhood Based Recommendation System</b>	<b>41</b>
3.1	Types of Neighborhood-based collaborative filtering . . . . .	42
3.1.1	User-based neighborhood collaborative filtering: . . . . .	42
3.1.2	Item-based neighborhood collaborative filtering: . . . . .	43
3.2	User-based VS Item-based Recommendation . . . . .	44
3.3	Rating Matrix . . . . .	45
3.4	Explicit and Implicit Ratings . . . . .	45
3.5	Similarity measures in collaborative filtering: . . . . .	46
3.5.1	Cosine similarity: . . . . .	46
3.5.2	Adjusted cosine vector: . . . . .	48
3.5.3	Pearson correlation: . . . . .	48
3.5.4	The Jaccard coefficient: . . . . .	49
3.5.5	Conditional probability-based similarity: . . . . .	49
3.5.6	Kendall's tau . . . . .	50
3.5.7	Euclidean distance . . . . .	50
3.5.8	Manhattan distance . . . . .	51

3.6	User-based Rating Prediction . . . . .	52
3.7	Item-based Rating Prediction . . . . .	52
3.8	Rating Normalization . . . . .	53
3.9	Neighborhood Selection . . . . .	54
3.9.1	Top-N Recommendation: . . . . .	54
3.9.2	Threshold filtering: . . . . .	55
3.9.3	Negative filtering: . . . . .	55
3.10	Dimensionality Reduction . . . . .	55
3.11	Matrix factorization techniques . . . . .	56
3.11.1	Singular Value Decomposition (SVD) . . . . .	57
3.11.2	Principal Component Analysis (PCA) . . . . .	58
3.11.3	Probabilistic Matrix Factorization (PMF) . . . . .	59
<b>4</b>	<b>Covariance Distance</b>	<b>60</b>
4.1	Predefined similarity measures . . . . .	61
4.2	Distance Covariance . . . . .	63
4.2.1	The Distance Covariance Function . . . . .	63
4.2.2	Estimation of distance covariance . . . . .	67
4.2.3	Asymptotic Distribution . . . . .	69
4.2.4	Generalisations and Modifications of the Distance Covariance . . . . .	69
4.3	The Distance Correlation Coefficient in High Dimensions . . . . .	70
4.4	Implimentation of distance covariance in recommendation systems . . . . .	71
4.4.1	Distance covariance in Item-Based Neighborhood . . . . .	71
4.4.2	Distance covariance in User-Based Neighborhood . . . . .	72
<b>5</b>	<b>Implementation of Restaurant RSs</b>	<b>73</b>
5.1	YASSIR EXPRESS . . . . .	74
5.1.1	The Application . . . . .	75
5.1.2	Flow . . . . .	75
5.1.3	Dataset . . . . .	77
5.2	Experimental Setup . . . . .	78
5.3	Data Processing . . . . .	79
5.3.1	Data Loading . . . . .	79
5.3.2	Data Preparation . . . . .	79
5.3.3	Data Exploration . . . . .	79
5.4	The Recommendations Systems Models . . . . .	83
5.4.1	Evaluation Metrics . . . . .	83
5.5	Neighborhood Collaborative Filtering . . . . .	84
5.5.1	User Neighborhood Collaborative Filtering . . . . .	84
5.5.2	Item Neighborhood Collaborative Filtering . . . . .	85
5.5.3	Similarity Metrics . . . . .	85
5.5.4	Results of the User-based model . . . . .	93
5.5.5	Results of the Item-based model . . . . .	99
5.5.6	Comparing the user-based and the item-based . . . . .	101
5.5.7	Conclusion . . . . .	103
5.6	Location-Based Recommendation System . . . . .	104
5.7	Popularity Based Recommendation Model . . . . .	107
<b>6</b>	<b>Appendix</b>	<b>109</b>
	<b>Bibliography</b>	<b>138</b>



# List of Figures

1.1	YASSIR EXPRESS:	13
1.2	YASSIR EXPRESS Services	14
2.1	The Long Tail: physical institutions can only provide what is popular, while online institutions can make everything available	21
2.2	Recommendation Systems Applications	22
2.3	Netflix Recommendation System	25
2.4	Amazon Recommendation System	26
2.5	Google Recommendation System	26
2.6	Google Recommendation System	27
2.7	Recommendation Systems Models	29
2.8	Hybrid Recommendation System	32
2.9	Accuracy metrics	34
3.1	User-based neighborhood collaborative filtering	42
3.2	The user-Based pipeline	43
3.3	Item-based neighborhood collaborative filtering	43
3.4	The item-Based pipeline	44
3.5	User Cosine Similarity Measure	46
3.6	Item Cosine Similarity Measure	47
3.7	Pearsons correlation	48
3.8	Jaccard coefficient	49
3.9	Matrix Factorization Techniques	56
3.10	Singular Value Decomposition (SVD)	57
5.1	Recommendation System Architecture	74
5.2	YASSIR EXPRESS Application	75
5.3	YASSIR EXPRESS User Experience	76
5.4	Steps performed to compare the recommendations systems	78
5.5	The states with the largest number of restaurants	79
5.6	The commune with the largest number of restaurants	80
5.7	The Rating Distribution	81
5.8	Box Plot Of The Rating	81
5.9	Flow chart Of The Neighborhood Collaborative Filtering	84
5.10	Similarities and Neighbors using Cosine Similarity metric	88
5.11	TOP Recommendation using Cosine similarity metric	91
5.12	TOP Recommendation using Euclidean Distance similarity metric	91
5.13	TOP Recommendation using Manhattan Distance similarity metric	91
5.14	MAE Evaluation results of the similarity measures using user-based	93
5.15	RMSE Evaluation results of the similarity measures using user-based	93
5.16	MSE Evaluation results of the similarity measures using user-based	94

5.17 Performance Comparison Between The Evaluation results of the similarity measures using user-based . . . . .	94
5.18 Top recommendation using the Distance Covariance similarity metric . . . . .	99
5.19 MAE Evaluation results of the similarity measures using item-based . . . . .	99
5.20 RMSE Evaluation results of the similarity measures using Item-based . . . . .	100
5.21 MSE Evaluation results of the similarity measures using Item-based . . . . .	100
5.22 Performance Comparison Between The Evaluation results of the similarity measures using item-based neighborhood . . . . .	100
5.23 Geographical Overview of Restaurants . . . . .	104
5.24 The table represents the longitude and latitude of the restaurants. . . . .	104
5.25 Silhouette Score Elbow for KMeans Clustering. . . . .	105
5.26 Top Recommendation Based on the Location-based recommendation system . .	106
5.27 Top Recommendation Based on the Popularity Model . . . . .	107
5.28 Evaluation Results of Recommendation Based on the Popularity Model . . . .	108

# List of Tables

2.1	Confusion Matrix . . . . .	36
3.1	User-Item Rating Matrix . . . . .	45
3.2	Example User-Item Rating Matrix . . . . .	47
4.1	TABLE OF SIMILARITIES . . . . .	62
5.1	Table summarizes the YASSIR EXPRESS dataset . . . . .	77
5.2	Table Of The Used Similarities Metrics . . . . .	86
5.3	Evaluation Results of the Similarity Measures on User-Based . . . . .	93
5.4	Results of the Similarity Measures on User-Based and Item-based First rating approach . . . . .	101
5.5	Results of the Similarity Measures on User-Based and Item-based second rating approach . . . . .	102
5.6	The Evaluation Results on the Popularity RSs . . . . .	108

---

## Executive Summary

**Chapter 1** This chapter provides presentation of the YASSIR EXPRESS company.

**Chapter 2** The aim of this chapter is to give an overview of the most important concepts used in recommendation systems, we provide what a recommendation system is and what it does, we identify the data and knowledge Sources used for a recommendation, we present the different models of RSs, then we end-up with the evaluation and challenges of the recommendation systems .

**Chapter 3** The state of the art, we identify the main models used on the recommendation systems , the neighborhood-based approach to collaborative filtering recommendation systems, We therefore focus on the presentation of neighborhood models, We give an overview of different similarity measures that are used on field of recommendation systems ,and we describe the process of neighborhood-based models.

**Chapter 4** This chapter represent the introductory of the distance covariance , As our main aim yield this distance can be used as new similarity measure on the neighborhood-based recommendation systems field ,we apply an appropriate modification of the distance covariance similarity measure.

**Chapter 5** This chapter presents the results and findings of implemented and the developed neighborhood-based models using the traditional similarity measures, We have designed an innovative proposal a new neighborhood-based approach using the distance covariance new similarity metric on the YASSIR EXPRESS dataset .However, the evaluation process of recommendation systems models performance were applied using the the Root Mean Squared Error (RMSE), Mean Squared Error (MSE) and the Mean Absolute Error (MAE), assuming that the lower they are the higher the quality of the recommendation.

# Chapter 1

## YASSIR EXPRESS

YASSIR <sup>1</sup> is a mobile application development company it was founded in 2017 by two doctoral students, El Mehdi Yettou and Nouredine Tayebi, two alumni from the polytechnic school of Algiers with extensive experience in the field of research at the Canadian University of Sherbrooke and the American University from Stanford, respectively. in order to facilitate the daily lives of citizens via an application that can be used either with your smartphone, tablet, or Pc. It is free to download via Google Play and the Apple store. The products that Yassir company own are Yassir Express <sup>2</sup> , Yassir Go <sup>3</sup> , Yassir Business <sup>4</sup> , Yassir market <sup>5</sup> . ***Yassir Go*** Initiator and leader of VTC in Algeria, and an innovative transport service. ***Yassir Business*** is a transportation solution for businesses. ***Yassir Gromadaire*** groceries application that markets wholesale food products. ***Yassir Express***



Figure 1.1: YASSIR EXPRESS:

Application developed by the company YASSIR, bringing together meal delivery services from restaurants partner, Yassir dark stores, and other services such as dry cleaning. Available on APP STORE about the smartphones under IOS and Google Play for Smartphones under

<sup>1</sup><https://yassir.com/en/home/>

<sup>2</sup><https://express.yassir.io/>

<sup>3</sup><https://app.yassir.io/auth>

<sup>4</sup><https://yassir.com/en/business-en/>

<sup>5</sup><https://market.yassir.io/>

Android.

Yassir Express offers a service by which E-consumers will be in touch with E-Suppliers to purchase products/services/Meals remotely presented on the application and benefit from delivery to their destination. In this context, it is recalled that Yassir Express acts as an intermediary. Its role is limited to hosting offers of products/services/meals from E-suppliers and putting them in contact with E-consumers.

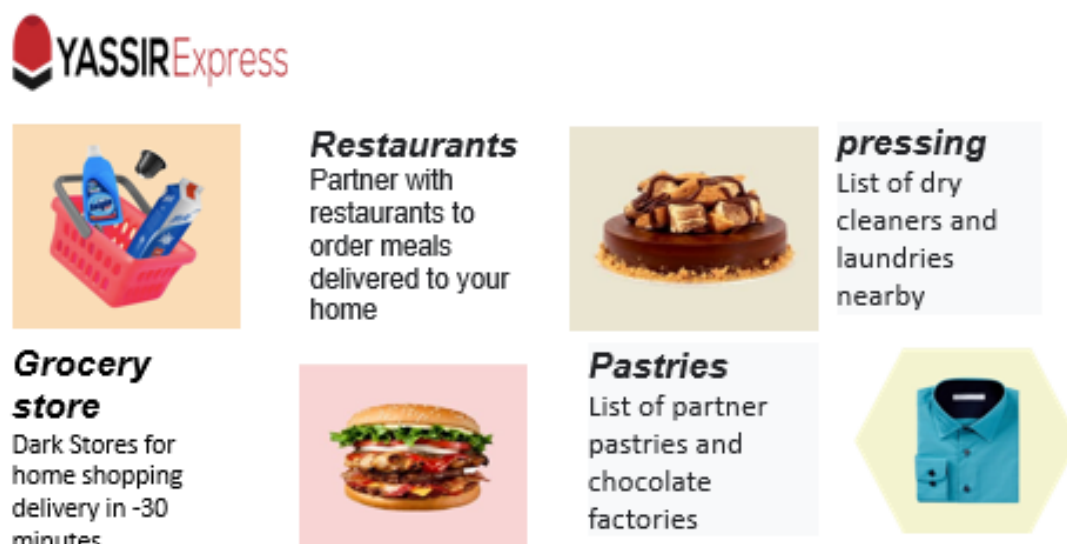


Figure 1.2: YASSIR EXPRESS Services

**E-consumer** means any natural or legal person who acquires, for consideration or free of charge, a good or service through the YASSIR EXPRESS Applications from an E-supplier for end-use.

**E-Supplier** means any natural or legal person who markets or offers the supply of goods or services or meals using electronic communications.

**E-payment** or electronic means of payment" means any payment instrument, authorized by the legislation in force, allowing its holder to make local or remote payments through an electronic system.

**User Content** means the textual, audio, and/or visual content and information communicated during your registration, including comments and feedback relating to the Services or Third-Party Services, as well as the content and information communicated during requests for assistance and on the occasion of participation in competitions and promotional operations.

**E-shop** designates the dedicated space within the YASSIR EXPRESS made available to the E-supplier to offer their products/services/menus for sale.

**Services** refers to the services available on the YASSIR EXPRESS applications and/or Website.

**Website** means the dynamic interface accessible at the addresses:..... any other website operated by YASSIR and made available as part of the Services.

**Prices** means the prices applicable to the Services, which will be expressly communicated to you via the Application and/or Website at the time of validation of the corresponding order.

**Access to the application and the service** The service is accessible, subject to the restrictions provided on the YASSIR EXPRESS application:

- To any natural person having total legal capacity by the law in force in the countries eligible to commit to these general conditions.
- To any legal person acting through a natural person having the legal capacity to contract in the name and on behalf of the legal person.

**Acceptance of terms and conditions** A checkbox materializes the acceptance of these general conditions in the registration form of the YASSIR EXPRESS application and / or website. This acceptance can only be whole and complete, and any conditional membership is considered null and void. E-consumers acknowledge that the use of the YASSIR EXPRESS application must only be for the Service provided for this purpose mentioned in the subject. The use of the application for a reason other than that mentioned in the subject exposes its author to penalties without prejudice to the legal actions that YASSIR may bring for the damage caused.

**Registration on the site and / or the application** Access to the Service requires the e-consumer to register on the YASSIR EXPRESS application by completing the form provided. The e-consumer must provide all the information marked as mandatory (Name, first name, address, and especially his telephone number). Any incomplete registration will not be validated. The e-consumer acknowledges that for the creation of his profile, YASSIR will be required to communicate his telephone number to one of his suppliers for authentication. After verification of the telephone number, the E-consumer's profile is automatically created ( The "profile"), thus giving him access to a personal space (hereafter: "the personal space") which allows him to manage his use of the Service, according to the form and technical means that YASSIR deems most appropriate to provide the said Service.

**Service Description** The Service consists of putting the e-consumer in direct contact with an E-supplier through the YASSIR EXPRESS application. This applies to all the eligible countries mentioned in article 2 hereof. Access to the Service requires that the e-consumer have a Smartphone and/or a device with internet access (3G and/or 4G connection). Any e-consumer wishing to place an order must connect to the application, once done, he will have to choose a product/meal and / or Service from one of the existing E-suppliers on the platform, all of the prices he has to pay are displayed on his screen, the e-consumer can then validate by pressing the "validate the order" button. The request will be sent instantly to the E-supplier.

**Product compliance** YASSIR EXPRESS will combine its best efforts so that the photographic representation of the products is as faithful as possible to the actual products; however, it is possible that the perception by the e-consumer of the photographic replica of the products does not correspond precisely to the product itself.

**Ordered** Each item chosen through a simple click is followed by a descriptive sheet, a photo, and a compatible accessory list, for which the E-supplier is solely responsible. With a simple click on the link or the Order button, the item of your choice is slipped into your basket. The latter contains the products/meals that the e-consumer plans to buy. It is possible at any time to access the basket (by clicking on the button representing a basket) and to modify its contents by adding or deleting an article.

- Validation of the order by the e-consumer: the e-consumer can consult the basket and click on Validate your order. Any order not validated cannot be honored.
- Availability: Product availability is displayed on the site/application on each product/meal sheet. If a product/meal is unavailable, the mention “Unavailable” will be displayed on the product/meal sheet.  
If, after the order, a product/meal becomes unavailable, YASSIR EXPRESS will inform the e-consumer by email of this unavailability. It will give the e-consumer the possibility of choosing between:
  - makes a pre-order of the product/meal as soon as the product is available, the order is deemed to have been validated.
  - cancels the order.
- Pre-order: is the opportunity to reserve the product/meal as soon as it is released and as soon as it is in stock/available. The order will be delivered to the e-consumer when the last pre-order product/meal is in stock or available.
- Cancel the order: The e-consumer will have the right to cancel his order provided that:
  - The cancellation is made within 10 minutes of its validation
  - The delivery person has not picked it up yet.
  - If the delivery person delivers the order to the destination, the e-consumer must pay the Service and delivery costs displayed beforehand on the YASSIR EXPRESS application even in the event of cancellation of this order

**Delivery** The Products/meals are delivered to the delivery address chosen by the e-consumer as indicated when ordering. The e-consumer undertakes to receive the Products/meals thus ordered. The receipt of the order is followed by an acknowledgment of receipt by the E-consumer, and a copy will be sent afterward..

- In the event of delivery to the workplace of the E-consumer, the deliverer, the e-consumer must take all necessary precautions so that the deliverer can deliver the order to him and the e-consumer can receive it. In the case otherwise, the e-consumer is obliged to pay the total amount of the order.
- YASSIR is responsible for the order of the e-consumer from the moment the delivery person picks it up from the E-supplier and delivers it to the destination indicated on the YASSIR EXPRESS application beyond this period and once the delivery has been confirmed. The e-consumer understands that YASSIR is not responsible for what may happen to the said product/meal delivered.
- The e-consumer acknowledges, for the smooth running of the service offered by YASSIR EXPRESS, that he will be required to communicate to the delivery person the following elements:
  - His identity document.
  - Order number.
  - His email address.



**Delivery time:** The e-consumer chooses the date and type of delivery according to the delivery slots available on the web application. The date is valid unless proven otherwise by agreement between the Parties. YASSIR EXPRESS and its partners will then implement their best efforts to meet the delivery times of the products ordered. If these deadlines cannot be met, the e-consumer will be automatically informed by any means, particularly by e-mail, SMS, or telephone.

**Delivery costs:** Delivery costs may vary depending on the basket amount, the place of delivery, the frequency of purchases by the e-consumer between two orders, and the day and time chosen. Delivery costs are calculated automatically and indicated when selecting the day and delivery window. In other words, the e-consumer can find out the amount of the delivery costs by clicking on the “Basket” icon on the line “delivery costs and service costs.”

**Data The E-consumer acknowledges and expressly accepts:** That the data collected on the YASSIR EXPRESS application and/or website and YASSIR’s computer equipment are proof of the reality of the operations carried out within the framework of these presents: That these data constitute the only mode of proof accepted between the parties, in particular for the calculation of the sums due to YASSIR. The e-consumer can access this data in his Personal Space.

#### **The obligations of YASSIR**

- Offer the e-consumer the free download of the YASSIR EXPRESS application.
- Provide the e-consumer with a presentation page on the YASSIR EXPRESS application based on his provided information.
- Make available to the e-consumer in the YASSIR EXPRESS Application a management service for his account, in particular, the history of his orders made.
- Deliver the order on time. Offer the e-consumer an evaluation and rating service to improve the service’s quality.
- Carry out a telephone follow-up to ensure the order’s delivery.
- Process order returns for products that arrive after the deadline.
- Process returns orders for products that do not comply with the request made by the e-consumer.

## Chapter 2

# An Introduction To Recommendation Systems

### 2.1 Introduction

The enormous amount of digital information generated by the increasing number of users on the internet has created the challenge of information overload. It becomes challenging, tedious, and time-consuming for users to retrieve the exact information on time from the web. Users need suggestions from friends or experts who have knowledge about the item. While this freedom of purchase has made online commerce a multibillion-dollar industry, A standard catalog of an online video on demand service such as Netflix <sup>1</sup> can have more than 100.000 titles. The e-commerce site Amazon <sup>2</sup> now contains several million product references. URLs of websites today are counted in billions or tens of billions, We are entering an era of huge databases (big data) <sup>3</sup> where the users cannot consider themselves to have an overview of what is available and what might be interesting to them it also made it more difficult to select the products that best fit their needs. This overwhelming size of data has shifted the focus of research community from simple information extraction to filtering of pertinent information this has increased the demand for recommendation systems, which provide suggestions of products to consumers. Recommendation systems were introduced in early 1990s as useful information filtering tools for guiding users in a personalized way to discover products or services they might be interested in from a large space of possible options.

Today, users are not even willing to bother to search for items, and instead, they prefer a more convenient way to find items matching their interests and preferences. Recommendation systems are heavily applied as personalized filters, and consequently, they are an integral part of users' daily lives. Whether users are watching movies, looking for new job positions, ordering food, or browsing an online catalog, a recommendation system is running in the background to generate personalized recommendations based on their users behavior and feedback. These include, for instance, rating, clicking, or reviewing items. A recommendation system is also essential from a business perspective since it has proven to increase a company's revenue. This statement is based on the observation that appropriate recommendations can strengthen user engagement and user satisfaction. As a result, users tend to spend more time consuming a service or purchasing more items. Recommendation systems have become an attractive research field to study both from industry and academia alike. In this regard, most research still focuses on improving recommendation quality as far as possible by developing new approaches.

---

<sup>1</sup><https://www.netflix.com/dz-en/>

<sup>2</sup><https://www.amazon.com/>

<sup>3</sup>Bigdata is a term used to describe a collection of data that is huge and yet growing exponentially with time.

The history of recommendation systems dates back to the year 1979 with relation to cognitive science [Rich, 1979], in which he uses stereotypes to perform a user modeling task . Recommendation systems gained prominence among other application areas such as approximation theory [Powell, 1981], information retrieval [Salton, Gerard, 1989], forecasting theories [Armstrong, 2001].

In the mid-1990s, recommendation systems became active in the research domain when the focus was shifted to recommendation problems by researchers that explicitly rely on user rating structure and also emerged as an independent research area.

The recommendation problem can be defined as estimating the response of a user to new items based on historical information stored in the system and suggesting to this user novel and original items for which the predicted response is high.

## 2.2 Recommendation Systems

is a specific type of information filtering technique that attempts to present information products (such as movies, music, news, books, restaurants...etc.) that are likely of interest to the user. which helps users navigate through large product assortments. They have great importance for the success of the e-commerce and IT industry nowadays and gradually have gained popularity in various applications RSs generate recommendations using various types of knowledge and data about users, the available items, and previous transactions stored in customized databases. The user can then browse the recommendations and may accept them or not and may provide them by implicit or explicit feedback.

All these user actions and feedback can be stored in the RSs<sup>4</sup> database and may be used for generating new recommendations. In recent years, the interest in recommendation systems has dramatically increased, as the following facts indicate: Recommendation systems play an important role in such highly rated Internet sites as Amazon.com, Netflix, YouTube, Yahoo<sup>5</sup>, Tripadvisor<sup>6</sup>, Last.fm<sup>7</sup>, and IMDb<sup>8</sup>.

Moreover, many media companies are now developing and deploying RSs as part of the services they provide to their subscribers. For example, Netflix is a streaming service that offers a wide variety of TV shows, movies, have made The Netflix Prize In 2006<sup>9</sup> Challenge A significant boost to research into recommendation systems was given when Netflix offered a prize of \$1,000,000 to the first person or team to beat their own recommendation algorithm, called CineMatch, by 10%. Initially there was a large explosion in the number of teams that submitted results, over 20,000 teams from over 150 countries registered in the first eight months After over three years of work, the prize was awarded in September 2009. The Netflix challenge consisted of a published dataset, giving the ratings by approximately half a million users on (typically small subsets of) approximately 17,000 movies. This data was selected from a larger dataset, and proposed algorithms were tested on their ability to predict the ratings in a secret remainder of the larger dataset. The information for each (user, movie) pair in the published dataset included a rating (1–5 stars) and the date on which the rating was made. The RMSE was used to measure the performance of algorithms. cinematic has an RMSE of approximately 0.95.

---

<sup>4</sup>RSs Recommendation Systems

<sup>5</sup><https://mail.yahoo.com/>

<sup>6</sup><https://www.tripadvisor.com/>

<sup>7</sup><https://www.last.fm/>

<sup>8</sup><https://www.imdb.com/>

<sup>9</sup>[https://www.wikiwand.com/en/Netflix\\_Prize](https://www.wikiwand.com/en/Netflix_Prize)

## 2.3 Data and Knowledge Sources

Recommendation systems are knowledge extraction systems that actively collect different types of information to make suggestions in order to build their recommendations. Data is primarily about the items to suggest and the users who will receive these recommendations. But, since the data and knowledge sources available for Recommendation systems can be very diverse, their use for making recommendations largely depends upon the recommendation techniques to be used.

as a general classification, data used refers to three kinds of objects:

items, users, and transactions, that is relations between the users and the items.

1. **Items :**  $\mathbf{T} = t_1, t_2, t_3, \dots, t_n$  Products and Services that a recommendation system recommends are referred to as items. There are items having low difficulty and value are news, web articles, e-books, DVDs and movies,....etc.

Items having high complexity and value are laptops, mobile phones, digital cameras, electrical appliances, PCs, Insurance policies, travel plans, financial investments, and jobs are considered as most complex items.

2. **Users:**  $\mathbf{U} = u_1, u_2, u_3, \dots, u_m$  Users of a recommendation system may have very diverse aims and features. In order to to personalize the recommendations, RSs exploit a range of information about the users. This information can be structured in various ways and again the selection of what information to model depends on the recommendation technique

therefore the system should exploit a variety of data about the users.

This data can be organized in several means it depending upon the recommendation approach.

3. **Transactions:** Any recorded communication between a recommendation system and the user is referred to as a transaction. These are logs that collect essential information generated at the time of interaction and these are valuable for the recommendation technique. The log may also have explicit feedback given by the user like ranking to the particular product. Ratings are the most common method of transaction data that a recommendation system gathers. There are varieties of forms that ratings could adopt as follows:

- **Binary rating** scale allowing users to assign items to two different classes (like/dislike) and a good example is YouTube that allows users to rate movies with either thump up or down.
- **Ordinal ratings**, such as “strongly agree, agree, neutral, disagree, strongly disagree” where users are asked to select the term that best indicates their opinion regarding an item (usually via questionnaire)
- **Unary rating**, by contrast, allows users to assign items only to a single class, which is positive in most cases and a prominent example is the Facebook’s “Like”-button. Purchased products in a web shop or clicked links on a new page can be implicit unary ratings and also in addition, unary rating can signal the purchase or observation of an item by users, or rating them positively. With the above cases, the absence of a rating indicates that there is no information relating the user to the item (perhaps a purchase was made somewhere else)

## 2.4 Functions of Recommendation Systems

The recommendation system offers suggestions to the user about a particular item that the user wants to use. Recommendation system plays different roles according to the user for example it's used by travel intermediary is usually used to increase the revenue like Expedia<sup>10</sup>. The following are different reasons to exploit RS technology by service providers:

- **Increase the revenue:**

The major objective of a commercial RSs is to increase the numbers of items that are sold. This function is most likely to be the most important in an industrial commercial RSs. The goal here is to actually sell more items than there would have been without any recommendations.

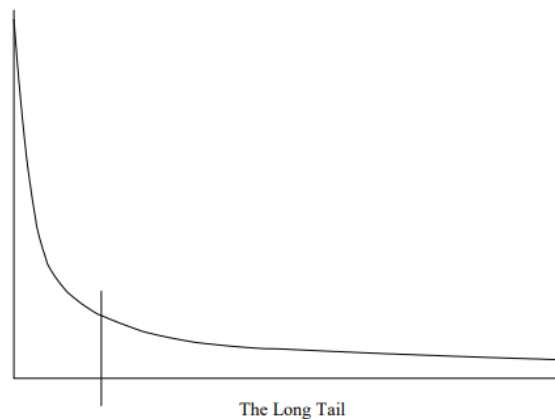


Figure 2.1: The Long Tail: physical institutions can only provide what is popular, while online institutions can make everything available

Recommendations are provided considering that suggested products and services meet the user's requirements. Non-commercial recommendations are used for similar objectives. Consider an example of a content writer who wants to increase the number of news readers on his site.

- **Increase diversity of items sold** The aim of this function is to invite users to select items that would remain unknown without recommendation. For example, in the case of a book RS (e.g Amazon bookstore), the service provider wants to be able to sell books from all its catalogue and not only the top 5 most popular ones.
- **User satisfaction:** The recommendation system helps in improving the experience of the user with the application or website. It provides interesting, significant, and, relevant recommendations as well as provides better human-computer interaction.
- **User loyalty:** A user always prefers to use a website or application which identifies its old users and treats him as a valuable user. It is a common feature of a recommendation system as it computes recommendations.
- **Better understanding of user needs:** The recommendation system is acting as an active learner of user's preferences by collecting explicitly or the predictions made by the system. The business holders may then re-use this information for improving the stock management or production of items.

<sup>10</sup><https://www.expedia.com/>

## 2.5 Recommendation Systems Applications in The industry

The recommendation system has been expanded and used in various service fields, Diverse applications in areas such as e-commerce, search, music , video and gaming apply similar techniques that leverage large volumes of data to better fulfill a user's needs. Nowadays Recommendation systems play a vital role in companies like Yahoo, Amazon Netflix, YouTube, TripAdvisor, Google, IMDB. The recommendation system was used were classified into six main categories:

Streaming Service, Social Network Service, Tourism Service, E-Commerce Service, Healthcare Service, Education Service, Academic Information Service.

The six main categories are divided based on the list of services that use a recommendation system with increasing users or increasing business value,

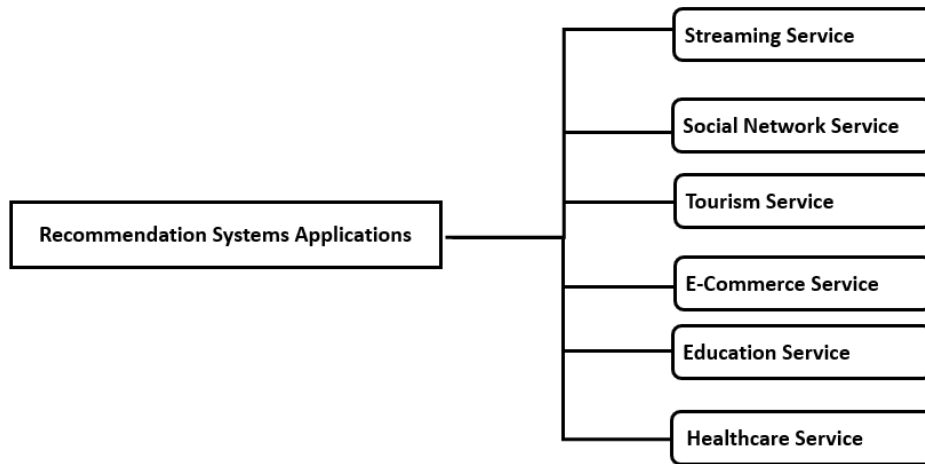


Figure 2.2: Recommendation Systems Applications

- **Streaming Service** In the past, video content, such as movies, were mainly consumed by users through TV or movie theaters. However, recently, a significant amount of video content is consumed through streaming platforms such as Netflix and YouTube. Audio content is also changing from downloading and consuming files to a user's local device to consuming content through streaming platforms such as Spotify <sup>11</sup>. Streaming services related to media content have been developed along with the recommendation system because it is necessary to reduce users' worries about choosing a vast amount of content and to provide content that is tailored to each user.
- **Social Network Service** Online social network services (SNS) such as Facebook<sup>12</sup>, Instagram<sup>13</sup>, Twitter<sup>14</sup>, and LinkedIn<sup>15</sup> are huge digital-based social exchanges where users can not only lifelog their daily life, hobbies, interests etc ..., but also provide a field of interaction with other users. The vast increase in the use of SNS was also accompanied by a vast increase in user-related data. It is possible to collect content information that users register with posts through SNS. In addition, user evaluation data can be collected, as well as rating data, these include various types of feedback data, such

<sup>11</sup><https://www.spotify.com/dz-fr/>

<sup>12</sup><https://www.facebook.com/>

<sup>13</sup><https://www.instagram.com/>

<sup>14</sup><https://twitter.com/home>

<sup>15</sup><https://www.linkedin.com/>

as likes and comments. The collected data are not only used for recommendations within SNS, but are also open to utilization in recommendation systems for other businesses.

- **Tourism Service** As the demand for travel has increased, recommendation systems have begun to be used in the tourism service field to recommend tourist destinations, route recommendations, and transportation methods.  
As the travel-related recommendation system uses situational data, such as review data and location data, user location, time, and weather, collected through SNS, research on recommendation systems using SNS has increased in the tourism service field.
- **E-Commerce Service** In recent years, with the development of digital platforms such as the web and applications, the form of consuming items has changed through e-commerce platforms such as Amazon, eBay<sup>16</sup>, and Alibaba<sup>17</sup>. E-commerce offers consumers many items and various options in the online environment and provides sellers with an easy way to sell. In particular, consumers were unable to go outside due to the lockdown measures due to COVID-19<sup>18</sup>, and as a result, they were unable to use physical Stores. Consequently, consumption using digital platforms increased exponentially.
- **Education Service** From the traditional form of education in the classroom or lecture hall, a new education trend, called Smart Learning, has formed through e-learning, in which learning is conducted through an online environment .  
Smart education has started to be gradually used in education due to the increase in the spread of various smart devices. Smart education can access vast digital resources and seamlessly provide personalized learning tailored to the needs, goals, talents, and interests of learners without time and space constraints. In this domain, recommendation systems help users to find, for instance, video courses and digital content by their current skills, and experiences. Therefore, this domain is predestined for recommendation systems due to the number of available subjects and different skill levels of the users. Providing users with the right and exciting content strengthens his or her confidence in these services. Famous services in this regard are Coursera <sup>19</sup>, Udemy <sup>20</sup> or Skillshare <sup>21</sup>.
- **Healthcare Service** the field of health recommendation systems that help users with professional treatment, the main goal is to provide suitable treatment methods according to the symptoms of various types of diseases and the stages of each disease. To this end, the health recommendation system analyzes the patient's information and the characteristics of the disease, offers an accurate diagnosis of the disease to the patient, and recommends an appropriate treatment according to the diagnosed disease.

---

<sup>16</sup><https://www.ebay.com/>

<sup>17</sup><https://www.alibaba.com/>

<sup>18</sup><https://www.who.int/emergencies/diseases/novel-coronavirus-2019>

<sup>19</sup><https://www.coursera.org/>

<sup>20</sup><https://www.udemy.com/>

<sup>21</sup><https://www.skillshare.com/>

## 2.6 Each Domain Has Its Unique Challenges

Today most e-commerce sites and applications are likely to have some recommendation engine powering their user experience.

Amazon is the first large company to be credited with having included a recommendation system at the core of their experience.

Amazon initially employed a simple item collaborative filtering approach, other retail companies such as eBay have followed the lead and incorporated recommendations in their experience.

News is also an area in which companies have applied recommendation approaches to personalize and focus on a user's interests. For example, Google News<sup>22</sup> was powered by some recommendations for news articles from the beginning. Yahoo! Has also invested in personalizing news and other web content. For news recommendation, some key challenges are freshness, where relevant articles may have a very limited period, and diversity, where there can be a large number of articles about the same topic.

Video recommendation has always been an active area of research, so it is not surprising that it is used in the industry to recommend various types of video-spanning movies and TV shows. For instance, recommendations have been an important component of the YouTube experience to help navigate the vast amounts of user-generated videos.

Music recommendation is also an active application area with interesting developments in the past years. Pandora, for instance, created a complete business model around the idea of creating personalized music stations. They created an approach that combined traditional collaborative filtering techniques with a curated approach called the Music Genome Project. Similarly, Apple's iTunes<sup>23</sup> application uses information about a user's music library to drive personalized mixes and playlists. More recently, Spotify started getting in the business of providing personalized music recommendations in its service. Social network companies have introduced several different recommendation avenues. Twitter, for example, introduced its Who to Follow recommendation algorithm to recommend new social connections. LinkedIn used a Survival Analysis approach to understand how likely a user is to change jobs. Google has also published work in recommendation systems for some social networks such as Orkut. Yahoo! has also worked on personalizing aspects such as comments on social sites.

---

<sup>22</sup><https://news.google.com/topstories?hl=en-US&gl=US&ceid=US:en>

<sup>23</sup><https://www.apple.com/itunes/>



## 2.7 Implementations Of Recommendation Systems In the Industry

### 2.7.1 Netflix Recommendation System

Netflix was founded in 1997 and initially provided a service to rent DVDs via email. In the following years, Netflix became a video streaming provider with approximately 182 million users across the world[[netflix](#)]. Nowadays, Netflix additionally evolved into a movie and series production company. Due to its business model of providing its customers with flexible monthly subscriptions to access their content and the growing number of competitors in the field of video streaming, Netflix aims to create high user engagement and retention. One major cornerstone in this regard is providing users with the content they are most interested in. In the past, customers went into their video rental store of trust and may be asked the owner for suggestions. Nowadays, the situation has changed users are not willing to spend much time searching for content on their own but rather receive interesting content automatically. What sounds like a simple task to do, indeed, is challenging because of the following reasons. First, the amount of possible movies and series to choose from is extended continuously and overwhelming list. And second, the number of users using the service is increasing from year to year. In this regard, Netflix's success is highly dependent on an accurate and well-engineered recommendation system. The Netflix lunch from 2006 – 2009 Netflix prize Offered a grand prize of \$ 1.000.000 USD to the team who could beat Netflix's own algorithm, Cinematch, by more than 10%, measured in RMSE.

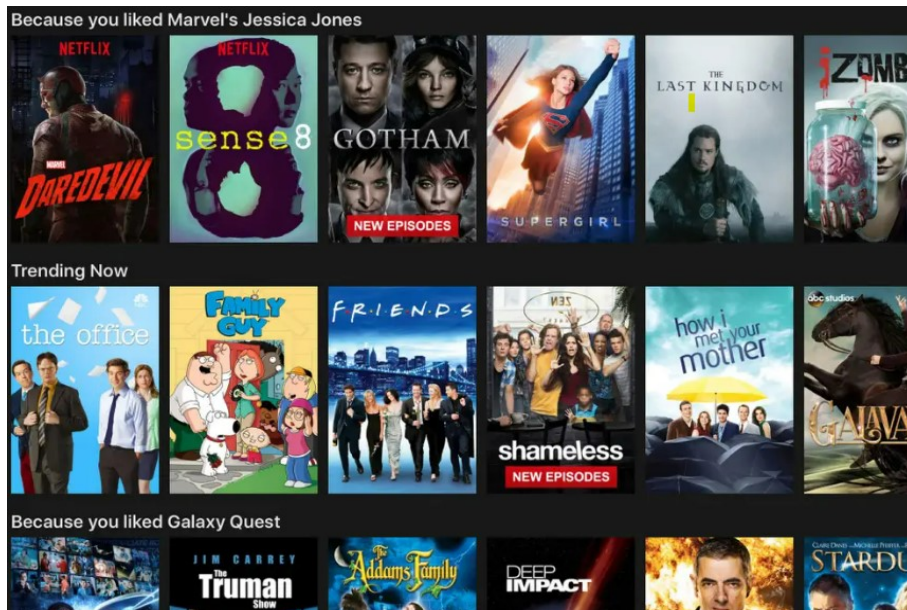


Figure 2.3: Netflix Recommendation System

### 2.7.2 Amazon Recommendation System

Amazon was also one of the pioneers in recommendation systems, especially in the commercial setting. During the early years, it was one of the few retailers that had the foresight to realize the usefulness of this technology. Originally founded as a book e-retailer, the business expanded to virtually all forms of products. Consequently, Amazon now sells virtually all categories of products such as books, CDs, software, electronics, and so on. The recommendations on Amazon are provided based on explicitly provided ratings, buying behavior, and browsing behavior. The ratings on Amazon are specified on a 5-star scale, with the lowest

rating being one star and the highest rating being five stars. The user-specific buying and browsing data can be easily collected when users are logged in with an account authentication mechanism supported by Amazon.

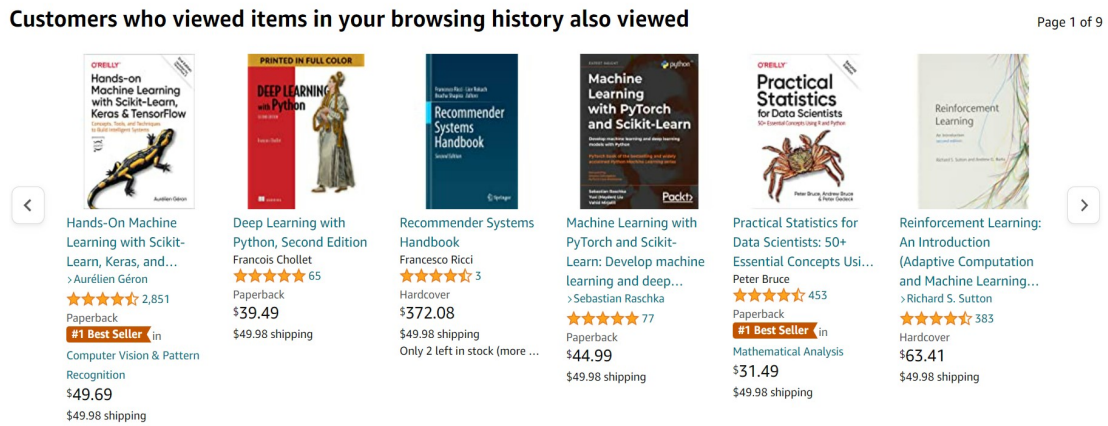


Figure 2.4: Amazon Recommendation System

### 2.7.3 Google News Personalization Recommendation System

System the Google News personalization system is able to recommend news to users based on their history of clicks. The clicks are associated with specific users based on identification mechanisms enabled by Gmail accounts. In this case, news articles are treated as items. The act of a user clicking on a news article can be viewed as a positive rating for that article.

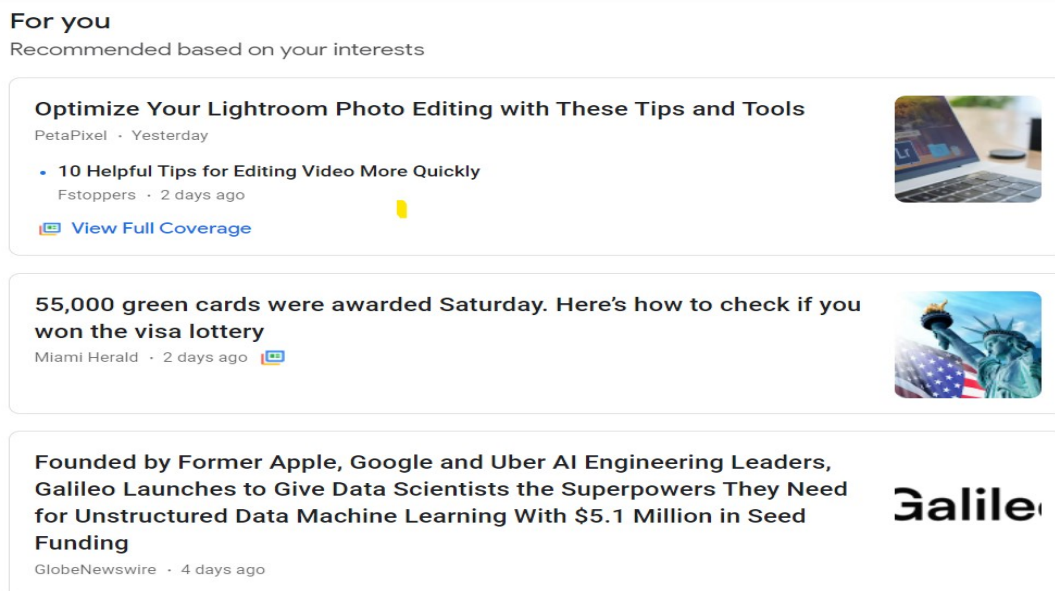


Figure 2.5: Google Recommendation System

### 2.7.4 Facebook Friend Recommendations System

Social networking sites often recommend potential friends to users in order to increase the number of social connections at the site. Facebook is one such example of a social networking Web site. This kind of recommendation has slightly different goals than a product recommendation. While a product recommendation directly increases the merchant's profit by facilitating product sales, an increase in the number of social connections improves the user's experience on a social network. This, in turn, encourages the growth of the social network. Social networks are heavily dependent on the growth of the network to increase their advertising revenues. Therefore, the recommendation of potential friends (or links) enables better growth and connectivity of the network. This problem is also referred to as link prediction in social network analysis. Such forms of recommendations are based on structural relationships rather than rating data.

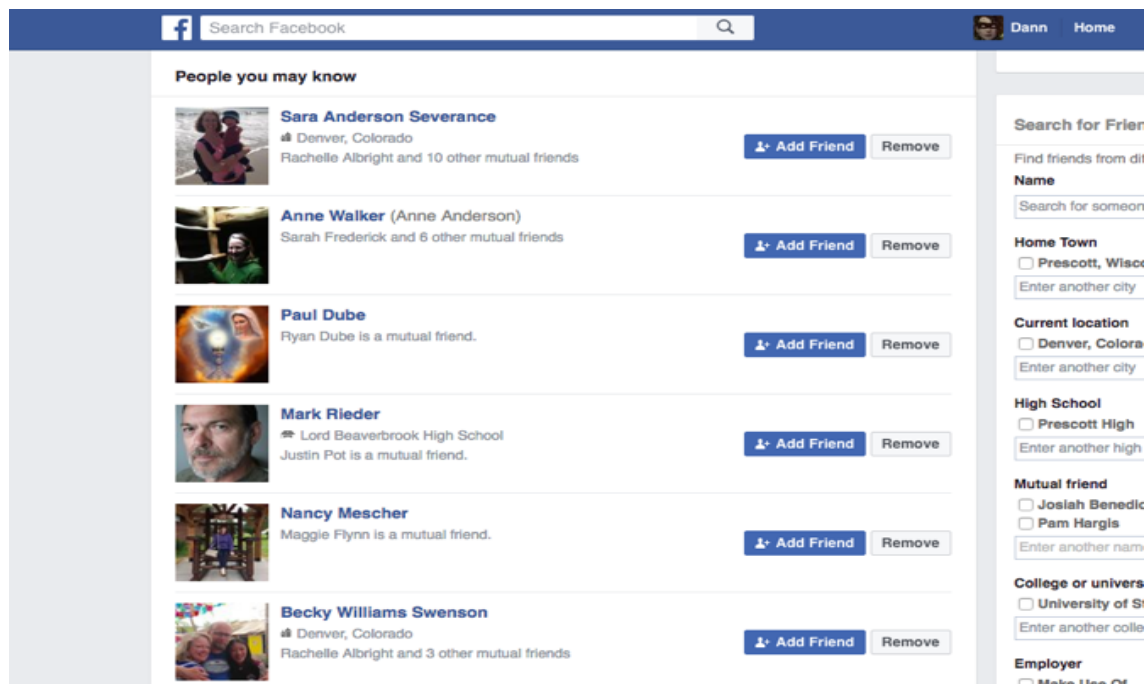


Figure 2.6: Google Recommendation System

## 2.8 Models of Recommendation Systems

There can be two types of recommendations **Personalized and Non-Personalized Recommendation Systems**.

Personalized Recommendation takes into consideration users previous history for rating and predicting items. On the other hand non-personalized recommendation systems recommend what is popular and relevant to all the users which can be a list of top-10 items for every new user.

### 2.8.1 Non-Personalized Recommendation System

In the case of Non-Personalized Recommendation Systems, the websites give a general item preference based on the feature evaluation of the items by different types of users. One such non-personalized recommendation is the news suggestion by online news agencies.

The non-personal recommendations do not have a direct impact on the users' buying habits and personal interests, so such types of systems have less concern in recommendation system research.

In other hands there is only one type of Non-Personalized Recommendation System method. Popularity-based recommendation systems .

1. **Popularity-Based Recommendation Systems** As the name suggests Popularity based recommendation system works with the trend. It basically uses the items which are in trend right now, the bestselling product or most popular product present in the market.

### 2.8.2 Personalized Recommendation System

The personalized recommendation systems rely on a user perspective and preferences, and the recommendation are based on the user behavior and interaction. The user interaction can be ratings of the items, searching and browsing history, items pages, add to cart, checkout, the time spent reviewing the items, clicks, likes ...etc.

As personalized recommendation systems analyze users' data, their purchases, rating, location, or profile, in that way, every user will get customized recommendations.

The personalized recommendation also considers the similarity with other users' behavior for better suggestions to increase their revenue.

1. **Collaborative-Based filtering:** Collaborative Filtering is an information filtering model that first appeared in the 1990s , the process of filtering or evaluating items through the opinions of other users. Collaborative Filtering is a model that constructs a user's preference database using the user's evaluation data to predict items that fit the user's taste, and then uses it for recommendation. Collaborative filtering There are two types of methods that are commonly used in collaborative filtering, Neighborhood-based Collaborative Filtering (Memory-Based) and Model-Based Collaborative Filtering

- (a) **Model-Based Filtering methods**

The recommendation system uses different machine learning and Artificial intelligence algorithms, for better prediction, taking into account the user's history, product features, and behaviors of other customers. Many of these cannot be computed directly, and have to be estimated using mathematical models. The model-based recommendation deals with the learning algorithms that examine the data and suggest the item. Different researchers have suggested a different

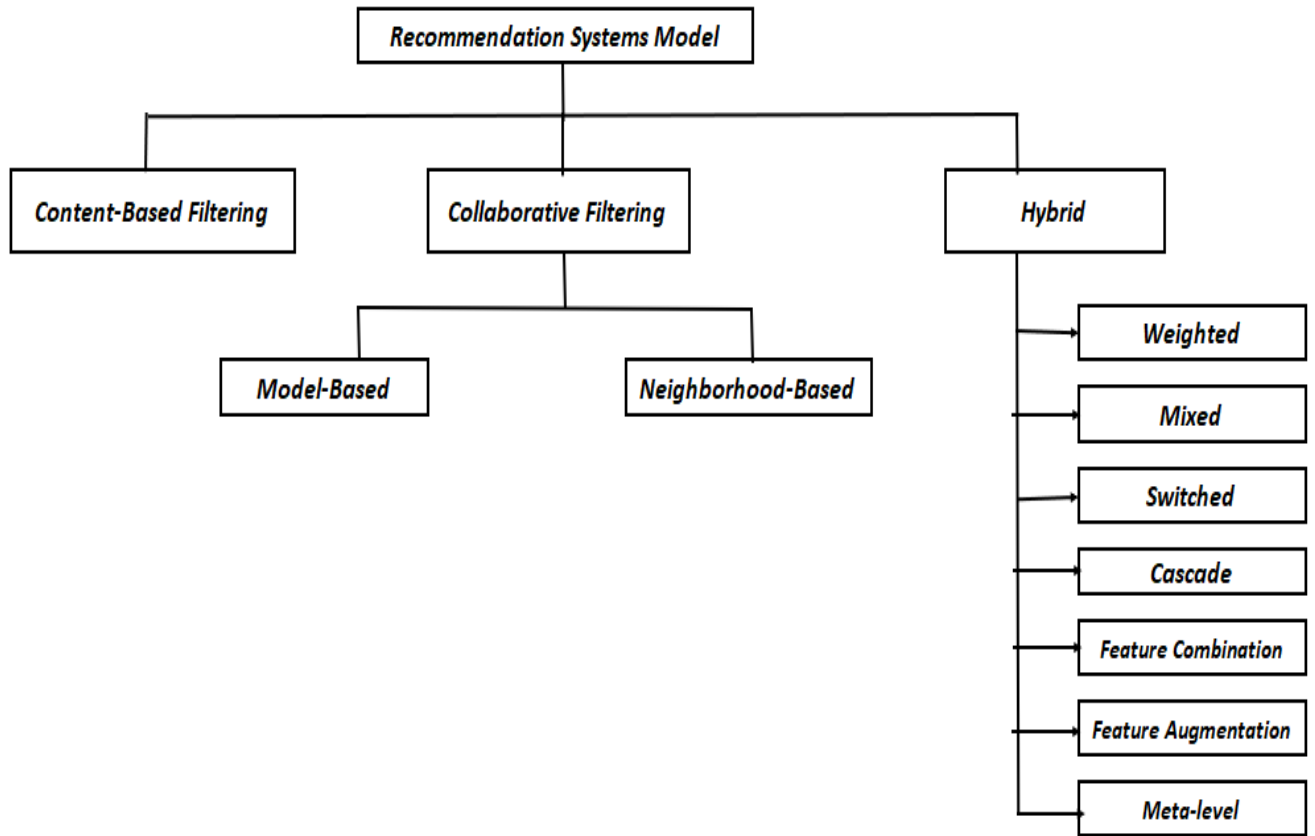


Figure 2.7: Recommendation Systems Models

mechanism to filter the suggestion. These mechanisms use nearest neighbor clustering, Bayesian probabilistic model, and neural networks for proper classification and refinement of information before providing suggestions to the customers. the main Advantages of Model-Based Collaborative Filtering are

- It helps to resolve the problem of sparsity and scalability.
- accuracy is better.

in addition the disadvantages of Model-Based are

- Implementation cost is high.
- There is trade-off between scalability and efficiency of model.
- Due to dimensionality reduction methods, it may loss valuable information

- (b) **Neighborhood-based methods** Neighborhood-based approach are also referred to as Memory-based collaborative filtering algorithms. These were among the earliest collaborative filtering algorithms, in which the ratings of user-item combinations are predicted on the basis of their neighborhoods. the Neighborhood-based approach is a Non-Parametric approach. The two methods that can be used to do this are called User-Based and item-based recommendation.

**Advantages of Neighborhood-based methods** the main advantages of Neighborhood-based methods are:

- **Simplicity:**

Neighborhood-based methods are intuitive and relatively simple to implement. In their simplest form, only one parameter (the number of neighbors used in the prediction) requires tuning.

- **Justifiability:**

Such methods also provide a concise and intuitive justification for the computed predictions. For example, in the item-based recommendation, the list of neighbor items, as well as the ratings given by the user to these items, can be presented to the user as a justification for the recommendation. This can help the user better understand the recommendation and its relevance, and could serve as a basis for an interactive system where users can select the neighbors for which greater importance should be given in the recommendation

- **Efficiency:**

One of the strong points of neighborhood-based systems is their efficiency. Unlike most model-based systems, they require no costly training phases, which need to be carried out at frequent intervals in large commercial applications. These systems may require pre-computing nearest neighbors in an offline step, which is typically much cheaper than model training, providing near-instantaneous recommendations. Moreover, storing these nearest neighbors requires very little memory, making such approaches scalable to applications having millions of users and items.

- **Stability:**

Another useful property of recommendation systems based on this approach is that they are little affected by the constant addition of users, items, and ratings, which are typically observed in large commercial applications. For instance, once item similarities have been computed, an item-based system can readily make recommendations to new users, without having to re-train the system. Moreover, once a few ratings have been entered for a new item, only the similarities between this item and the ones already in the system need to be computed.

**Disadvantages of Neighborhood-based**

- These systems are fully dependent on the rankings provided by the users.
- These systems are not able to handle the sparse data which result in performance degradation.
- Recommendation cannot be made for new customers and new products.
- System is not scalable.

## 2. Content-Based:

The system learns to recommend items that are similar to the ones that the user liked in the past. The similarity of items is calculated based on the features associated with the compared items. the similarity in the rating history of the users In content-based recommendation systems, the descriptive attributes of items are used to make recommendations. The term “content” refers to these descriptions.

In Content-Based methods, the ratings and user behavior of users are combined with the content information available in the items. For example, consider a situation where John has rated the movie Terminator highly, but we do not have access to the ratings of other users. However, the item description of Terminator contains similar genre keywords as other science fiction movies, such as Alien and Predator. In such cases, these movies can be recommended to John. In Content-Based methods, the item descriptions, which are labeled with ratings, are used as training data to create a user-specific classification or regression modeling problem. For each user, the training documents correspond to the descriptions of the items she has bought or rated. The class (or dependent) variable corresponds to the specified ratings or buying behavior. These training documents are used to create a classification or regression model, which is specific to the user at hand (or active user). This user-specific model is used to predict whether the corresponding individual will like an item for which her rating or user behavior is unknown.

Advantages of Content-Based Recommendation technique such as :

- **User Independence**

Content-Based Recommendation system builds a user profile only based upon the rating or purchase done by the user in the past. No neighbor is considered for building the profile of the user who has the same interest as of user.

- **Transparency**

Explanation facility of a content-based recommendation system is transparent to the user which means it provides explanations of the recommendations.

- **New Item**

It does not suffer from the first-rate problem which means if an item is not rated by any user, it is still able to recommend that item to the user.

And also disadvantages of content-Based Recommendation are:

- **Limited Content Analysis**

One of the shortcomings of a content-based recommendation system is limited content associated with the item in terms of a number of features and types. Proper differentiation cannot be done between the item’s user likes and items user dislikes if available data is insufficient. Representation sometimes is able to capture only certain aspects of user choice but not all. For example, Web pages, feature extraction techniques from the text completely overlook visual qualities and additional multimedia information.

- **Over-Specialization**

Content-based recommendation system does not have any essential method to explore something unpredicted. The system can recommend only those items which result in a high score while matching with the user profile

- **New User**

To make a recommendation system to learn about user preferences, sufficient ratings need to be collected. System is not able to provide reliable recommendations to the new users as no past data is available

### 3. Hybrid Recommendation Systems:

Hybrid Recommendation System has the potential to tackle both content and collaborative RSs in its term and combine the efforts of both to produce a system from their strategic advantage. Both filtering models feature limitations because the Content-Based Filtering model relies on metadata about the user's item, and Collaborative Filtering relies on the user's item rating, A Hybrid recommendation model was proposed to solve the limitations of both recommendation filtering models and improve the recommendation performance. The Hybrid Recommendation model is mainly designed to solve the sparsity Problem, the main goal of most studies dealing with the Hybrid recommendation model is to compensate for the lack of rating data by integrating the information of the Content-Based Filtering and Collaborative Filtering models.

The Hybrid recommendation model have in seven types: Weighted Hybridization, Switching Hybridization, Cascade Hybridization, Mixed Hybridization, Feature Combination, Feature-Augmentation, and Meta-level.

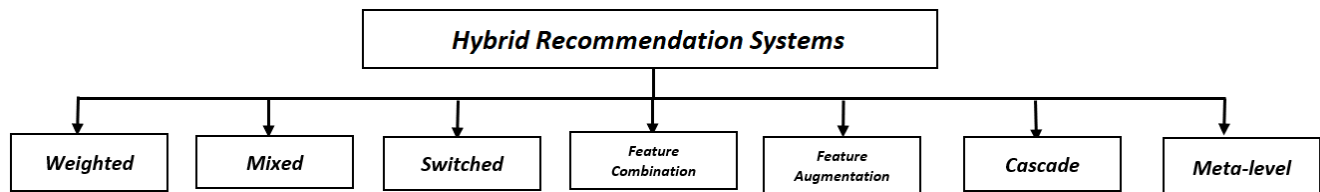


Figure 2.8: Hybrid Recommendation System

(a) **Weighted:**

The score of different recommendation components are combined numerically.

(b) **Mixed:**

Recommendations from different recommenders are presented together.

(c) **Switched:**

In this, system switches among various recommending methods depending on some conditions. For instance, a Colaborative-Based filtering method may shift to the content-based recommendation if the collaborative filtering method doesn't offer sufficient reliable results.

(d) **Feature Combination:**

Features derived from different knowledge sources are combined together and given to a single recommendation algorithm.

(e) **Feature Augmentation:** One recommendation technique is used to compute a feature or set of features, which is then part of the input to the next technique.

(f) **Cascade:**

Recommenders are given strict priority, with the lower priority ones breaking ties in the scoring of the higher ones.

(g) **Meta-level:**

One recommendation technique is applied and produces some sort of model, which is then the input used by the next technique.

4. **Demographic:** This type of system recommends items based on the demographic profile of the user, such as age, gender, Location, education, etc. for classifying groups



of users. For example, if it is given in a site that 70% of people of your age use this product and getting benefits. Many commercial sites and industries have adopted this strategy as it is not too complex and the implementation is also simple. In a Demographic-Based Recommendation System, proper market research is needed along with the cultural value of citizens in the specified region, accompanied by a short survey report to accumulate data for categorization.

5. **Community-based recommendation system:** The main idea in the community-based recommendation is to recommend to the user based on the preference of his friends, following the epigram ‘Tell me who your friends are, and I will tell you who you are’.
6. **Contextual recommendation system:** In this kind of recommendation system, the results of recommendation vary according to the context of the user. For example, in the temporal context, the clothes recommended in summer vary totally from those recommended in winter. For a social context, an example recommending a restaurant for a Saturday night with friends, which would vary from a restaurant for a lunch during the week with co-workers
7. **Knowledge-based system:** Recommend items based on specific domain knowledge about how certain item features meet users’ needs and preferences. Knowledge-Based Recommendation Systems are beneficial in the context of items that are not often purchased. Examples include items such as real estate, automobiles, tourism requests, financial services, or expensive luxury goods. In such cases, sufficient ratings may not be available for the recommendation process. As the items are bought rarely, and with different types of detailed options, it is difficult to obtain a sufficient number of ratings for a specific instantiation of the item at hand. A particular item may have attributes associated with it that correspond to its various properties, and a user may be interested only in items with specific properties. For example, cars may have several makes, models, colors, engine options, and interior options, and user interests may be regulated by a particular combination of these options. Thus, in these cases, the item domain tends to be complex in terms of its various properties, and it is hard to associate sufficient ratings with the large number of combinations at hand. Such cases can be addressed with knowledge-based recommendation systems.
8. **Location-Based Recommendation:**

Location-based recommendation systems have witnessed an increasing interest in recent years because of the increasing prevalence of mobile phones and other GPS-enabled devices. Location-based recommenders incorporate the location of users to provide relevant and precise recommendations. These can be a point of interest recommendations, such as restaurants, events in nearby locations, or posts and local trends on social media. For example, a traveling user may wish to determine the closest restaurant based on her previous rating history for other restaurants. Generally, the recommendation of places always has a location aspect built into it. An example of such a system is Foursquare which recommends various areas such as restaurants or nightlife venues. Two types of the spatial locality are common to such systems.

- (a) **User-Specific Locality:** The geographical location of a user has an essential role in her preferences.
- (b) **Item-specific locality:** The geographical location of an item (e.g., restaurant) might have an impact on the relevance of the item, depending on the current location of the user. Users are generally not willing to travel very far from their current location. This type of locality is referred to as a travel locality.

## 2.9 Evaluating Recommendation Systems

In order to evaluate the performance of the Recommendation Systems models, The quality of RSs may be measured in many ways, and different metrics can be applied.

The simplest way to evaluate the performance of a recommendation system is:

### 2.9.1 Accuracy

is one of the most fundamental measures through which recommendation systems are evaluated. Three types of accuracy may be considered:

**Prediction-based metrics , Decision-based metrics , Rank-based metrics.**

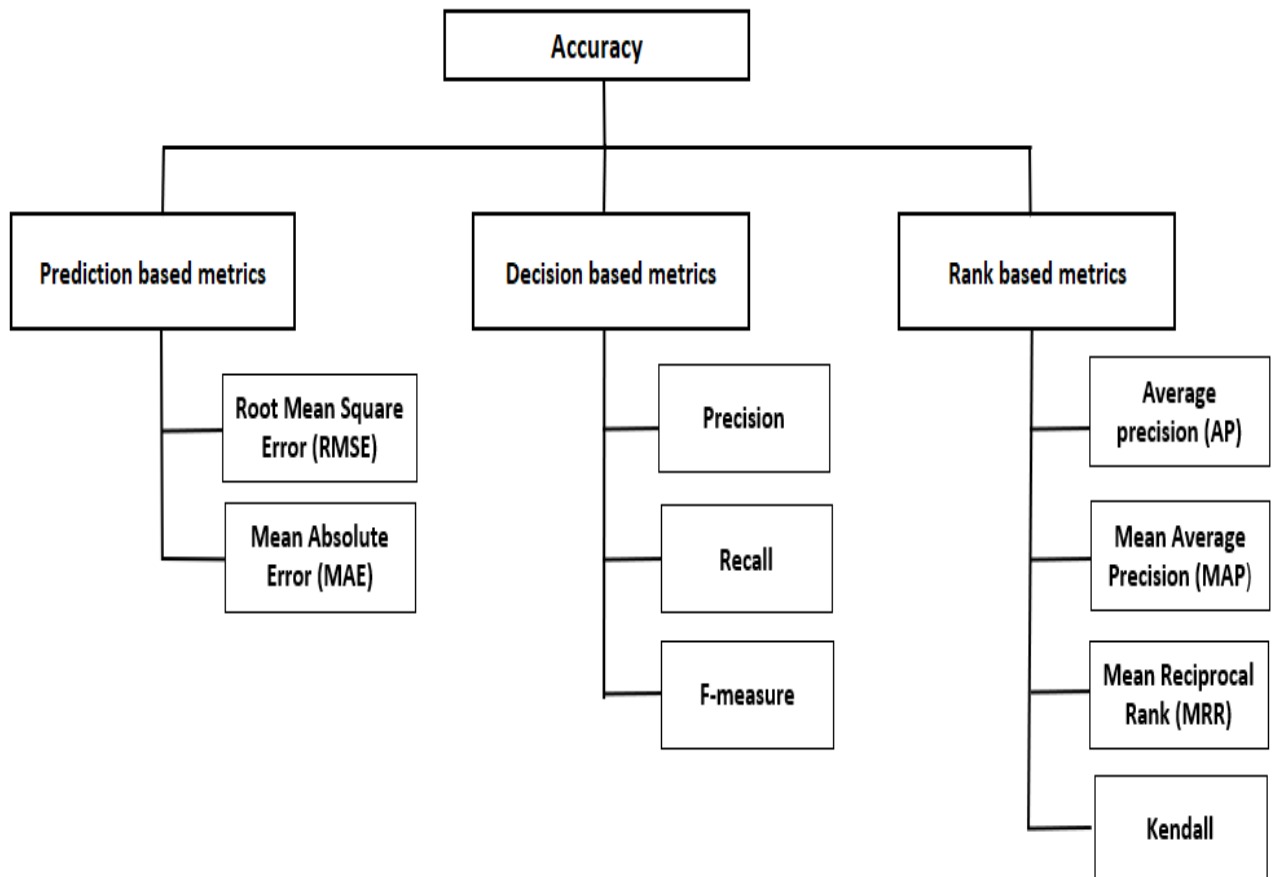


Figure 2.9: Accuracy metrics

### Prediction Based Metrics

Used to measure the difference between predicted rating  $\hat{r}_{ui}$  and the rating given by the user  $r_{ui}$ .

- **Root Mean Square Error (RMSE)**

it evaluates the difference between the ratings predicted  $\hat{r}_{ui}$  by the RS and the ratings given by the user  $r_{ui}$ . It also has become popular in the recommendation system field with the Netflix Challenge (Bell et al., 2007), reducing the RMSE amounts to increase the relevance and precision of the recommendations.

The RMSE is defined as follows :

$$RMSE = \sqrt{\frac{\sum_{(u,i,r) \in R} (\hat{r}_{ui} - r_{ui})^2}{|R|}}$$

Where  $R$  be the set of ratings ,normally a Test Set.

- **Mean Absolute Error (MAE)**

The Mean Absolute Error (MAE) Similar to  $RMSE$  is a metric that is used to detect the accuracy of the system by comparing the predicted ratings against the actual ratings of the items. The average absolute difference between the estimated and the user's true rating is termed as the mean absolute error. The relationship between the metrics measure and the system's performance is inverse, which implies better performance of the system ,The lower the  $MAE$ , the more accuracy the recommendation system predicts rating.

$$MAE = \frac{\sum_{(u,i,r) \in R} |\hat{r}_{ui} - r_{ui}|}{|R|}$$

Where  $R$  denotes the amount of user  $u$  rated items.,normally a Test Set,  $r_{ui}$  determines the actual rating that user  $u$  rates item  $i$  and  $\hat{r}_{ui}$  denotes the predicted rating of item  $i$  for user  $u$ .

### Decision based metrics

Decision-based metrics evaluates the top-N recommendations for a user. There are four different cases to take into account:

**Confusion Matrix** A confusion matrix is a table that is used to visualizes and summarizes the performance of a classification algorithm

**True Positive (TP)** The system recommends an item the user is interested in.

**False Positive (FP)** The system recommends an item the user is not interested in.

**True Negative (TN)** The system does not recommend an item the user is not interested in.

**False negative (FN)** The system does not recommend an item the user is interested in.

	Relevant	Not Relevant
Recommended	True Positive	False Positive
Not Recommended	False Negative	True Negative

Table 2.1: Confusion Matrix

- **Precision** Precision defines the segment of suggested items that are identical to the users' preferences in the testing dataset, It is also known as positive predictive value and measures the relevancy of the result. Assuming that, in the system, items ratings belong to 3–5 are positive ratings and 1 and 2 are negative ratings. Thus, when the item's real rating exists within 3–5 and the item's predicted rating is 3–5, then we can conclude that it is like the true positive (TP), whereas, when the item's real rating belongs to 1 and 2, but the item's predicted rating belongs to 3–5, then we can declare this as false positive (FP). Thus, the precision is computed as the ratio of the number of relevant recommended items to the number of all the recommended items, Higher the precision better is a recommendation .

$$Precision = \frac{TP}{TP + FP}$$

Where **TP+FP** refers to the total suggested items .

- **Recall**

The recall is the mean quantity of items of the testing dataset that exists among the ranked list from the training dataset. It also called sensitivity and measures the amount of truly relevant results that are identified. Based on the above precision's hypothesis, an item's real rating exists within 3–5 and the item's predicted rating belongs to 3–5, then, we can call it a true positive (TP), on the other hand, when the item's predicted rating belongs to 1–2, then it is classified as a false negative (FN) ,the equation of the recall computation as the ratio of number of relevant recommended items to the total number of relevant items.

$$Recall = \frac{TP}{TP + FN}$$

- **F-measure:**

It has been seen that, when the precision and recall cannot anticipate decent results, then  $F_\beta - measure$  is used as a weighted harmonic mean of the precision and recall to ensure better evaluation of the test system

The general formula of F-measure is given as follows :

$$F_{\beta} - measure = \frac{(\beta^2 + 1)(Precision.Recall)}{\beta^2.Precision + Recall}$$

$\beta$  is a parameter that controls a balance between Precision and Recall.

The most common  $F - measure$  is  $F_1 - measure$ , where  $\beta = 1$ .

$F_1 - measure$  the standard F-Measure, the harmonic mean of precision and recall.

$$F_1 - measure = \frac{2 \times (Precision + Recall)}{(Precision + Recall)}$$

### Rank-based metrics

A rank metric measures the ability of a RSs to estimate the correct order of items concerning the user's preference, which is called the measurement of rank correlation in statistics.

We need to be careful with these metrics as it makes no sense to compute them if a user appears only once in the test dataset. Rank-based metrics are sometimes hard to interpret

- **Average precision (AP)**

Average Precision (AP) equals:

$$AP = \frac{\sum_{k=1}^n (P(k) \times rel(k))}{NumberofRelevantItems}$$

Where  $P(k)$  is Precision at top-k, and  $rel(k)$  is an indicator function equaling 1 if the item at rank is a relevant item, and zero otherwise,  $n$  is the number of items in the recommendation list.

- **Mean Average Precision (MAP)**

Mean Average Precision (MAP) is defined as the precision (i.e. the percentage of relevant items among the first  $k$  recommendations) at the position  $k$  in the ranked results

$$MAP = \frac{\sum_{q=1}^Q AP(q)}{Q}$$

where  $Q$  is the whole number of the users involved in this recommendation system.

- **Mean Reciprocal Rank (MRR)**

The mean reciprocal rank (MRR) is a popular evaluation metric that is used produce a list of ranked items for queries by calculate the average of reciprocal of the rank in which the first correct item was retrieved at each prediction. It will match the recommended list with the predicted recommended list. It will give a higher score if the match is found at the top of the list. The MRR is computed as follows:

$$MRR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{rank_i}$$

Where  $rank_i$  is the rank of  $i$ th item.

- **Kendall-  $\tau$**

Kendall Rank correlation is a non-parametric test that measures the strength of dependence between two variables. If we consider two samples, a and b, where each sample size is n, we know that the total number of pairings with a b is  $n(n-1)/2$ . The following formula is used to calculate the value of Kendall rank correlation:

$$\tau = \frac{n_c - n_d}{\frac{1}{2}n(n-1)}$$

Where  $n_c$  Number of concordant.  $n_d$  Number of discordant.

### 2.9.2 Non Accuracy Metrics

**Coverage:** As the prediction accuracy of a recommendation system, especially in collaborative filtering systems, in many cases grows with the amount of data, some algorithms may provide recommendations with high quality, but only for a small portion of the items where they have huge amounts of data. There are two types of coverage, which are referred to as user-space coverage and item-space coverage, respectively.

**Novelty:** The novelty of a recommendation system evaluates the likelihood of a recommendation system to give recommendations to the user that they are not aware of, or that they have not seen before]. Unseen recommendations often increase the ability of the user to discover important insights into their likes and dislikes that they did not know previously. This is more important than discovering items that they were already aware of but have not rated.

**Confidence :** The recommendation can be defined as the system's trust. Estimating ratings is an inexact process that can vary significantly with the specific training data. Furthermore, the algorithmic methodology might also considerably impact the predicted ratings. This always leads to uncertainty in the user about the accuracy of the predictions. Many recommendation systems may report ratings together with confidence estimates. For example, a confidence interval on the range of predicted ratings may be provided. In general, recommendation systems that can accurately recommend smaller confidence intervals are more desirable because they bolster the user's trust in the system. Two algorithms that use the same method for reporting confidence can measure how well the predicted error matches these confidence intervals. For example, if two recommendation systems provide 95% confidence intervals for each rating, one can measure the absolute width of the intervals reported by the two algorithms. The algorithm with the smaller confidence interval width will win if both algorithms are correct (i.e., within the specified intervals) at least 95% of the time on the hidden ratings. If one of the algorithms falls below the required 95% accuracy, then it automatically loses. Unfortunately, if one system uses 95% confidence intervals and another uses 99% confidence intervals, it is impossible to compare them meaningfully. Therefore, it is possible to use such systems only by setting the same confidence level in both cases.

**Trust:** in the recommendation system is How much the user trusts the recommendation, this aspect varies between users. Some users trust recommendations that offer them new items, while others prefer having some of their preferred items in the recommendation list in order to trust the system.

**Diversity :** The notion of diversity implies that the set of proposed recommendations within a single recommended list should be as diverse as possible, usually users prefer systems with high diversity as they can explore new horizons of items. for example If all three restaurants (eg. French, Thai, Italian, Indian, and Chinese) are of a particular cuisine and contain similar dishes, then there is little diversity in the recommendations.

**Serendipity :** is a measure of the level of surprise in successful recommendations. In other words, recommendations need to be unexpected, as serendipitous recommendation

helps the user find a surprisingly interesting item he might not have otherwise discovered. for example, Consider the case where a particular user frequently orders from Indian restaurants. The recommendation of a new Pakistani restaurant to that user might be novel if that user has not eaten at that restaurant earlier. However, such a recommendation is not serendipitous because it is well known that Indian and Pakistani food is almost identical. On the other hand, if the recommendation system suggests a new Ethiopian restaurant to the user, then such a recommendation is serendipitous because it is less obvious.

**Robustness and Stability :** The robustness is the ability to make good recommendations in the presence of fake information, As more people rely on recommendation systems to guide them through the item space, influencing the system to change the rating of an item may be profitable to an interested party, For example, an owner of a hotel may wish to boost the rating for their hotel. This can be done by injecting fake user profiles that rate the hotel positively or by injecting fake users that rate the competitors negatively. Such attempts to influence the recommendation are typically called attacks. An attack occurs when an agent tries to influence the behavior of the recommendation system. To measure the actual influence of an attack on the outcome of a recommendation system, Robustness measures the shift in the overall accuracy before and after an attack; the stability measure expresses the attack-induced change of the ratings predicted for the attacked items.

**Scalability :** In recent years, the sizes of the data sets continue to increase over time. As a result, it has become increasingly essential to design recommendation systems that can perform effectively and efficiently in the presence of large amounts of data. Scalability is typically measured by experimenting with growing data sets, showing how the speed and resource consumption behaves as the task scales up, A variety of measures are used for determining the scalability. Such as Training time: Most recommendation systems require a training phase, which is separate from the testing phase, Therefore, as long as the training time is of the order of a few hours, it is quite acceptable in most real settings. And Prediction time Once a model has been trained, it is used to determine the top recommendations for a particular user. Furthermore, the computation increases as the number of users increases. User-based methods work fine for thousands of users, but scalability gets to be a problem when we have a million users. This is where the importance of scalability has become particularly great in recent years because of the increasing importance of the “big data” paradigm.

**Utility :** Many e-commerce websites employ a recommendation system in order to improve their revenue by, e.g., enhancing cross-sell. In such cases the recommendation systems can be judged by the revenue that it generates for the website, Utility can be measured cleanly from the perspective of the recommendation system owner.

**Privacy:** Providing personal data to the recommendation systems may increase the system performance but may lead to problems of data privacy and security. Users are reluctant to feed data into recommendation systems that suffer from the data privacy issue

## 2.10 Challenges in Recommendation Systems

This section explores the main challenges that face the Recommendation Systems.

### 2.10.1 Cold Start

Cold starts can be classified into two distinct subsets: item cold starts and user cold starts. Whenever a new item is displayed on an the web application <sup>24</sup>, it goes through the item cold start, and this means that there are no rating due to the lack of user interaction. If there are not enough user interactions, RSs can't retrieve the user preferences, which leads to the RSs being unable to display any recommendation. The cold-start behavior occurs when a user creates an account for the first time and does not have any item preferences or history available to base recommendations. The solution to the problem can be provided by incorporating the global relevance where the system can retrieve the item based on its popularity.

### 2.10.2 Sparsity

The sparsity problem in recommendation systems remains that only a few values in the user-item matrix are known, while most items are missing. In many large-scale applications, both the number of items and the number of users are large, which leads to the sparsity problem, where the user-item matrix becomes a sparse matrix<sup>25</sup>. in which most of the rating are zero . Generally, where the majority of the users do not rate most of the items, consequently, the rating matrix becomes very sparse. Due to this, the data sparsity problem arises that declines the chances of finding a set of users with similar ratings.

### 2.10.3 Scalability

To provide a recommendation is challenging when the recommendation systems deals with huge datasets. Since some recommendation algorithms work extremely well with small datasets, but its efficiency may not work for large datasets.

---

<sup>24</sup>**A web application** (or web app) is application software that runs on a web browser.

<sup>25</sup>**Sparse matrix** is a matrix in which most of the elements are zero



## Chapter 3

# Neighborhood Based Recommendation System

The collaborative filtering approach recommends items based on the opinions of other users. Typically, by computing the similarity of users, a set of “nearest neighbor” users whose known preferences correlate significantly with a given user. Preferences for unseen items are predicted for the user based on a combination of the preferences known from the nearest neighbors. Thus, in this approach, users share their preferences regarding each item that they consume so that other users of the system can better decide which items to consume (Herlocker et al. 1999). The collaborative filtering approach is the most successful and widely adopted recommendation technique to date.

As mentioned, the collaborative filtering approach utilizes user preferences to generate recommendations. Several different techniques have been proposed for collaborative filtering recommendations, including neighborhood-based, model-based ,latent factors , Due to space limitation, we will only review the neighborhood-based collaborative filtering techniques since they are the most prevalent algorithms used in collaborative filtering for recommendation. The neighborhood-based CF algorithm use the entire user-item database to generate a prediction. Every user is part of a group of people with similar interests. The neighborhood-based CF algorithm, a prevalent memory-based CF algorithm, uses the following steps:

- **Neighborhood Similarity** Computing the similarities between users or items.
- **Produce a prediction  $\hat{r}_{ui}$  of user  $u$  to item  $i$ .**
- **Select the neighborhood of the active user.**

### 3.1 Types of Neighborhood-based collaborative filtering

There are two main type of Neighborhood-based collaborative filtering

- User-based neighborhood collaborative filtering.
- Item-based neighborhood collaborative filtering.

#### 3.1.1 User-based neighborhood collaborative filtering:

User-based neighborhood collaborative filtering is the most successful method for building recommendation systems to date, and is extensively used in many commercial, method first seek who shared the same rating pattern with the target user and then use the ratings of the similar users to estimate the predictions and then recommendation.

This method calculate the rating for a yet unrated item of the active user, average the ratings of the nearest neighbors about this particular item. In order to generate more accurate predictions, rating values of neighbor are assigned with weights according their similarity to the target user.

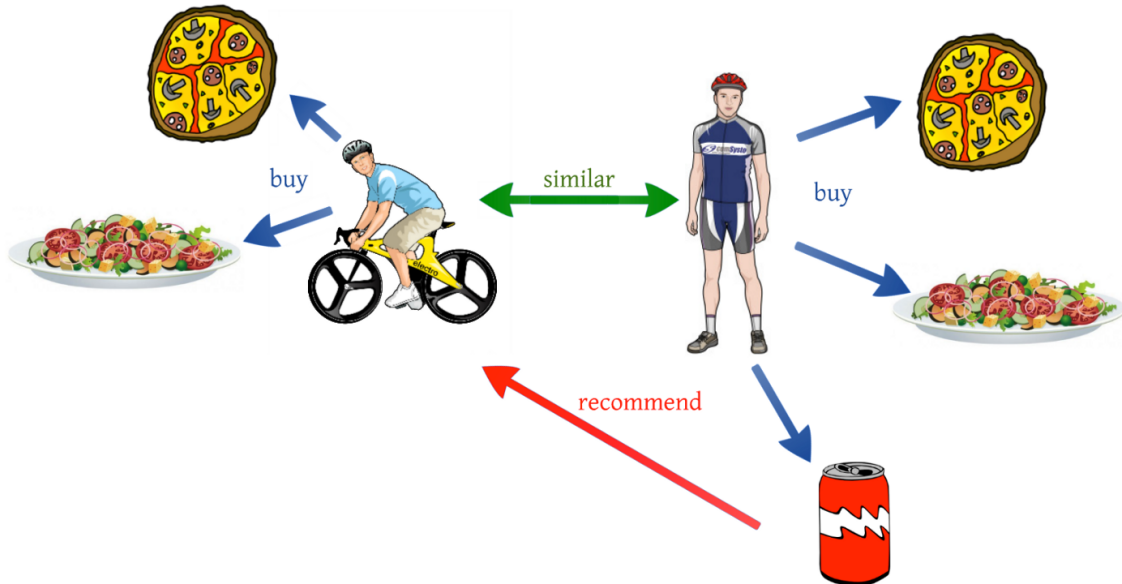


Figure 3.1: User-based neighborhood collaborative filtering

The main idea behind user-based filtering is that if we are able to find users that have bought similar items in the past, they are more likely to buy similar items in the future too. This method for generating the more precise prediction, weights allocate to the values of neighbor based their similarity to the active user.

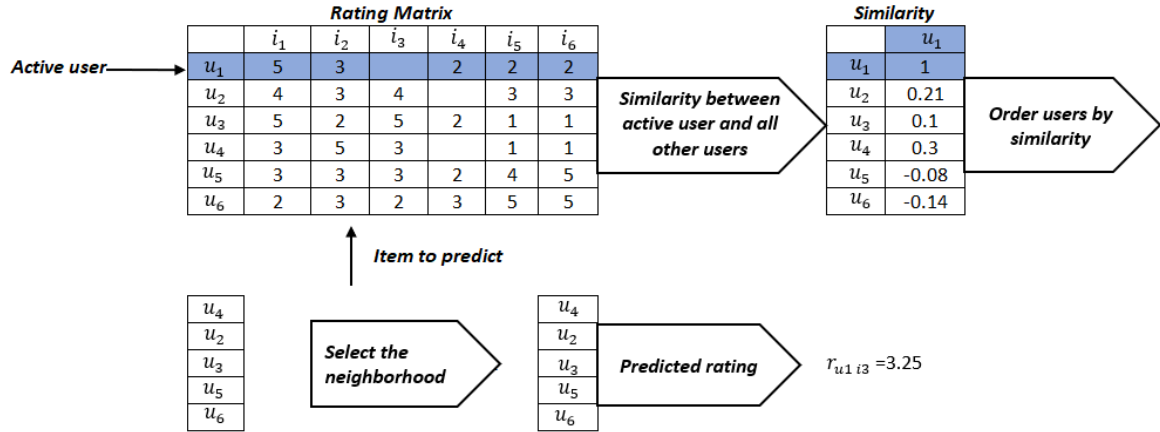


Figure 3.2: The user-Based pipeline

### 3.1.2 Item-based neighborhood collaborative filtering:

Item-based neighborhood collaborative filtering is the transpose of user-based nearest neighbor algorithms that produce predictions based on similarities between items.

An item-based method exploits the similarity among the items.

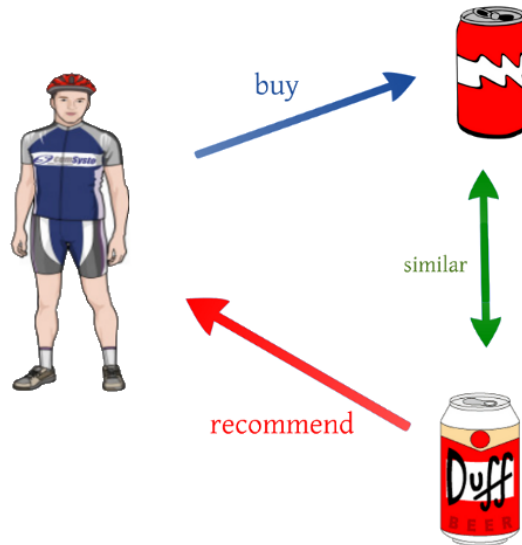


Figure 3.3: Item-based neighborhood collaborative filtering

This method looks into the set of items that a user has rated to (the ones the user already buys) and computes the similarity among the target item (to decide whether is worth to recommend it to the user or not).

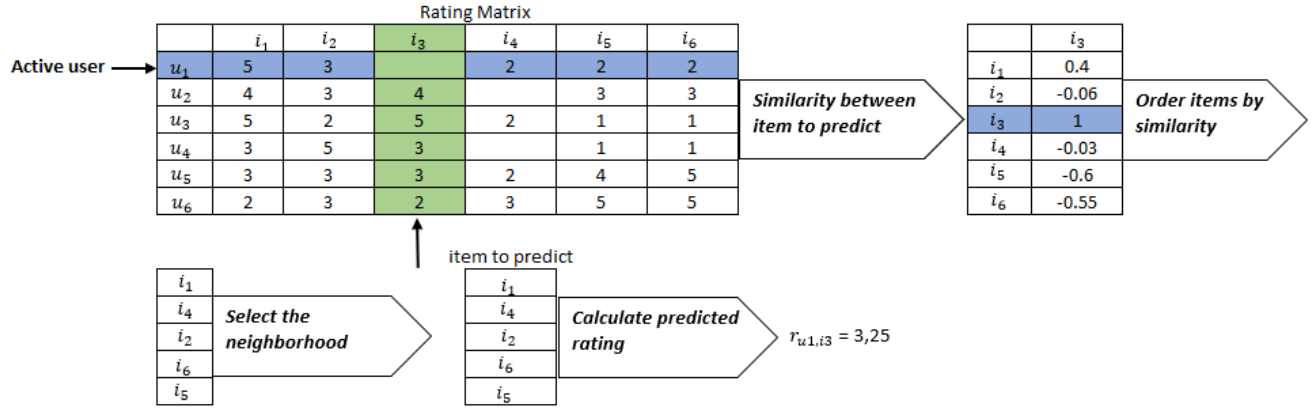


Figure 3.4: The item-Based pipeline

### 3.2 User-based VS Item-based Recommendation

When choosing between the implementation of a user-based and an item-based neighborhood recommendation system, five criteria should be considered:

- Accuracy:** The accuracy of neighborhood recommendation methods depends mostly on the ratio between the number of users and items in the system. As the similarity between two users in user-based methods, which determines the neighbors of a user, is normally obtained by comparing the ratings made by these users on the same items, an item-based method usually computes the similarity between two items by comparing ratings made by the same user on these items. In general, in cases where the number of users is much greater than the number of items, such as large commercial systems like Amazon.com, item-based methods can therefore produce more accurate recommendations. Likewise, systems that have fewer users than items, e.g., a research paper recommender with thousands of users but hundreds of thousands of articles to recommend, may benefit more from user-based neighborhood methods.
- Efficiency:** The memory and computational efficiency of recommender systems also depends on the ratio between the number of users and items. Thus, when the number of users exceeds the number of items, as is most often the case, item-based recommendation approaches require much less memory and time to compute the similarity weights (training phase) than user-based ones, making them more scalable.
- Stability:** The choice between a user-based and an item-based approach also depends on the frequency and amount of change in the users and items of the system. If the list of available items is fairly static in comparison to the users of the system, an item-based method may be preferable since the item similarity weights could then be computed at infrequent time intervals while still being able to recommend items to new users.
- Justifiability:** An advantage of item-based methods is that they can easily be used to justify a recommendation. Hence, the list of neighbor items used in the prediction, as well as their similarity weights, can be presented to the user as an explanation of the recommendation. By modifying the list of neighbors and/or their weights, it then becomes possible for the user to participate interactively in the recommendation process. User-based methods, however, are less amenable to this process because the active user does not know the other users serving as neighbors in the recommendation.

- **Serendipity:** In item-based methods, the rating predicted for an item is based on the ratings given to similar items. Consequently, recommendation systems using this approach will tend to recommend to a user items that are related to those usually appreciated by this user. user-based approaches are more likely to make serendipitous recommendations This is particularly true if the recommendation is made with a small number of nearest-neighbors.

### 3.3 Rating Matrix

User/Item	$i_1$	$i_2$	...	$i_j$	...	$i_n$
$u_1$	$r_{11}$	$r_{12}$	...	$r_{1j}$	...	$r_{1n}$
$u_2$	$r_{21}$	$r_{22}$	...	$r_{2j}$	...	$r_{2n}$
$u_3$	$r_{31}$	$r_{32}$	...	$r_{3j}$	...	$r_{3n}$
.	.	.	...	.	...	.
.	.	.	...	.	...	.
.	.	.	...	.	...	.
$u_i$	$r_{i1}$	$r_{i2}$	...	$r_{ij}$	...	$r_{in}$
.	.	.	...	.	...	.
.	.	.	...	.	...	.
.	.	.	...	.	...	.
$u_m$	$r_{m1}$	$r_{m2}$	....	$r_{mj}$	...	$r_{mn}$

Table 3.1: User-Item Rating Matrix

Let  $R$  be a user-item matrix of size  $m \times n$ , where  $n$  is the total number of the items, while  $m$  is the total number of the users, each row represents a user, and each column represents one single item, each cell would represent the rating  $r_{ui}$  given by user  $u$  to the item  $i$  or null if the user did not rate the item yet, In most cases, this matrix is sparse because users do not normally rate all the items in the data set.

### 3.4 Explicit and Implicit Ratings

For the purpose of gathering the preference of users, the system shall collect ratings from the user's interaction. The ratings can be either explicit or implicit.

The explicit ratings, the most prominent form of such feedback in the literature is the user-provided ratings, e.g., on a [1-5] scale often displayed as "stars". On the implicit ratings, we do not ask users to give any ratings, and we just observe their behavior the implicit ratings can be gathered using the user's feedback and interaction, An example of this is keeping track of what a user clicks on in the online New York Times.

Another implicit rating is what the users buy. Nonetheless, explicit ratings are not always readily available as many users may be reluctant to rate items. Implicit ratings are often readily available since it is mainly concerned with modeling implicit behavior such as user clicks. Many recommendation systems are centered on implicit ratings, which indirectly reflect users' opinions through observing user behavior.

There are diverse forms of implicit ratings on the Web pages: purchase history, browsing history, watches, and even mouse movements, the clicks on the link, time spent looking at a page, and repeated visits referring a page to others. Also, of example on Music players: what the person plays skipping tunes number of times a tune is played, and others like, save, share. IMDB also collects ratings ranging from one to ten stars for movies, and YouTube provides the thumbs-up and thumbs-down buttons for users to show their preferences.

It is apparent that gathering explicit ratings requires users to indicate their interests proactively. For example, a user that purchased many books by the same author probably likes that author. Note that implicit ratings are inherently noisy. We can only guess their preferences and true motives. A user who watches a movie does not necessarily indicate a positive view of that movie. Explicit ratings struggle mainly from the users who are lazy and do not rate items, and also, they may lie or give only partial information.

### 3.5 Similarity measures in collaborative filtering:

Similarity computation between items or users is a critical step in Neighborhood-based collaborative filtering algorithms.,

For item-based CF algorithms, the basic idea of the similarity computation between item  $i$  and item  $j$  is first to work on the users who have rated both of these items and then apply a similarity computation to determine the similarity  $s_{ij}$ .

For a user-based CF algorithm, we first calculate the similarity,  $s_{uv}$  between the users  $u$  and  $v$  who have both rated the same items. There are many different methods to compute similarity or weight between users or items.

#### 3.5.1 Cosine similarity:

Usually, the cosine similarity metric is used to estimate the similarity between two users or items. This technique presents a user as a vector of ratings rated by himself and an item as a vector of ratings rated by the set of users. The cosine between two vectors representing two users (or items) indicates the similarity value between each other. A value close to 1 indicates that it exists a strong correlation between the two variables.

A value close to 0 indicates that there is no correlation (independent variables). the cosine measure for users and items, calculating the Cosine Vector (CV) (or Vector Space) similarity between these vectors indicate the distance of them to each other :

**Cosine user similarity**

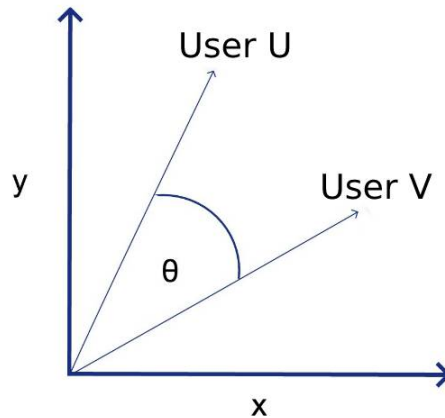


Figure 3.5: User Cosine Similarity Measure

Cosine similarity is measured by looking at the angles between the rating vectors, the smaller the angle between the users' vectors, the more similar the users.

$$\text{cosine}(u, v) = \frac{\sum_{i \in I_{uv}} r_{ui} r_{vi}}{\sqrt{\sum_{u \in I_u} r_{ui}^2} \sqrt{\sum_{v \in I_v} r_{vi}^2}}$$

where  $I_u$  and  $I_v$  denote the sets of items rated by users  $u$  and  $v$ , respectively, and  $I_{uv}$  denotes the set of items commonly rated by both  $u$  and  $v$ .  $r_{ui}$  and  $r_{vi}$  are the ratings values on item  $i$  given by users  $u$  and  $v$ , respectively

For example, this is the rating matrix.

	$i_1$	$i_2$
$u_1$	3	5
$u_2$	4	1
$u_3$	2	5

Table 3.2: **Example User-Item Rating Matrix**

We will look at content is to see the rows of the rating matrix as vectors in space, and then look at the angle between them. the we plot the data iton the coordinate system ( This works for more than two items as well) From the angles between the vectors, it's easy to see that  $u_1$  and  $u_3$  are much closer than  $u_2$  and  $u_1$  when it comes to ratings. Because of that, you can assume that the preferences of  $u_3$  and  $u_1$  are more similar.

#### Cosine item similarity

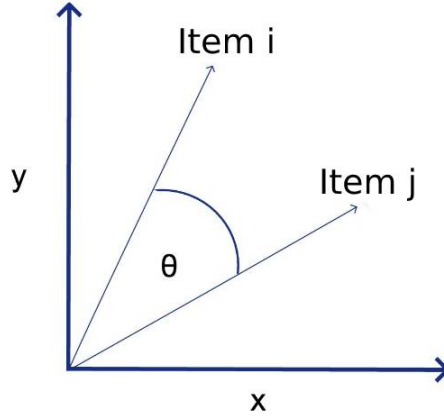


Figure 3.6: Item Cosine Similarity Measure

$$\text{cosine}(i, j) = \frac{\sum_{u \in U_{ij}} r_{ui} r_{uj}}{\sqrt{\sum_{u \in U_i} r_{ui}^2} \sqrt{\sum_{u \in U_j} r_{uj}^2}}$$

Where  $U_i$  and  $U_j$  denote the sets of users who rated the items  $i$  and  $j$ , respectively, and  $U_{ij}$  represents the set of users who rated both items  $i$  and  $j$ .  $r_{ui}$  and  $r_{uj}$  are the rating values assigned by the same user  $u$  on items  $i$  and  $j$ , respectively

A shortcoming of this measure is that it does not examine the differences in the mean and variance of the ratings made by user's  $u$  and  $v$ .

### 3.5.2 Adjusted cosine vector:

The Adjusted cosine measure calculates the correlation value between two users.

$$ACosine(u, v) = \frac{\sum_{i \in I_{uv}} (r_{ui} - \bar{r}_i)(r_{vi} - \bar{r}_i)}{\sqrt{\sum_{i \in I_{uv}} (r_{ui} - \bar{r}_i)^2} \sqrt{\sum_{i \in I_{uv}} (r_{vi} - \bar{r}_i)^2}}$$

Where  $I_{uv}$  denotes the set of items commonly rated by both  $u$  and  $v$ .  $\bar{r}_i$  denotes the average ratings on  $i$ .  $r_{ui}$  and  $r_{vi}$  denote, respectively, the ratings of the user  $u$  and  $v$  on the item  $i$ . The Adjusted cosine measure calculates the correlation value between two items.

$$ACosine(i, j) = \frac{\sum_{u \in U_{ij}} (r_{ui} - \bar{r}_u)(r_{uj} - \bar{r}_u)}{\sqrt{\sum_{u \in U_{ij}} (r_{ui} - \bar{r}_u)^2} \sqrt{\sum_{u \in U_{ij}} (r_{uj} - \bar{r}_u)^2}}$$

Where  $U_{ij}$  represents the set of users who rated both items  $i$  and  $j$ .  $\bar{r}_u$  denotes the average ratings by  $u$ . We mention that  $r_{ui}$  and  $r_{uj}$  are the ratings of user  $u$  on items  $i$  and  $j$ , respectively.

### 3.5.3 Pearson correlation:

This measure was proposed by Karl Pearson (Pearson, 1895) to measure linear relationships and became widely used in statistics field.

Pearson's correlation, measures the linear correlation between two vectors of ratings. PCC formula returns a value between -1 and 1, where:

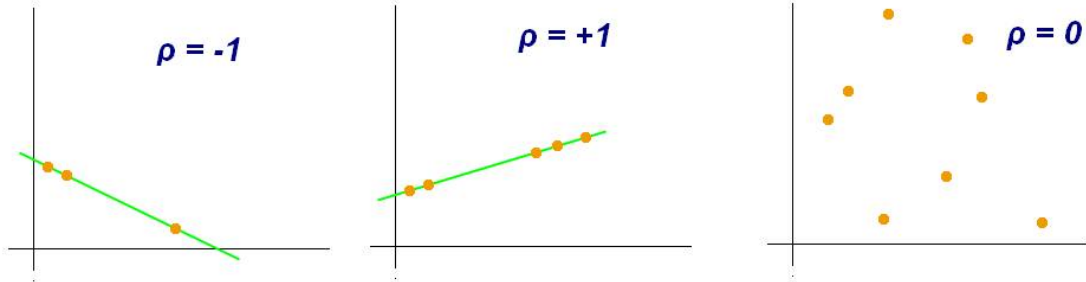


Figure 3.7: Pearsons correlation

- 1 indicates a strong positive correlation,
- 1 indicates a strong negative correlation
- 0 indicates no correlation at all

Pearson's correlation (PC) is a well-known metric that compares ratings, where the effects of mean and variance have been eliminated, is the Pearson Correlation (PC) similarity calculates the similarity between two users  $u$  and  $v$ .

$$PCC(u, v) = \frac{\sum_{i \in I_{uv}} (r_{ui} - \bar{r}_u)(r_{vi} - \bar{r}_v)}{\sqrt{\sum_{i \in I_{uv}} (r_{ui} - \bar{r}_u)^2} \sqrt{\sum_{i \in I_{uv}} (r_{vi} - \bar{r}_v)^2}}$$

Where  $I_{uv}$  denotes the set of items commonly rated by both  $u$  and  $v$ .  $\bar{r}_u$  and  $\bar{r}_v$  denote the average ratings of the users  $u$  and  $v$  on item  $i$  in  $I_{uv}$ , respectively.  $r_{ui}$  and  $r_{vi}$  are ratings of users  $u$  and  $v$  on the same item  $i$ . this Formula calculates the similarity between two items  $i$  and  $j$ :



$$PCC(i, j) = \frac{\sum_{u \in U_{ij}} (r_{ui} - \bar{r}_i)(r_{uj} - \bar{r}_j)}{\sqrt{\sum_{u \in U_{ij}} (r_{ui} - \bar{r}_i)^2} \sqrt{\sum_{u \in U_{ij}} (r_{uj} - \bar{r}_j)^2}}$$

let  $U_{ij}$  denote the set of users who provided a rating for both  $i$  and  $j$ , i.e  $u(i, j) = u(i) \cap u(j)$ ,  $U_{ij}$  represents the set of users who rated both items  $i$  and  $j$ ,  $\bar{r}_i$  and  $\bar{r}_j$  denote the average ratings on  $i$  and  $j$  in  $U_{ij}$ , respectively.  $r_{ui}$  and  $r_{uj}$  are ratings of user  $u$  on items  $i$  and  $j$ .

### 3.5.4 The Jaccard coefficient:

The Jaccard index, denoted by  $J$ , computes the similarity of two sets. The Jaccard coefficient between two finite sets is defined as the intersection's cardinality divided by the union's cardinality. It measures the ratio of the number of elements shared between the two sets to the total number of elements in both sets.  $J$  index takes a value between 0 and 1, the closer the index to 1, the more similar the two vectors. Calculates the Jaccard index for users similarity of two vectors  $u$  and  $v$ , while  $u$  and  $v$  can be users (set of ratings assigned by the same user).

$$J(u, v) = \frac{|u \cap v|}{|u \cup v|}$$

Calculates the Jaccard index for items similarity of two vectors  $i$  and  $j$ , while  $i$  and  $j$  can be items (set of ratings assigned by the same item).

$$J(i, j) = \frac{|i \cap j|}{|i \cup j|}$$

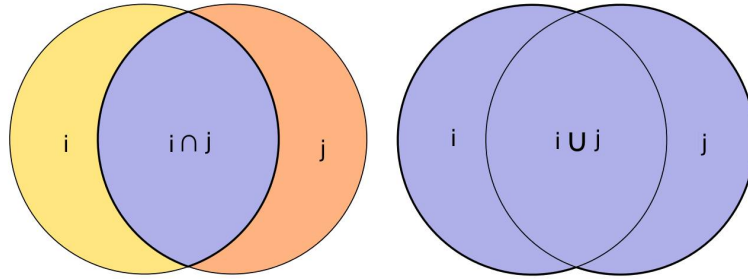


Figure 3.8: Jaccard coefficient

### 3.5.5 Conditional probability-based similarity:

[Karypis (2001)] proposed the conditional probability-based metric as a similarity metric for item-based CF Top-N recommendations.

The similarity between two items  $i$  and  $j$  is simply the probability of purchasing (rating) one given that the other has already been purchased. Thus the probability of purchasing  $a$  given that  $b$  has been purchased is determined as the number of users that purchased both items divided by the total number of users that purchased  $b$ . Note that this metric gives asymmetric similarities since  $(P(i|j) \neq P(j|i))$ . The similarity of  $i$  to  $j$  is given by:

$$Sim(i, j) = P(i|j) = \frac{freq(i, j)}{freq(j)}$$

According to the [Deshpande and Karypis (2004)], one of the shortcomings of an asymmetric metric is that each item tends to have high conditional probabilities concerning the

most favored items. To reduce this shortcoming, the following form of the conditional probability is presented by the followed formulas :

$$Sim(a, b) = P(a|b) = \frac{\sum_{u: R_{u,b} > 0} R_{u,b}}{Freq(a), (Freq(b)^\alpha)}$$

where,  $\alpha \in [0, 1]$  and  $Freq(a)$  indicates the number of users that have a transaction on item  $i$  in the training data and  $R(u, b)$  is the  $(u, b)$  element in the normalized user-item matrix.

### 3.5.6 Kendall's tau

As Spearman's rank correlation coefficient, Kendall's Tau ( $\tau$ ) evaluates statistical monotone relationships based on the ranks of the data. The value returned by  $\tau$  ranges from **-1** (as the rank of one variable increases, the other one decreases) to **1** (the ranks of both variables increase together). At the same time, **0** indicates no relationship between the two variables. This measure is mainly based on counting the concordant pairs (ordered in the same way) and discordant pairs (ordered differently). Kendall's tau for computing the association strength between two vectors of ratings is defined as:

$$\tau = \frac{c - d}{c + d}$$

Where  $c$  is the number of concordant pairs, and  $d$  is the number of discordant pairs.

### 3.5.7 Euclidean distance

The Euclidean distance denoted by  $d$ , from a user  $u$  to a user  $v$  (or from an item  $i$  to an item  $j$ ) is the length of a line segment between the two users (or items) in the Euclidean space. In practical terms, each user is represented by its Cartesian coordinates with respect to the basis of items (the same thing for an item which is represented with respect to the basis of users) and the distance between two users (or two items) is the absolute value of the numerical difference of their coordinates. The Euclidean distance ( $d$ ) formula representing the correlation between two users  $u$  and  $v$  is as follows:

$$d(u, v) = \sqrt{\sum_{i \in I_{uv}} (r_{vi} - r_{ui})^2}$$

Where  $I_{uv}$  denotes the set of items commonly rated by both  $u$  and  $v$ ,  $r_{ui}$  and  $r_{vi}$  denote the rating of the user  $u$  and  $v$ , respectively, on the item  $i$ . This formula provides the Euclidean distance between two items  $i$  and  $j$ .

$$d(i, j) = \sqrt{\sum_{u \in U_{IJ}} (r_{uj} - r_{ui})^2}$$

$U_{ij}$  denotes the set of users who rated both items  $i$  and  $j$ ,  $r_{ui}$  and  $r_{uj}$  denote the rating of the user  $u$  on items  $i$  and  $j$ , respectively. The Euclidean distance should be normalized to become a similarity measure.

### 3.5.8 Manhattan distance

The Manhattan distance, also known as city blocks and taxicab, between two vectors is equal to the one-norm of the distance between the vectors. To adapt this measure to the RSs domain, we need to represent a user by its Cartesian coordinates with respect to the basis of items (the same thing for an item). The Manhattan distance ( $d_1$ ) between two users  $u$  and  $v$  is as follows:

$$d_1(u, v) = \sum_{i \in I_{uv}} (|r_{vi} - r_{ui}|)$$

$I_{uv}$  denotes the set of items commonly rated by both  $u$  and  $v$ ,  $r_{ui}$  and  $r_{vi}$  denote the rating of the user  $u$  and  $v$ , respectively, on the item  $i$ . the Euclidean distance between two items  $i$  and  $j$ .

$$d_1(i, j) = \sum_{u \in U_{ij}} (|r_{uj} - r_{ui}|)$$

$U_{ij}$  represents the set of users who rated both items  $i$  and  $j$ ,  $r_{ui}$  and  $r_{uj}$  denote the rating of the user  $u$  on items  $i$  and  $j$ , respectively.

### 3.6 User-based Rating Prediction

User-based neighborhood recommendation methods predict the rating  $r_{ui}$  of a user  $u$  for a new item  $i$  using the ratings given to  $i$  by users most similar to  $u$ , called nearest-neighbors. Suppose we have for each user  $vu$  a value  $S_{uv}$  representing the similarity between  $u$  and  $v$ .

The k-nearest-neighbors (k-NN) of  $u$  is one of the simplest forms of machine learning algorithms mostly used for classification, denoted by  $N(u)$ , are the  $k$  users  $v$  with the highest similarity  $S_{uv}$  to  $u$ . However, only the users who have rated item  $i$  can be used in the prediction of  $r_{ui}$  and we instead consider the  $k$  users most similar to  $u$  that have rated  $i$ . We note set of neighbors as  $N_i(u)$  The rating  $r_{ui}$  can be estimated as the average rating given to  $i$  by these neighbors :

$$\widehat{r_{ui}} = \frac{1}{|N_i(u)|} \sum_{v \in N_i(u)} r_{vi}$$

As this equation does not take into account the fact that the neighbors can have different levels of similarity. A common solution to this problem is to weigh the contribution of each neighbor by its similarity to  $u$ . However, if these weights do not sum to 1, the predicted ratings can be well outside the range of allowed values. Consequently, it is customary to normalize these weights such that the predicted rating becomes :

$$\widehat{r_{ui}} = \frac{\sum_{v \in N_i(u)} S_{uv} r_{vi}}{\sum_{v \in N_i(u)} |S_{uv}|}$$

$|s_{uv}|$  is used instead of  $s_{uv}$  v because negative weights can produce ratings outside the allowed range.

### 3.7 Item-based Rating Prediction

Item-based approaches look at ratings given to similar items. Denote by  $N_u(i)$  the items rated by user  $u$  most similar to item  $i$ . The predicted rating of  $u$  for  $i$  is obtained as a weighted average of the ratings given by  $u$  to the items of  $N_u(i)$ :

$$\widehat{r_{ui}} = \frac{\sum_{j \in N_u(i)} S_{ij} r_{uj}}{\sum_{j \in N_u(i)} |S_{ij}|}$$

### 3.8 Rating Normalization

Assigning a rating  $r_{ui}$  to an item  $i$  for user  $u$ , each user has its own personal scale, some users might be reluctant to give high/low scores to items they liked/disliked. Two of the most popular rating normalization schemes that have been proposed to convert individual ratings to a more universal scale are mean-centering and Z-score.

- **Mean-centering** In user-based recommendation, a raw rating  $r_{ui}$  is transformation to a mean-centered one  $h(r_{ui})$  by subtracting to  $r_{ui}$  the average  $\bar{r}_u$  of the ratings given by user  $u$  to the items in  $I_u$ .

$$h(r_{ui}) = r_{ui} - \bar{r}_u$$

Using this approach the user-based prediction of a rating  $r_{ui}$  is obtained as

$$\widehat{r_{ui}} = \bar{r}_u + \frac{\sum_{v \in N_i(u)} s_{uv}(r_{vi} - \bar{r}_v)}{\sum_{v \in N_i(u)} |s_{uv}|}$$

The Mean-centering normalization technique is most often used in item-based, where The item-mean-centered normalization of  $r_{ui}$  is given by

$$h(r_{ui}) = r_{ui} - \bar{r}_i$$

$\bar{r}_i$  corresponds to the mean rating given to item  $i$  by user in  $U_i$

Using this approach the item-based prediction of a rating  $r_{ui}$  is obtained as:

$$r_{ui} = \bar{r}_i + \frac{\sum_{j \in N_u(i)} s_{ij}(r_{uj} - \bar{r}_j)}{\sum_{j \in N_u(i)} |s_{ij}|}$$

- **Z-score normalization** This is very similar to the above method, just that we also divide by the standard deviation in this case. This method provides a slight lift in the accuracy when compared to mean centering. in addition, the z-score is the signed number of deviations from the mean indicating that a datum is above the mean if positive and below the mean otherwise.

In user-based methods, the normalization of a rating  $r_{ui}$  divides the user-mean-centered rating by the standard deviation  $\sigma_u$  of the ratings given by user  $u$ :

$$h(r_{ui}) = \frac{r_{ui} - \bar{r}_u}{\sigma_u}$$

A user-based prediction of rating  $r_{ui}$  using g this normalization approach would therefore be obtained as :

$$\widehat{r_{ui}} = \bar{r}_u + \sigma_u \frac{\sum_{v \in N_i(u)} s_{uv}(r_{vi} - \bar{r}_v)/\sigma_v}{\sum_{v \in N_i(u)} |s_{uv}|}$$

The z-score normalization of  $r_{ui}$  in item-based methods divides the item- mean-centered rating by the standard deviation of ratings given to item  $i$ :

$$h(r_{ui}) = \frac{r_{ui} - \bar{r}_i}{\sigma_i}$$

The item-based prediction of rating  $r_{ui}$  would then be

$$\widehat{r_{ui}} = \bar{r}_i + \sigma_i \frac{\sum_{j \in N_u(i)} s_{ij}(r_{uj} - \bar{r}_j)/\sigma_j}{\sum_{j \in N_u(i)} |s_{ij}|}$$

## 3.9 Neighborhood Selection

The number of nearest-neighbors to select have a severe impact on the quality of the recommender system. We call it a neighborhood because we talk about items with a small distance between them. In this section, I'll cover three ways to define and calculate such a neighborhood: Top-N filtering, Threshold filtering, and Negative filtering.

### 3.9.1 Top-N Recommendation:

For each user or item, only a list of the  $N$  nearest-neighbors and their respective similarity weight is kept. To avoid problems with efficiency or accuracy,  $N$  should be chosen carefully. Thus, if  $N$  is too large, an excessive amount of memory will be required to store the neighborhood lists, and predicting ratings will be slow. On the other hand, selecting a too-small value for  $N$  may reduce the coverage of the recommendation method, which causes some items to be never recommended.

#### User-Based Top-N Recommendation

User-based recommendation systems compute the top- $N$  recommended items for that user as follows. First, they identify the  $k$  most similar user in the database.

Once this set of the  $k$  most similar has been discovered, we identify the set  $C$  of items purchased by the group and their frequency. Using this set, user-based techniques then recommend the  $N$  most frequent items in  $C$  that are not already in the set of items that have already been purchased by the user for which we want to compute the top- $N$  recommendations.

#### Item-Based Top-N Recommendation

Item-based top- $N$  recommendation algorithms that use item-to-item similarity to compute the relations between the items, for each item, the  $k$  most similar items are computed, and their corresponding similarities are recorded. Now, for each user that has purchased a set  $T$  of items that the customer has already purchased, we want to compute the top- $N$  recommendations items as follows.

First, we identify the set  $C$  of the candidate recommended items by taking the union of the  $k$  most similar items for each item  $i \in T$ , and removing from the union any items that are already in  $T$ . Then, for each item  $c \in C$  we compute its similarity to the set  $T$  as the sum of the similarities between all the items  $i \in T$  and  $c$ , using only the  $k$  most similar items of  $i$ . Finally, the items in  $C$  are sorted in non-increasing order with respect to that similarity, and the first  $N$  items are selected as the top- $N$  recommended set.

### 3.9.2 Threshold filtering:

Instead of keeping a fixed number of nearest-neighbors, this approach keeps all the neighbors whose similar weight has a magnitude greater than a given threshold  $s_{min}$ . While this is more flexible than the previous filtering technique, as only the most significant neighbors are kept, the correct value of  $s_{min}$  may be challenging to determine.

### 3.9.3 Negative filtering:

In general, negative rating correlations are less reliable than positive ones. Intuitively, this is because a strong positive correlation between two users is a good indicator of their belonging to a common group. However, although a negative correlation may indicate membership in different groups, it does not tell how diverse these groups are or whether these groups are compatible with other categories of items. While negative correlations provide no significant improvement in the prediction accuracy,

The three filtering approaches are not mutually exclusive and can be combined to fit the needs of the recommendation systems.

## 3.10 Dimensionality Reduction

The dimension reduction phase transforms the original user-item matrix into a lower-dimensional space to address the sparsity and scalability problems often encountered in collaborative filtering recommendation scenarios. The original representation of the input data to a collaborative filtering system is an  $n \times m$  user-item matrix  $R_{n \times m}$ . This representation may potentially pose sparsity and scalability problems for collaborative filtering systems [B.Sarwar et al. 2001]. In practice, when a large set of items is available, users may have rated or chosen a very low percentage of items, resulting in a very sparse user preference matrix.

Consequently, a collaborative filtering recommendation system may be unable to make recommendations for a particular user.

When building collaborative filtering, The user-item matrix or the rating matrix usually contains very large numbers of users or items, we are dealing with millions of users rating millions of items. This may be problematic during data exploration and analysis processes.

Dimensionality reduction is one of the oldest approaches used to respond to this problem to improve performance, speed up calculations, and avoid the curse of dimensionality. Dimensionality reduction aims to select or extract an optimal subset from the matrix.

The main aims of dimensionality reduction are as follows:

- Easy visualize and understand data.
- Reduce the required storage space.
- Reduce learning and use times.

### 3.11 Matrix factorization techniques

Matrix factorization is a subset of dimensionality reduction in which a data user-item matrix is reduced to the product of many low-rank matrices.

Recommendation systems mainly base their suggestions on rating data of items, which are often placed in a matrix with one representing users and the other representing items of interest, the ratings are given explicitly by users creating a sparse user-item rating matrix because an individual user is likely to rate only a small fraction of the items that belong to the item set, i.e., many values in the rating matrix are null since all users do not rate all items.

And due to the high-dimensional data of the users and items,, MF was proposed to deal with these problems. Matrix factorization is the decomposition of the initial matrix into two or smaller metrics to reduce required storage space and processing time.

Another challenging issue with this user-item rating matrix is the scalability of data (i.e., the large number of possible registered users or inserted items), which may affect the time performance of a recommendation algorithm we can deal with all the mentioned challenges by applying matrix decomposition methods (also known as factorization methods).

Matrix factorization techniques gained popularity during the Netflix challenge because they are fast and accurate.

Matrix Factorization decomposes a matrix into a product of Three matrices, it's a powerful technique to find the hidden structure behind the data.

The main popular three methods are:

- Singular Value Decomposition (SVD)
- Principal Component Analysis (PCA)
- Probabilistic Matrix Factorization (PMF)

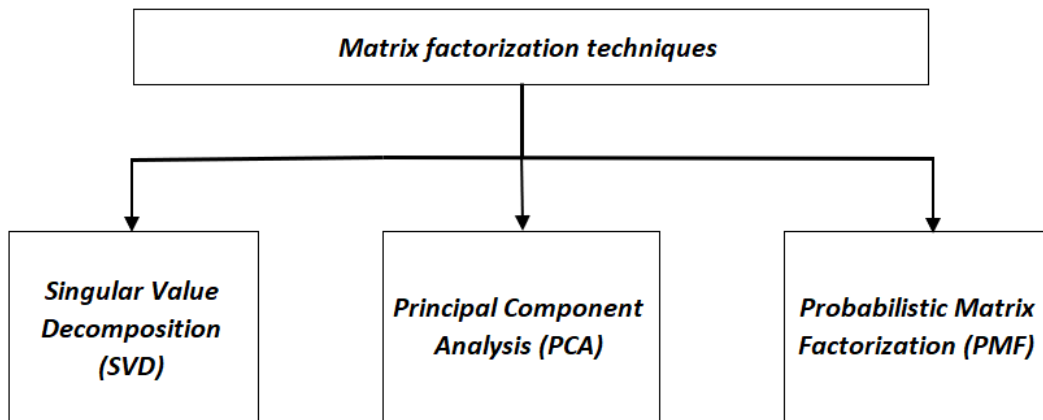


Figure 3.9: Matrix Factorization Techniques



### 3.11.1 Singular Value Decomposition (SVD)

It is the most popular matrix factorization method that can be applied in recommendation systems to solve sparsity and scalability problems. The Singular Value Decomposition (SVD) is used as a dimensionality reduction technique also to produce a low-dimensional representation of the original user-item matrix.

It is nothing more than decomposing vectors into orthogonal axes by transforming data from a high-dimensional space into a low-dimensional space. The key an issue in an SVD decomposition is to find a lower-dimensional feature space thus can be expressed in the following way is so that the low-dimensional

$$SVD(A) = U \times S \times V^T$$

Where  $A$  is the  $m \times n$  matrix with rank  $r$ ,  $U$  and  $V$  are  $m \times m$  and  $n \times n$  are orthogonal

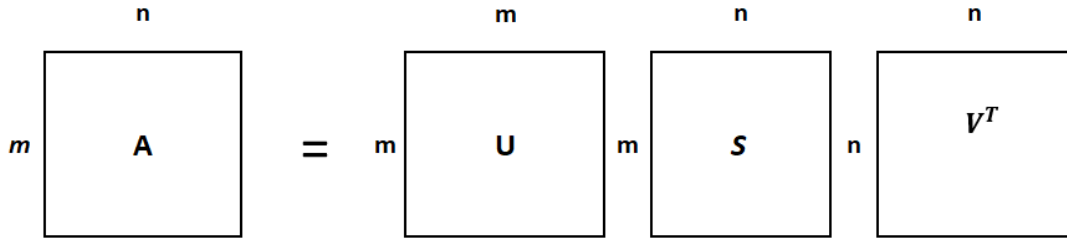


Figure 3.10: Singular Value Decomposition (SVD)

matrices respectively.  $S$  is a diagonal matrix  $m \times n$  singular orthogonal matrix with non-negative elements and

the  $m \times m$  matrix  $U$  is called orthogonal where  $U^T U = I_{m \times m}$ , Also matrix  $V$  is orthogonal where  $V^T V = I_{n \times n}$ .

The initial  $r$  diagonal elements of  $S$ ,  $(s_1, s_2, s_3, \dots, s_r)$  are called the singular values of matrix  $A$  have the property that  $s_i > 0$  and  $s_1 \geq s_2 \geq s_3 \geq \dots \geq s_r$ .

Accordingly, the first  $r$  columns of  $U$  are called the left singular vectors (the eigenvectors of  $AA^T$ ) and The first  $r$  columns of  $V$  are called the right singular vectors (the eigenvectors of the  $A^T A$ ).

If we focus only on these  $r$  nonzero singular values, the effective dimensions of the SVD matrices  $U$ ,  $S$  and  $V$  will become  $m \times r$ ,  $r \times r$  and  $r \times n$  respectively.

In SVD, particularly useful in the case of Recommendation Systems, is that it can provide the best low-rank approximation of the original matrix  $A$ . By retaining the first  $k \ll r$  singular values of  $S$  and discarding the rest, which can be translated as keeping the  $k$  largest singular values, based on the fact that the elements in  $S$  are sorted.

The resulting diagonal matrix  $S$  is termed  $S_k$ . Matrices  $U$  and  $V$  should be also reduced accordingly.  $U_k$  is produced by removing  $r - k$  columns from matrix  $U$ .  $V_k$  is produced by removing  $r - k$  rows from matrix  $V$ . Matrix  $A_k$  which is defined as:

$$A_k = U_k \times S_k \times V_k^T$$

Stands for the closest linear approximation of the original matrix  $A$  with reduced rank  $k$ . Once this transformation is completed, users and items can be represented as points in the  $k$ -dimensional space.

[B.Sarwar et al. 2001] Applied the Matrix Factorization model SVD to reduce the dimensionality of a rating matrix. Using the MovieLens dataset, they selected 943 users to form

a  $943 \times 1682$  matrix that is 95.4% sparse (each user's on average rates 5% of the 1682 movies). They first fill missing values using user and movie rating averages and then apply SVD. For this large dataset Item-Based technique provide optimal accuracy with significantly faster and high-quality online recommendations than the user-user (k-Nearest-Neighbor) method.

We define the original user-item matrix  $R$ , of size  $m \times n$ , which includes the ratings of  $m$  users on  $n$  items  $r_{ij}$ . In order to eliminate all missing rating  $r_{ui}$  we replace them by Computing the average of a  $\bar{r}_u$  then we normalized matrix  $R$ . Then the normalized matrix  $R$

$$R = U \times S \times V^T$$

Perform the dimensionality reduction step by keeping only  $k$  diagonal entries from matrix  $S$  to obtain a  $k \times k$  matrix,  $S_k$ . Similarly, matrices  $U_k$  and  $V_k$  of size  $m \times k$  and  $k \times n$  are generated.

The "reduced" user-item matrix,  $R$ , is obtained by

$$R = U_k \times S_k \times V_K$$

### 3.11.2 Principal Component Analysis (PCA)

Principal component analysis (PCA) is also the powerful dimensionality reduction technique that transforms high-dimensions data into lower dimensions while retaining as much information as possible and is a particular realization of the matrix factorization approach.

PCA allows to obtain an ordered list of components that account for the largest amount of the variance from the data in terms of least square errors, The amount of variance captured by the first component is larger than the amount of variance on the second component and so on. We can reduce the dimensionality of the data by neglecting those components. PCA is a statistical technique used to reduce the dimensionality of the data for computational efficiency consisting of many correlated variables. It retains the variation present in the dataset up to the maximum extent and uses the ideas of variances and covariances. PCA is a multi-variate analysis technique that transforms a set of correlated variables into a set of orthogonal components. Although correlated components are required to reproduce the total variability of the original variables, much of this variability can be often explained by a small number  $k$  of the principal components (PC's). The  $k$  PC's can then replace the original variables, thus compressing the original data set into a smaller one that consists of the  $k$  PC's .

[Goldberg et al] Proposed an approach to use Principal Component Analysis (PCA) in the context of an online joke recommendation system. Their system, known as Eigentaste's, In Eigentaste they addressed sparseness using universal queries, ensuring that all users rate a common set of k-items. So, the resulting sub-matrix will be dense and directly compute the square symmetric correlation matrix and then do linear projection using the Principle Component Analysis (PCA), a closely-related factor analysis technique first described by Pearson in 1901. Like SVD, PCA reduces the dimensionality of the matrix by optimally projecting highly correlated data along with a smaller number of orthogonal dimensions.

### 3.11.3 Probabilistic Matrix Factorization (PMF)

Probabilistic Matrix Factorization (PMF) is a probabilistic linear model with Gaussian observation noise. The Probabilistic Matrix Factorization (PMF) models the user preference matrix as a product of two lower-rank user and item matrices. Suppose we have  $N$  users and  $M$  items. Let  $R_{ij}$  be the rating value of user  $i$  for item  $j$ ,  $U_i$  and  $V_j$  represent  $D$ -dimensional user-specific and item-specific latent feature vectors respectively.

The conditional distribution over the observed ratings  $R \in \mathbb{R}^{N \times M}$  and the prior distributions over  $U \in \mathbb{R}^{D \times N}$  and  $V \in \mathbb{R}^{D \times M}$  are given by :

$$P(R|U, V, \sigma^2) = \prod_{i=1}^N \prod_{j=1}^M [N(R_{ij}|U_i^T V_j, \sigma^2)]^{I_{ij}}$$

Where  $N(x|u, \sigma^2)$  is the probability density function of the Gaussian distribution with mean  $u$  and variance  $\sigma^2$ ,  $I_{ij}$  is the indicator variable that is equal to 1 if user  $i$  rated item  $j$  and equal to 0 otherwise.

We also place zero-mean spherical Gaussian priors[1, 11] on user and item feature vectors:

$$P(U|\sigma_u^2) = \prod_{i=1}^N N(U_i|0, \sigma_u^2 I)$$

$$p(V|\sigma_v^2) = \prod_{j=1}^M N(V_j|0, \sigma_v^2 I)$$

## Chapter 4

# Covariance Distance

The objective of a recommendation systems is to provide users with personalized recommendations while selecting an item among a set of items (restaurants, books, etc...). Collaborative filtering is one of the most successful approaches for an online store to make personalized recommendations through its recommendation systems. The neighborhood-based is one of the effective used model the critical component of this method is the mechanism of finding similarities between items or users in the item-based neighborhood or user-based neighborhood, respectively.

In addition, item-based and user-based neighborhood collaborative filtering algorithms produce predictions based on similarities between items and users, using the user's ratings the calculation of similarities metrics, Studies, and real-world implementations so far has relied on traditional distance and vector similarity measures, mainly although the most widely used similarity measures such as cosine, pearson correlation, and adjusted cosine have been proved to be successful in many studies in terms of their effectiveness in the recommendation systems, those similarity metrics have been mentioned in the previous chapter, they are limited to be used.

This chapter provides an explanation of the distance covariance a new similarity measure, that we tested in term of improving recommendation systems performance.

## 4.1 Predefined similarity measures

The definitions of the main similarity measures that have been used for neighborhood-based recommendation systems are summarized in the table . The first one, the cosine measure, looks at the angle between two rating vectors where a smaller angle implies greater similarity. The second one, Adjusted Cosine Similarity, has proven more efficient for the item-based neighborhood [B. Sarwar et al. 2001]. the third one, Pearson’s correlation, measures the linear correlation between two rating vectors.

Measure	User Similarity Measure	Item similarity Measure
<b>Cosine</b>	$\text{cosine}(u, v) = \frac{\sum_{i \in I_{uv}} r_{ui} r_{vi}}{\sqrt{\sum_{i \in I_u} r_{ui}^2} \sqrt{\sum_{i \in I_v} r_{vi}^2}}$	$\text{cosine}(i, j) = \frac{\sum_{u \in U_{ij}} r_{ui} r_{uj}}{\sqrt{\sum_{u \in U_i} r_{ui}^2} \sqrt{\sum_{u \in U_j} r_{uj}^2}}$
<b>Adjusted cosine</b>	$\text{ACosine}(u, v) = \frac{\sum_{i \in I_{uv}} (r_{ui} - \bar{r}_i)(r_{vi} - \bar{r}_i)}{\sqrt{\sum_{i \in I_{uv}} (r_{ui} - \bar{r}_i)^2} \sqrt{\sum_{i \in I_{uv}} (r_{vi} - \bar{r}_i)^2}}$	$\text{ACosine}(i, j) = \frac{\sum_{u \in U_{ij}} (r_{ui} - \bar{r}_u)(r_{uj} - \bar{r}_u)}{\sqrt{\sum_{u \in U_{ij}} (r_{ui} - \bar{r}_u)^2} \sqrt{\sum_{u \in U_{ij}} (r_{uj} - \bar{r}_u)^2}}$
<b>Pearson correlation</b>	$\text{PCC}(u, v) = \frac{\sum_{i \in I_{uv}} (r_{ui} - \bar{r}_u)(r_{vi} - \bar{r}_v)}{\sqrt{\sum_{i \in I_{uv}} (r_{ui} - \bar{r}_u)^2} \sqrt{\sum_{i \in I_{uv}} (r_{vi} - \bar{r}_v)^2}}$	$\text{PCC}(i, j) = \frac{\sum_{u \in U_{ij}} (r_{ui} - \bar{r}_i)(r_{uj} - \bar{r}_j)}{\sqrt{\sum_{u \in U_{ij}} (r_{ui} - \bar{r}_i)^2} \sqrt{\sum_{u \in U_{ij}} (r_{uj} - \bar{r}_j)^2}}$
<b>The Jaccard coefficient</b>	$J(u, v) = \frac{ u \cap v }{ u \cup v }$	$J(i, j) = \frac{ i \cap j }{ i \cup j }$
<b>Euclidean distance</b>	$d(u, v) = \sqrt{\sum_{i \in I_{uv}} (r_{vi} - r_{ui})^2}$	$d(i, j) = \sqrt{\sum_{u \in U_{ij}} (r_{uj} - r_{ui})^2}$
<b>Manhattan distance</b>	$d_1(u, v) = \sum_{i \in I_{uv}} ( r_{vi} - r_{ui} )$	$d_1(i, j) = \sum_{u \in U_{ij}} ( r_{uj} - r_{ui} )$

Table 4.1: TABLE OF SIMILARITIES

## 4.2 Distance Covariance

Classical dependence measures such as Pearson correlation, Spearman's  $\rho$ , and Kendall's  $\tau$  can detect only monotonic or linear dependence to overcome these limitations, [Szekely et al.2007] Proposed distance covariance and its derived correlation.

The distance covariance is a weighted  $L_2$  distance between the joint characteristic function and the product of marginal distributions. This measure can detect the presence of a dependence structure when this sample size is large enough. There has been considerable recent interest in measuring dependence by employing the concept of distance covariance, a new and appealing measure of dependence for random variables.

The distance covariance is an effective measure of dependence between sets of multivariate random variables based on the assumption that the data are *i.i.d* (independent and identically distributed). However, as this assumption is often violated in many practical problems, proposed to extend the notion of distance covariance to the case of dependent data.

### 4.2.1 The Distance Covariance Function

The distance covariance function as a new measure of dependence between multivariate random vectors, say  $X$  and  $Y$ , of arbitrary, not necessarily equal in dimensions, let note  $p$  and  $q$  be the dimension of  $X$  and  $Y$ , respectively.

A crucial feature of the distance covariance is that it equals zero if and only if  $X$  and  $Y$  are mutually independent. Hence the distance covariance is sensitive to arbitrary dependencies, this is in contrast to the classical covariance, which is generally capable of detecting only linear dependencies. where the distance covariance have been applied to multivariate data to detect arbitrary types of non-linear associations between sets of variables.

The definition of distance covariance relies on the joint characteristic function of  $X$  and  $Y$ , which is denoted by

$$\phi_{X,Y}(t, s) = E[\exp^{i(t'X + s'Y)}]$$

Where  $(t, s) \in R^{p+q}$ , and  $i^2 = -1$ .

The corresponding marginal characteristic function of multivariate random vector  $X$  and  $Y$  are denoted respectively by

$$\phi_X(t) = E[\exp(it'X)]$$

$$\phi_Y(s) = E[\exp(is'Y)]$$

The distance covariance function is defined as the nonnegative square root of a weighted  $L_2$  distance between the joint and the product of the marginal characteristic functions of two random vectors, namely

$$V^2(X, Y) = \int_{R^{p+q}} |\phi_{X,Y}(t, s) - \phi_X(t)\phi_Y(s)|^2 \omega(t, s) dt ds \quad (4.1)$$

Where  $\omega(., .) : R^{p+q} \rightarrow R$  is a weight function for which the previous integral exists .

The distance covariance function  $V^2(X, Y)$  leads to the definition of the distance correlation function between  $X$  and  $Y$  , which is the positive square root of

$$R^2(X, Y) = \begin{cases} \frac{V^2(X, Y)}{\sqrt{V^2(X, X)V^2(Y, Y)}}, & V^2(X, X)V^2(Y, Y) > 0 \\ 0, & V^2(X, X)V^2(Y, Y) = 0. \end{cases}$$

The previous display shows that the distance correlation function is a coefficient analogous to coefficient. the coefficient measures the linear relationship between  $X$  and  $Y$  and can be zero even when the variables are dependent,

The choice of the weight function  $\omega(.,.)$  is crucial. Choosing the non-integrable weight function of the form

$$\omega(t, s) = (c_p c_q |t|_p^{1+p} |s|_q^{1+q})^{-1}$$

where  $c_d = \pi^{(1+d)/2} / \Gamma((1+d)/2)$ , with  $\Gamma(.)$  denoting the Gamma function.  $R^2(X, Y)$  is scale and rotation invariant. In particular  $\omega(t, s)$  is the only weight function such that  $V^2(X, Y)$  is rigid motion invariant and scale equivariant. However, the choice of  $\omega(t, s)$  yields to computational advantages.

Denote by  $|x|_p$  The Euclidean norm <sup>1</sup> of  $x$  in  $R^p$ .

The distance covariance  $V^2(X, Y)$  between  $X$  and  $Y$  will be defined as the following equation :

$$V^2(X, Y) = \int_{R^{p+q}} \frac{|\phi_{(X,Y)}(t, s) - \phi_X(t)\phi_Y(s)|^2}{c_p c_q |t|_p^{1+p} |s|_q^{1+q}} dt ds$$

As the  $V^2(X, Y)$  is nonnegative, it follows that  $V^2(X, Y) \geq 0$ . Furthermore  $V^2(X, Y) < \infty$  whenever  $X$  and  $Y$  have finite first moments.

An advantage of this equation is that it directly implies one of the most important properties of the distance covariance, the characterization of independence. For all  $X$  and  $Y$ ,  $V^2(X, Y) = 0$  if and only if  $X$  and  $Y$  are independent.

### Proof

If  $X$  and  $Y$  are independent then  $\phi_{X,Y}(t, s) = \phi_X(t)\phi_Y(s)$  for all  $s$  and  $t$ , hence  $V^2(X, Y) = 0$ . If  $X$  and  $Y$  are not independent then the functions  $\phi_{X,Y}(t, s)$  and  $\phi_X(t)\phi_Y(s)$  are not identical. Since characteristic functions are continuous then there exists an open set  $A \subseteq R^p \times R^q$  such that  $|\phi_{X,Y}(t, s) - \phi_X(t)\phi_Y(s)|^2 > 0$  for all  $(t, s) \in A$ , hence in the equation  $V^2(X, Y) > 0$ .

The main properties of  $V^2(X, Y)$ ,  $R^2(X, Y)$ , when  $\omega(t, s)$  is used, are given by the following:

- If  $E(|X|_p + |Y|_q) < \infty$ , then the distance correlation,  $R(X, Y)$ , satisfies  $0 \leq R(X, Y) \leq 1$ , and  $R(X, Y) = 0$  if and only if  $X$  and  $Y$  are independent.
- $R(X, Y)$  is invariant under orthogonal transformations of the form  $(X, Y) \rightarrow (\alpha_1 + b_1 C_1 X, \alpha_2 + b_2 C_2 Y)$ , where  $\alpha_1, \alpha_2$  are arbitrary vectors,  $b_1, b_2$  are arbitrary nonzero numbers and  $C_1, C_2$  are arbitrary orthogonal matrices.
- If  $p = q = 1$  and  $X$  and  $Y$  have standard normal distributions with  $r = \text{Cov}(X, Y)$ , then  $R(X, Y) = |r|$  and

$$R^2(X, Y) = \frac{r \arcsin r + \sqrt{1 - r^2} - r \arcsin r/2 - \sqrt{4 - r^2} + 1}{1 + \pi/3 - \sqrt{3}}$$

- When  $E(|X|_p^2 + |Y|_q^2) < \infty$  then

$$V^2(X, Y) = E|X - X'|_p |Y - Y'|_q + E|X - X'|_p |Y - Y''|_q - 2E|X - X'|_p |Y - Y''|_q$$

with  $(X', Y')$  and  $(X'', Y'')$  being independent copies of  $(X, Y)$

<sup>1</sup>The Euclidean norm ( $L_2$  norm) of a vector  $x \in R^p$  is defined by  $|x|_p = \sqrt{x_1^2 + x_2^2 + x_3^2 + \dots + x_p^2}$



**Theorem**

Suppose that  $(X_1, Y_1), \dots, (X_4, Y_4)$  are independent, identically distributed (i.i.d.)

$$V^2(X, Y) = E[|X_1 - X_2| \cdot |Y_1 - Y_2| - 2|X_1 - X_2| \cdot |Y_1 - Y_3| + |X_1 - X_2| \cdot |Y_3 - Y_4|] \quad (4.2)$$

**Proof**

Note that  $\langle \cdot, \cdot \rangle$  the corresponding inner product.

$$\begin{aligned} |\phi_{(X,Y)}(t, s) - \phi_X(t)\phi_Y(s)|^2 &= (\phi_{(X,Y)}(t, s) - \phi_X(t)\phi_Y(s)) \overline{(\phi_{(X,Y)}(t, s) - \phi_X(t)\phi_Y(s))} \\ &= E[e^{i\langle t, X_1 - X_2 \rangle + i\langle s, Y_1 - Y_2 \rangle} - 2e^{i\langle t, X_1 - X_2 \rangle + i\langle s, Y_1 - Y_3 \rangle} + e^{i\langle t, X_1 - X_2 \rangle + i\langle s, Y_3 - Y_4 \rangle}] \end{aligned}$$

As this last expression is real, any term of the form  $e^{iz}$ ,  $z \in \mathbb{R}$  can be replaced by  $\cos z$

$$\begin{aligned} V^2(X, Y) &= \int_{\mathbb{R}^{p+q}} \frac{|\phi_{(X,Y)}(t, s) - \phi_X(t)\phi_Y(s)|^2}{c_p c_q |t|_p^{1+p} |s|_q^{1+q}} dt ds \\ c_p c_q V^2(X, Y) &= \int_{\mathbb{R}^{p+q}} \frac{A_{12}(t, s) - 2A_{13}(t, s) + A_{34}(t, s)}{|t|_p^{1+p} |s|_q^{1+q}} dt ds \end{aligned} \quad (4.3)$$

Where, for each  $(j, k)$ ,

$$A_{jk}(t, s) = E(\cos(\langle t, X_1 - X_2 \rangle + \langle s, Y_j - Y_k \rangle)) \quad (4.4)$$

Replacing  $t$  by  $-t$  in (4.3) we also obtain

$$c_p c_q V^2(X, Y) = \int_{\mathbb{R}^{p+q}} \frac{A_{12}(t, -s) - 2A_{13}(t, -s) + A_{34}(t, -s)}{|t|_p^{1+p} |s|_q^{1+q}} dt ds \quad (4.5)$$

And by adding (4.3) and (4.5) we find that

$$c_p c_q V^2(X, Y) = \int_{\mathbb{R}^{p+q}} \frac{B_{12}(t, s) - 2B_{13}(t, s) + B_{34}(t, s)}{|t|_p^{1+p} |s|_q^{1+q}} dt ds \quad (4.6)$$

Where for each  $(j, k)$ ,

$$B_{jk}(t, s) = \frac{1}{2}(A_{jk}(t, s) + A_{jk}(t, -s)) \quad (4.7)$$

On applying to  $A_{jk}(t, s)$  (4.4) and  $B_{jk}(t, s)$  (4.7) the trigonometric identity,

$$\cos(x + y) + \cos(x - y) = 2 \cos x \cos y$$

We deduce that

$$B_{jk}(t, s) = E(\cos(\langle t, X_1 - X_2 \rangle \langle s, Y_j - Y_k \rangle)) \quad (4.8)$$

For  $j, k \in 1, 2, 3, 4$  we apply to Eq (4.8)

$$\begin{aligned} & \cos \langle t, X_1 - X_2 \rangle \cos \langle s, Y_j - Y_k \rangle = \\ & (1 - \cos \langle t, X_1 - X_2 \rangle)(1 - \cos \langle s, Y_j - Y_k \rangle) \\ & - 1 + \cos \langle t, X_1 - X_2 \rangle + \cos \langle s, Y_j - Y_k \rangle \end{aligned} \quad (4.9)$$

Then we obtain

$$\begin{aligned} c_p c_q V^2(X, Y) = & \int_{R^{p+q}} (E[(1 - \cos \langle t, X_1 - X_2 \rangle)(1 - \cos \langle s, Y_1 - Y_2 \rangle) \\ & - 2E[(1 - \cos \langle t, X_1 - X_2 \rangle)(1 - \cos \langle s, Y_1 - Y_3 \rangle)] \\ & + E[(1 - \cos \langle t, X_1 - X_2 \rangle)(1 - \cos \langle s, Y_3 - Y_4 \rangle)]] \frac{dt ds}{|t|_p^{1+p} |s|_q^{1+q}} \end{aligned} \quad (4.10)$$

Which is obtained by decomposing all summands on the right-hand side using Eq.(4.9) and observing that all terms which are not of the form  $E(\cos \langle t, X_i - X_j \rangle \cos \langle s, Y_l - Y_k \rangle)$  cancel each other. By applying the Fubini-Tonelli Theorem and the linearity of expectation and integration, we obtain

$$\begin{aligned} c_p c_q V^2(X, Y) = & E \int_{R^{p+q}} [(1 - \cos \langle t, X_1 - X_2 \rangle)(1 - \cos \langle s, Y_1 - Y_2 \rangle) \\ & - 2(1 - \cos \langle t, X_1 - X_2 \rangle)(1 - \cos \langle s, Y_1 - Y_3 \rangle) \\ & + (1 - \cos \langle t, X_1 - X_2 \rangle)(1 - \cos \langle s, Y_3 - Y_4 \rangle)] \frac{dt ds}{|t|_p^{1+p} |s|_q^{1+q}} \end{aligned} \quad (4.11)$$

### 4.2.2 Estimation of distance covariance

The empirical distance covariance and correlation measures are functions of the double centred distance matrices of the samples. Suppose that  $(X_i, Y_i)$ ,  $i = 1, \dots, n$  is a random sample from the joint distribution of the random vectors  $X$  and  $Y$ . Based on this sample, consider the  $n \times n$  Euclidean pairwise distance matrices with elements

$$(a_{ij}) = (|X_i - X_j|_p)$$

and

$$(b_{ij}) = (|Y_i - Y_j|_q)$$

These matrices are double centred so that their row and column means are equal to zero. In other words, let

$$\begin{aligned} A_{ij} &= a_{ij} - \bar{a}_{i.} - \bar{a}_{.j} + \bar{a}_{..} \\ B_{ij} &= b_{ij} - \bar{b}_{i.} - \bar{b}_{.j} + \bar{b}_{..} \end{aligned}$$

Where  $\bar{a}_{i.} = (\sum_{j=1}^n a_{ij})/n$ ,  $\bar{a}_{.j} = (\sum_{i=1}^n a_{ij})/n$ ,  $\bar{a}_{..} = (\sum_{i,j=1}^n a_{ij})/n^2$ . Similarly, we define the quantities  $\bar{b}_{i.}$ ,  $\bar{b}_{.j}$  and  $\bar{b}_{..}$ . The sample distance covariance is defined by the square root of the statistic

$$\hat{V}^2(X, Y) = \frac{1}{n^2} \sum_{i,j=1}^n A_{ij} B_{ij}$$

We compute the squared sample distance correlation  $\hat{R}^2(X, Y)$

$$\hat{R}^2(X, Y) = \begin{cases} \frac{\hat{V}^2(X, Y)}{\sqrt{\hat{V}^2(X, X) \hat{V}^2(Y, Y)}}, & \hat{V}^2(X, X) \hat{V}^2(Y, Y) > 0 \\ 0, & \hat{V}^2(X, X) \hat{V}^2(Y, Y) = 0. \end{cases}$$

Furthermore, note that by multiplying out  $\hat{V}^2(X, Y)$  we obtain the alternative expression

$$\hat{V}^2(X, Y) = \frac{1}{n^2} \sum_{i,j=1}^n a_{ij} b_{ij} + \frac{1}{n^4} \sum_{i,j=1}^n a_{ij} \sum_{i,j=1}^n b_{ij} - \frac{2}{n^3} \sum_{i,j,k=1}^n a_{ij} b_{jk}$$

The main properties of sample distance covariance and correlation functions are the following:

- The estimators  $\hat{V}^2(X, Y)$  and  $\hat{R}^2(X, Y)$  are both strongly consistent, that is, they both converge almost surely to their population  $V^2(X, Y)$  and  $R^2(X, Y)$  respectively, as  $n$  tends to infinity, provided that  $E(|X|_p) < \infty$  and  $E(|Y|_q) < \infty$ .
- $\hat{V}^2(X, Y) \geq 0$  the equality holds when  $X$  and  $Y$  are independent.
- $0 \leq \hat{R}(X, Y) \leq 1$ .
- $\hat{R}(X, \alpha + bXC) = 1$ , where  $\alpha$  is a vector,  $b$  is a nonzero real number and  $C$  is an orthogonal matrix.

$\hat{V}^2(X, Y)$  is not an unbiased estimator of  $V^2(X, Y)$ , where we have

$$E[\hat{V}^2(X, Y)] = \frac{(n-1)(n-2)^2}{n^3} V^2(X, Y) + \frac{2(n-1)^2}{n^3} \gamma - \frac{(n-1)(n-2)}{n^3} \alpha \beta$$

Where

$$\alpha = E|X - X'|_p$$

,

$$\beta = E|Y - Y'|_q$$

and

$$\gamma = E[|X - X'|_p |Y - Y'|_q]$$

$X'$  (respectively  $Y'$ ) denotes an independent of  $X$  (respectively  $Y$ ). An unbiased estimator of squared distance covariance proposed by Székely and Rizzo (2014) is given by

$$\hat{V}_U(X, Y) = \frac{1}{n(n-3)} \sum_{i \neq j} \tilde{A}_{ij} \tilde{B}_{ij}$$

For  $n > 3$ , where the so-called  $U$ -centred matrices,  $\tilde{A}_{ij}$  have the additional property that  $E[\tilde{A}_{ij}] = 0$  for all  $i, j$  and are defined by

$$\tilde{A}_{ij} = \begin{cases} a_{ij} - \frac{1}{n-2} \sum_{l=1}^n a_{il} - \frac{1}{n-2} \sum_{k=1}^n a_{kj} + \frac{1}{(n-1)(n-2)} \sum_{k,l=1}^n a_{kl}, & i \neq j \\ 0, & i = j \end{cases}$$

### 4.2.3 Asymptotic Distribution

Under the null hypothesis of independence between the random vectors  $X$  and  $Y$ .

$$n\hat{V}^2(X, Y) \xrightarrow{n \rightarrow \infty} \sum_{j=1}^{\infty} \lambda_j Z_j \quad (4.12)$$

In distribution, where  $Z_j$  are independent standard normal variables and  $\lambda_j$  are eigenvalues that depend on the joint distribution of the random vectors  $(X, Y)$ , provided that  $E[|X|_p] < \infty$  and  $E[|Y|_q] < \infty$ .

Moreover, the quadratic form

$$Q := \sum_{j=1}^{\infty} \lambda_j Z_j$$

Satisfies  $E[Q] = E[|X - X'|_p |Y - Y''|_q]$  On the other hand, if  $X$  and  $Y$  are not independent, then

$$n\hat{V}^2(X, Y) \xrightarrow{n \rightarrow \infty} \infty \quad (4.13)$$

The equations (4.12) and (4.13) imply a test for independence based on the statistic  $nV^2$ . While Székely et al. (2007) derive an asymptotic test with asymptotic significance level at most  $\alpha$  they point out that this significance level can be quite conservative for many distributions.

### 4.2.4 Generalisations and Modifications of the Distance Covariance

The  $\alpha$  distance covariance function,  $V^{(\alpha)}(X, Y)$  is a generalisation of the distance covariance function (4.1) for a constant  $\alpha$  that lies in the interval  $[0, 2]$ , when  $\omega(t, s)$  is used.

$$V^{2(\alpha)}(X, Y) = \frac{1}{C_{(p, \alpha)} C_{(q, \alpha)}} \int_{R^{p+q}} \frac{|\phi_{(X, Y)}(t, s) - \phi_X(t) \phi_Y(s)|^2}{|t|_p^{\alpha+p} |s|_q^{\alpha+q}} dt ds$$

Where  $C(p, \alpha)$  and  $C(q, \alpha)$  are given by

$$C(d, \alpha) = \frac{2\pi^{d/2} \Gamma((1 - \alpha/2))}{\alpha 2^\alpha \Gamma((d + \alpha)/2)}$$

The estimator of  $V^{2(\alpha)}(X, Y)$  is computed by defining a distance matrix with elements  $a_{ij} = |X_i - X_j|_p^\alpha$  and  $b_{ij} = |Y_i - Y_j|_q^\alpha$

$V^{2(\alpha)}(X, Y)$  Extends the theory of distance covariance because  $V^2(X, Y)$  is obtained a special case when  $\alpha = 1$ .

### 4.3 The Distance Correlation Coefficient in High Dimensions

As stated in asymptotic Tests, distance covariance provides a test for independence for random  $X, Y$  of arbitrary dimension  $p$  and  $q$ , respectively. Distance correlation is often used when needed as measure for independence than applying a formal test. While this appears to work well for univariate random variables, some authors have reported that the direct interpretation of distance correlation is questionable in high dimensions. [Dueck et al.(2014)] showed that for fixed  $q$  and standard multivariate normal random vectors  $(X, Y) \in \mathbb{R}^{p+q}$ .

$$\lim_{p \rightarrow \infty} R(X, Y) = 0 \quad (4.14)$$

Irrespective of the dependence structure between  $X$  and  $Y$  On the other hand, Székely et al.(2013) proved that for fixed  $n$ ,

$$\lim_{p, q \rightarrow \infty} \widehat{R}(X, Y) = 1 \quad (4.15)$$

Provided that the coordinates of  $X$  and  $Y$  are i.i.d. and the second moments of  $X$  and  $Y$  exist. Comparison of (4.14) and (4.15) reveals that there are actually two different effects, which complicate the interpretation of distance correlation in high dimension.

Equation (4.14) shows that the population distance correlation is close to 0 when comparing a low-dimensional and a high-dimensional vector.

Equation (4.15) shows that the sample distance correlation is close to 1 when comparing two high-dimensional vectors. These two problems may occur simultaneously, when comparing two high-dimensional vectors with substantially distinct dimensions.

In particular, (4.14) shows that we can find sequences  $(p_k)_{k \in \mathbb{N}}, (q_k)_{k \in \mathbb{N}}$  such that  $p_k \rightarrow \infty, q_k \rightarrow \infty$  and

$$\begin{aligned} \lim_{k \rightarrow \infty} R(X, Y) &= 0 \\ \lim_{k \rightarrow \infty} \widehat{R}(X, Y) &= 1 \end{aligned}$$

Where for any  $k$ ,  $X$  and  $Y$  are  $p_k$  and  $q_k$ -dimensional standard normal vectors and  $n$  is fixed. When  $p = q$ , the population distance correlation  $R(X, Y)$  can attain both the values 0 ( $X$  and  $Y$  independent) and 1 ( $Y = \alpha + \beta CX$  with  $\alpha \in \mathbb{R}^p, \beta \in \mathbb{R}, C$  orthogonal)

Equation (4.15) we can apply an unbiased estimator for the distance covariance as described in Székely et al. (2013). This has the further advantage that transformation of this estimator converges to a Student- $t$  distribution with  $n(n-3)/2$  degrees of freedom, which enables a distance correlation  $t$ -test for independence in high dimensions. When  $p$  and  $q$  differ (but are both high-dimensional), the  $t$ -test may still be applied, however, a direct interpretability is critical (recall (4.14)).

## 4.4 Implimentation of distance covariance in recommendation systems

Distance covariance is recently introduced dependency measures between random vectors , the implimentation of distance covariance in user-based neighborhood or item-based neighborhood recommendation systems . where we had used to compute the estimators of the squared distance covariance  $\hat{V}^2(i, j)$  in item-based and  $\hat{V}^2(u, v)$  in user-based .

### 4.4.1 Distance covariance in Item-Based Neighborhood

The empirical distance covariance for Item-based neighborhood  $\hat{V}^2(r_i, r_j)$  with a simple form. Suppose that we have  $n$  observations of  $r_i$  denotes the rating vecteur of item  $i, r_j$  denotes the rating vecteur of item  $j$  :

$$\hat{V}^2(r_i, r_j) = \frac{1}{n^2} \sum_{u,v=1}^n a_{uv} b_{uv} + \frac{1}{n^4} \sum_{u,v=1}^n a_{uv} \sum_{u,v=1}^n b_{uv} - \frac{2}{n^3} \sum_{u,v,k=1}^n a_{uv} b_{uk}$$

$(r_{ui}, r_{uj})$  ,  $i = 1, \dots, n$  and  $j = 1, \dots, n$  is a random sample from the joint distribution of the random vectors  $r_i$  and  $r_j$ . Based on this sample, consider the  $n \times m$  Euclidean pairwise distance matrices with elements

$$(a_{uv}) = (|r_{ui} - r_{vj}|_p)$$

and

$$(b_{uv}) = (|r_{uj} - r_{vj}|_q)$$

These matrices are double centred so that their row and means are equal to zero. In other words, let

$$A_{uv} = a_{uv} - \bar{a}_{u.} - \bar{a}_{.v} + \bar{a}_{..}$$

$$B_{uv} = b_{uv} - \bar{b}_{u.} - \bar{b}_{.v} + \bar{b}_{..}$$

where  $\bar{a}_{u.} = (\sum_{v=1}^n a_{uv})/n$ ,  $\bar{a}_{.v} = (\sum_{u=1}^n a_{uv})/n$ ,  $\bar{a}_{..} = (\sum_{u,v=1}^n a_{uv})/n^2$

Similarly, we define the quantities  $\bar{b}_{u.}$  ,  $\bar{b}_{.v}$  and  $\bar{b}_{..}$ . The sample distance covariance is defined by the square root of the statistic

$$\hat{V}^2(r_i, r_j) = \frac{1}{n^2} \sum_{u,v=1}^n A_{uv} B_{uv}$$

$$\hat{V}^2(r_i, r_j) = \frac{1}{n^2} \sum_{u,v=1}^n a_{uv} b_{uv} + \frac{1}{n^4} \sum_{u,v=1}^n a_{uv} \sum_{u,v=1}^n b_{uv} - \frac{2}{n^3} \sum_{u,v,k=1}^n a_{uv} b_{uk}$$

#### 4.4.2 Distance covariance in User-Based Neighborhood

The empirical distance covariance for User-based neighborhood  $\hat{V}^2(r_u, r_v)$  with a simple form. Suppose that we have  $m$  observations of  $r_u$  denotes the rating vecteur of user  $u$ ,  $r_v$  denotes the rating vecteur of user  $v$  :

$$\hat{V}^2(r_u, r_v) = \frac{1}{m^2} \sum_{i,j=1}^m a_{ij}b_{ij} + \frac{1}{m^4} \sum_{i,j=1}^n a_{ij} \sum_{i,j=1}^m b_{ij} - \frac{2}{m^3} \sum_{i,j,k=1}^m a_{ij}b_{ik}$$

$(r_{ui}, r_{vi})$ ,  $u = 1, \dots, m$  and  $v = 1, \dots, m$  is a random sample from the joint distribution of the random vectors  $r_u$  and  $r_v$ . Based on this sample, consider the  $n \times m$  Euclidean pairwise distance matrices with elements

$$(a_{ij}) = (|r_{ui} - r_{uj}|_p)$$

and

$$(b_{ij}) = (|r_{vi} - r_{vj}|_q)$$

These matrices are double centred so that their row and means are equal to zero. In other words, let

$$A_{ij} = a_{ij} - \bar{a}_{i.} - \bar{a}_{.j} + \bar{a}_{..}$$

$$B_{ij} = b_{ij} - \bar{b}_{i.} - \bar{b}_{.j} + \bar{b}_{..}$$

Where  $\bar{a}_{i.} = (\sum_{j=1}^m a_{ij})/m$ ,  $\bar{a}_{.j} = (\sum_{i=1}^m a_{ij})/m$ ,  $\bar{a}_{..} = (\sum_{i,j=1}^m a_{ij})/m^2$

Similarly, we define the quantities  $\bar{b}_{i.}$ ,  $\bar{b}_{.j}$  and  $\bar{b}_{..}$ . The sample distance covariance is defined by the square root of the statistic

$$\hat{V}^2(r_u, r_v) = \frac{1}{m^2} \sum_{i,j=1}^n A_{ij}B_{ij}$$

$$\hat{V}^2(r_u, r_v) = \frac{1}{m^2} \sum_{i,j=1}^m a_{ij}b_{ij} + \frac{1}{m^4} \sum_{i,j=1}^m a_{ij} \sum_{i,j=1}^n b_{ij} - \frac{2}{m^3} \sum_{i,j,k=1}^n a_{ij}b_{ik}$$



## Chapter 5

# Implementation of Restaurant Recommendation Systems

### Introduction

Recommendation systems help users deal with data overload by recommending items based on their preferences. Currently, recommendation systems play an essential role in most commercial applications. These systems help increase users' satisfaction and content providers' revenue by capturing users' personalized preferences. The most common approach is Neighborhood collaborative filtering, which employs users' interactions and interests to define their preferences among items.

Nowadays, eating out at restaurants has become one of the needs of today's society due to the lifestyle. A large number of restaurant choices and a lack of information about the restaurant becomes an obstacle to users' needs in choosing a restaurant. A restaurant is a business that sells various menus of food and drinks to its users, and the bigger the business, the more menus are offered. Consequently, choosing an excellent appropriate restaurant to order from it can be challenging and time-consuming for the users since the choice of the available restaurants provided is overwhelming. However, users want a restaurant with a good reputation and fit their tastes, so restaurant ratings from other users are required in the restaurant recommendation process. Recommendation systems help users tackle the problem of finding restaurants that suit their preferences by generating a personalized list of recommendations that might interest them by learning through their previous ratings. In this case, A restaurant recommendation system is becoming a necessity following the trend of eating out at restaurants. The task of such a system is to generate a Top-N list of restaurants that may interest a user, in which the user's previous rating behavior.

The case study is the YASSIR EXPRESS web application where we apply our recommendation systems models, and we implemented a widely used approach user-Neighborhood collaborative filtering and the Item-Neighborhood collaborative filtering, One of the main components of Neighborhood collaborative filtering is the similarity measure used to determine the set of similar users or items. Several similarity metrics have been proposed, with different performances in terms of accuracy of recommendations, the calculation of similarities has relied on traditional distance and vector similarity measures such as Cosine, Adjusted Cosine Euclidean, and Manhattan distance which, however, have been seldom questioned in terms of their effectiveness in the recommendation systems. We will present a new similarity metric distance covariance function as a new measure of dependency between two random vectors multidimensional (not necessarily equal) to generate the list of recommendations to a target user. We have Also implemented popularity recommendation systems and location-based recommendation systems.

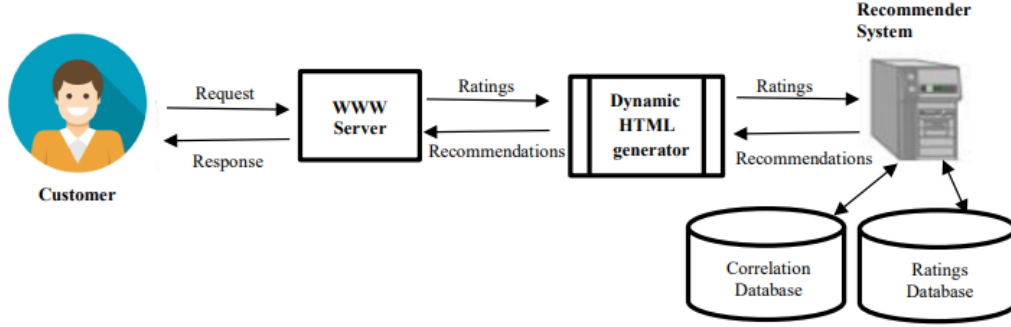


Figure 5.1: Recommendation System Architecture

## 5.1 YASSIR EXPRESS

Yassir is a company with a headquarters in Palo Alto, Silicon Valley, USA, and a second in Algiers, Algeria. This changes the daily life of the services provided in Algeria and Africa through innovative digital mobile solutions and also by creating a bridge between Silicon Valley and the Algerian ecosystem.

It was founded in 2017 by two doctoral students, El Mehdi Yettou and Nouredine Tayebi, two alumni from the polytechnic school of Algiers with extensive experience in the field of research at the Canadian University of Sherbrooke and the American University from Stanford, respectively.

A vision to prove the model for start-ups that are overgrowing and have both economic and social impact on their environment is also possible in Algeria.

YASSIR EXPRESS Released in 2018, its application was developed by the company Yassir, bringing together meal delivery services from restaurants partner, Yassir dark stores, and other services such as dry cleaning.

YASSIR EXPRESS has established itself in the two neighboring countries, Tunisia and Morocco, available on the APP store for the smartphones under ios and Google Play for Smartphones under Android. YASSIR EXPRESS offers a service by which e-consumers will be in touch with e-suppliers to purchase products, services, and meals remotely presented on the application and benefit from delivery to their destination. In this context, it is recalled that Yassir Express acts as an intermediary. Its role is limited to hosting offers of products, services, and meals from e-suppliers and putting them in contact with e-consumers. The opportunities behind a recommendation system for a YASSIR EXPRESS are several indeed, it would allow a better user experience, and They provide an effective way of driving traffic and increasing average order value, which increases the company's revenue.

### 5.1.1 The Application

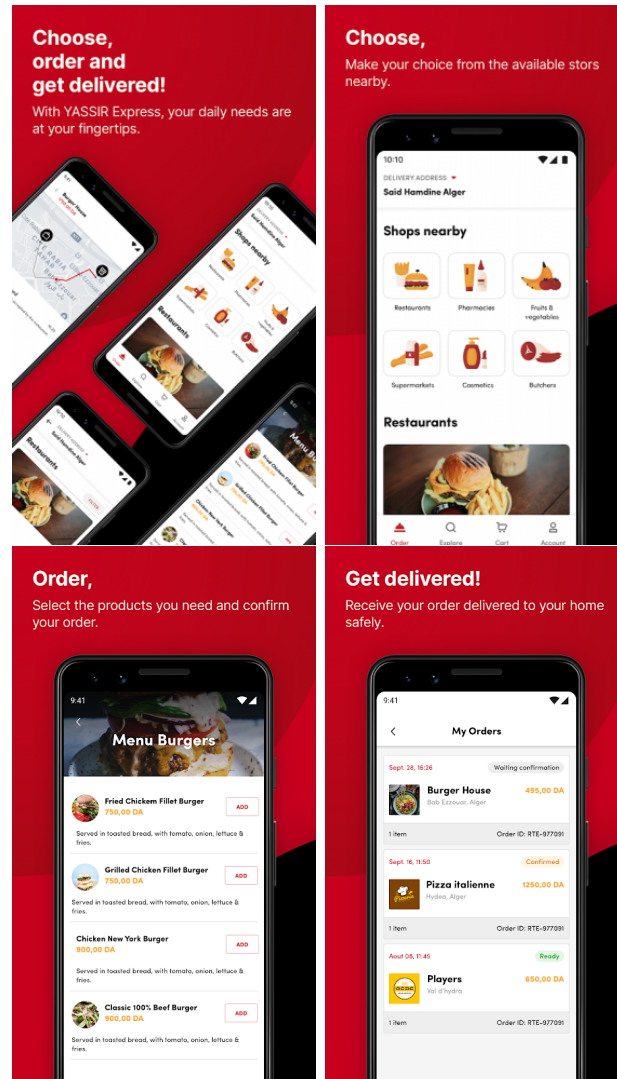


Figure 5.2: YASSIR EXPRESS Application

### 5.1.2 Flow

In this section, we will see how the application works more in detail, trying to follow the user from switching on the app to choosing a restaurant. The user journey starts at the home page where the user adds his/her location position, then the user can choose between a chosen or selected set of restaurants based on the restaurant near to his location, next the user can pick the available dish the from he have select, As we can see on the following figure [5.1.2](#)

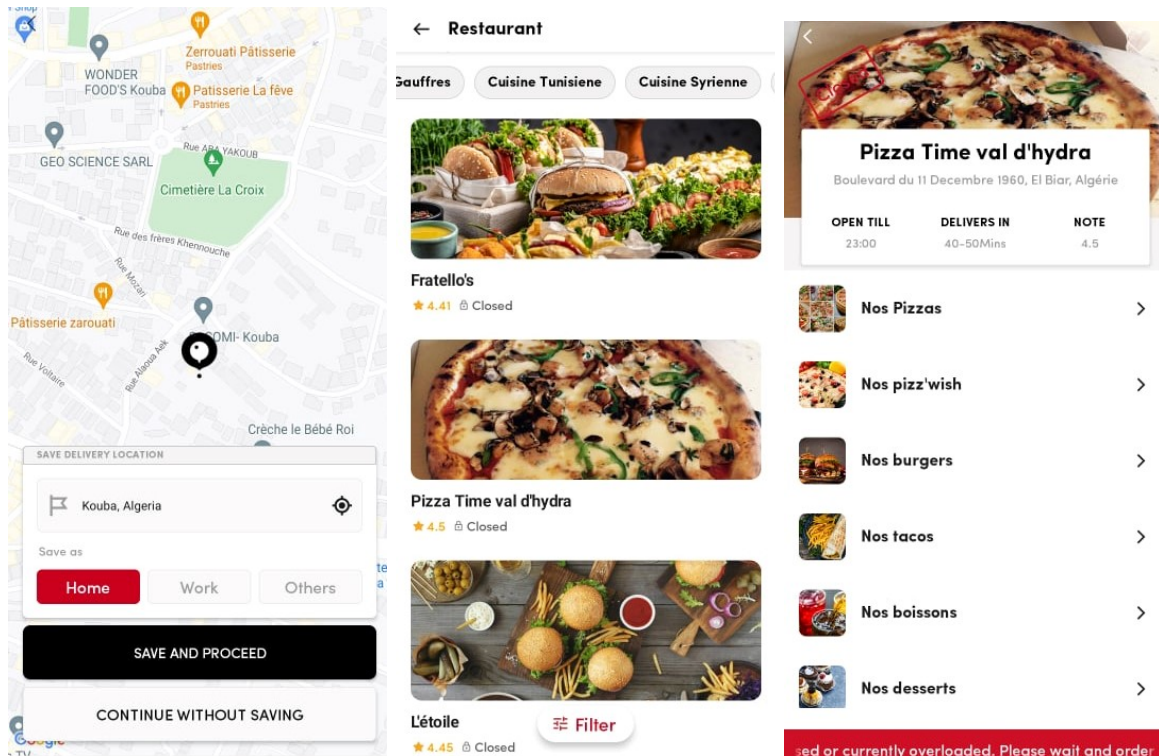


Figure 5.3: YASSIR EXPRESS User Experience

## Frameworks and libraries variety



**Google Colab** Google Colaboratory, or “Google Colab” for short, is a product from Google Research. Colab allows anybody to write and execute arbitrary python code in other words, Colab is a web IDE for python through the browser and is especially well suited to machine learning, data analysis, and education. Colab is a hosted Jupyter notebook service that requires no setup to use while providing access free of charge to computing resources, including GPUs.

**Python** Python is a high-level, interpreted, general-purpose programming language. Its design philosophy emphasizes code readability with the use of significant indentation. Python is dynamically-typed and garbage-collected. It supports multiple programming paradigms, including structured (particularly procedural), object-oriented and functional programming. Due to its comprehensive standard library, it is often described as a "batteries included" language. Python is a programming language preferred for programming due to its vast features, applicability, and simplicity. Python is the best-fit machine learning due to its independent platform and popularity in the programming community.

**Python Libraries** Python Libraries are vital in the preparation of a suitable programming environment. It offers a reliable environment that reduces software development time significantly. A library includes prewritten code that developers can use to speed up coding when working on complex projects. Here are the main Python Libraries that we have used on our project.

- **pandas** It is a software library written for the Python programming language for data manipulation and analysis. In particular, it offers data structures and operations for manipulating numerical tables.
- **NumPy** It is a package that defines a multi-dimensional array object and associated fast math functions that operate on it. It also provides simple routines for linear algebra and fft and sophisticated random-number generation.
- **Seaborn** It is a Python data visualization library based on matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics.
- **Scikit-learn** (formerly scikits.learn and also known as sklearn) is a free software machine learning library for the Python programming language. It features various classification, regression and clustering algorithms, including support-vector machines, random forests, gradient boosting, k-means .
- **Plotly** is Python graphing library makes interactive, publication-quality graphs. Examples of how to make line plots, scatter plots, area charts, bar charts, error bars, box plots, histograms...
- **Surprise** is a Python scikit for building and analyzing recommender systems .
- **SciPy** in Python is an open-source library used for solving mathematical, scientific, engineering, and technical problems. It allows users to manipulate the data and visualize the data using a wide range of high-level Python commands.
- **Matplotlib** is a comprehensive library for creating static, animated, and interactive visualizations in Python.

### 5.1.3 Dataset

The dataset used for training the recommendation system was the YASSIR EXPRESS, it was composed of one table with all the information about users and restaurants collected by the web application.

YASSIR EXPRESS comprises **104358** users with information on **1069** restaurants (items), the dataset was available in the form of two CSV files<sup>1</sup>, then we did the merge. For the reason of anonymity, we are under obligation to not disclose the dataset's real numbers, names, times, and values.

<i>Dataset</i>	<i>Users</i>	<i>Items</i>
<i>Yassir Express</i>	<i>104358</i>	<i>1069</i>

Table 5.1: Table summarizes the YASSIR EXPRESS dataset

A dataset contains information about the users and the restaurants where all the main information about restaurants are collected, for instance:

- **order id**: The ID of the order.
- **user id**: The ID of the user.
- **restaurant**: The ID of the restaurant.
- **request date**: The request time of the order.

<sup>1</sup>CSV (comma-separated values) file is a text file in which information is separated by commas.

- **delivery date:** The delivery time of the order.
- **pick up date:** The pick up date date of the order.
- **restaurant accepted date:** The time that the restaurant accepted the order.
- **restaurant address lat:** The restaurant latitude address .
- **restaurant address lng:** The restaurant longitude address.
- **commune res:** The restaurant commune address.
- **wilaya res:** The restaurant wilaya address.

## 5.2 Experimental Setup

We followed the process described in figure 5.2 with the aim to apply an experimental comparative evaluation result between the similarity metrics using different evaluation metrics. In Step 1, Yassir Express dataset was selected. In In step 2, the data loading was performed.the data Preparation then exploration. In Step 3 recommendations models were implemented. Finally, in Step 4, we evaluated the recommendation models using RMSE, MSE, and MAE metrics.

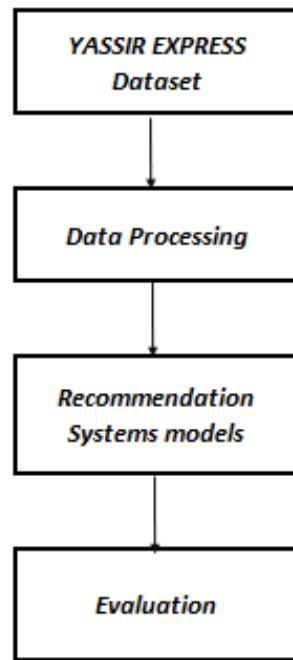


Figure 5.4: Steps performed to compare the recommendations systems

## 5.3 Data Processing

### 5.3.1 Data Loading

The data loading was from two CSV files, the resulting format is a Pandas dataframe structure, which can be considered a table representation of the data.

### 5.3.2 Data Preparation

This module's functions provide easy manipulation and filtering of the data built on the Python library and its data frame structure. First, we removed the missing values (NaN), then transformed the different data structures into appropriate data types such as time and the string datatype. For instance, convert user and item information in numerical kinds and dates represented by strings into a corresponding date.

### 5.3.3 Data Exploration

To explore our dataset, we will see how many unique commune and states are in the dataset, and then we will display most of the ten cities and states with the largest number of restaurants.

```
data['wilaya_res'].value_counts().sort_values(ascending=False)
data['commune_res'].value_counts().sort_values(ascending=False)
```

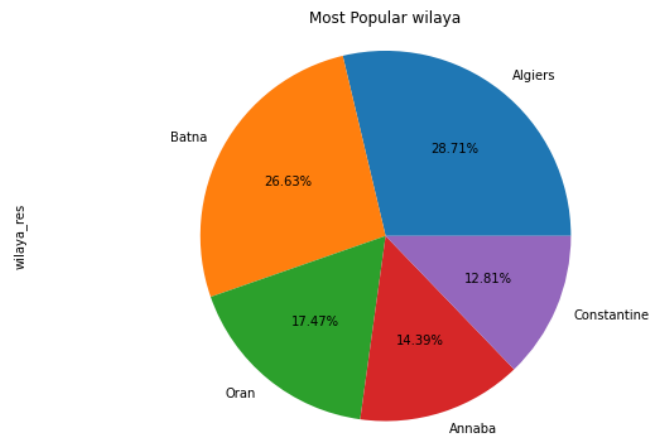


Figure 5.5: The states with the largest number of restaurants

This figure 5.3.3 displays the Algerian states that has the highest number of restaurants.

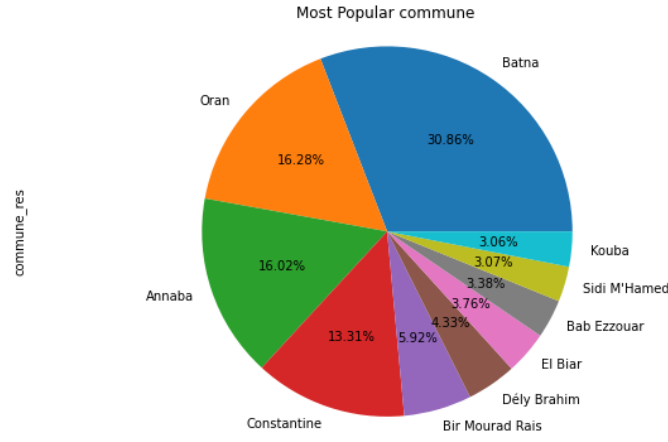


Figure 5.6: The commune with the largest number of restaurants

This figure indicate that Batna is the commune that have the largest number of restaurants.

### User Rating

The YASSIR EXPRESS dataset doesn't contain any user rating, this problem occurs when the users are lazy or don't rate items. Due to this situation, we defined an two implicit rating, as we have previously mentioned, the implicit rating can be gathered using the user's interaction to deduce the user's preferences.

**The First implicit rating function** that outputs is the rating  $r_{ui}$  [1-5] of the user  $u$  to item  $i$ . We have considered the duration per minute between the delivery time and request time as the implicit rating that we will work with. In contrast, the delivery duration is the most crucial feature that the users focus on in the food delivery services. Moreover, the users do not order from a restaurant that takes a long time to deliver. Furthermore, there are other features like the taste and flavor of the delivered dishes, hygiene, price, etc. as our dataset doesn't contain any of that information, we can't add these features on our implicit rating computation.

```
data.delivery_date = pd.to_datetime(data.delivery_date)
data.request_date = pd.to_datetime(data.request_date)
data["diff"] = (data["delivery_date"] - data["request_date"])
               .astype('timedelta64[m]')
```

The k-means clustering algorithm was applied to cluster the delivery duration by selecting the k=5 values the number of clusters, then we were created the rating scalar by assing for each cluster his own rating value.

We plotted the distribution of the ratings (that we have created). where the total of each rating number is calculated. This plot shows 5.3.3 that a rating range of 2 to 4 is available for this dataset and most restaurants have a rating of 1.



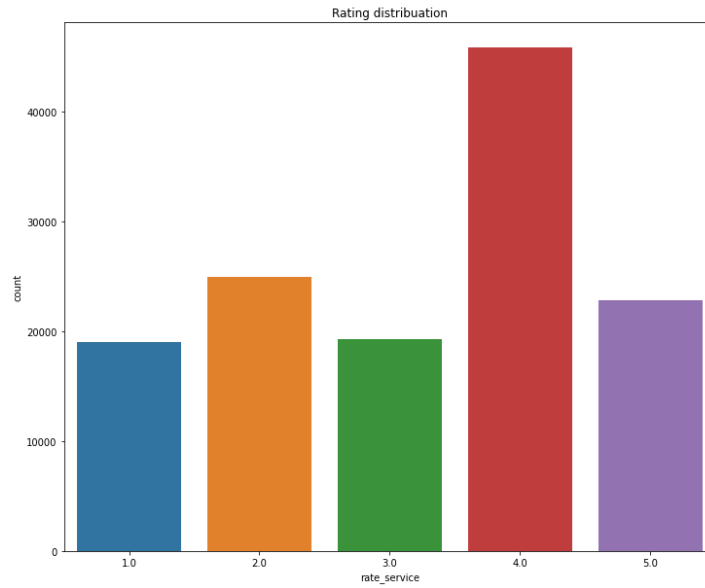


Figure 5.7: The Rating Distribution

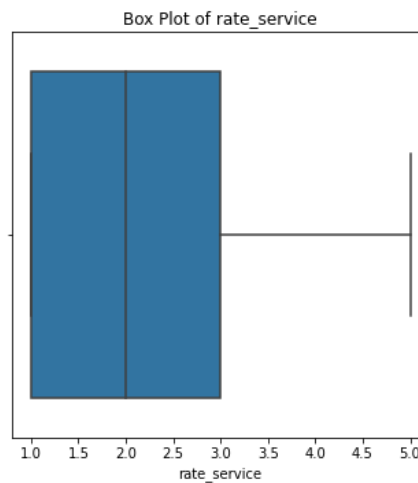


Figure 5.8: Box Plot Of The Rating

The box plot displays the rating distribution 5.3.3 is right-skewed, the horizontal line (sometimes called the "whiskers") goes from 1 to 5 (from minimum to maximum), while the box extends from 1 (= Q1) to 3 (= Q3) with a middle vertical line at 2 . where the 25% of the rating are between 2 and 1 while 75% of our rating are from 2 to 5 and the , the median is the average value from a set of rating shown by the line that divides the box into two parts it's equal to 2.

**The Second implicit rating function** to go Further, we gathered all the users that have purchased more the one time from the same restaurant, and we computed the mean of the delivery duration between the restaurants that the user has been purchased more than once, then we divided by the number of times that the user has been purchased from the same restaurant more than once.

```
# count the number of restaurant that has been purchased  
#from the same user  
data.groupby(['user_id', 'restaurant'])["diff1"].count()  
# the mean of the delivery duration between  
# the same restaurant that has been purchased from the same user  
data.groupby(['user_id', 'restaurant'])["diff1"].mean()  
# rating formula  
data.groupby(['user_id', 'restaurant'])["diff1"].mean() /  
data.groupby(['user_id', 'restaurant'])["diff1"].count()
```

**Important** all the followed models have been calculated using the first rating approach

## 5.4 The Recommendations Systems Models

Recommendation systems are one of the most commonly used practical systems in data science. In this section, we will focus on the implementation, and describe in detail the recommendations that we built.

### 5.4.1 Evaluation Metrics

Before applying a recommendation systems in real life applications, it is important to measure the effectiveness of a RSs. To evaluate the performance of the RS, we intend to utilize each time a different similarity measure. Therefore, we evaluate the performance of the RS using 3 well-known evaluation metrics which are widely used for evaluating regression models :

- **Root Mean Square Error (RMSE)**

it evaluates the difference between the ratings predicted  $\hat{r}_{ui}$  by the RS and the ratings given by the user  $r_{ui}$ .

$$RMSE = \sqrt{\frac{\sum_{(u,i,r) \in R} (\hat{r}_{ui} - r_{ui})^2}{|R|}}$$

- **Mean Absolute Error (MAE):**

it evaluates the difference between the ratings predicted by the RS and the ratings given by the users. It returns a positive value.

$$MAE = \frac{\sum_{(u,i,r) \in R} |\hat{r}_{ui} - r_{ui}|}{|R|}$$

Where  $R$  denotes the amount of user  $u$  rated items, normally a Test Set,  $r_{ui}$  determines the actual rating that user  $u$  rates item  $i$  and  $\hat{r}_{ui}$  denotes the predicted rating of item  $i$  for user  $u$ . The lower the MAE, the more accuracy the recommendation system predicts rating.

- **Mean squared error (MSE)**

The MSE is used to evaluate recommendation systems by computing the average value of all absolute value differences between the predicted ratings  $\hat{r}_{ui}$  and the user's true ratings  $r_{ui}$ . The MSE can be computed using the following equation

$$MSE = \frac{\sum_{(u,i,r) \in R} (\hat{r}_{ui} - r_{ui})^2}{|R|}$$

## 5.5 Neighborhood Collaborative Filtering

Today, recommendation systems play a vital role in the acceleration of searches by internet users to find what they are interested in. Among the strategies proposed for recommendation systems, Neighborhood collaborative filtering (also known as memory-based systems) has received due attention regarding its simplicity and efficiency.

Neighborhood-based collaborative filtering models directly exploit the feedback from the users. We can distinguish either user-based or item-based approaches within Neighbourhood-based. The critical factor for the success of this strategy returns to the similar calculation methods that affect the accuracy of its recommendations. In this section, we will implement the user-Neighborhood and the item-Neighborhood collaborative filtering. We carried out an experimental comparative evaluation result between the similarity metrics using different evaluation metrics, MSE, RMSE, and MAE.

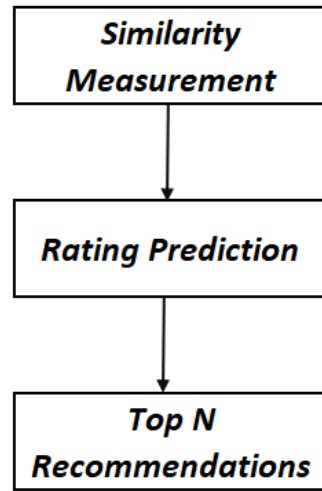


Figure 5.9: Flow chart Of The Neighborhood Collaborative Filtering

### 5.5.1 User Neighborhood Collaborative Filtering

User Neighborhood CF is based on the main idea that users who have an interest in the same items and similar ratings will thus have similar preferences. In contrast, the recommendations are based on user feedback from similar users (known as neighbors), where they give a similar rating behavior. The recommendation system is able to predict if a user might be interested in an unseen item. In the user-based recommendation, we need to find  $k$  most similar users and make a prediction of a target user,  $k$ -NN computes a specific neighborhood for each user consisting of the top  $k$  users with higher similarity to the target user and the Top- $N$  recommendation.

The process of a typical User Neighborhood is generally divided into three steps:

1. **Similarity measurement** (between users ) The recommendation system computes the similarities between the active user and users who have rated the same items. Neighborhood Similar users are regrouped in a subset, namely the « neighborhood » of the active user.
2. **Rating prediction** (using ratings of similar users ) User-based prediction of a rating  $r_{ui}$  is obtained as:

$$\widehat{r_{ui}} = \bar{r}_u + \frac{\sum_{v \in N_i(u)} s_{uv}(r_{vi} - \bar{r}_v)}{\sum_{v \in N_i(u)} |s_{uv}|}$$

3. **The Top- $N$  recommendations** Are obtained by selecting the  $N$  items which provide the most satisfaction to the user according to the prediction by decreasing the order of predicted rating.3.9.1

We applied an experimental comparative evaluation result between the similar items as the critical step in the proposed user-based recommendation is to determine the similarity between the users, In implementing this model, we applied different similarity metrics such as the Cosine, Euclidean Distance, and Manhattan Distance using different evaluation metrics, MSE, RMSE, and MAE.

### 5.5.2 Item Neighborhood Collaborative Filtering

Item k-nearest Neighborhood CF method predicts the user preference based on the similarity with the k nearest items. While the item-based determine the similarities among the various items and then used them to identify the set of items to be recommended. First,the method is used to compute the similarity between the items, and then combine these similarities in order to compute the similarity between a set of items and a candidate recommender item.

The process of a typical Item Neighborhood is generally divided into three steps as the user Neighborhood :

1. **Similarity computation** (among the items). The recommendation system computes the similarities between the items , Neighborhood Similar items are regrouped in a subset namely the « neighborhood » of items. .
2. **Rating prediction** (using ratings of similar items ). item-based prediction of a rating  $r_{ui}$  :

$$r_{ui} = \bar{r}_i + \frac{\sum_{j \in N_u(i)} s_{ij}(r_{uj} - \bar{r}_j)}{\sum_{j \in N_u(i)} |s_{ij}|}$$

3. **Top- $N$  recommendation.** Are obtained by choosing the  $N$  items which provide most satisfaction to the user according to prediction by decreasing order of predicted rating. 3.9.1

The experimental comparative evaluation result between the similar items as the critical step in the proposed item-based recommendation in determining the similarity between the items. In implementing this model, applied different similarity metrics as the Cosine, Adjusted Cosine, Euclidean, Manhattan Distance, and Distance Covariance Distance.

### 5.5.3 Similarity Metrics

In this section, the similarity metrics used in neighborhood-based collaborative filtering are described.

**Cosine** This metric first transfers inputs to the vector space and then utilizes the angle between two rate vectors as the degree of the similarity of two users,or items.

**Adjusted Cosine** This metric is extensions of the Cosine metric which are defined similar to Pearson correlation. The Adjusted cosine measure calculates the correlation value between two users.

**Euclidean Distance** The euclidean distance compute the correlation between two users (or two items) is the absolute value of the numerical difference of their coordinates.

**Manhattan Distance** The Manhattan Distance it's the one norm of the distance the measure the similarity between two users (or two items)

**Covariance Distance** Since we want to have a similarity measure, we take its complement to one. Hence, the more ratings are close, the more the distance is short. And therefore the more users or items are considered to be alike, Covariance Distance It's new similarity metrics that have been used for the first time in the Recommendation systems Domaine. Distance covariance is recently introduced. dependency measures between random vectors, where we will use the empirical distance covariance to compute the similarity between the two users(or two items).

Measure	User Similarity Measure	Item similarity Measure
Cosine	$cosine(u, v) = \frac{\sum_{i \in I_{uv}} r_{ui} r_{vi}}{\sqrt{\sum_{i \in I_u} r_{ui}^2} \sqrt{\sum_{i \in I_v} r_{vi}^2}}$	$cosine(i, j) = \frac{\sum_{u \in U_{ij}} r_{ui} r_{uj}}{\sqrt{\sum_{u \in U_i} r_{ui}^2} \sqrt{\sum_{u \in U_j} r_{uj}^2}}$
Adjusted cosine	$ACosine(u, v) = \frac{\sum_{i \in I_{uv}} (r_{ui} - \bar{r}_i)(r_{vi} - \bar{r}_i)}{\sqrt{\sum_{i \in I_{uv}} (r_{ui} - \bar{r}_i)^2} \sqrt{\sum_{i \in I_{uv}} (r_{vi} - \bar{r}_i)^2}}$	$ACosine(i, j) = \frac{\sum_{u \in U_{ij}} (r_{ui} - \bar{r}_u)(r_{uj} - \bar{r}_u)}{\sqrt{\sum_{u \in U_{ij}} (r_{ui} - \bar{r}_u)^2} \sqrt{\sum_{u \in U_{ij}} (r_{uj} - \bar{r}_u)^2}}$
Covariance Distance	$\hat{V}^2(u, v) = \frac{1}{m^2} \sum_{i,j=1}^m a_{ij} b_{ij} + \frac{1}{m^4} \sum_{i,j=1}^n a_{ij} \sum_{i,j=1}^m b_{ij} - \frac{2}{m^3} \sum_{i,j,k=1}^m a_{ij} b_{ik}$	$\hat{V}^2(i, j) = \frac{1}{n^2} \sum_{u,v=1}^n a_{uv} b_{uv} + \frac{1}{n^4} \sum_{u,v=1}^n a_{uv} \sum_{u,v=1}^n b_{uv} - \frac{2}{n^3} \sum_{u,v,k=1}^n a_{uv} b_{uk}$
Euclidean distance	$d(u, v) = \sqrt{\sum_{i \in I_{uv}} (r_{vi} - r_{ui})^2}$	$d(i, j) = \sqrt{\sum_{u \in U_{ij}} (r_{uj} - r_{ui})^2}$
Manhattan distance	$d_1(u, v) = \sum_{i \in I_{uv}} ( r_{vi} - r_{ui} )$	$d_1(i, j) = \sum_{u \in U_{ij}} ( r_{uj} - r_{ui} )$

Table 5.2: Table Of The Used Similarities Metrics

**Important:** Notice that due to the computational limitations (RAM Saturation problem) of a large number of the users in our dataset, we limited our user neighborhood implementation model by applying only the Euclidean distance, Cosine, and Manhattan Distance similarity metrics.

### Implementation of the User Neighborhood

- We identify the set  $G_u$  of  $k$  most similar users.  $G_u$  is the group users similar to the active user  $u$ . The similarity between two users  $u$  and  $v$  can be measured by the similarity metrics.
- Find the set  $C$  of candidate items, purchased by the group and not purchased by the active user  $u$ .
- Candidate items have to be the most frequent items purchased by the group. Aggregate ratings of users in  $G_u$  to make predictions for user  $u$  on items he has not already purchased.
- The Top- $N$  recommendations are obtained by choosing the  $N$  items which provide most satisfaction to the user according to prediction.

*Step 1. Identify  $G_u$  the set of  $k$  users similar to an active user  $u$*

Defined the function the compute the User-Item matrix  $R$ . Each row of the rating matrix  $R$  represents ratings of user  $u$  on all items of the database. Missing ratings are replaced with 0,  $U$ th row of the rating matrix  $R$ . Ratings of user  $u$  on all items in the database,

```
def ratings_matrix(data1):
    return csr_matrix(pd.crosstab(data1.userid, data1.itemid,
        data1.rating,aggfunc=sum).fillna(0).values)
R = ratings_matrix(data1)
```

We create a nearest neighbors model with sklearn through the function `createmodel()`. This function creates and fit a nearest neighbors model with user's ratings,  $n\_neighbors = 51$  neighbors, define the number of neighbors to return. With  $k = 50$  neighbors,  $|G_u| = 51$  as  $G_u$  contains 50 similar users added to the active user  $u$ . That is why  $n\_neighbors = 51$ .

```
def create_model(rating_matrix, metric):
    model = NearestNeighbors(metric=metric, n_neighbors=51,
        algorithm='brute')
    model.fit(rating_matrix)
    return model
```

Then we defined `nearestneighbors()` function that returns the knn users for each user.

```
def nearest_neighbors(rating_matrix, model):
    '''-rating_matrix: rating matrix of shape (nbusers, nbitems)
    -model: nearest neighbors model
    -similarities :Distances of the neighbors from the referenced user.
    -neighbors :Neighbors of the referenced user in decreasing order
        of similarities.'''
    similarities, neighbors = model.kneighbors(rating_matrix)
    return similarities[:, 1:], neighbors[:, 1:]
```

*Similarities and Neighbors using Cosine Similarity metric*

```
model_cos = create_model(rating_matrix=R,metric='cosine')
similarities_cos, neighbors_cos = nearest_neighbors(R, model_cos)
similarities_cos, neighbors_cos
```

```
(array([[0.0513167, 0.0513167, 0.0513167, ..., 0.0513167, 0.0513167,
        0.0513167],
       [0.      , 0.      , 0.      , ..., 0.      , 0.      ,
        0.      ],
       [0.      , 0.      , 0.      , ..., 0.      , 0.      ,
        0.      ],
       ...,
       [1.      , 1.      , 1.      , ..., 1.      , 1.      ,
        1.      ],
       [0.      , 0.      , 0.      , ..., 0.      , 0.      ,
        0.      ],
       [1.      , 1.      , 1.      , ..., 1.      , 1.      ,
        1.      ]]), array([[33754, 27427, 40713, ..., 46026, 1712, 1064],
       [ 8749, 43319, 33058, ..., 53374, 27442, 39403],
       [56230, 39487, 30542, ..., 48147, 30951, 48002],
       ...,
       [46377, 46376, 46375, ..., 46382, 46392, 46406],
       [69029, 43427, 36259, ..., 54814, 37488, 56931],
       [46377, 46376, 46375, ..., 46382, 46392, 46406]]))
```

Figure 5.10: Similarities and Neighbors using Cosine Similarity metric

*Similarities and Neighbors using Euclidean Distance Similarity metric*

```
model_euclidean=
create_model(rating_matrix=R,metric='euclidean')
similarities_euclidean, neighbors_euclidean=
nearest_neighbors(R,model_euclidean)
similarities_euclidean, neighbors_euclidean
```

*Similarities and Neighbors using Manhattan Distance Similarity metric*

```
model_manhattan= create_model(rating_matrix=R,metric='manhattan')
similarities_manhattan, neighbors_manhattan =
nearest_neighbors(R,model_manhattan)
similarities_manhattan, neighbors_manhattan
```

**Step 2. Find candidate items**

The set  $C$  of candidate items are the most frequent ones purchased by users in  $G_u$  for an active user  $u$  and not purchased by  $u$ . Function *findcandidateitems()* find items purchased by these similar users as well as their frequency. Note that the frequency of the items in the set  $C$  can be computed by just counting the actual occurrence frequency of that items.

**Find candidate items with Cosine similarity metric**

```
def find_candidate_items_cos(userid):
    """
    Find candidate items for an active user
    :param userid : active user
    :param neighbors : users similar to the active user
    :return candidates : top 50 of candidate items
    Gu_items : frequent items of  $G_u$  in decreasing order of frequency.
    activeitems : items already purchased by the active user.
    candidates :frequent items of  $G_u$  not purchased by the active user  $u$ .
    """
    user_neighbors = neighbors_cos[userid]
```



```

activities = data1.loc[data1.userid.isin(user_neighbors)]
# sort items in decreasing order of frequency
frequency = activities.groupby('itemid')['rating'].count()
.reset_index(name='count').sort_values(['count'], ascending=False)
Gu_items = frequency.itemid
active_items = data1.loc[data1.userid == userid].itemid.to_list()
candidates_cos =
np.setdiff1d(Gu_items, active_items, assume_unique=True)[:50]
return candidates_cos

```

We applied the same function to get the *Find candidate items with Euclidean Distance similarity metric* and *Find candidate items with Manhattan Distance similarity metric*

### Step 3. Rating prediction

Predict the score of  $u$  on a candidate item  $i$ .

Rating prediction with the cosine similarity metric

```

def predict_cos(userid, itemid):
    """
    :param
    - userid : user id for which we want to make prediction.
    - itemid : item id on which we want to make prediction.
    :return
    - r_hat_cos : predicted rating of user userid on item itemid. """
    user_similarities = similarities_cos[userid]
    user_neighbors = neighbors_cos[userid]
    # get mean rating of user userid
    user_mean = mean[userid]
    # find users who rated item 'itemid'
    iratings = np_ratings[np_ratings[:, 1].astype('int') == itemid]
    # find similar users to 'userid' who rated item 'itemid'
    suri = iratings[np.isin(iratings[:, 0], user_neighbors)]
    # similar users who rated current item (surci)
    normalized_ratings = suri[:, 4]
    indexes =
    [np.where(user_neighbors == uid)[0][0]
    for uid in suri[:, 0].astype('int')]
    sims = user_similarities[indexes]
    num = np.dot(normalized_ratings, sims)
    den = np.sum(np.abs(sims))
    if num == 0 or den == 0:
    return user_mean
    r_hat_cos= user_mean +
    np.dot(normalized_ratings, sims) / np.sum(np.abs(sims))
    return r_hat_cos

```

We applied the same function to get the predict with Euclidean Distance similarity metric and predict with Manhattan Distance similarity metric.

Rating prediction for the active user on each candidate item with Cosine similarity metric

```
def user2userPredictions_cos(userid, pred_path):  
    """  
    Make rating prediction for the active user on each  
    candidate item and save in file prediction.csv  
    :param  
    - userid : id of the active user  
    - pred_path : where to save predictions  
    """  
    # find candidate items for the active user  
    candidates_cos= find_candidate_items_cos(userid)  
    # loop over candidates items to make predictions  
    for itemid in candidates_cos:  
        # prediction for userid on itemid  
        r_hat_cos = predict_cos(userid, itemid)  
        # save predictions  
        with open(pred_path, 'a+') as file:  
            line = '{}',{},{ }\n'.format(userid, itemid, r_hat_cos)  
            file.write(line)
```

We applied the same function to get the *user2userPredictions with Euclidean Distance similarity metric* and *user2userPredictions with Manhattan Distance similarity metric*.

#### *Step 4. Top-N recommendation*

Function `user2userRecommendation()` reads predictions for a given user and return the list of items in decreasing order of predicted rating.

```
def user2userRecommendation_cos(userid):  
    saved_predictions = 'predictions_cos.csv'  
    predictions =  
    pd.read_csv(saved_predictions, sep=',',  
    names=['userid', 'itemid', 'predicted_rating'])  
    predictions = predictions[predictions.userid==userid]  
    df =  
    predictions.sort_values(by=['predicted_rating'], ascending=False)  
    df.userid = df.userid.tolist()  
    df.itemid =df.itemid.tolist()  
    df = pd.merge(df, data2, on='itemid', how='inner')  
    df = df.drop_duplicates(subset=['itemid'])  
    return df
```

We applied the same function to get the *user2userRecommendation with Euclidean Distance similarity metric* and *user2userRecommendation with Manhattan Distance similarity metric*. *Let us make top recommendation for a given user.*

```
user2userRecommendation_cos(244)
```

	userid	itemid	predicted_rating
<b>0</b>	244	358	1.439946
<b>2221</b>	244	84	1.267857
<b>2322</b>	244	668	0.625000
<b>2444</b>	244	662	0.125000
<b>2675</b>	244	409	0.101636

Figure 5.11: TOP Recommendation using Cosine similarity metric

```
user2userRecommendation_euclidean(244)
```

	userid	itemid	predicted_rating
<b>0</b>	244	477	2.725000
<b>234</b>	244	206	1.925000
<b>445</b>	244	358	1.553571
<b>2666</b>	244	625	1.375000
<b>2677</b>	244	309	1.325000

Figure 5.12: TOP Recommendation using Euclidean Distance similarity metric

```
user2userRecommendation_manhattan(244)
```

	userid	itemid	predicted_rating
<b>0</b>	244	477	2.725000
<b>234</b>	244	358	1.553571
<b>2455</b>	244	309	1.325000
<b>2559</b>	244	119	0.125000
<b>2612</b>	244	662	0.125000
<b>2843</b>	244	444	0.125000

Figure 5.13: TOP Recommendation using Manhattan Distance similarity metric

*Evaluation of user-based recommendation system with the Manhattan distance similarity metric*

We split the dataset with a ratio of 75%, 25% into training set and test set ,We use the

training set to obtain the training model, and then apply the model to test set to test the performance of our models.

```
raw_examples, raw_labels = get_examples(data1, labels_column='rating')
# train test split
(x_train, x_test), (y_train, y_test) =
train_test_split(examples=raw_examples, labels=raw_labels)
def evaluate_cos(x_test, y_test):
    preds = CFU_MAE_cos =
    np.sum(np.absolute(y_test - np.array(preds))) / x_test.shape[0]
    print('>> CFU_MAE_cos: ' + str(CFU_MAE_cos))
    CFU_MSE_cos = mean_squared_error(y_test, np.array(preds))
    print('>> CFU_MSE CFU_MSE_cos: ' + str(CFU_MSE_cos))
    CFU_RMSE_cos = math.sqrt(CFU_MSE_cos)
    print('>> CFU_RMSE_cos' + str(CFU_RMSE_cos))
    evaluate_cos.list1=[CFU_MAE_cos,CFU_MSE_cos,CFU_RMSE_cos]
    return evaluate_cos.list1
evaluate_cos(x_test, y_test)
```

We applied the same function to evaluate the *User-based model using Euclidean Distance similarity metric* and *User-based model using Manhattan Distance similarity metric*.

## 5.5.4 Results of the User-based model

	Similarity Measure	RMSE	MAE	MSE
User-Based	Euclidean Distance	<i>0.90187</i>	<i>0.64881</i>	<i>0.81338</i>
	Manhattan Distance	<i>0.89382</i>	<i>0.64353</i>	<i>0.79891</i>
	Cosine	<i>0.93452</i>	<i>0.69218</i>	<i>0.87333</i>

Table 5.3: Evaluation Results of the Similarity Measures on User-Based

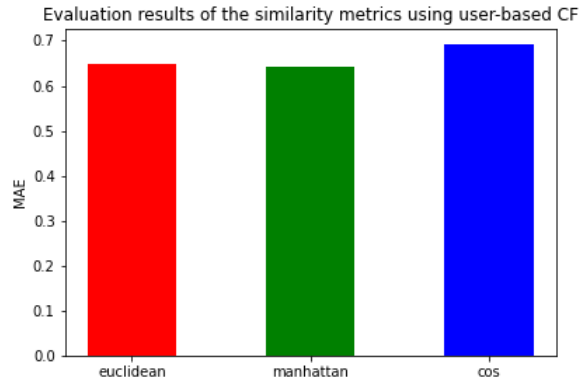


Figure 5.14: MAE Evaluation results of the similarity measures using user-based

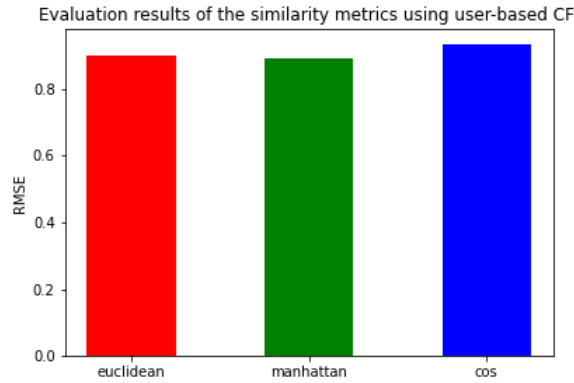


Figure 5.15: RMSE Evaluation results of the similarity measures using user-based

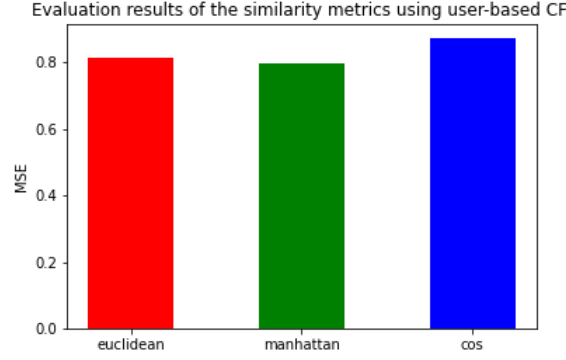


Figure 5.16: MSE Evaluation results of the similarity measures using user-based

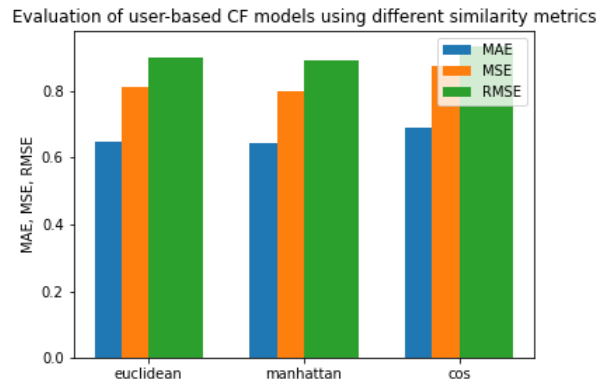


Figure 5.17: Performance Comparison Between The Evaluation results of the similarity measures using user-based

We applied user-based recommendation system using three similarity measures cosine, mahattan and euclidean distance , which was evaluated using RMSE, MSE, and MAE.

In this section, we present the obtained findings, We perform user-based Neighborhood. The experimental study was carried out using different similarity measures, and After applying the user-based CF models, we obtained the results shown in 5.3 In the table and the graphs, we present the values reached by each similarity measure in terms of RMSE, MSE, and MAE. As shown in Table 5.3, Manhattan distance similarity provides the best result for the three evaluation metrics RMSE, MSE, and MAE, while the lower those evaluation metrics' values are, the better our model's accuracy. For instance, Manhattan distance gets the best values 0.64353(MAE), 0.79891(MSE), 0.89382(RMSE).

The main conclusion that we can deduce from this part of the experimental study is that the Manhattan distance measure outcores the other measures it was applied in the user-based approach with the YASSIR EXPRESS dataset. At the same time different results may be obtained by using another dataset.

### Implementation of the Item Neighborhood

- First identify the  $k$  most similar items for each item in the dataset and record the corresponding similarities. To compute similarity between two items, we applied the Cosine, Euclidean Distance, Manhattan Distance, Covariance distance, Adjusted Cosine.
- Finding candidate items for user  $u$ , where the set  $I_u$  of items already rated by user  $u$ , then we take the union of similar items as  $C$  for all items in  $I_u$ , then we exclude from the set  $C$  all items in  $I_u$ , to avoid recommending to a user items he has already purchased.
- Find similarity between each candidate item and the set  $I_u$ .
- Rank candidate items according to their similarities to  $I_u$ .
- Top- $N$  recommendations for a given user.

#### *Similarities and Neighbors with covariance distance Similarity metric*

```
def cov_dist(x, y):
    return dcor.distance_covariance(x, y, method='AVL')

def dist_it(np_ratings, nb_items):
    similarities_cov_it = np.zeros(shape=(nb_items, nb_items))
    similarities_cov_it.fill(-1)
    def _progress(count):
        sys.stdout.write
        ('\rComputing similarities. Progress status: %.1f%%' %
         (float(count / nb_items)*100.0))
        sys.stdout.flush()
    items = sorted(data1.itemid.unique())
    for i in items[:-1]:
        for j in items[i+1:]:
            scores =
            np_ratings[(np_ratings[:, 1] == i) | (np_ratings[:, 1] == j), :]
            vals, count = np.unique(scores[:, 0], return_counts = True)
            scores = scores[np.isin(scores[:, 0], vals[count > 1]), :]
            if scores.shape[0] > 2:
                x = scores[scores[:, 1].astype('int') == i, 2]
                y = scores[scores[:, 1].astype('int') == j, 2]

            m1 = max(len(x), len(y))
            x = np.concatenate((x, np.zeros(m1-len(x))))
            y = np.concatenate((y, np.zeros(m1-len(y))))
            w = dcor.distance_covariance(x, y)
            similarities_cov_it[i, j] = w
            similarities_cov_it[j, i] = w
            _progress(i)
            _progress(nb_items)
    #get neighbors by their neighbors in decreasing
    #order of similarities
    neighbors_cov_it =
    np.flip(np.argsort(similarities_cov_it), axis=1)
    #sort similarities in decreasing order
    similarities_cov_it = np.flip(np.sort(similarities_cov_it), axis=1)
```

```
return similarities_cov_it , neighbors_cov_it
```

```
nb_items = data1.itemid.nunique()
similarities_cov_it, neighbors_cov_it =
dist_it(np_ratings, nb_items=nb_items)
```

We computed the similarities and the neighbors of the Cosine, Adjusted Cosine, Euclidean, and the manhattan similarity metrics in the same way we previously calculated them on the user-based.

*Finding candidate items for user  $u$  using the covariance distance similarity metric*

```
def candidate_items_cov_it(userid):
    """
    :param userid : user id for which we wish to find candidate items
    :return : I_u, candidates
    """
    #1. Finding the set I_u of items already rated by user userid
    I_u = np_ratings[np_ratings[:, 0] == userid]
    I_u = I_u[:, 1].astype('int')
    #2. Taking the union of similar items for all
    # items in I_u to form the set of candidate items
    c_cov = set()
    for iid in I_u:
        #add the neighbors of item iid in the set of candidate items
        c_cov.update(neighbors_cov_it[iid])
    c_cov = list(c_cov)
    #3. exclude from the set C all items in I_u.
    candidates_cov_it = np.setdiff1d(c_cov, I_u, assume_unique=True)
    return I_u, candidates_cov_it

#i_{u} number of items purchased by user u
test_user = uencoder.transform([1])[0]
i_u, u_candidates_cov_it = candidate_items_cov_it(test_user)
```

*Find similarity between each candidate item and the set  $I_u$  using Distance Covariance Similarity metric*

```
def similarity_with_Iu_cov_it(c_cov, I_u):
    """:param c : itemid of a candidate item
        :param I_u : set of items already purchased by a given user
        :return w : similarity between c and I_u """
    w = 0
    for iid in I_u :
        #get similarity between itemid and c, if c
        #is one of the k nearest neighbors of itemid
        if c_cov in neighbors_cov_it[iid] :
            w = w +
            similarities_cov_it[iid, neighbors_cov_it[iid] == c_cov][0]
    return w
```

*Rank candidate items according to their similarities using Covariance Distance Similarity metric*



---

```
def rank_candidates_cov(candidates_cov_it, I_u):
    """rank candidate items according to their similarities with i_u
    :param candidates : list of candidate items
    :param I_u : list of items purchased by the user
    :return ranked_candidates : dataframe of
    candidate items, ranked in descending order
    of similarities with I_u"""
    #list of candidate items mapped to their
    #corresponding similarities to I_u
    sims =
    [similarity_with_Iu_cov_it(c_cov, I_u) for c_cov in candidates_cov_it ]
    candidates_cov_it = iencoder.inverse_transform(candidates_cov_it )
    mapping = list(zip(candidates_cov_it , sims))
    ranked_candidates_cov = sorted(mapping,
    key=lambda couple:couple[1], reverse=True)
    return ranked_candidates_cov
```

Now that we defined all functions necessary to build our item Top-N recommendation, let's define function that makes top- recommendations for a given user.

```
def topn_recommendation_cov(userid, N=50):
    #find candidate items
    I_u, candidates_cov_it = candidate_items_cov_it(userid)
    #rank candidate items according to their similarities with I_u
    ranked_candidates_cov = rank_candidates_cov( candidates_cov_it, I_u)
    # get the first N row of ranked_candidates to
    build the top N recommendation list
    topn_cov = pd.DataFrame(ranked_candidates_cov[:N],
    columns=['itemid', 'similarity_with_Iu_cov_it'])
    topn_cov = pd.merge(topn_cov, data2, on='itemid', how='inner')
    topn_cov = topn_cov.drop_duplicates(subset=['itemid'])
    return topn_cov
```

Before recommending the previous list to the user, we can go further and predict the ratings the user would have given to each of these items, sort them in descending order of prediction and return the reordered list as the new top N recommendation list.

#### *Rating prediction using distance covariance similarity metric*

```
def predict_cov(userid, itemid):
    """
    Make rating prediction for user userid on item itemid
    :param userid : id of the active user
    :param itemid : id of the item for which we are making prediction
    :return r_hat : predicted rating
    """
    #Get items similar to item itemid
    #with their corresponding similarities
    item_neighbors = neighbors_cov_it[itemid]
    item_similarities= similarities_cov_it[itemid]
    #get ratings of user with id userid
    uratings = np_ratings[np_ratings[:, 0].astype('int') == userid]
    #similar items rated by item the user of i
```

```

siru = uratings[np.isin(uratings[:, 1], item_neighbors)]
scores = siru[:, 2]
indexes =
[np.where(item_neighbors == iid)[0][0]]
for iid in siru[:,1].astype('int'):
sims =
item_similarities[indexes]
dot = np.dot(scores, sims)
som = np.sum(np.abs(sims))
if dot == 0 or som == 0:
return mean[userid]
return dot / som

```

*Predict ratings the user would have given to the previous top-N recommended list with the Distance Covariance similarity metric*

```

def topn_prediction_cov(userid):
    """
    :return topn :
    initial topN recommendations returned by the function item2item_topN
    :return topn_predict :
    topN recommendations reordered according to rating predictions
    """
    #make top N recommendation for the active user
    topn_cov = topn_recommendation_cov(userid)
    #get list of items of the top N list
    itemids = topn_cov.itemid.to_list()
    predictions = []
    #Make prediction for each item in the top N list
    for itemid in itemids:
        r = predict_cov(userid, itemid)
        predictions.append((itemid,r))
    predictions =
    pd.DataFrame(predictions, columns=['itemid','prediction'])
    #merge the predictions to top N_list and rearrange
    #the list according to predictions
    topn_predict_cov =
    pd.merge(topn_cov, predictions, on='itemid', how='inner')
    topn_predict_cov =
    topn_predict_cov.sort_values(by=['prediction'],ascending=False)
    return topn_cov, topn_predict_cov

topn_cov, topn_predict_cov = topn_prediction_cov(userid=test_user)
topn_predict_cov

```

	itemid	similarity_with_Iu_cov_it	prediction
	8	373	0.088735
	5	342	0.094360
	23	583	0.072364
	48	587	0.060425
	17	116	0.075117

Figure 5.18: Top recommendation using the Distance Covariance similarity metric

### 5.5.5 Results of the Item-based model

We split the dataset with a ratio of 75%, 25% into training set and test set ,We use the training set to obtain the training model, and then apply the model to test set to test the performance of our models.

Evaluation of the item-based recommendation systems with Covariance distance ,Euclidean Distance,Manhattan Distance,Cosine,Adjusted Cosine Similarity metrics was applied using the same function as we did on the user-based and we obtained this results

	Similarity Measure	RMSE	MAE	MSE
Item-Based	Euclidean Distance	<i>0.93329</i>	<i>0.676635</i>	<i>0.87103</i>
	Manhattan Distance	<i>0.93345</i>	<i>0.67690</i>	<i>0.87133</i>
	Cosine	<i>0.94761</i>	<i>0.67799</i>	<i>0.89796</i>
	Adjusted Cosine	<b>2.40673</b>	<b>1.67612</b>	<b>5.79237</b>
	Covariance Distance	<b>3.14617</b>	<b>2.25624</b>	<b>9.89840</b>

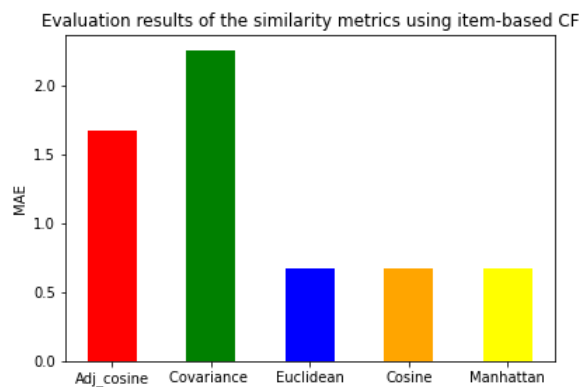


Figure 5.19: MAE Evaluation results of the similarity measures using item-based

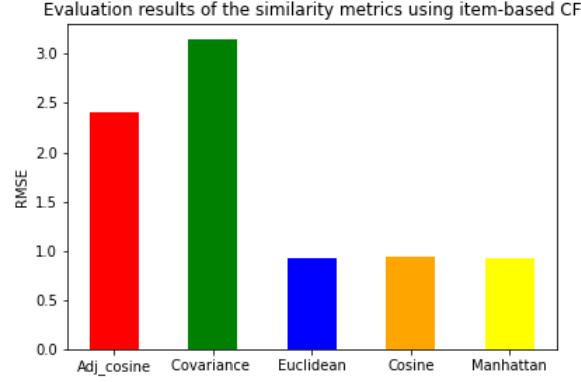


Figure 5.20: RMSE Evaluation results of the similarity measures using Item-based

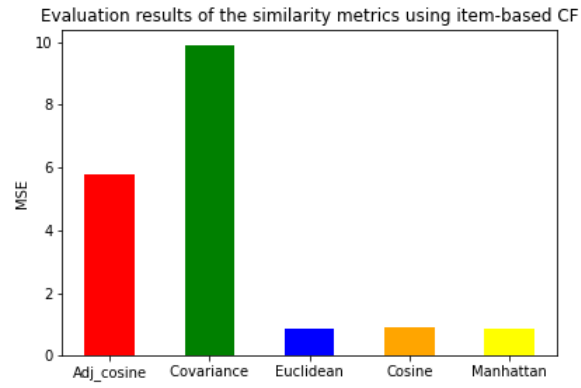


Figure 5.21: MSE Evaluation results of the similarity measures using Item-based

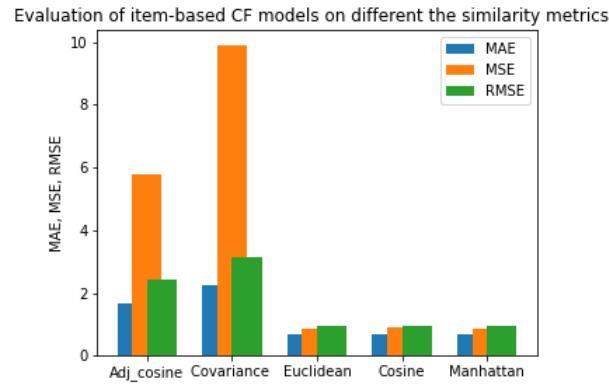


Figure 5.22: Performance Comparison Between The Evaluation results of the similarity measures using item-based neighborhood

We applied Item-based neighborhood, which was evaluated using RMSE, MSE, and MAE. In this section, we present the obtained findings.

We performed item-based Neighborhood. The experimental study was carried out using different similarity measures Cosine, Adjusted Cosine, Manhattan Distance, Euclidean Distance, and the Covariance Distance we. Got the results are shown in 5.5.5 In the Table and the graphs that present the values reached by each similarity measures in terms of RMSE, MSE, MAE.

As displayed in the Table, the Euclidean distance similarity provides the best result for the three evaluation metrics RMSE, MSE, and MAE . While the lower those evaluation metrics values are the better our models accuracy. For instance, Euclidean distance gets the best values:

0.676635 (MAE), 0.87103 (MSE), 0.93329 (RMSE).

In contrast, the worst similarity metric was the distance covariance in terms of effectiveness of the model :

2.25624 (MAE), 9.89840 (MSE), 3.14617 (RMSE).

The main conclusion that can be deduced from this part of the experimental study is that the Euclidean distance measure outcores the other measures if it is applied in the Items-based approach on the YASSIR EXPRESS dataset, while we may obtain different results by using another dataset.

### 5.5.6 Comparing the user-based and the item-based

We present the obtained findings, we perform two types of experiments, using two models user-based and item-based using the first rating approach .

	Similarity Measure	RMSE	MAE	MSE
User-Based	Euclidean Distance	<i>0.90187</i>	<i>0.64881</i>	<i>0.81338</i>
	Manhattan Distance	<i>0.89382</i>	<i>0.64353</i>	<i>0.79891</i>
	Cosine	<i>0.93452</i>	<i>0.69218</i>	<i>0.87333</i>
Item-Based	Euclidean Distance	<i>0.93329</i>	<i>0.676635</i>	<i>0.87103</i>
	Manhattan Distance	<i>0.93345</i>	<i>0.67690</i>	<i>0.87133</i>
	Cosine	<i>0.94761</i>	<i>0.67799</i>	<i>0.89796</i>
	Adjusted Cosine	<b>2.40673</b>	<b>1.67612</b>	<b>5.79237</b>
	Covariance Distance	<b>3.14617</b>	<b>2.25624</b>	<b>9.89840</b>

Table 5.4: Results of the Similarity Measures on User-Based and Item-based First rating approach

The comparison between the user-based and item-based recommendation systems using the different similarity measures the evaluation results in terms of their effectiveness using the **RMSE,MSE,MAE** . The evaluation results shown on the table 5.4 manhattan distance has the best results on the user-based are RMSE , MAE and MSE are 0.89382,0.64353 and 0.79891, respectively, in contrast, the performance of the user-based model was the lowest using the Cosine similarity metric. On the other hand, we can observe that the item-based delivers the best results using. Euclidean distance RMSE, MAE and MSE 0.93329, 0.676635 and 0.87103

Nevertheless, the manhattan and euclidean distance the evaluated measures provide close results, meanwhile, the performance of the item-based model was the lowest while using the covariance distance similarity metric .

As shown in 5.4 we can remark that it is evidently clear that user-based provides results better than those offered by the item-based for all similarity metrics in terms of their performance .

We present the obtained findings, we perform two types of experiments, user-based and item-based using the second rating approach .

	Similarity Measure	RMSE	MAE	MSE
User-Based	Euclidean Distance	<i>3.10530</i>	<i>2.77726</i>	<i>3.22170</i>
	Manhattan Distance	<i>3.09143</i>	<i>2.76914</i>	<i>3.19123</i>
	Cosine	<i>3.23437</i>	<i>2.92977</i>	<i>3.52367</i>
Item-Based	Euclidean Distance	<i>3.33196</i>	<i>2.96546</i>	<i>3.77413</i>
	Manhattan Distance	<i>3.33244</i>	<i>2.96580</i>	<i>3.77541</i>
	Cosine	<i>3.2835</i>	<i>2.84632</i>	<i>1.64744</i>
	Adjusted Cosine	<i>5.18968</i>	<i>4.72174</i>	<i>12.17406</i>
	Covariance Distance	<i>5.82345</i>	<i>5.25817</i>	<i>16.61884</i>

Table 5.5: Results of the Similarity Measures on User-Based and Item-based second rating approach

We compare the user-based and the item-based recommendation systems using the evaluation results in terms of their effectiveness using the RMSE,MSE,MAE.

We show the results provided by the different similarity measures using the item-based and the user-based provided by the second rating formula. , The evaluation results shown on the table 5.5 manhattan distance has the best results on the user-based are RMSE , MAE and MSE are 3.09143, 2.76914 and 3.19123, respectively, in contrast, the performance of the user-based model was the lowest using the Cosine similarity metric. On the other hand, we can observe that the item-based delivers the best results using Cosine RMSE, MAE and MSE 3.2835 , 2.84632, 1.64744

Nevertheless, the manhattan and euclidean distance the evaluated measures provide close results, meanwhile, the performance of the item-based model was the lowest while using the covariance distance similarity metric with RMSE, MAE and MSE 5.82345, 5.25817, 16.61884. As shown in 5.5 we can remark that it is evidently clear that user-based provides results better than those offered by the item-based for all similarity metrics in terms of their performance. We can determine the user-based and the item-based recommendation systems using the evaluation results in terms of their effectiveness using the RMSE, MSE, and MAE by using the First and second rating approaches that we present their findings on the 5.4 and 5.5 respectively that the performance of the first rating approach are the best, in addition, the second approach has provided the worst recommendation of the items with arent relevant to the users preferences.

### 5.5.7 Conclusion

In this section, we conducted a comparative study of important clues related to the similarity metrics and their relation commonly used in the field of collaborative filtering recommendation systems to determine the effectiveness of our new similarity metric the covariance distance. In fact, we conducted an experimental study for the user-based and the item-based where each time a different similarity metrics to prove that the performance of the Recommendation Systems is strongly dependent on the used similarity metric where for the user-based we used the cosine ,manhattan and the euclidean distance ,we used a set of well-known evaluation metrics RMSE, MSE, and MAE, and on the item-based we applied the cosine , adjusted cosine, manhattan, euclidean, and the covariance distance. Findings showed that user-based and item-based don't share the same similarity metric that offers the best performance to the system and our new similarity metric, the distance covariance, has been the worst in terms of performance and effectiveness of the item neighborhood by using the first and the second rating approach. Indeed, the manhattan distance similarity gives the best performance for the user-based similarity, using the two rating approach,

The Euclidean distances for the item-based are on the first rating approach. Meanwhile, the Cosine similarity has the best performance on the second rating approach. As shown in previous sections, a bad choice of a similarity metric will undoubtedly lead to lower performance and vice versa. Also, we can deduce that there is no specific similarity metric suitable for both user-based and item-based approaches. Some measures perform better in one approach than the other. Moreover, this study demonstrates that similarity measure choice is not the only challenge for CF-based RSs.

## 5.6 Location-Based Recommendation System

### *Geographical Overview of Restaurants*

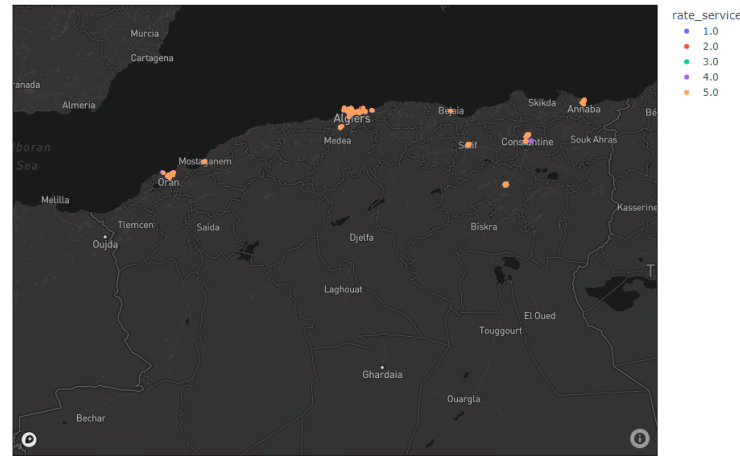


Figure 5.23: Geographical Overview of Restaurants

#### ***K-Means Clustering***

Based on the above visualization, we can see that there are areas where there is a high density of Restaurants while there are also other areas that have fewer Restaurants. Depending on where a user is currently located, it would make sense to recommend restaurants that are close to his location regardless of where he is in a city. Hence, one way where we can group restaurants together based on geographical location Coordinates latitude and longitude is using the K-Means Clustering Algorithm.

The **K-Means** algorithm predicts the cluster where the user is located in and pulls out this cluster's top 10 restaurants and recommends them to him.

	restaurant_address_lng	restaurant_address_lat
0	-0.653765	35.675536
1	-0.653765	35.675536
2	3.185681	36.722018
3	3.185681	36.722018
4	3.185681	36.722018

Figure 5.24: The table represents the longitude and latitude of the restaurants.

The silhouette Method is a method that find the optimal number of clusters and interpretation and validation of consistency within clusters , The silhouette method computes silhouette coefficients of each point that measure how much a point is similar to its own cluster compared to other clusters ,by providing a graphical representation . The silhouette value ranges between  $[-1, 1]$ , where a high value indicates that the object is well matched to its own cluster and poorly matched to neighboring clusters.



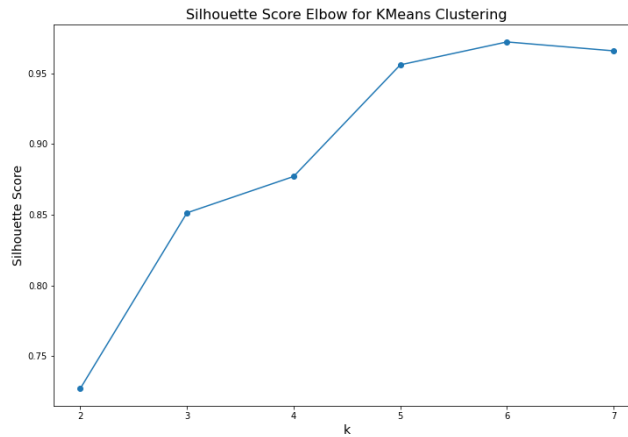


Figure 5.25: Silhouette Score Elbow for KMeans Clustering.

The Silhouette Score reaches its global maximum at the optimal  $k=6$  as we can see in graphical representation of the Silhouette Score Elbow for KMeans Clustering.

#### ***Location-Based Recommendation Function***

```
# Creating Location-Based Recommendation Function
def location_based_recommendation(df,
restaurant_address_lat, restaurant_address_lng):
#Predict the cluster for longitude and latitude provided
init_cluster =
kmeans.predict(
np.array([restaurant_address_lng, restaurant_address_lat]).
reshape(1,-1))[0]
print("This restaurant belongs to init_cluster:", init_cluster)
return
df[df['init_cluster']==init_cluster].iloc[0:]
[['restaurantname', 'restaurant_address_lat',
'restaurant_address_lng','rate_service',
'restaurant_Score','init_cluster']]

# Top 10 Recommended Restaurants based on Location of a given user
df= location_based_recommendation( top_restaurants,user1.latitude,
user1.longitude)
```

#### ***Haversine distance***

Haversine distance is the Great Circle Distance (GCD) between two geographic coordinates. A GCD is incidentally the shortest distance between the two coordinates(long/lat). This is almost similar to a Euclidean distance, except that we are accounting for the spherical nature of the Earth (yes, we are generalizing the Earth as a sphere with a radius of 6378000 KM to make our lives easier)

```
def loc_rec(user_latitude,user_longitude,display=10,distance=3):
df=location_based_recommendation( top_restaurants, user1.latitude,
user1.longitude)
df = df.drop_duplicates()
df['Distance']=df.apply(lambda x :Distance(user1.latitude, user1.longitude,
x.restaurant_address_lat, x.restaurant_address_lng),axis=1)
```

```

df=df[df['Distance']<=distance]
return df.head(display).reset_index(drop=True)

loc_rec(36.726,3.0827)
#latitude and longitude <location> of a given user

```

	restaurant_address_lat	restaurant_address_lng	rate_service	restaurant_Score	init_cluster	Distance
0	36.735729	3.101273	5.0	144	1	1.98
1	36.740225	3.089829	5.0	117	1	1.71
2	36.738411	3.093935	5.0	59	1	1.71
3	36.733574	3.108903	5.0	11	1	2.49
4	36.718308	3.091916	5.0	177	1	1.19
5	36.728701	3.076883	5.0	841	1	0.60
6	36.729527	3.084235	5.0	419	1	0.42
7	36.747394	3.084140	5.0	24	1	2.39
8	36.737815	3.094307	5.0	123	1	1.67
9	36.728669	3.077295	5.0	122	1	0.57

Figure 5.26: Top Recommendation Based on the Location-based recommendation system

## 5.7 Popularity Based Recommendation Model

Popularity based recommendation system uses the items that are in trend right now ,It is a very fundamental type of recommendation systems which gives recommendations based on the popularity of the items ,As mentioned in section 2.8.1. For example, the most popular item will be recommended first.It ranks items based on its popularity i.e. the rating count. If a item is highly rated then it is most likely to be ranked higher and hence will be recommended. As it is based on the items popularity,that is why it's can't be personalized and hence same set of items will be recommended for all the users.The implemented Popularity recommendation systems we selected 70% of the dataset as the training set and the remaining 30% as the testing test.

### The Recommendation Based on the Popularity Model

	userid	itemid	rating_count	Rank
<b>446</b>	265	452	12369	1.0
<b>367</b>	265	373	6687	2.0
<b>442</b>	265	448	4183	3.0
<b>293</b>	265	299	4087	4.0
<b>146</b>	265	149	3215	5.0
<b>355</b>	265	361	3076	6.0
<b>643</b>	265	652	3065	7.0
<b>632</b>	265	641	2816	8.0
<b>478</b>	265	484	2724	9.0
<b>257</b>	265	262	2177	10.0

Figure 5.27: Top Recommendation Based on the Popularity Model

This table represents the ten recommended items based on the popularity model for the userid 265. Our model depends on the Rank variable, which ranks the items based on the highest rating count, the rating count variable counts the number of users that have order from the itemid (the restaurant) . In this section, we will plot the evaluation metrics results. The evaluation metrics are **RMSE,MSE, and MAE ,Recall ,Precision.**

## Results

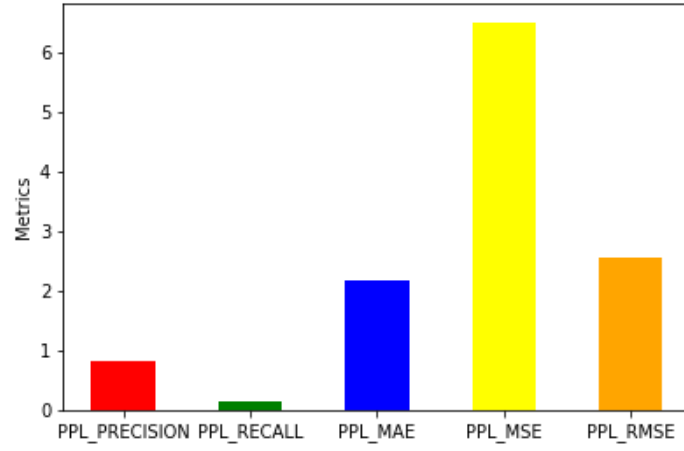


Figure 5.28: Evaluation Results of Recommendation Based on the Popularity Model

Model	RMSE	MSE	MAE	Recall	Precision
Popularity RSs	2.54954	6.50016	2.19050	0.16666	0.83334

Table 5.6: The Evaluation Results on the Popularity RSs

To evaluate the effectiveness of the popularity recommendation systems ,we test the performance of our model using MAE ,RMSE,MSE.

The evaluation results are shown in 5.6 that the precision and the recall of our model is about 83% and 16% respectively, RMSE(2.54954), MAE(2.19050),MSE(6.50016).





# Bibliography

- [1] User modeling via stereotypes. *Cognitive Science*, 328–353.
- [2] Approximation theory and methods M. J. D. (Michael James David) Powell, 1981.
- [3] Automatic text processing : the transformation, analysis, and retrieval of information by computer Salton, Gerard ,1989 .
- [4] Principles of forecasting: a handbook for researchers and practitioners JS Armstrong , 2001
- [5] Recommender systems : the textbook Charu C. Aggarwal.
- [6] Practical recommender systems Kim Falk.
- [7] Hands-On Recommendation Systems with Python : Start Building Powerful and Personalized, Recommendation Engines with Python. Rounak. Banik 2018.
- [8] Recommender Systems Handbook 3rd ed. 2022 Edition: Ricci, Francesco,Rokach, Lior, Shapira.
- [9] Building Recommender Systems with Machine Learning and AI: Help people discover new products and content with deep learning, neural networks, and machine learning recommendations. by Frank Kane .
- [10] Recommender Systems: An Introduction 1st Edition by Dietmar Jannach.
- [11] Statistical Methods for Recommender Systems 1st Edition by Deepak K. Agarwal.
- [12] Matrix and Tensor Factorization Techniques for Recommender Systems (SpringerBriefs in Computer Science) 1st ed. 2016 Edition .by Panagiotis Symeonidis, Andreas Zioupos.
- [13] Recommender Systems: Algorithms and Applications 1st Edition, Kindle Edition by P. Pavan Kumar (Editor), S. Vairachilai (Editor), Sirisha Potluri (Editor), Sachi Nandan Mohanty (Editor)
- [14] Number of netflix paying streaming subscribers worldwide from 3rd quarter 2011 to 1st quarter 2020.<https://www.statista.com/statistics/250934/quarterly-number-of-netflix-streaming-subscribers-worldwide>. Last accessed: 2020-12-09
- [15] Burke, R. Hybrid recommender systems: Survey and experiments. User Model. User-Adapt. Interact. 2002, 12, 331–370.
- [16] "Item-based collaborative filtering recommendation algorithms" , Published in the Proceedings of the 10th International Conference on World Wide Web, WWW 2001,Sarwar Badrul,Karypis George,Karypis George,Karypis George.
- [17] MEASURING AND TESTING DEPENDENCE BY CORRELATION OF DISTANCESPDF Gábor J Székely, Maria L Rizzo et al.35, 6, 2007

- [18] Dueck, J., Edelmann, D., Gneiting, T. & Richards, D. The affinity invariant distance correlation. *Bernoulli*, official journal of the Bernoulli Society for Mathematical Statistics and Probability 2014-11-04, 2305–2330.
- [19] Probabilistic Matrix Factorization, Ruslan Salakhutdinov, Andriy Mnih.
- [20] Sage: Recommender engine as a cloud service Sage: Recommender engine as a cloud, Ronen, Royi Koenigstein, Noam Ziklik, Elad Sitruk, Mikael Yaari, Ronen Haiby-Weiss, Neta
- [21] Eigentaste: A Constant Time Collaborative Filtering Algorithm Goldberg Ken, Roeder Theresa, Gupta Dhruv, Perkins Chris
- [22] Item-based collaborative filtering recommendation algorithms, Sarwar Badrul, Karypis George, Konstan Joseph, Riedl John.
- [23] Item-based top-N recommendation algorithms Deshpande Mukund, Karypis George
- [24] Big Data Recommender Systems: Algorithms, Architectures, Big Data, Security and Trust (2019), Osman Khalid Samee U. Khan Albert Y. Zomaya.
- [25] Collaborative filtering based on iterative principal component analysis Dohyun Kim-Bong Jin Yum *Expert Systems with Applications* (2005)
- [26] Matrix and Tensor Factorization Techniques for Recommender Systems Panagiotis Symeonidis Andreas Zioupos *Book* (2016)
- [27] Non-Personalized Recommender Systems and User-based Collaborative Recommender Systems, Anil Poriya, Neev Patel, Rekha Sharma, Tanvi Bhagat
- [28] A hybrid recommendation system considering visual information for predicting favorite restaurants, Chu Wei Ta, Tsai Ya Lun, November 2017
- [29] Item-based collaborative filtering recommendation algorithms, Sarwar Badru, Karypis George, Konstan Joseph, Riedl John.
- [30] Root mean square error (RMSE) or mean absolute error (MAE)? -Arguments against avoiding RMSE in the literature, Chai, T. Draxler, R. R.
- [31] Scientific paper recommendation: A survey, Bai Xiaomei, Wang Mengyang, Lee Ivan, Yang Zhuo, Kong Xiangjie, Xia Feng.
- [32] A recommender system based on collaborative filtering using ontology and dimensionality reduction techniques Mehrbakhsh Nilashi, Othman Ibrahim, Karamollah Bagherifard.
- [33] A Probabilistic View of Neighborhood-Based Recommendation Methods Jun Wang, Qiang Tang, *IEEE International Conference on Data Mining Workshops, ICDMW* (2016).
- [34] Application of dimensionality reduction in recommender system-a case study B Sarwar, George Karypis, J Riedl, *ACM WebKDD 2000 Web Mining for ECommerce Workshop* (2000).
- [35] A Survey of Recommendation Systems: Recommendation Models, Techniques, and Application Fields Hyeyoung Ko, Suyeon Lee, Anna Choi. *Electronics* (Switzerland) (2022)



- [36] Estimating probabilities in recommendation systems Mingxuan Sun, Guy Lebanon, Paul Kidwell *Journal of the Royal Statistical Society. Series C: Applied Statistics* (2012)
- [37] Setting goals and choosing metrics for recommender system evaluations Gunnar Schröder, Maik Thiele, Wolfgang Lehner *CEUR Workshop Proceedings* (2011)
- [38] Implicit feedback techniques on recommender systems applied to electronic books Edward Rolando Núñez-Valdéz Juan Manuel Cueva Lovelle, Carlos Enrique Montenegro Marín *Computers in Human Behavior* (2012)
- [39] Restaurant Recommendation System Based on User Ratings with Collaborative Filtering Achmad Arif, Munaji Andi Wahju, Rahardjo Emanuel *IOP Conference Series: Materials Science and Engineering* (2021)
- [40] Recommending based on implicit feedback Dietmar Jannach Lukas Lerche Markus Zanker *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* (2018)
- [41] A hybrid online-product recommendation system: Combining implicit rating-based collaborative filtering and sequential pattern analysis Keunho Choi Donghee Yoo[...] Yongmoo Suh *Electronic Commerce Research and Applications* (2012)
- [42] A recommender system based on implicit feedback for selective dissemination of ebooks Edward Rolando Núñez-Valdez David Quintana[...] Enrique Herrera-Viedma *Information Sciences* (2018)
- [43] A new similarity measure for collaborative filtering to alleviate the new user cold-starting problem Hyung Jun Ahn *Information Sciences* (2008)
- [44] A new similarity measure for collaborative filtering based recommender systems Achraf Gazdar Lotfi Hidri *Knowledge-Based Systems* (2020)
- [45] Restaurant Recommender System Using User-Based Collaborative Filtering Approach: A Case Study at Bandung Raya Region Alif Azhar Fakhri Z. K.A. Baizal Erwin Budi Setiawan *Journal of Physics: Conference Series* (2019)
- [46] A sub-one quasi-norm-based similarity measure for collaborative filtering in recommender systems Shan Jiang Shu Chong Fang[...] John E. Lavery *Information Sciences* (2019)
- [47] A similarity measure based on Kullback–Leibler divergence for collaborative filtering in sparse data Jiangzhou Deng Yong Wang[...] Younghee Park *Journal of Information Science* (2019)
- [48] An item–item collaborative filtering recommender system using trust and genre to address the cold-start problem Mahamudul Hasan Falguni Roy *Big Data and Cognitive Computing* (2019)
- [49] The distance correlation t -test of independence in high dimension Gábor J. Székely-Maria L. Rizzo *Journal of Multivariate Analysis* (2013)
- [50] Distance covariance for discretized stochastic processes Herold Dehling Muneya Matsui[...] Laleh Tafakori Bernoulli (2020)
- [51] A Basic Treatment of the Distance Covariance Dominic Edelmann Tobias Terzer Donald Richards *Sankhya B* (2021)

- 
- [52] An Updated Literature Review of Distance Correlation and Its Applications to Time Series Dominic EdelmannKonstantinos FokianosMaria Pitsillou International Statistical Review (2019)
- [53] Similarity measures for Collaborative Filtering-based Recommender Systems: Review and experimental comparison Fethi Fkih Journal of King Saud University - Computer and Information Sciences (2021)