

01/06

Research Methodology \Rightarrow Collecting literature
for literature review

Interest :- Content Recommendation

Movie Recommendation \Rightarrow Start collecting
literature

(Need to check what models are used nowadays
for Movie Recommendation).

* I will also review about Netflix , as it is
currently famous for movie database and
recommendation. (Need one section for that
in literature review)

* Check for database available :- which one will be
feasible !

* Thinking to recommend movies based on

Plan this later... } $\begin{cases} \hookrightarrow \text{genre} \\ \hookrightarrow \text{cast} \\ \hookrightarrow \text{year} \\ ?? \end{cases}$

02/06

⇒ Finalized Movie Recommendation System.

Various literatures P:- (Need to select best few)
for review

1. Movie GRIN: A Movie Recommender System

- ✓ Uses hybrid recommendation approach
↳ with the help of machine learning
and cluster analysis

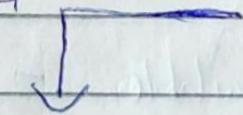
One of the approach used to hybrid
P :- collaborative approach. user preference
↓
(i/p)

Machine learning ⇒ SVM
model



Prediction (o/p)
of movie title

further
i/p



K-Means Clustering



Cluster of similar movies.

2. MovieMender : Movie Recommendation System

→ uses hybrid approach → collaborative flat filtering + content-based

★ Takes user-ratings matrix as database using web crawler.

↳ Fill it using content-based predictor

⇒ Recommendations are provided by collaborative filtering

3. An Improved approach for Movie Recommendation System :- Hybrid Approach

X ✓ (Again!!)

Content - based + collaborative filtering

★ SVM as classifier & genetic algorithm for optimization

↳ Classified movies from SVM are clustered using k-Means for recommendation purpose.

05/06

My Project idea:- (as of now)

- ★ Movie Recommendation based on some criterion (genre, title, year)
- ★ Hybrid approach (for sure)!!
- ★ Need to check which classifier/method can be used for approach.

* If time persists ?? will try for UI.

Review Research continued :- Definitely good paper

4. Netflix Recommender System :- for review

- ✓ ↳ Good paper about different recommendation algorithms used by Netflix.
- ↳ Can inspire for 2-3 ideas (for sure)
- ↳ Gives different perspective of how any recommendation can be built i.e., what users want as good recommendations?

5. Personalized Movie Recommendation based on Deep Learning.

↳ collaborative filtering based approach

↳ uses deep learning models

↳ Seq2Seq model based on

Complex Research

LSTM RNN

Paper ; but good one

as it uses deep neural net approach

6. Movie Recommendation System using Sentiment

Analysis from Microblogging Data

↳ hybrid approach

Tweets (User reviews) \rightarrow Database } unique set

\Rightarrow Sentiment analysis by VADER employed to provide recommendations.

7. Movie REC \rightarrow hybrid approach

X

II

* using cumulative wt. of different attributes

* Uses k-means algorithm.

10/06

Literature Review → Selected Papers

1. Movie CFN! → hybrid approach
(use traditional model)
2. Netflix Recommender System → gives idea about different different recommendation techniques
3. Movie Recommendation using microblogging data → different kind of dataset (tweets)
↳ leads to unique approach
4. Personalized Movie Recommendation based on deep learning → Neural Net.
(Deep Learning approach).

Review Paper Structure :- 1. Abstract : (Summary)

2. Introduction : About Recommendation System :— * general context
* types of approach
↳ Further, a crisp intro. about movie recommender recommendations.

3. Literature Review → structure in Order of each paper.

but start with 'Netflix Recommender' paper
→ It gives a good base.

Then all 3 papers.

4. Conclusion : (whole paper in 1 para).

↓ Final Review Paper Structure.

11/06 → Change of plans for review
paper → Not including deep
learning paper (too complex
and taking space).

* Cannot remove anyone of other 2 papers
as they are uniquely setting the tone
of review.

⇒ Recommendation techniques :- Recommender
~~Content-Based Filtering~~ Systems
Content-Based :- ~~Filtering~~ → Hand book,
User-Ratings matrix :- ↴ 2nd ed.
matrix of m users and n items

Each row \Rightarrow users ; Each column \Rightarrow items
(i, j) cell \Rightarrow rating given by user i to item j .

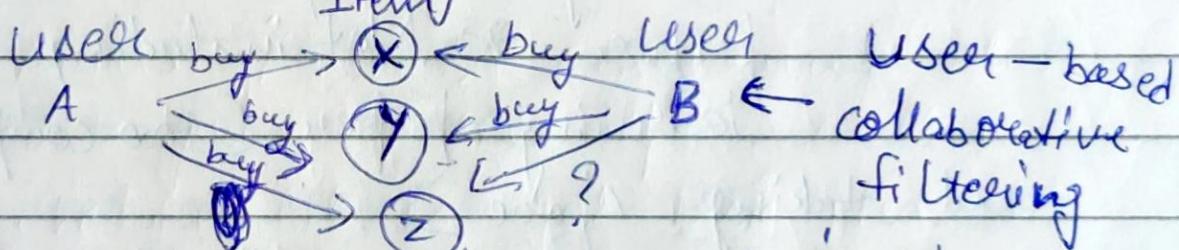
* This matrix is usually sparse as users have not rated most movies.

\Rightarrow To predict those empty values \Rightarrow

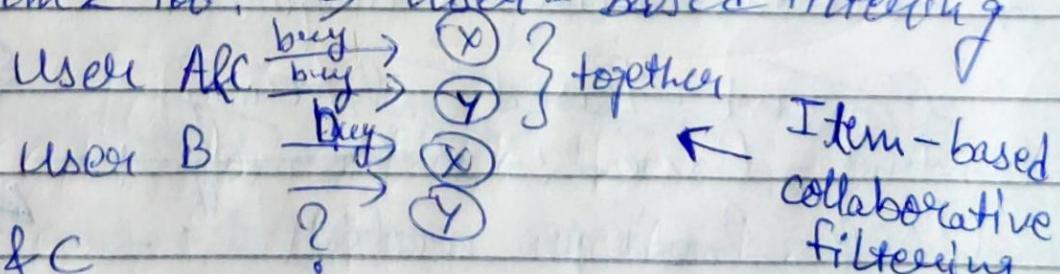
To recommend movie

\hookrightarrow type of content-based

\Rightarrow Collaborative filtering



How likely will be that user B is going to buy item Z, given Z is bought by A too ? Since both users have bought similar items X, Y in past, it is highly likely that user B will buy item Z too. \Rightarrow User-based filtering



If user A & B buy item X & Y together, then there are high chance that both items will be

bought by B ! \Rightarrow item-based filtering

15/06

\Rightarrow Idea checked before:— create a movie recommendation system based on genre, year, cast.



* This demands similarity between two movies in order to get us good recommendation
 \hookrightarrow Metadata-based Recommender.

* Need to check techniques which uses similarity measure:— vector representation

Count Vectorizer

represents count of each word in vocabulary according to its frequency, in form of one-hot vector.



If word present in doc, then +1, else 0.

TF-IDF vectorizer

assigns weight along with frequency to each word.

TF \Rightarrow term frequency in a document

IDF \Rightarrow No. of documents that contain, term

Weight \Rightarrow $\langle 0, 1 \rangle$

Most of the related work used this.

⇒ TF-IDF vectorizer uses cosine similarity for checking similarity b/w 2 texts/values.

$$\text{cosine} = \frac{\mathbf{a} \cdot \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|}$$



$$\|\mathbf{a}\| \|\mathbf{b}\|$$

Higher the cosine similarity, more similar two ~~text~~ text will be.

→ This will make pretty good model!

* Advantage :— * I understand this concept well.

* Can be helpful to get basic idea of content-based (metadata) movie recommendation.

* Fulfill my idea for recommending movies based on genre, ~~title~~, year, cast provided by user.

Need to check the limitation of this method.

* Will help in NLP too (currently ongoing feature)

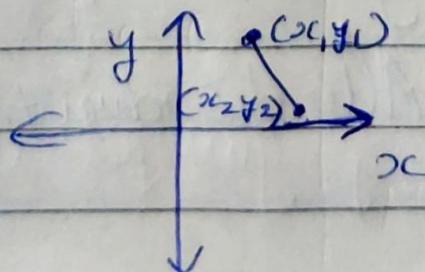
20/05

Various Collaborative - filtering techniques :-

- ★ Similarity measures
- ★ Supervised learning
- ★ Clustering
- * Remember user - ratings matrix.
 ↳ Need to fill the empty cells of matrix.
- * Empty cells i.e. rating can be predicted by finding users who has similar taste
 ↳ We can calculate check which user has similar taste by measuring distance
 - 2 measures
 - popularity used
 - Euclidean distance &
 - Pearson correlation.

⇒ Euclidean distance :- dist. b/w 2 datapts. plotted in n-dimensional space.

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$



Plotting of
n-dimensional
space

- ⇒ Pearson correlation :- ranges b/w -1 to 1
- 1 ⇒ negative correlation
(totally opposite)
 - 0 ⇒ no relation
 - +1 ⇒ positive correlation
(very similar)

★ Better than euclidean as it checks for both similarity & dissimilarity.

~~30/06~~ ⇒ Research Proposal ⇒ not perfect plan selected.

★ Need to give any one above learned method in detail.

★ Will be discussed with Dr. T if any further changes required.

★ Structure of proposal ⇒ Abstract (summary)

⇒ Intro. & Background ⇒ Detailed intro. about movie recommendation system & various approaches.

↳ Research question (most imp. part) ⇒ All disadvantages of basic model.

- ⇒ Research Methodology ⇒ Methods used; approach used.
- ⇒ Schedule (whole schedule of research with Gantt chart)
- ⇒ Summary (Conclusion of Research Proposal)

Research Proposal ⇒ My Hybrid Approach

Content-based

Collaborative-

Disadvantage :- cold-start problem
 ↳ recommendations cannot be made for new users.

Research methodology
 methods → cosine similarity for checking proposed similarity b/w users & movies

↳ TF-IDF vectorizer to give weightage to different metadata of movie

* My main approach still stands with the idea that user will be able to get movie recommendations based on genre/year/cast (any detail i/p user gives).

10/07

Searching for Dataset

- ↳ checking other literature
- ↳ Kaggle
- ↳ githubs

⇒ Checked various literatures related to movie recommendation systems :-

★ Most of them are using 2 datasets (~~open-sourced~~) Not sure!!

1. MovieLens dataset (grouplens.org/datasets/movielens)

2. TMDb dataset (themoviedb.org)

⇒ I will be using MovieLens dataset as it has multiple variants available in various sizes. Also, tmdb dataset need to have api.

★ MovieLens dataset can be directly used in colab using just wget & unzip.

11/07 MovieLens dataset (Different variants)

* I will be using only 4 variants, either because ~~either~~ it's used in many literatures and results can be compared ~~or~~ ~~it's~~

Dataset
can
be
changed!

1. ml-latest-small.zip

↳ smallest dataset

↳ does not have other metadata such as tags.

2. ml-latest.zip

↳ largest dataset

↳ let's see if it fits first!

3. ml-20m.zip

↳ not too large; not too small

↳ 20M movie ratings

↳ popularly used in most research works.

↳ stable dataset

4. ml-25m.zip

↳ too large dataset

↳ let's see if it fits first!

↳ stable dataset

15/07 Going through datasets for check learning the data!

⇒ Exploratory analysis ⇒ will help for data pre-processing!

Important files -

Movies - 25M

→ features.csv

→ movies.csv

→ tags.csv

Movies - latest - small

→ features.csv

→ movies.csv

★ User ~~details~~ ^{ratings} :- Only unique userId provided ratings for each movie

★ Movies.csv ⇒ movieId, genre, title

I guess these files will be enough for research project implementation.

⇒ Selected both dataset as while ~~testing~~ ^{Checking} similarity between movies / users, large dataset will give us hard time !!

16/07 Final Report \Rightarrow Introduction
 \downarrow

I will write a brief introduction about Recommendation System.

\hookrightarrow different types of recommendation systems

\Rightarrow General overview of Movie recommendation system \Rightarrow what input/output is expected?

Related Work: Thinking to write regarding Netflix \Rightarrow That will be helpful as I have thought of research project structure.

It has resemblance to Netflix Recommendation algorithms: (resemblance as in algorithms used for recommendations).

+ abstract

This¹ will constitute starting of Final Report !!

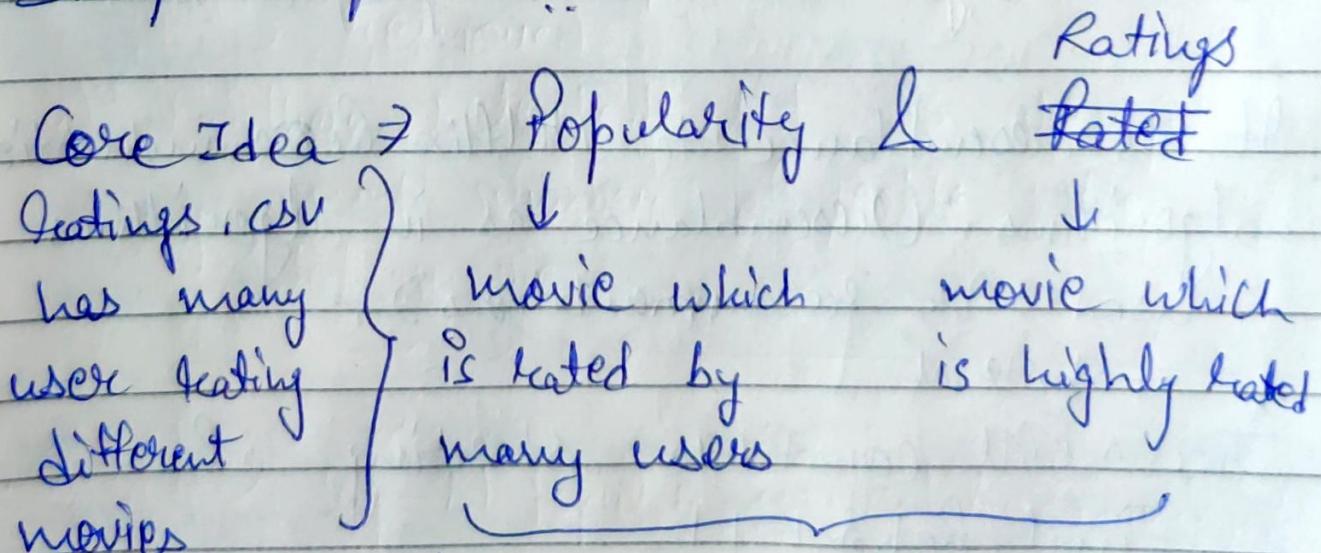
20/07 Started coding part for my research project!

General Idea :- To provide any new user a set of movie recommendations

- ★ User need to provide simple inputs
 - ↳ number of recommendations, n
 - ↳ year-range; from-year, to-year
 - ↳ genre

- ★ These inputs can be easily given by any user who wants to check out movies.

Simple inputs !!



If both Popularity & Ratings are weighted accordingly, then ...

a good movie recommendations can be provided.

★ Since this will utilize ^{release} year and genre of movies, this is our ~~content~~ = content-based recommendation system.

★ Data pre-processing :-

↳ ratings.csv and movies.csv has been used.

↳ timestamp from ratings.csv has been removed

↳ Since release year is provided with movie title, it has been extracted separately and placed in new column \rightarrow year

↳ ~~year~~ year } $\xrightarrow{\text{int}}$
Typecasting genre } $\xrightarrow{\text{datatype}}$ { str
 $\xrightarrow{\text{str}}$

★ Since there is year-range, few if-else has been placed so user does not give wrong input \Rightarrow genre must be from given list

★ from-year < to-year

21/07 Top - N Popular Movies

- ⇒ Data preprocessing completed!
- ⇒ input bounds provided so user does not feed wrong input
- ⇒ Core idea implementation:-

- ★ Count of users who rated each movie has been extracted by checking ratings provided to each movie grouped by movieId
 - ★ Movies are filtered ⇒ movies which has less ratings than avg rating and less count of ~~people~~ user who rated it ⇒ removed
 - ★ Further, selected list is arranged by most number of people who rated movie.
-
- ↳ Top-N Popular movies recommender system completed !!

22/07 Top-N Rated Movies

- Data pre processing completed (same)
- Input bounds provided so user does not feed wrong input (same)
- Core idea implementation :-

★ Count of user who rated each movie will be calculated by grouping using movieId.

★ Average rating of movie will be calculated once the total count of users are calculated.

⇒ Movies will be filtered ⇒ same as $\not\in$ Top-N Popular; movies having less than avg. rating and avg. count will be filtered out.

⇒ Finally, selected movies will be arranged as per highest rating provided by users.

↳ Top-N Rated Movies

Recommender System completed!!

25/07 ⇒ Final Report ⇒ Methodology

↳ Hybrid approach

↳ Since, using both content-based
and collaborative filtering approaches

⇒ I will mention both approaches

in perspective of movie recommender
systems.

Content - Based



uses metadata
of movie (year,
genre) for
recommendations

user need to
have basic idea
about movies.

⇒ Simple I/p ⇒ year, genre

collaborative filtering



user ratings of
input movie to
check similarity with
other movies;

(other users who rated
same input movie –
what they watched) as
recommendations.

- ⇒ Also mention both similarity measures
- ↳ using Pearson Correlation Coefficient
- ↳ using Cosine Similarity.

This constitutes methodology section
of Final Report!

Note:- Dataset description need to be included
in Implementation Section

30/07 Collaborative - Filtering approach

↳ user need to provide movie

↓

movie ratings will be extracted
which will give us users

↓

Other movies rated by these
users will help us to get
selected movies for recommendations

Calculated using similarity
score

31/07

Top-N Similar Movies by Cosine Similarity

⇒ Same data preprocessing
as content based movie
recommender.

★ User-Rating matrix
is calculated which
gives details of
each movie title
rated by each user

★ Cosine similarity
is measured between
each movie ratings
provided by each
user.

★ Total count of ~~few~~
users who rated movie
is also considered, where
only top 0.01% is selected

Top-N Similar Movies
by Pearson Correlation

⇒ Same data preprocessing
as content based movie
recommender.

★ Here, also user-
rating matrix will
be created.

★ coldWith function
used with method =
"Pearson" to check
similarity between
movie ratings.

★ Here also, total
count of users who
rated movie (only top
0.01%) is selected

01/08 Implementation Section of Final Report

→ Dataset description :— Details about dataset used ; little bit description about main/important .csv files.

→ Data Preprocessing :— Preprocessing of data before implementing functions

for content-based and collaborative filtering recommendations.

→ Implementation of content-based movie recommender systems

 → Top-N Popular Movie Recomm. System

 → Top-N Rated Movie Recomm. System

→ Implementation of collaborative filtering movie recommender systems

 → Top-N Similar (by Pearson Correlation) Movie Recomm. System

 → Top-N Similar (by Cosine Similarity) Movie Recomm. System

05/08

Evaluation for Recommender Systems

⇒ Considering user-surveys as movies recommended will be evaluated effectively by humans.

★ 10 Random users (family & friends)

⇒ Testing of Content-Based Movie Recommender System ⇒ Number of recommendations

are kept same fixed ⇒ 10

Users	From year	To year	genre	Popular	Rated
1	2011	2015	Action	10/10	9/10
2	1991	2000	Romance	7/10	7/10
3	2005	2015	Mystery	10/10	8/10
4	2001	2015	Fantasy	9/10	6/10
5	2011	2019	Sci-Fi	10/10	10/10
6	2005	2015	Animation	10/10	5/10
7	1994	2019	War	8/10	8/10
8	2011	2019	Adventure	10/10	9/10
9	2011	2019	Horror	8/10	6/10
10	2005	2010	Thriller	9/10	8/10

* All year-range & genre is fixed by each user to evaluate ^{Top-N} popular & rated movies

* Observations :- * Top-N popular Movies

Recommender was effective! 7 users out of 10 ~~that liked~~ liked atleast 9 movies out of 10 recommendations.

* Top-N Rated was not effective compared to Top-N Popular \Rightarrow Only 3 users out of 10 liked atleast 9 movies out of 10 recommendations.

06/08

\Rightarrow Testing of Collaborative filtering Recommender movie

User	From	To	movie	Cosine	Pearson
1	2005	2015	Inception	9/10	10/10
2	2000	2018	Ratatouille	9/10	7/10
3	1994	2004	Ocean's 11	10/10	10/10
4	2000	2018	Gladiator	9/10	8/10
5	2011	2018	Conjuring	9/10	10/10
6	2012	2018	Deadpool	9/10	9/10
7	2008	2018	Avengers	9/10	8/10
8	1994	2000	Fight Club	10/10	9/10
9	1991	2000	Pulp Fiction	9/10	9/10
10	2001	2018	Lakeland	9/10	8/10

- ★ Both similar movies recommended did good job.
- ★ However, movies recommended using cosine similarity of ~~will~~ be was very effective.
- ★ Also, the user survey helped to avoid cold-start problem, since users found at least 2 new movies from recommendations.
- ★ Few users actually watched them and responded that they liked the ~~new user~~ movies recommended by recommender system.

⇒ Results & Discussion in Final Report

↳ Need to plot ~~above~~ previous mentioned results.

- ★ Barplot will help to visualize the effectiveness of recommender system.
- ★ Result will be in form of how many movies did user liked from the recommendations!

10/08 ➔

Few ^{final} sections and touch-up

of reports left ➔ Abstract, Conclusion,

→ GitHub repository is created
and code notebooks will be uploaded

⇒ Presentation Structure :-

1 min. Introduction :- About Movies Recommender System

2 mins. Methods :- All 4 notebooks explained

1 min. Results :- Results of user-survey presented

1 min. Discussion :- Comparison & Discussion from above results.

~5 mins

4 Deliverables need to be submitted
by 15th August ➔ Report, Repository,
Presentation and Research Journal !!

15/08

D - DAY

Conclusions :- Good Research work

Basic concepts of Movie Recommender system is studied and researched well.

Next ?: {Thinking to increase complexity by bringing deep learning models.

- Learn about 'Surprise' library
- Will try playing with GUI!

★ Thanks to Dr. T for effectively distributing Research work. Able to do research from scratch!

Final Deliverables

- ⇒ Presentation → Completed
- ⇒ Final report → Completed
- ⇒ Journal → Completed
- ⇒ Others Repository → Completed

X