

Recommender Systems

SIR ZAIN

Project 10%.

Finals 50%.

S1 , S2 30%.

Assigments 5%.

Quiz 5%.

have to use any 2 techniques
show rating predictions from both

such koy purchase X
recommended purchase ✓



Information Overload

↳ too much info



million flavours of icecream

HOW TO CHOOSE

TOP 5



TOO MANY OPTIONS

MAKE IT LESSER

EASIER TO CHOOSE

Implicit

- ↳ MORE USEFUL
- ↳ many people just scroll, don't like ↳
- ↳ is collective ??
- ↳ likes is collective
- ↳ has no dislike

PRO
↳ harder to collect
CON
↳ not everyone rates

↳ -ve comments are implicit ratings

↳ useful



Genre

↳ normally rated high by active user

1. Similarity b/w items
2. Similarity b/w users



Explicit IMDb

↳ on a scale
↳ 1 star, i don't like it
↳ 5 star, i love it

PRO
↳ easy to collect
CON
↳ not everyone rates

↳ assumes no rating is positive

↳ not useful

Recommender

- ↳ A suggestion to help user in decision making

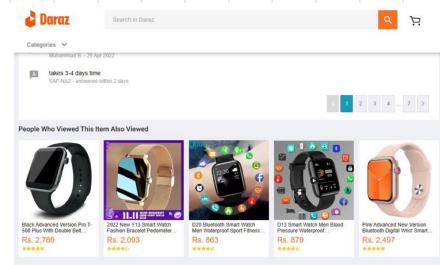
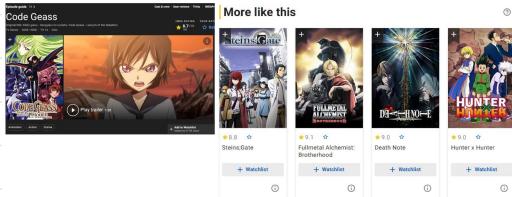
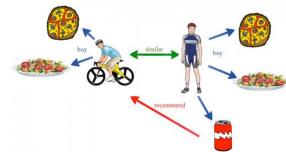


Adv of Recommendation

- ↳ To increase sales
- ↳ To improve user experience
- ↳ To maximise productivity

Example

If **cyclist A** always orders **pizza** and **salad** at our food joint, and **cyclist B** eats **pizza** and **salad** as well. Then recommend **cola** to **cyclist A** if **cyclist B** had been ordering it lately



Recommender system (RS)

- ↳ helps match users with items
- ↳ It eases info overload
- ↳ helps improve user experience
- ↳ helps increase sales
- ↳ suggest novel items

How does RS work?

- ↳ It bases recommendation on
- ↳ Past behavior pattern of that user
- ↳ Similarity with other users
- ↳ Item similarity
- ↳ context

$$f: U \times I \rightarrow R$$

↓ ↓ ↓ → ordered list
 users items recommended items

- For each user u , what we want to do is to choose the item i that maximizes f

$$S_u = \underset{i}{\operatorname{argmax}} f(u, i)$$

PARADIGMS OF RECOMMENDER SYSTEM

1. Collaborative based

a. Collaborative Filtering (CF)

1. User based

↳ find users similar to me
↳ then recommend me what they like

2. Item based

↳ recommend me an item
similar to the ones I normally like



2. Content based approach

↳ doesn't req. other users data

↳ only users own history patterns

↳ minimal collaboration



3. Demographic approach

↳ e.g. Spotify



4. Social Tagging → aka connection based

↳ connection based

mutual friends suggested

on social media
↳ insta
↳ facebook

friend liked Pages suggested

similarity based X
trust based X → moi si

based on what people in
a user's connection like

e.g. Fb group recommendation

5. Trust based approach

↳ normally in news recommendation

e.g. Stack Overflow
social media posts

↳ reliability is important

mujey achha laga hai
TUM BHI LELO

6. Hybrid approach

↳ combination of 2 or > of

the above mentioned approaches

* collaborative and similarity
are SAME THING

Collaborative Filtering (CF)

- ↳ uses similarities to recommend items to active user
- ↳ USES similarity, ratings
- ↳ It has 2 approaches

1. User based CF

- ↳ find users similar to me
- ↳ then recommend me what they like
 - ↳ if most similar user rated it high then it should be recommended
- ↳ users are vectors
- ↳ similarity b/w users

PROS

- ↳ more diverse recommendations
- ↳ better choice if no of items > no of users

CONS

- ↳ not stable
 - e.g. likeness changes over time - I rated something high but I don't like it as much anymore
 - after a year user preference might change
- ↳ can't provide indepth analysis of a user

↳ needs ratings!
↳ no features
only raw ratings

2. Item based CF

- ↳ recommended me an item similar to the ones I normally like
- ↳ items are vectors
- ↳ similarity b/w items

PROS

- ↳ provides more accurate recommendations
- ↳ better choice if no of users > no of items
- ↳ is more stable
- ↳ can provide indepth analysis on users

CONS

- ↳ is prone to shifting attacks
 - ↳ major attack
 - ↳ sub mil key more to gain rating artificially
- ↳ provides much less diversity than user based CF
 - ↳ campaign to pay to give high or low rating

similarity is measured using

VECTOR SPACE

- ↳ interactive Table / Matrix as vectors

TERMINOLOGIES OF CF

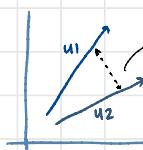
- ↳ Active User
 - the user that's given recommendation
- ↳ Active item
 - deciding whether it should be recommended or not
- ↳ Interactive Table / Matrix
- ↳ Metric
 - to measure similarity between users
- ↳ Method
 - for subset consisting of closest neighbour

INPUT

	I ₁	I ₂	I ₃	I ₄
U ₁	4	3	2	
U ₂	3	2	5	
U ₃	1	2	3	
U ₄	2	4		

Rows (all) → user
cols (all) → item

VECTOR SPACE MODEL



Distance measure

- ↳ Euclidean
- ↳ Cosine
- ↳ Manhattan
- ↳ Correlation
- ↳ Raw distance

1.1 User based CF APPROACH

Q) U_1 is active user

Predict for I_5 of U_1 ,

STEPS

1. Choose a K

↳ taking $K=2$

2. Measure Similarity

1. Euclidean - U_2, U_3

2. Cosine - U_2, U_3

3. Pearson - U_2, U_4 → uses mean

4. Adj Cosine Similarity → uses mean

no of users similar ↗ hyperparameter
we choose

Active user	I ₁ I ₂ I ₃ I ₄ I ₅				
	U ₁	5	3	4	4
U ₂	3	1	2	3	3
U ₃	4	3	4	3	5
U ₄	3	3	1	5	4
U ₅	1	5	5	2	1

dim = 5
only 4 compared as
we don't count active item

no of users = 5
only 4 compared as
we don't count active user

3. Choose Prediction Method

1. Average

2. Weighted Average Item 5 of $U_1 = 3$

3. Mean Adjusted Weighted Avg → uses Mean

4. Significance Weighting

STEP 1. CHOOSE K

↳ Take $K=2$

STEP 2. MEASURING USER SIMILARITY

$$\text{Similarity} = 1 - \text{Distance} \longrightarrow \text{anti-distance}$$

1. Find distance b/w AC and all other users

2. Find k nearest neighbour of AC

↳ Some similarity measures

1. Euclidean

good to start from
not the best

↳ Lesser distance = most similar

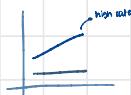
$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

e.g. $U_1 = (3, 10)$ $U_2 = (7, 5)$
 $d = \sqrt{(7-3)^2 + (5-10)^2}$
 $= 5$

CONS

↳ doesn't consider patterns

e.g. Pattern
 U_2 rates low
 U_3 rates high \rightarrow doesn't take this into consideration



↳ doesn't consider slope, only distance

	I ₁	I ₂	I ₃	I ₄	I ₅	$d_{U_1 \text{ with } U_2}$	Similarity with U_2
Active User	5	3	4	4	?	0	
✓ U_2	3	1	2	3	3	3.60	-2.6 $\rightarrow 1-3.6$
✓ U_3	4	3	4	3	5	1.41	-0.41 $\rightarrow 1-1.41$
U_4	3	3	1	5	4	3.74	-2.74 $\rightarrow 1-3.74$
U_5	1	5	5	2	1	5	-4 $\rightarrow 1-5$

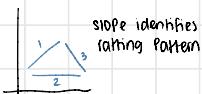
for K=2, U_2 and U_3

2. COSINE SIMILARITY

sparse matrix as
not everyone rates

↳ considers slope and distance

- ↳ low rater
- ↳ high rater
- ↳ indifferent



↳ treats each item separately

$$\cos\theta = \frac{A \cdot B}{\|A\| \cdot \|B\|}$$

	I ₁	I ₂	I ₃	I ₄	I ₅	$\cos\theta$ similarity with U_2
Active User	5	3	4	4	?	
✓ U_2	3	1	2	3	3	0.97
✓ U_3	4	3	4	3	5	0.99
U_4	3	3	1	5	4	0.89
U_5	1	5	5	2	1	0.79

for K=2, U_2 and U_3

CON

↳ can't measure correlation
b/w 2 users

User 2 likes something else
User 3 likes something else

3. Pearson Correlation Coefficient (r)

↳ measures magnitude and orientation b/w data points

↳ strength and r/s is given by -1 to 1

• -ve → strong -ve correlation → if U₁ likes this U₂ hates it → inverse

• 0 → no correlation

• +ve → strong +ve correlation → if U₁ likes this U₂ likes it just as much

↳ treats each item as a whole collection

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \cdot \sqrt{\sum (y_i - \bar{y})^2}}$$

CON

↳ if uniform distribution

the line will be flat

there might be a user that rated everyone 3

$\sqrt{(3-3)^2} = 0 \text{ D.O.D.}$ → UNDEFINED

↓ SOLUTION

↳ add another item and give it another value → kinda like something → but not for negative

↓ How many times will we do this

↓ too many calculations!

CF

↳ high resource

↳ high processing power

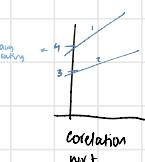
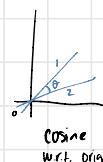
↳ high storage

↳ will work for < 100 users

↳ uses magnitude

↳ what is my general rating pattern →

↳ which is most similar user



	I ₁	I ₂	I ₃	I ₄	I ₅	Mean	r
Active User	U ₁	5	3	4	4	?	4
✓	U ₂	3	1	2	3	3	2.4
	U ₃	4	3	4	3	5	3.8
✓	U ₄	3	3	1	5	4	3.2
	U ₅	1	5	5	2	1	2.8

$$\frac{(5-4)(3-2.4) + (3-4)(1-2.4)}{\sqrt{(5-4)^2 + (3-4)^2} \cdot \sqrt{(4-4)^2 + (1-2.4)^2}} = \dots$$

0.85

$$\frac{1+5+5+2+1}{5} = 3.2$$

0.70

-0.76

for k=2, U₂ and U₄

* Exam advice

↳ make subtraction column with mean
↳ many people make this mistake

Pearson Correlation Coefficient: Issues

- Underlying assumption is that users dislike what they rated below average
- This is not true in practice (we rate only what we liked or highly disliked)
- The correlation flattens in case of uniformly distributed ratings

$$sim(a, b) = \frac{\sum_{p \in P} (r_{a,p} - \bar{r}_a)(r_{b,p} - \bar{r}_b)}{\sqrt{\sum_{p \in P} (r_{a,p} - \bar{r}_a)^2} \cdot \sqrt{\sum_{p \in P} (r_{b,p} - \bar{r}_b)^2} + \epsilon}$$

!!! Will be zero in case of uniform rating !!!

*interaction table must never change

STEP 3. Prediction Method

* example below are based on Using Pearson correlation as a similarity measure

1. Average

	I_1	I_2	I_3	I_4	I_5	
<u>Active User</u>	U_1	5	3	4	4	U
	U_2	3	1	2	3	3
	U_4	3	3	1	5	4

switch up

$\frac{3+4}{2} = 3.5$

CON

Outliers effect

$$\underline{1} \div 3.5$$

E1 ①) \cup_1 is $A \cup$, find $1s$

	Item 1	Item 2	Item 3	Item 4	Item 5	Mean	Pearson Correlation Similarity with User 1
User 1	5	3	4	4	?	4	1
User 2	3	1	2	3	3	2.4	= 0.85
User 3	4	3	4	3	5	3.8	= 0
User 4	3	3	1	5	4	3.2	= 0.70
User 5	1	5	5	2	1	2.8	= -0.76

K=2

2. Weighted Average

L multiply with AI

$$R_U = \left(\sum_{u=1}^n R_u \right) / n$$

e1
27

	I ₁	I ₂	I ₃	I ₄	I ₅	r
Active User	U ₁	5	3	4	4	3
	U ₂	3	1	2	3	3
	U ₄	3	3	1	5	4
P _{I₅}	$\frac{5+3}{5+3+1+5} = \frac{8}{14} = \frac{4}{7}$		$\frac{3+1}{3+1+2+3} = \frac{4}{10} = \frac{2}{5}$		$\frac{4+3}{4+3+3+4} = \frac{7}{14} = \frac{1}{2}$	
	$3 \times 0.85 + 4 \times 0.7$		$10.85 + 10.7$		$= 3.45$	
	Weight of I ₅		Weight of I ₂		Weight of I ₄	
	Weight of I ₁		Weight of I ₃		Weight of I ₅	

RIS: 3 → Should we recommend
or not?

Depends on cutoff
YOU CHOOSE \uparrow
taking 3 as cut off
hyper parameter based on domain knowledge
business decision

$$R_{31} = \frac{(7 \cdot 0.894) + (6 \cdot 0.939)}{|0.894| + |0.939|} \approx 6.49$$

$$R_{36} = \frac{(4 \cdot 0.894) + (4 \cdot 0.939)}{|0.894| + |0.939|} = 4$$

Conclusion
↳ incorrect assumption based
on correlation

E2 Q) U_3 is AU , find I_1 and I_2

	Item 1	Item 2	Item 3	Item 4	Item 5	Item 6	Mean	Pearson Correlation Similarity with User 3
User 1	7	6	7	4	5	4	5.5	0.894
User 2	6	7	?	4	3	4	4.8	0.939
User 3	?	3	3	1	1	?	2	1
User 4	1	2	2	3	3	4	2.5	-1
User 5	1	?	1	2	3	3	2	-0.817

K₂2

3. Mean Adjusted Ratings

↳ independent of whether user is generous or not

↳ mean centered Prediction function

↳ removes bias

$$R_U = \overline{r}_a + \frac{\sum_{b \in N} sim(a, b) * (\overline{r}_{b,p} - \overline{r}_b)}{\sum_{b \in N} sim(a, b)}$$

E2
二

	I ₁	I ₂	I ₃	I ₄	I ₅	I ₆	Mean	r
U ₁	7	6	7	4	5	4	5.5	0.894
U ₂	6	7	?	4	3	4	4.8	0.939
• U ₃	?	3	3	1	1	?	2	

$$R_{31} = 2 + \frac{(1.5 \cdot 0.894) + (1.2 + 0.939)}{0.894 + 0.939} \approx 3.35 = 3$$

mean of active user

$$R_{36} = 2 + \frac{(-1.5 * 0.894) + (-0.8 * 0.939)}{|0.894| + |0.939|} \approx 0.86 \Rightarrow 1$$

SOLUTION

Consider the given interaction matrix for Crunchyroll website users.

					Mem
User 1	9	7	2	8	6.5
User 2	?	7	5	7	6.33
User 3	10	4	9	5	7
User 4	8	5	8	6	6.75

- a) Find the Cosine similarity between all users. $\begin{matrix} 1 & 1 & 1 & 1 & 1 & 1 \end{matrix}$

a) Cosine Similarity

$$U_1, U_2 = 0.958$$

$$U_2, U_3 = 0.881$$

$$U_1, U_3 = 0.839$$

$$U_2, U_4 = 0.943$$

$$U_1, U_4 = 0.833$$

- b) Predict the missing rating $R(U_2, I_1)$ using mean-centered prediction function for $k=2$.

b) $k=2$, closest are U_1, U_4

$$R_{U_2, I_1} = 6.33 + \frac{(9-6.5)0.958 + (8-6.75)(0.943)}{10.9581 + 10.943} = 8.21 \approx 8$$

Item based collaborative filtering

- ↳ same as user based but similarity b/w items
- ↳ items are vectors
- ↳ more stable

"Things don't change as much as people do."
— Made-up quote

Q) V_3 is Active User, Predict for I_1 and I_6

~~skipped~~ Adjusted cosine similarity

$$sim(\vec{a}, \vec{b}) = \frac{\sum_{u \in U} (r_{u,a} - \bar{r}_u)(r_{u,b} - \bar{r}_u)}{\sqrt{\sum_{u \in U} (r_{u,a} - \bar{r}_u)^2} \sqrt{\sum_{u \in U} (r_{u,b} - \bar{r}_u)^2}}$$

almost always used for item based CF

doesnt matter which in rows or columns
can transpose for ownself

	I_1	I_2	I_3	I_4	I_5	I_6	Mean
U_1	7	6	7	4	5	4	5.5
U_2	6	7	7	4	3	4	4.8
U_3 Active User	7	3	3	1	1	?	2
U_4	1	2	2	3	3	4	2.5
U_5	1	?	1	2	3	3	2

for I_1 , calculate Adj cosine with all other items

Adj cosine

$$(1,3) = \frac{(1.5)(1.5) + (-1.5)(-0.5) + (-1)(-1)}{\sqrt{1.5^2 + -1.5^2 + -1^2} \cdot \sqrt{1.5^2 + -0.5^2 + -1^2}} = 0.912$$

$$(1,4) = \frac{(1.5)(-1.5) + (1.2)(-0.5) + (-1)(0.5)}{\sqrt{1.5^2 + 1.2^2 + 1^2} \cdot \sqrt{-1.5^2 + -0.5^2 + 0.5^2}} = -0.824$$

$$(1,2) = \frac{(1.5)(0.5) + (1.2)(2.2) + (-1)(-0.5)}{\sqrt{1.5^2 + 1.2^2 + 1^2} \cdot \sqrt{0.5^2 + 2.2^2 + 0.5^2}} = 0.735$$

do me rest

↓ after subtracting mean

	I_1	I_2	I_3	I_4	I_5	I_6	Mean
U_1	1.5	0.5	1.5	-1.5	0.5	-1.5	5.5
U_2	1.2	2.2	?	-0.5	-1.8	-0.8	4.8
U_3	7	1	1	-1	1	?	2
U_4	-1.5	-0.5	-0.5	0.5	0.5	1.5	2.5
U_5	-1	?	-1	0	1	1	2

use only dim commonly present

$K=2$, so I_2 and I_3

~~skipped~~ Prediction function

↳ Using weighted avg

$$R_{31} = \frac{(3)(0.735) + (3)(0.912)}{|0.735| + |0.912|} = 3 \rightarrow \text{Predicted rating of } I_1 \text{ of } V_3$$

Do the same steps for I_6

Collaborative Filtering

- ↳ no distribution b/w dependent and independent variables
- ↳ similar to missing value analysis but with a much larger matrix

	I ₁	I ₂	I ₃	I ₄
U ₁	5	?	?	3
U ₂	5	1	?	?
U ₃	3	?	2	5

↓ SOLUTION

Active User

VS

Classification

- ↳ distribution b/w dependent and independent variables

Significance Weighting

- ↳ tweak the score based on no of items rated together by these users
- ↳ uses discount factor
- ↳ kick in when no of common ratings b/w 2 users < a particular threshold β

optional
but recommended

□ The reliability of any similarity function $sim(u, v)$ between two users u and v is often affected by the number of common ratings between u and v i.e. $|I_u \cap I_v|$

□ When the two users have only a small number of ratings in common, the similarity function $sim(u, v)$ should include a **discount factor** to de-emphasize the importance of that particular user pair

$$\text{DiscountSim}(u, v) = \text{Sim}(u, v) \cdot \frac{\min(|I_u \cap I_v|, \beta)}{\beta}$$

$\beta \rightarrow$ threshold → hyperparameter value

	I ₁	I ₂	I ₃	I ₄	I ₅	I ₆	I ₇ , I ₈	R	Similarity
U ₁	7	6	7	4	5	4	4	0.894	1
U ₂	6	7	7	4	3	4	3	0.939	0.7
Active User U ₃	7	3	3	1	1	?			
U ₄	1	2	2	3	3	4	4	-1	
U ₅	1	?	1	2	3	3	3	-0.817	

Common Items

taking $B=4$

$U_3, U_1 = 0.894 \times \frac{\min(4, 4)}{4} \rightarrow 1$

$U_3, U_2 = 0.939 \times \frac{\min(3, 4)}{4} \cdot 0.7 \rightarrow 0.7$

Collaborative Filtering Issues

CONS

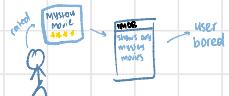
- ↳ highly dependent on data
- ↳ not good for sparse matrix → less data
- ↳ not everything is rated

↳ Serendipity

↳ basically randomness → so no overfitting

- ↳ expand user's taste into neighbouring areas
- ↳ measure of recommending something new
- ↳ personalization is indirectly proportional to serendipity
large value of K → high serendipity

it is good to recommend something new to the user



↳ Cold Start User

- ↳ new user with no info
- ↳ hard to recommend
- SOLUTION** ↳ make them rate popular items in the beginning

Cold Start item

- ↳ items with no rating
- ↳ so no similarity can be calculated
- SOLUTION** ↳ randomly pop up those items and make user rate it

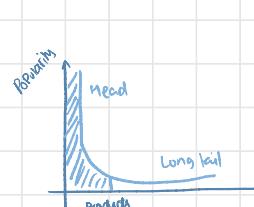
	CoV	CoI
U-CF	X cold start user wants to rate	X cold start item
I-CF	X no common items	X no user rated so cold start item

↳ Long tail

- ↳ a large no. of items will be unrated
- ↳ often leads to sparsity → not having enough data to make predictions
- ↳ products with less views

↓ SOLUTION

- can never be fully resolved → but can work with CF
- ↳ prioritize unrated movies more



↳ Scaling

- ↳ NOT scalable
- ↳ as if too many users or items the too high resource which is hard to provide

SOLUTION

↳ Use offline training

- ↳ Pre-computed calculations can be stored offline → save CPU resource

- ↳ matrix factorization → save storage

CF

- ↳ high resource
- ↳ high processing power
- ↳ high storage

will work for < 100 users

* recommender system needs to be simple and resource req. is low → as too much data

CRITERIA TO JUDGE RECOMMENDER SYSTEMS

MEMORY Based

- ↳ uses entire data everytime rating is predicted

User based CF

new data calculated from old time

Content based RS

- ↳ match user with items that are similar to what they highly rated in the past
- ↳ history based → uses your own past ratings
- ↳ uses item features, not similarity b/w users
 - ↳ there are no other users
- ↳ good for personalized recommendations

↳ serendipity
high personalization
features

Model Based

- ↳ uses data once to create a model then makes new prediction w/o using the entire data again

Item based CF

→ calc similarity b/w items and store them offline

Content based RS

→ pre stores traversed data

Collaborative based RS

- ↳ uses similarity, using ratings

amazon

USES

- ↳ movie, blogs, websites
news articles, wiki pages



Content

- ↳ explicit features of the item

e.g. for movie

↳ Genre: Action/Superhero

↳ Actor: Robert Downey Jr.

↳ Year: 2019

□ It can also be textual content (title, description, table of content, etc.)

□ Several techniques to compute the distance between two textual documents

□ Can use NLP techniques to extract content features

□ Can be extracted from the signal itself (audio, image)

Vector Space Model (VSM)

- ↳ each item is saved as a vector of its features
- ↳ in n-dimensional space
- ↳ similarity: angle b/w 2 vectors
- ↳ every user has a user profile vector

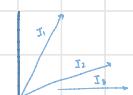


dim: no. of features

Content Based



dim: no. of users



dim: no. of items

Collaborative Filtering

Content-based Filtering APPROACH 1

Steps

1. User Profile

- ↳ One hot vector \times Input user ratings
- ↳ Normalise = $\frac{\text{E column}}{\text{E total}}$

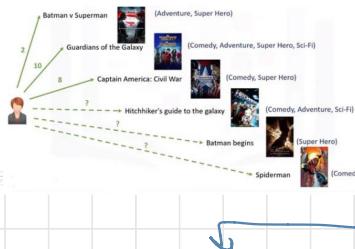
2. Recommend

1. Weighted Avg

- ↳ One hot vector \times User Profile
- ↳ Row Sum

2. Cosine Similarity

Q) features \rightarrow 1 : Genre



◻ We first create a user profile based on already rated items (and attributes)

		Comedy	Adventure	Super Hero	Sci-Fi
		0	2	2	0
Input User Ratings		10	10	10	10
		8	0	8	0
		14/10	19/10	20/10	10/10
		0.3	0.2	0.33	0.16

Normalized $\frac{\text{E column}}{\text{E total}}$

◻ Which one of the yet unrated movies should we be recommending?

	Comedy	Adventure	Super Hero	Sci-Fi
M1	1	1	0	1
M2	0	0	1	0
M3	1	0	1	0

	Comedy	Adventure	Super Hero	Sci-Fi
User Profile	0.3	0.2	0.33	0.16

	G ₁	G ₂	G ₃	G ₄	Weighted Avg
M ₁	0.3	0.2	0	0.16	0.66 \rightarrow Recommend this movie
M ₂	0	0	0.33	0	0.33
M ₃	0.3	0	0.33	0.16	0.63

Row sum $0.3+0.2+0+0.16 = 0.66$

*can use cosine similarity as well

Q) Using cosine similarity

◻ Task: Find Cosine similarity between user profile and movie matrix and compare with previous results.

0.09 0.04 0.0756

if someone wants rating, make in scale of 0 and 10

$$0.66 \approx 6.6 \approx 7$$

$$0.33 \approx 3.3 \approx 3$$

$$0.63 \approx 6.3 \approx 6$$

	G ₁	G ₂	G ₃	G ₄	Weighted Avg
M ₁	0.3	0.2	0	0.16	0.76
M ₂	0	0	0.33	0	0.64
M ₃	0.3	0	0.33	0.16	0.86 \rightarrow Recommend this movie

?

Question 2 (CLO: 1)

10 points

For the given data, create user profile and item profiles. Then use these profiles to predict ratings for Item 2 and Item 4 with content-based filtering.

User	Item 1	Item 2	Item 3	Item 4
Features	f_1, f_2	f_1, f_3	f_2, f_3	f_2
Rating	5		3	

Note: The allowed values for ratings are 1, 2, 3, 4 and 5.

Rated items User Ratings User Profile → Normalise User Profile

$$\begin{array}{|c|c|c|} \hline F_1 & F_2 & F_3 \\ \hline I_1 & 1 & 1 \\ \hline I_3 & 0 & 1 \\ \hline \end{array} \times \begin{array}{|c|} \hline R \\ \hline 5 \\ \hline 3 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline F_1 & F_2 & F_3 \\ \hline I_1 & 5 & 5 \\ \hline I_3 & 0 & 3 \\ \hline \end{array}$$

Total: 16

$$\begin{array}{|c|c|c|} \hline F_1 & F_2 & F_3 \\ \hline 0.3 & 0.5 & 0.2 \\ \hline \end{array}$$

Unrated items Normalized User Profile Low Sum Weighted Avg

$$\begin{array}{|c|c|c|} \hline F_1 & F_2 & F_3 \\ \hline I_2 & 1 & 0 \\ \hline I_4 & 0 & 1 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline F_1 & F_2 & F_3 \\ \hline 0.3 & 0.5 & 0.2 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline F_1 & F_2 & F_3 \\ \hline I_2 & 0.3 & 0 \\ \hline I_4 & 0 & 0.5 \\ \hline \end{array} = \frac{0.5}{0.5}$$

either I_2 or I_4 can be recommended

Content-based Filtering APPROACH 2

↳ Text recommendation → TF-IDF → Making Features

$$TF_{tk,dj} = \frac{freq(tk)}{\max(ti, dj)}$$

↑ row wise

just raw frequency
doesn't give imp of term

$$IDF_{tk} = \log_{10} \frac{N}{n_k}$$

↑ column wise

the more the term occurs
the lesser the imp

STEPS

1. TF
2. IDF
3. $TF \times IDF$
4. Magnitude = $\sqrt{\sum \text{of row}^2}$

$$5. \text{Normalise} = \frac{TF \cdot IDF}{\text{Magnitude}}$$

6. Cosine Similarity

Q) BLOG 1 is AI

* if active item not given
choose any by yourself

BLOG	Business	Invention	Hospital	TFB	TFI	TFH
Active Item 1	31	27	0	1.14	0.87	0
2	29	61	3	0.47	2.10	0.04
3	52	99	6	0.52	1.90	0.06
4	10	29	5	0.34	2.9	0.11
5	13	41	7	0.31	3.15	0.11
Doc Freq	2500	1200	800			

IDF	1.3	1.61	1.79
	\downarrow $\log_{10} \frac{50000}{2500}$	\downarrow $\log_{10} \frac{50000}{1200}$	\downarrow $\log_{10} \frac{50000}{800}$

Corpus contains

- N → 50000 blogs
- 2500 → Business
- 1200 → Invention
- 800 → hospital

↓ TF x IDF

0.47x13

BLOG	Business	Invention	Hospital	Magnitude
Active Item	1.48	1.40	0	$2.03 \rightarrow \sqrt{1.48^2 + 1.4^2 + 0^2}$
2	0.61	3.38	0.07	3.43
3	0.67	3.05	0.10	3.12
4	0.44	4.66	0.30	4.69
5	0.40	5.07	0.30	5.01
Doc freq	2500	1200	800	
IDF	1.3	1.61	1.79	

Normalise = Divide by magnitude

using length so scale is 0-1

BLOG	Business	Invention	Hospital	Cosine Similarity
Active Item	0.75	0.68	0	
2	0.17	0.98	0.02	0.79
3	0.21	0.91	0.03	0.82
4	0.09	0.99	0.06	0.74
5	0.07	0.99	0.05	0.73
Doc freq	2500	1200	800	
IDF	1.3	1.61	1.79	

BLOG 2 and 3 will be recommended

PROS

- ↳ Scalable new users → doesn't need too many resources
- ↳ highly personalised → max personalization
- ↳ partially solves cold start item → doesn't solve for cold start user
- ↳ recommendations are explained



you liked this
hence it is recommended

CONS

- ↳ hard to construct features → if features are not readily available nice images
- ↳ too many calculation
- ↳ little to no serendipity
- ↳ suffers from cold start user problem

DESIGNING A RECOMMENDATION SYSTEM

Probabalistic CF

aka Naive Bayes CF
doesn't have vector same → rating are features
so more accurate than normal CF

Naive Bayes classifier

↳ a supervised multi class classification algorithm

$$\text{Prediction} = \frac{\text{Posterior}}{\text{likelihood} \times \text{Prior}}$$
$$P(Y|X) = \frac{P(X|Y) \times P(Y)}{P(X)}$$

if features stay constant then remove it

Naive Bayes Collaborative Filtering

↳ Probabilities rated based on Bayes Rule

↳ assumes ratings are independent → can't use this also → imp if this is false

↳ this approach can be both

1. user based
2. item based

1. User based NB CF

$$P(r_i=y) = \frac{(\text{no of users who rated item } i \text{ as } y) + \alpha}{(\text{no of users who rated item } i) + \beta}$$

Prior
likelihood
smoothing → so value isn't 0

$$P(r_j=k | r_i=y) = \frac{(\text{no of users who rated item } j \text{ as } k \text{ and item } i \text{ as } y) + \alpha}{(\text{no of users who rated item } j \text{ and rated item } i \text{ as } y) + \beta}$$

item
smoothness

$$\beta = R \times \alpha$$

no of possible ratings
hyperparameter

Q) How probable is the rating 1 for item 5?

	I ₁	I ₂	I ₃	I ₄	I ₅	
User	U ₁	1	3	3	2	?
User	U ₂	2	4	2	2	4
User	U ₃	1	3	3	5	1
User	U ₄	4	5	2	3	3
User	U ₅	1	1	5	2	1

↓ X column
→ Y rows

Take out for all items → rating → 1
→ 2
→ 3
→ 4
→ 5

*if item not rated
don't consider the column

*if no rating given,
then dont count that user

* if β not given
 $d = 0.002$
 $\beta = \alpha \times \text{possible ratings}$

$$P(x|r_{i_5} = 1) = P(r_{i_1}: 1 | r_{i_5} = 1) \times P(r_{i_2}: 3 | r_{i_5} = 1) \times P(r_{i_3}: 3 | r_{i_5} = 1) \times P(r_{i_4}: 2 | r_{i_5} = 1)$$

likelihood for r_{i_1}
 likelihood for r_{i_2}
 likelihood for r_{i_3}
 likelihood for r_{i_4}
 likelihood for r_{i_5}

I: as 1
 I: as 1
 I: as 1
 I: as 1

2 + 0.01
 2 + 0.05
 2 + 0.05
 2 + 0.05

x
 x
 x
 x

1 + 0.01
 1 + 0.05
 1 + 0.05
 1 + 0.05

2 + 0.05
 2 + 0.05
 2 + 0.05
 2 + 0.05

= 0.117

$$P(X|r_{i_5}=2) = P(r_{i_1}=1|r_{i_5}=2) \times P(r_{i_2}=3|r_{i_5}=2) \times P(r_{i_3}=3|r_{i_5}=2) \times P(r_{i_4}=2|r_{i_5}=2)$$

$$\frac{0+0.01}{0+0.05} \times \frac{0+0.01}{0+0.05} \times \frac{0+0.01}{0+0.05} \times \frac{0+0.01}{0+0.05} = 0.0016$$

$$P(X|r_{i_5}=3) = P(r_{i_1}=1|r_{i_5}=3) \times P(r_{i_2}=3|r_{i_5}=3) \times P(r_{i_3}=3|r_{i_5}=3) \times P(r_{i_4}=2|r_{i_5}=3)$$

no of users who failed based on receiving reward
 in as 1

how many 3's in 5

$\frac{1}{4,2,1}$
 $\frac{1}{3,3}$
 $\frac{1}{1,4}$

$\frac{0+0.01}{1+0.05} \times \frac{0+0.01}{1+0.05} \times \frac{0+0.01}{1+0.05} \times \frac{0+0.01}{1+0.05} = 0.00000008^2$

$$P(X|r_{15}=4) = P(r_{11}=1|r_{15}=4) \times P(r_{12}=3|r_{15}=4) \times P(r_{13}=3|r_{15}=4) \times P(r_{14}=2|r_{15}=4) = 0.00000083$$

$$\therefore \frac{0+0.01}{1+0.05} \times \frac{0+0.01}{1+0.05} \times \frac{0+0.01}{1+0.05} \times \frac{1+0.01}{1+0.05}$$

$$P(X|r_{15}=5) = P(r_{11}=1|r_{15}=5) \times P(r_{12}=3|r_{15}=5) \times P(r_{13}=3|r_{15}=5) \times P(r_{14}=2|r_{15}=5) =$$

$$\frac{0+0.01}{1+0.05} \times \frac{0+0.01}{1+0.05} \times \frac{0+0.01}{1+0.05} \times \frac{0+0.01}{1+0.05} = 0.0016$$

Active User	Active Item				
	I ₁	I ₂	I ₃	I ₄	I ₅
U ₁	1	3	3	2	?
U ₂	2	4	2	2	4
U ₃	1	3	3	5	1
U ₄	4	5	2	3	3
U ₅	1	1	5	2	1

* On Ans, don't write Posterior Probability
but the class label

If multiply numerator by a weight → overfitting
If multiply denominator by a weight → underfitting

* If prior 0, then no need
to calculate its likelihood

Prior

no of users
who rated
I₅ on 1

$$P(r_{is}=1) = \frac{2+0.01}{4+0.05} = 0.496$$

no of users
who rated I₅

Posterior

$$P(r_{is}=1|x) = 0.117 \times 0.496 \rightarrow \text{gives max}$$

$$= 0.058$$

likelihood of

Prior of r_{is}=1

Hence
U₁ I₅ = 1 rating

$$P(r_{is}=2) = \frac{0+0.01}{4+0.05} = 0.002 , P(r_{is}=1|x) = 0.0016 \times 0.002$$

$$= 0.0000032$$

$$P(r_{is}=3) = \frac{1+0.01}{4+0.05} = 0.249 , P(r_{is}=1|x) = 0.000000082 \times 0.249$$

$$= 0.00000002$$

$$P(r_{is}=4) = \frac{1+0.01}{4+0.05} = 0.249 , P(r_{is}=1|x) = 0.000000083 \times 0.249$$

$$= 0.00000002$$

$$P(r_{is}=5) = \frac{0+0.01}{4+0.05} = 0.002 , P(r_{is}=1|x) = 0.0016 \times 0.002$$

$$= 0.0000032$$

2. Item based NB CF

no of users
FOCUS

$$P(r_{ui}=y) = \frac{(\text{no of items that user } i \text{ has given a rating } y) + \alpha}{(\text{no of total items the user has rated}) + \beta}$$

likelihood

$$P(r_j=k|r_i=y) = \frac{(\text{no of items that both user } j \text{ and user } i \text{ have rated as } y) + \alpha}{(\text{no of items that user } j \text{ has given a rating } y \text{ and }) + \beta}$$

for which user i has given a rating

AU ki existing ratings

Q) How probable is the rating 1 for item 5?

Tell how many decimal places you rounded off

Active User	I ₁	I ₂	I ₃	I ₄	I ₅
U ₁	1	3	3	2	?
U ₂	2	4	2	2	4
U ₃	1	3	3	5	1
U ₄	4	5	2	3	3
U ₅	1	1	5	2	1

likelihood
↑

$$P(r_{uj}=1 | r_{u1}=1) = P(r_{u2}=1 | r_{u1}=1) \times P(r_{u3}=1 | r_{u1}=1) \times P(r_{u4}=1 | r_{u1}=1) \times P(r_{u5}=1 | r_{u1}=1)$$

$$= \frac{0+0.01}{1+0.05} \times \frac{1+0.01}{1+0.05} \times \frac{0+0.01}{1+0.05} \times \frac{1+0.01}{1+0.05} \cdot 0.0000839$$

no. of rated
U₃, U₅ rated 1

 no. of items rated
U₃ rated 1 and
U₅ rated something
no. of items rated
U₃, U₅ rated 1

$$P(r_{uj}=2 | r_{u1}=2) = P(r_{u2}=2 | r_{u1}=2) \times P(r_{u3}=2 | r_{u1}=2) \times P(r_{u4}=2 | r_{u1}=2) \times P(r_{u5}=2 | r_{u1}=2)$$

$$= \frac{1+0.01}{1+0.05} \times \frac{0+0.01}{1+0.05} \times \frac{0+0.01}{1+0.05} \times \frac{1+0.01}{1+0.05} \cdot 0.0017624$$

$$P(r_{uj}=3 | r_{u1}=3) = P(r_{u2}=3 | r_{u1}=3) \times P(r_{u3}=3 | r_{u1}=3) \times P(r_{u4}=3 | r_{u1}=3) \times P(r_{u5}=3 | r_{u1}=3)$$

$$= \frac{0+0.01}{2+0.05} \times \frac{2+0.01}{2+0.05} \times \frac{0+0.01}{2+0.05} \times \frac{0+0.01}{2+0.05} \cdot 0.0000001$$

no. of rated
U₃, U₅ rated 3

 no. of items rated
U₃ rated 3 and
U₅ rated something
no. of items rated
U₃, U₅ rated 3

$$P(r_{uj}=4 | r_{u1}=4) = P(r_{u2}=4 | r_{u1}=4) \times P(r_{u3}=4 | r_{u1}=4) \times P(r_{u4}=4 | r_{u1}=4) \times P(r_{u5}=4 | r_{u1}=4)$$

$$= \frac{0+0.01}{0+0.05} \times \frac{0+0.01}{0+0.05} \times \frac{0+0.01}{0+0.05} \times \frac{0+0.01}{0+0.05} \cdot 0.006$$

$$P(r_{uj}=5 | r_{u1}=5) = P(r_{u2}=5 | r_{u1}=5) \times P(r_{u3}=5 | r_{u1}=5) \times P(r_{u4}=5 | r_{u1}=5) \times P(r_{u5}=5 | r_{u1}=5)$$

$$= \frac{0+0.01}{0+0.05} \times \frac{0+0.01}{0+0.05} \times \frac{0+0.01}{0+0.05} \times \frac{0+0.01}{0+0.05} \cdot 0.006$$

	I ₁	I ₂	I ₃	I ₄	I ₅	
Active User	U ₁	1	3	3	2	?
	U ₂	2	4	2	2	4
	U ₃	1	3	3	5	1
	U ₄	4	5	2	3	3
	U ₅	1	1	5	2	1

Prior
↑

$$P(U_1=1) = \frac{1}{4} \xrightarrow{\substack{\text{no of items} \\ \text{that AU rated 1}}} = 0.25$$

↓
NO. OF ITEMS
AU rated

$$, P(r_{U_1, i_5=1} | x) : 0.0000839 \times 0.25 \\ = 0.000,020,9$$

$$P(U_1=2) = \frac{1}{4} = 0.25$$

$$, P(r_{U_1, i_5=2} | x) : 0.0011624 \times 0.25 \xrightarrow{\substack{\text{MAX} \\ \text{Hence} \\ U, i_5=2 \text{ rating}}}$$

$$P(U_1=3) = \frac{2}{4} = 0.5$$

$$, P(r_{U_1, i_5=3} | x) = 0.000,000,1 \times 0.5 \\ = 0.000,000,0$$

$$P(U_1=4), \frac{0}{4} \xrightarrow{?}$$

$$, P(r_{U_1, i_5=4} | x) = 0.006 \times 0 \xrightarrow{?} \\ = 0$$

$$P(U_1=5) = \frac{0}{4}$$

$$, P(r_{U_1, i_5=5} | x) = 0.006 \times 0 \\ = 0$$

NAIVE BAYES CF

PROS

- ↳ more accurate recommendations
- ↳ can provide ranking of predicted ratings
- ↳ can provide confidence level for a prediction
↳ handles sparse data well
↳ simple implementation

how likely
rating 2
can come

→ be 90% sure
prediction is 4

CONS

- ↳ can be computationally intractable
- ↳ serendipity can't be controlled
- ↳ independence b/w rating is required

Provide brief answers for the following and justify:

- a) Why is serendipity so important in recommender systems? How do we control serendipity in item-based Collaborative Filtering?

Answer: By serendipity we mean recommending items to the user such that his taste is expanded. We can control serendipity in item-based collaborative filtering by adjusting the value of hyperparameter k (as in number of items to be considered in the neighborhood).

- b) How do we handle cases where the ratings are allowed to be real-valued in Naïve Bayes Collaborative Filtering?

Answer: By assuming the ratings to be normally distributed and using the Gaussian or z-distribution to estimate the likelihood.

- c) Let us say that we have a photography website. The site contains both user-created and AI-generated images. Should our recommender system be giving more importance to user-created images over the AI-generated ones for recommendation? Provide a reason.

Answer: Although there is no definite answer to this open issue yet, one strong argument is to give preference to AI generated images as often user feedback is important in performance evaluation and improvement of the model.

- d) Can content-based recommender systems solve cold-start problem if a new item is added?

Answer: Yes. It can solve cold-start problem if a new item is added since this new item can also be recommended to a user based on feature similarity.

- e) If an active user has k yet unrated items in a system of m total items, how many vectors would be present in the vector space model at one time?

Answer: $k+1$ (one user vector and k item vectors).

History Based Techniques

- ↳ uses historical data to make recommendations
 - ↳ collaborative Filtering
 - ↳ Content based Filtering
- ↳ These techniques
 - ↳ require data
 - ↳ are not adaptable
 - ↳ doesn't ask user for preferences

Knowledge based RS

- ↳ exploits user requirements
- ↳ Used when
 - ↳ customers have explicit requirements
 - ↳ it's difficult to get range of a specific type of item
 - ↳ when ratings are time sensitive
- ↳ Application Domain
 - ↳ expensive items
 - ↳ not frequently purchased
 - ↳ low ratings
 - ↳ e.g. car names
 - ↳ Time spans important → technological product
 - ↳ Exploits requirements of user → buying property

Knowledge-based Systems: Example

- ↳ It has 2 types
 - 1. Constraint Based Systems
 - ↳ explicitly defined by conditions
 - 2. Rule Based Systems
 - ↳ similarity to specific requirements
- If you're looking for a house or a car online. You input price, how many rooms, how much total floor space, etc., and the website returns a list of houses based on those constraints
 - Issue:
 - Is this recommendation or simple query search?

Knowledge base: Example

```

 $V_C = \{kl_c: [\text{expert, average, beginner}] \dots \text{ /* level of expertise */}$ 
 $wr_c: [\text{low, medium, high}] \dots \text{ /* willingness to take risks */}$ 
 $id_c: [\text{shortterm, mediumterm, longterm}] \dots \text{ /* duration of investment */}$ 
 $aw_c: [\text{yes, no}] \dots \text{ /* advisory wanted? */}$ 
 $ds_c: [\text{savings, bonds, stockfunds, singleshares}] \dots \text{ /* direct product search */}$ 
 $sl_c: [\text{savings, bonds}] \dots \text{ /* type of low-risk investment */}$ 
 $av_c: [\text{yes, no}] \dots \text{ /* availability of funds */}$ 
 $sh_c: [\text{stockfunds, singlshares}] \dots \text{ /* type of high-risk investment */} \}$ 

```

```

 $V_{PROD} = \{name_p: [\text{text}] \dots \text{ /* name of the product */}$ 
 $er_p: [1..40] \dots \text{ /* expected return rate */}$ 
 $ri_p: [\text{low, medium, high}] \dots \text{ /* risk level */}$ 
 $mniv_p: [1..14] \dots \text{ /* minimum investment period of product in years */}$ 
 $inst_p: [\text{text}] \dots \text{ /* financial institute */} \}$ 

```

Knowledge base: Example

$C_R = \{CR_1: wr_c = \text{high} \rightarrow id_c \neq \text{shortterm},$
 $CR_2: kl_c = \text{beginner} \rightarrow wr_c \neq \text{high}\}$

$C_F = \{CF_1: id_c = \text{shortterm} \rightarrow mniv_p < 3,$
 $CF_2: id_c = \text{mediumterm} \rightarrow mniv_p \geq 3 \wedge mniv_p < 6,$
 $CF_3: id_c = \text{longterm} \rightarrow mniv_p \geq 6,$
 $CF_4: wr_c = \text{low} \rightarrow ri_p = \text{low},$
 $CF_5: wr_c = \text{medium} \rightarrow ri_p = \text{low} \vee ri_p = \text{medium},$
 $CF_6: wr_c = \text{high} \rightarrow ri_p = \text{low} \vee ri_p = \text{medium} \vee ri_p = \text{high},$
 $CF_7: kl_c = \text{beginner} \rightarrow ri_p \neq \text{high},$
 $CF_8: sl_c = \text{savings} \rightarrow name_p = \text{savings},$
 $CF_9: sl_c = \text{bonds} \rightarrow name_p = \text{bonds} \}$

$C_{PROD} = \{C_{PROD_1}: name_p = \text{savings} \wedge er_p = 3 \wedge ri_p = \text{low} \wedge mniv_p = 1 \wedge inst_p = A;$
 $C_{PROD_2}: name_p = \text{bonds} \wedge er_p = 5 \wedge ri_p = \text{medium} \wedge mniv_p = 5 \wedge inst_p = B;$
 $C_{PROD_3}: name_p = \text{equity} \wedge er_p = 9 \wedge ri_p = \text{high} \wedge mniv_p = 10 \wedge inst_p = B\}$

Conceptual Goals of Recommender Systems

Approach	Conceptual Goal	Input
Collaborative	Give me recommendations based on a collaborative approach that leverages the ratings and actions of my peers/myself.	User ratings + community ratings
Content-based	Give me recommendations based on the content (attributes) I have favored in my past ratings and actions.	User ratings + item attributes
Knowledge-based	Give me recommendations based on my explicit specification of the kind of content (attributes) I want.	User specification + item attributes + domain knowledge