

(\*) Stemming and Lemmatization :-

→ Process of reducing inflected words to their word STEM.

### Stemming.

history → histori  
historical → historical

finally → final  
final → final  
finalized → finalized

### Lemmatization

history → history  
historical → historical

finally → final  
final → final  
finalized → finalized

gone  
goes → go  
going

#### NOTE:-

- (1) Lemmatization the converted words will have meaning whereas in stemming in some of the cases it is not possible.
- (2) Lemmatization takes more time than stemming.

Q. where we can use stemming and lemmatization?

→ stemming :- sentiment Analysis, spam classifier.

Reason :- In this cases we require only base word. we are having the words in our dataset.

Lemmatization :- chatBot, Question Answer App

Reason :- The response that human gets from app should be meaningful.

NOTE:- (1) Lemmatization is slower and accurate.  
(2) Stemming is faster but inaccurate.

BEST

## Practical implementation of stemming:-

```

import nltk
from nltk.stem import PorterStemmer
from nltk.corpus import stopwords
paragraph = "The quick brown fox jumps over the lazy dog"
sentences = nltk.sent_tokenize(paragraph)
stemmer = PorterStemmer()
    
```

**convert paragraph into sentences**

e.g.: of, them, from the, etc...  
these words do not play an imprtnt role in +ve sentiment analysis.

stemmer = PorterStemmer()

# Now we will remove stop words from each sentences...  
# then after removing the stop words we will stemming  
# each word using stemmer.

for i in range(len(sentences)):

# convert each sentence into words--

words = nltk.word\_tokenize(sentences[i])

converting sentences into words

# Stemming the words--

words = [stemmer.stem(word) for word in words]

stemming the words which are not stopwords.

if word not in set(stopwords.words('english'))]

# join the words

sentences[i] = ' '.join(words)

## Problem with stemming:-

→ Produced intermediate representation of the word may not have any meaning.

Example: fina, histori, etc.

## (\* Practical implementation of lemmatization :-

responsible for  
Lemmatization

```
import nltk
from nltk.stem import WordNetLemmatizer
from nltk.corpus import stopwords
```

```
paragraph = " ".join(sentences)
sentences = nltk.sent_tokenize(paragraph)
```

```
# creating object of WordNetLemmatizer
```

```
lemmatizer = WordNetLemmatizer()
```

### # Lemmatizing

1. creating words from sentences.
2. lemmatizing words which are not stopwords.
3. joining the words to form the sentences.

III

```
for i in range(len(sentences)):
```

Responsible for converting  
sentences to word

```
# 1
words = nltk.word_tokenize(sentences[i])
```

```
words = [lemmatizer.lemmatize(word) for word in words]
```

Responsible for  
lemmatizing the  
words which  
are not stopwords.

```
# 3. Join the words...
```

```
sentences[i] = " ".join(words)
```

Ques.

When to use lemmatization and when to use stemming?

→ If the ultimate goal is to understand the meaning of the text, lemmatization is the better choice.

If the goal is simply to extract features from the text, stemming can be more efficient option.

## (\*) Bag of Words :-

sent-1 → He is a good boy.

sent-2 → she is a good girl.

sent-3 → Boy & girl are good

stop words

lowering sentence

sent-1 → good boy

sent-2 → good girl

sent-3 → boy girl  
good.

words →	good	boy	girl
freqn →	3	2	2

} sort the frequency in desc. order if not sorted.

Binary Bag of words

BOW

vectors

f3

	good	boy	girl
Sent1	1	1	0
Sent2	1	0	1
Sent3	1	1	1

## \* Problems with Bag of words :-

- (1) It does not tell us about the importance of the words / features.
- (2) Not preferable for larger datasets.
- (3) Does not take into account the order of words in a document.
- (4) Treats all words as equally important.
- (5) computationally expensive.

## \* Practical implementation of Bag of words:-

```
import nltk
```

```
paragraph = " "
```

```
# cleaning the text...
```

```
import re
```

```
from nltk.corpus import stopwords ← for removing stop  
keywords.
```

```
from nltk.stem import PorterStemmer
```

```
from nltk.stem import WordNetLemmatizer
```

```
stemmer = PorterStemmer() ← for stemming
```

```
lemmatizer = WordNetLemmatizer() ← for lemmatizing
```

```
# converting paragraph into sentences...
```

converting para to sent.

```
sentences = nltk.sent_tokenize(paragraph)
```

```
corpus = [] ← for storing cleaned sentences.
```

```
for i in range(len(sentences)) :
```

```
    review = re.sub('[^a-zA-Z]', ' ', sentences[i])
```

```
    review = review.lower()
```

```
    review = review.split() ← for converting sent to list of words.
```

```
    review = [stemmer.stem(word) for word in review if word not in
```

```
    set(stopwords.words('english'))]
```

↑ Stemming the words other than  
stopwords.

```
review = ' '.join(review)
```

```
corpus.append(review)
```

converting into  
counted  
vector

```
# creating the Bag of words Model...
```

```
from sklearn.feature_extraction.text import CountVectorizer
```

```
cv = CountVectorizer()
```

```
cv = CountVectorizer() → for converting into matrix form.
```

```
X = cv.fit_transform(corpus).toarray()
```

# (\*) TF-IDF (Term Frequency - Inverse Document Frequency)

sent<sub>1</sub> → good boy  
 sent<sub>2</sub> → good girl  
 sent<sub>3</sub> → boy girl good

words	good	boy	girl
freq <sup>n</sup>	3	2	2

$\text{TF} = \frac{\text{No. of repetition of words in a sentence}}{\text{No. of words in a sentence}}$

TF

	Term .. Frequency		
	sent <sub>1</sub>	sent <sub>2</sub>	sent <sub>3</sub>
good	4/2	4/2	4/3
boy	4/2	0	4/3
girl	0	4/2	4/3

$\text{IDF} = \log \left( \frac{\text{No. of Sentences}}{\text{No. of sentences containing that word}} \right)$

TF-IDF

	f <sub>1</sub>	f <sub>2</sub>	f <sub>3</sub>
sent 1	good	boy	girl
sent 2	0	0	4/2 log(3/2)
sent 3	0	4/3 log(3/2)	1/3 log(3/2)

IDF

words	IDF
good	$\log(3/3) = 0$
boy	$\log(3/2) =$
girl	$\log(3/2)$

## \* Practical Implementation of TF-IDF

```
import nltk  
paragraph = " "  
# cleaning  
import re  
from nltk.stem import PorterStemmer  
from nltk.stem import WordNetLemmatizer  
from nltk.corpus import stopwords  
stemmer = PorterStemmer()  
lemmatizer = WordNetLemmatizer()
```

we can use any  
of them as per  
the requirements

### # convert para to sent...

```
sentences = nltk.sent_tokenize(paragraph)  
corpus = [] ← for storing cleaned sentences...
```

```
for i in range(len(sentences)):
```

```
    review = re.sub('[^a-zA-Z]', ' ', sentences[i])
```

```
    review = review.lower()
```

```
    review = review.split() ← conv. sent to words...
```

```
    review = [stemmer.stem(word) for word in review]
```

stemming

```
    if word not in  
        set(stopwords.words('english'))]
```

```
    review = '/'.join(review)
```

```
    corpus.append(review)
```

### # creating the TF-IDF model...

```
from sklearn.feature_extraction.text import TfidfVectorizer  
tfid = TfidfVectorizer() → convert into TF-IDF table/matrix  
X = tfid.fit_transform(corpus).toarray()
```

(\*) WORD EMBEDDINGS (Feature Representation).

	Boy	Girl	King	Queen	Apple	Mango
Gender	-1	1	-0.92	0.93	0.0	0.1
Royal	0.01	0.02	0.95	0.96	-0.02	0.01
Age	0.03	0.04	0.7	0.6	0.92	0.93
⋮	⋮	⋮	⋮	⋮	⋮	⋮
<u>300 dim</u>	[	[	[	[	[	[

[300x1000] ↗ Low Dimension and Dense Matrix.

## (\*) Recommender Systems (K)

### \* Key Problems:-

- (1) Gathering "known" ratings for matrix
- (2) Extrapolate unknown ratings from the known ones  
↳ interested in what you like instead of what you don't like.
- (3) Evaluating extrapolation methods.  
↳ How to measure performance of recomm. model.

### (I) Gathering ratings:-

#### a) Explicit:

↳ ask users to rate items.

Problem: only small fraction of users leave ratings and reviews.

Doesn't scale  
↳ so, data not sufficient or excellent.

#### b) Implicit:

↳ learn ratings from user actions.

e.g. purchase → high rating

But we cannot capture the low rating using this method.

NOTE: Most of the recomm. system uses combination of explicit and implicit ratings.

When explicit ratings are available they use them and supplement implicit rating with them and use it.

### (II) Extrapolating utilities :-

key problem: matrix  $U$  is sparse as most of users haven't leave ratings and reviews.

#### Cold start problem :-

↳ New items have no ratings.

↳ New user have no history.

## Approaches to Recommender System

Content Based

Collaborative

Latent Factor

### (\*) Content Based Recommendation System :-

- Objective :- Recommend items  $\rightarrow$  customer (X) similar items to previous items rated highly by customer (X)

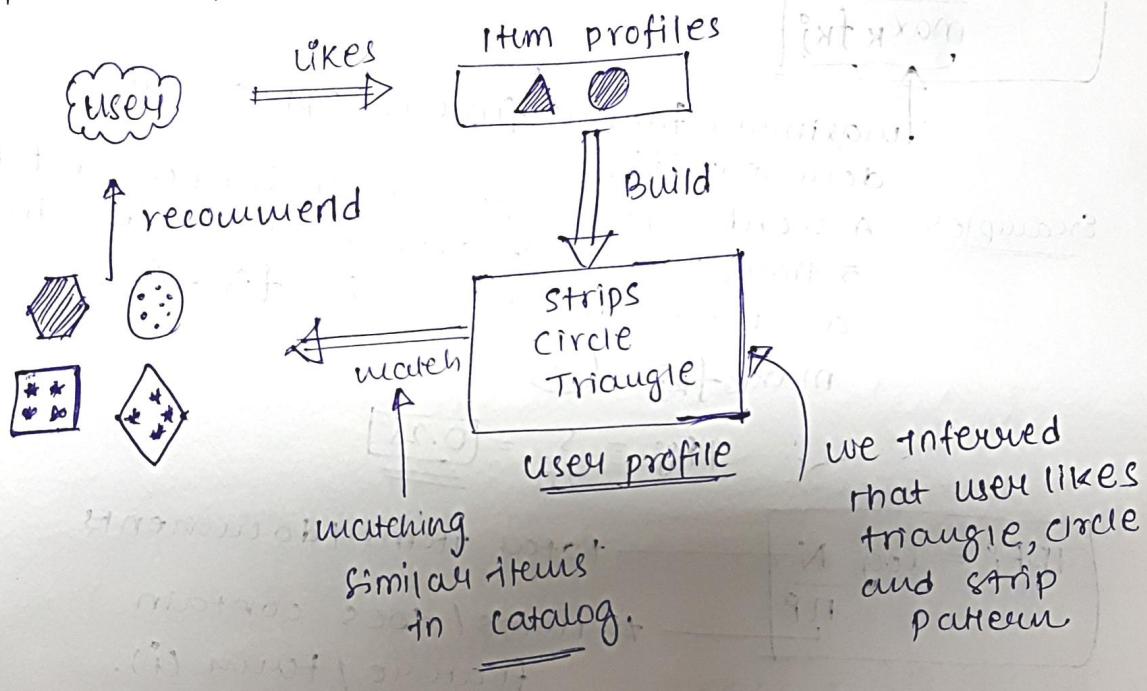
Example :- movies

↳ same genre

Blogs

↳ same author, content, etc..

- Firstly, we will find set of items user likes using implicit or explicit data.



## (=) Item Profile :-

↳ for each item create item profile.

Think item profile as a vector.

↳ ~~each~~ one entry for each feature. Image  
↳ it might be boolean or real valued.

set of features  
of items.

Ex:- movies

↳ author, genre.

Image

↳ tags, metadata

Hence, Profile  $\vec{x}$ )

↳ set of important words in item document

How to get them?

↳ using TF-IDF.

## (#) TF-IDF :-

$f_{ij}^o$  = frequency of feature  $(i)$  in item/doc  $(j)$

$$TF_{ij}^o = \frac{f_{ij}^o}{\max_k f_{kj}^o}$$

maximum time a feature appear in document  $(j)$ .

Example: A word apple appears in a document  $(j)$  5 times. But it appears 20 times in another document. So  $f_{ij}^o = 5$ .

$$\max_k f_{kj}^o = 20$$

$$TF_{ij}^o = \frac{5}{20} = 0.25$$

$$IDF_i^o = \log \frac{N}{n_i^o}$$

total items / documents.

# items / docs contain feature / term  $(i)$ .

TF-IDF score:  $w_{ij} = TF_{ij}^o \times IDF_i^o$

NOTE:- Need to Normalize TF to discount for "longer" document.

The more the common term is → larger the → lower the IDF;

NOTE:- IDF gives lower weights to more common words and higher weights to less common or rarer words.

⇒ create TF-IDF vector, sort it and put some threshold.  
⇒ Doc profile :- set of words with highest TF-IDF scores, together with their scores.

(=) User Profile :-

- ↳ user has rated items with profile  $i_1, i_2, \dots, i_n$ .
- ↳ simple average of rated item profiles.
- ↳ problems :- doesn't take into account the likeness of item.  
It takes all items to be equally likely.

Solution :- Take weighted average.

Variant :- Normalize weights using average rating of user.

Example: Boolean utility matrix

items → movies

feature → "Actor"

item profile → vector of 0 or 1 for each actor.

X has watched 5 movies.

2 consisting of Actor A. → weight = 2/5 = 0.4.

3 consisting of Actor B. → weight = 3/5 = 0.6.

3

Example :- Star ratings.

Actor A's movies rated 3 and 5

Actor B's movies rated 1, 2 and 4

(\*) Normalization helps to capture the idea that some are the +ve ratings and some of them are -ve ratings.

(\*) users are diff. their perspective might be diff. some consider 4 as average rating and some consider 4 as +ve rating (high).

⇒ Normalizing ratings by subtracting user's mean:-

$$\text{weight} = (0+2)/2 = \underline{\underline{1}}$$

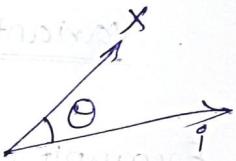
Actor "A" :- 0, +2

$$\text{Actor "B" :- } -2, -1, +1 \quad \text{weight} = \underline{\underline{-2/3}}$$

(=) Making Predictions :-

Taking pair of user profile ( $X$ ) and item profile ( $i$ ) and finding cosine similarity betn them.

$$U(X, i) = \cos(\theta) = (X, i) / (||X|| \cdot ||i||)$$



(\*) Smaller the ( $\theta$ )

higher the similarity.

NOTE: cosine distance is the actually angle ( $\theta$ ) and cosine similarity is the angle ( $180 - \theta$ ). For our convenience, we use  $\cos(\theta)$  as our similarity measure and call it cosine similarity...

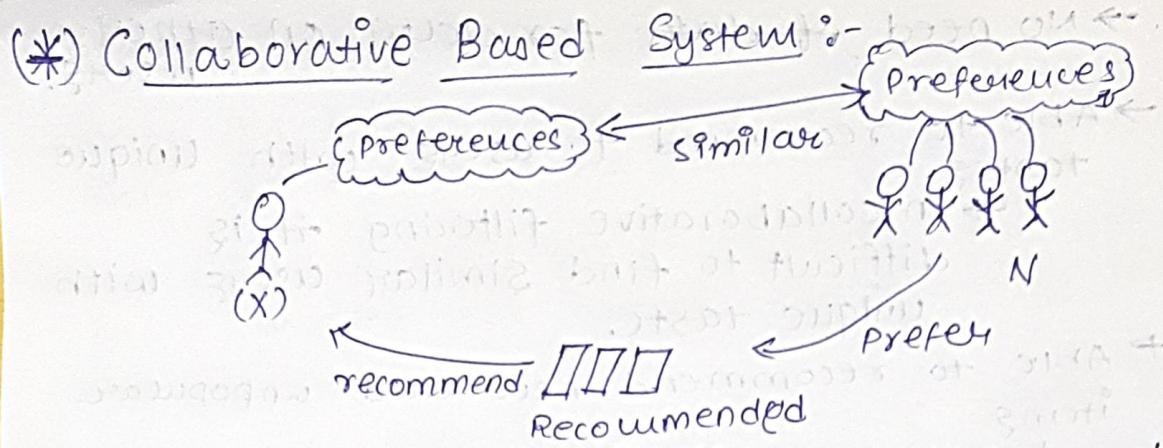


- No need of data for data on other users.
- Able to recommend to users with unique tastes
  - ↳ in collaborative filtering it is difficult to find similar users with unique taste.
- Able to recommend new and unpopular items
  - ↳ No first-rater problem.

- Explanations for recommended items
  - ↳ we can explain user that based on what feature the given item is recommended...



- Finding the appropriate features is hard.
- overspecialization
  - ↳ never recommends items outside user's content profile.
  - ↳ people might have multiple interests
  - ↳ unable to exploit quality judgements of other users.
- cold-start problem for new users.
  - ↳ How to build user-profile?



- ⇒ Find set N of other users whose rating is X's rating.
- ⇒ Estimate X's ratings based on ratings of users in N.

Example:-

	HP1	HP2	HP3	TW	SW1	SW2	SW3
A	4			5	1		
B	5	5	4				
C				2	4	5	
D		3					3

⇒  $\gamma_A, \gamma_B$  are the respective rating vectors for users A and B.

⇒ we need similarity metric  $\text{sim}(A, B)$ .

⇒ Here we can see  $\text{sim}(A, B) > \text{sim}(A, C)$ .

OPTION-1 :- Jaccard Similarity.

$$\text{sim}(A, B) = |\gamma_A \cap \gamma_B| / |\gamma_A \cup \gamma_B|$$

$$\text{sim}(A, B) = 4/5 \quad \} \quad \text{sim}(A, B) < \text{sim}(A, C)$$

$$\text{sim}(A, C) = 2/4 \quad \}$$

It counters the notion that we have captured intuitively.

Problem:- ignores rating values...

- Option-2:- cosine similarity :-

$$\text{sim}(A, B) = \cos(\gamma_A, \gamma_B)$$

$$\downarrow \\ \underline{0.38} \quad \text{sim}(A, C) = \underline{0.32}$$

$$\text{sim}(A, B) > \text{sim}(A, C)$$

But not by too much.

$\Rightarrow$  we had filled empty ratings with zero.

Problem :- treats missing ratings as negative.

correlation

- Option 3:- centred cosine: (also known as Pearson)

$\rightarrow$  normalize the rating by subtracting row mean.

$$\text{avg. rating of } A = 10/3 = \underline{3.33}$$

$$\text{avg. rating of } B = 14/3$$

$$\text{u of } \underline{C} = 11/3$$

$$\text{u of } D = 6/2 = 3$$

⋮

Modified:

	HPI	HP2	HP3	TW	SW1	SW2	SW3
A	2/3			5/3	-7/3		
B	4/3	4/3	-2/3	-5/3	4/3	4/3	
C							0
D	0						

sum of  
ratings  
row wise  
= 0

zero becomes  
avg. rating of  
each user.

$$\text{Now, } \text{sim}(A, B) = \cos(\gamma_A, \gamma_B) = \underline{0.09}$$

$$\text{sim}(A, C) = \underline{-0.56}$$

$$\therefore \text{sim}(A, B) > \text{sim}(A, C)$$



captures intuition better  
 ↳ missing ratings treated as "average".  
 ↳ Handles "tough raters" and "easy raters".

## $\Rightarrow$ Rating Predictions :-

- Let  $\underline{r_x}$  be the user  $x$ 's ratings.
- Let  $\underline{N}$  be the set of  $K$  most similar to  $\underline{x}$  users, who have also rated item  $i$ .
- Prediction for user  $x$  and item  $i$
- set  $N$  consist only those users who have rated item  $(i)$ .

Option 1 :- Average Rating from the neighbours.

$$r_{xi} = \frac{1}{K} \sum_{y \in N} r_{yi}$$

Here, all users (neighbours) given equal importance.

Option 2 :- weighted Average

$$r_{xi} = \frac{\sum_{y \in N} s_{xy} r_{yi}}{\sum_{y \in N} s_{xy}}$$

$\boxed{\text{sim}(x,y)}$

## (\*) ITEM-ITEM Collaborative Filtering :-

- 1) For item  $i$ , find other similar items.
  - 2) Estimate rating for item  $i$  based on ratings for similar items.
- (\*) can use the same similarity metrics and prediction functions as in user-user model.

$$r_{xi} = \frac{\sum_{j \in N(i; x)} S_{ij} \cdot r_{xj}}{\sum_{j \in N(i; x)} S_{ij}}$$

rating of user  $x$  on item  $j$

similarity of items  $i$  and  $j$

set of items rated by  $x$  similar to  $i$ .

Example:- Item-Item CF ( $|N|=2$ ) users. No. of neighbours = 2.

	1	2	3	4	5	6	7	8	9	10	11	12	Sim( $i, m$ )
1	1		3		?	5			5		4		1.00
2		5	4			4		4		2	1	3	-0.18
3	2	4		1	2		3		4	3	5		0.41
4		2	4	3	5	2		4			2		-0.10
5			4		4								
6	1		3		3			2			2	5	-0.31
correlation as 1.00													0.59

Here we use Pearson cosine similarity:

- (1) subtract the mean rating  $m_i$ , for each movie  $i$
- $$m_i = (1+3+5+5+4)/5 = \underline{\underline{3.6}}$$
- $$\text{row 1} = [-2.6, 0, -0.6, 0, 0, 1.4, 0, 0, 1.4, 0, 0, 4, 0]$$

- (2) compute the cosine similarities between rows...

$$S_{13} = 0.41 \quad S_{16} = 0.59$$

Predict by taking weighted average :-

$$\hat{Y}_{15} = \frac{(0.41 * 2 + 0.59 * 3)}{(0.41 + 0.59)} \stackrel{\text{Q}}{=} [2.6]$$

### (\*) Item-Item vs User-User :-

→ in practice item-item outperforms user-user in many use cases. like movie, books, etc.

Why?

↳ item-item are "simple", than user.

↳ item belongs to small set of "genres", users have varied ~~tests~~ tastes.

↳ item similarity more meaningful than user similarity.

### (\*) Implementing Collaborative Filtering :-

(=) complexity :-

→ most expensive step is finding K-most similar user or

items  $\rightarrow O(|U|)$   $|U| \leftarrow$  size of utility matrix.

→ To expensive to do at runtime

↳ could pre-compute

↳ Naive pre-computation takes time

$O(n \cdot |U|)$   $n = \# \text{ items / users}$

Ways :-

→ near-neighbour search in high dimensions.

→ clustering → we can restrict the in cluster searching to save time.

→ Dimensionality reduction

## $\Rightarrow$ Pros/cons of collaborative filtering :-

- works for any kind of item. } it is tough to find right  
    ↳ no feature selection needed. } set of features
- Cold Start Problem. Need enough no. of users to find a match
- Sparsity. - user/ratings matrix is sparse.
  - Hard to find users that have rated same items.
- First rater Problem :-
  - cannot recommend an unrated item.
  - New items, Esoteric items.
- Popularity Bias :-
  - Tend to recommend popular items.

## $\Rightarrow$ Hybrid Methods :-

$\Rightarrow$  Add content based methods to collaborative Filtering

- Item profiles for new item problems.
- Demographics to deal with new user problem.

$\Rightarrow$  Implement two or more different recommenders and combine predictions.

- using a linear model.
- Ex:- global baseline + collaborative filtering.

## $\Rightarrow$ Global Baseline Estimates :-

$\rightarrow$  Estimate Hiren's rating for movie KGF2

Problem :- Hiren has not rated any movie similar to KGF2.

### Global Baseline Approach :

- Mean movie rating :- 3.7 stars.
- KGF2 is 0.5 stars above avg.
- Hiren's rating is 0.2 below the avg.
- Base Line estimate :-  $3.7 + 0.5 - 0.2 = \boxed{4 \text{ stars}}$

(\*) Haven't used any information about the movies Hiren rated, similar to KGF2.

### (≡) Combining Global Baseline with CF :-

⇒ Global Baseline estimates

- Hiren will give KGF2 4 stars.

⇒ Local Neighbourhood (CF/NN):-

- Hiren didn't like related movie KGF.

- Rated it 1 star below the avg. rating

⇒ Final estimate.

- Hiren will rate KGF2  $4 - \frac{1}{3} = \underline{\underline{3 \text{ stars}}}$

#### • common Practice:-

- Define similarities  $s_{ij}$  of items  $i$  and  $j$

- select the  $K$  nearest neighbours  $N(i; n)$

↳ items most similar to  $i$  and rated by user  $x$ .

- Estimate rating  $\hat{r}_{xi}$  as the weighted avg:

$$\hat{r}_{xi} = b_{xi} + \frac{\sum_{j \in N(i; x)} s_{ij} (\hat{r}_{xj} - b_{xj})}{\sum_{j \in N(i; x)} s_{ij}}$$

Baseline estimator  
for  $\hat{r}_{xi}$

$$b_{xi} = \bar{u} + b_x + b_i$$

overall  
mean  
movie  
rating

$b_x$  = rating deviation of user  $x$

$$= (\text{avg. rating of user } x) - \bar{u}$$

$b_i$  = rating deviation of movie ( $i$ )

## (\*) Evaluating Recommender Systems :-

- Take the piece of data from the user-item rating matrix, and consider as test set.
- compare the predictions against withheld known ratings
- Root mean squared Error.

$$\sqrt{\frac{\sum_{(x_i, i) \in T} (r_{xi} - r_{xi}^*)^2}{N}}$$

$$N = |T|$$

$r_{xi}^*$  ← actual.

$r_{xi}$  ← predicted.

## (\*) Problems with Error Measures :-

- Narrow focus on accuracy sometimes misses
  - ↳ Prediction diversity.
  - ↳ Prediction content.
  - ↳ order of predictions.
- we care only to predict high ratings (in practice)
  - ↳ RMSE might penalize a method that does well for high ratings and badly for others.

Alternative: precision at top K.

↳ % of predictions in the user's top K withheld ratings.