



Non personalized
(TOP-Pop)

Personalized

Content Based
find item with
some features
(CBF)

Collaborative
Based
(CRS)

Context
aware
Hybrid

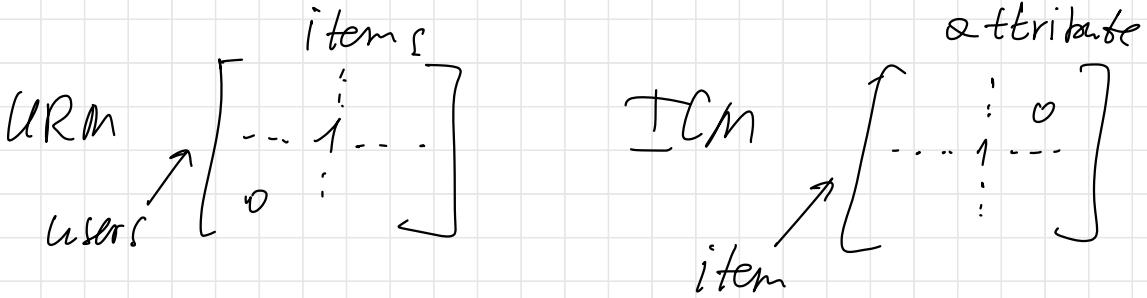
User Based
Memory Based

Item Based
Model Based
SLIM

Matrix
factorization
(MF)

DL,
Factoriza-
tion
machine

I didn't write metrics formulas
here



{0, 1} - implicit parameter (interested or not)

[0 - 5] - explicit parameter, (how user evaluated recommendation)

Dot product

$$\vec{i} \cdot \vec{j} = \sum_k i_k \cdot j_k = s_{ij}$$

↑
similarity

Cosine similarity

$$\frac{\vec{i} \cdot \vec{j}}{\|\vec{i}\| \cdot \|\vec{j}\|}$$

normalised

$$\|\vec{i}\| \cdot \|\vec{j}\| \rightarrow \cos \theta$$

θ -angle between

Shrink term

$$\frac{\vec{i} \cdot \vec{j}}{\|\vec{i}\| \cdot \|\vec{j}\|}$$

vector \vec{i}, \vec{j}

if interested only in one item/attribute

\Rightarrow will be not relevant (some kind
Content Based Filtering (ICM) of regulariza-
tion)

item $\begin{bmatrix} - & 0,2 \\ 0,1 & - \end{bmatrix}$] Similarity matrix
for item
similar, with empty
diagonal

Rating prediction based on Item Sim Matrix

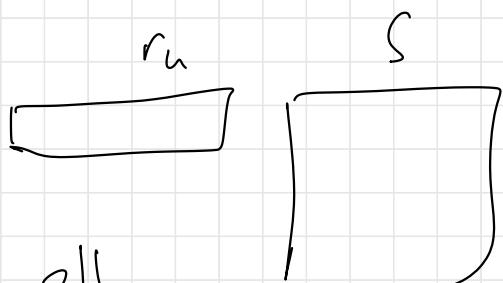
normalised

$$r_{ui} = \frac{\sum_j (r_{uj} \cdot s_{ij})}{\sum_j s_{ij}}$$

user u put rating
for item j

from similarity Matr.

$$\tilde{r}_u = r_u \cdot s$$



$$\tilde{R} = R \cdot S \quad - \text{for all users}$$

$$r_{ui} = \sum_{j=KNN(i)} r_{uj} s_{ij}$$

only if j in

$$\sum_{j=KNN(i)} s_{ij}$$

to K nearest
to i

TF - IDF

Term frequency - inverse
document
frequency)

$$TF_{a,i} = \frac{N_{a,i}}{N_i} \text{ often } = 1$$

$N_{a,i}$ - appearances of attribute a in item i

N_i - number of attributes of item i

$$IDF_a = \log \frac{N_{ITEMS}}{N_a}$$

N_{ITEMS} - total number
of items

N_a - number of
items with
attribute a

Collaborative filtering (URM)

User Based similarity

$$\frac{\bar{r}_u \cdot \bar{r}_o}{|\bar{r}_u| \cdot |\bar{r}_o| + t}$$

rating of user g to item i

Cosine Similarity

$$S_{uv} = \frac{\sum_{i \in KNN(u)} r_{ui} \cdot s_{uo}}{\sqrt{\sum_{i \in KNN(u)} s_{uo}^2}}$$

where $\begin{bmatrix} - & 0,5 \\ - & \dots \\ 0,5 & \dots \\ - & \end{bmatrix}$

But when we use explicit rating we have take into account, that different users put their rating differently \Rightarrow remove bias effect.

$$\bar{r}_u = \frac{\sum_i r_{ui}}{N_u + t}$$

$$\bar{r}_o = \frac{\sum_i r_{oi}}{N_o + t}$$

from Cosine to Pearson correlation

$s_{uo} = \frac{\sum_i (r_{ui} - \bar{r}_u)(r_{oi} - \bar{r}_o)}{\sqrt{\sum_i (r_{ui} - \bar{r}_u)^2 \cdot \sum_i (r_{oi} - \bar{r}_o)^2} + t}$

How much user u like item i compared to his average opinion

Prediction

$$r_{ui} - \bar{r}_u = \frac{\sum_{\sigma \in KNN(u)} (r_{\sigma i} - \bar{r}_{\sigma}) s_{\sigma u}}{\sum_{\sigma \in KNN(u)} s_{\sigma u}} \Rightarrow$$

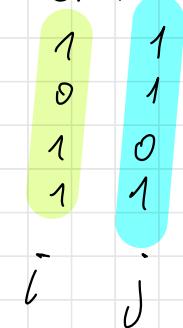
$$r_{ui} = \frac{\sum_{\sigma \in KNN(u)} (r_{\sigma i} - \bar{r}_{\sigma}) s_{\sigma u}}{\sum_{\sigma \in KNN(u)} s_{\sigma u}} + \bar{r}_u$$

you like
i more or
less

user average
opinion

Item based Collaborative filtering

Item similarity (URM)

user []

$$\begin{matrix} & \text{Item} \\ \text{user} & \left[\begin{matrix} 1 & 1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{matrix} \right] \\ & i \quad j \end{matrix}$$

$$s_{ij} = \frac{\sum_u r_{ui} \cdot r_{uj}}{\sqrt{\sum_u r_{ui}^2 \cdot \sum_u r_{uj}^2}}$$

$$r_{ui} = \frac{\sum_{j \in KNN(i)} r_{uj} \cdot s_{ij}}{\sum_{j \in KNN(i)} s_{ij}}$$

prediction

$$\hat{r}_u = r_u \cdot s$$

$$\hat{R} = R \cdot S$$

for explicit ratings

$$\bar{r}_u = \frac{\sum_i r_{ui}}{N_u + t}$$

user bias

$$s_{ij} = \frac{\sum_u (r_{ui} - \bar{r}_u)(r_{uj} - \bar{r}_u)}{\sqrt{\sum_u (r_{ui} - \bar{r}_u)^2 \cdot \sum_u (r_{uj} - \bar{r}_u)^2}}$$

same

from cosine to
adjusted cosine

Comparison User Based and Item Based

User Based

$$r_{ui} = \frac{\sum_{j \in KNN(u)} r_{uj} \cdot s_{u,j}}{\sum_{j \in KNN(u)} s_{u,j}}$$

similarity between

users: cosine similarity

Item Based

$$r_{ui} = \frac{\sum_{j \in KNN(i)} r_{uj} \cdot s_{i,j}}{\sum_{j \in KNN(i)} s_{i,j}}$$

similarity between

items: cosine

similarity

as Content based, but
without knowing attributes

For explicit ratings

$$r_{ui} = \frac{\sum_{j \in KNN(u)} (r_{uj} - \bar{r}_u) \tilde{s}_{u,j}}{\sum_{j \in KNN(u)} \tilde{s}_{u,j}}$$

Pearson correlation

coefficient (remove user bias)

$$r_{ui} = \frac{\sum_{j \in KNN(i)} (r_{uj} - \bar{r}_i) \tilde{s}_{i,j}}{\sum_{j \in KNN(i)} \tilde{s}_{i,j}}$$

Cosine adjusted
(remove user bias)

Model Based / MEMORY Based

Model : $\begin{cases} f(URM) & \text{Collaborative} \\ f(ICM) & \text{Content Based} \end{cases}$

estimating ratings $g(\text{model}, \text{user profile})$

recommendation only for users from data set
user profile $\in URM \rightarrow$ Memory Based

user profile $\notin URM \rightarrow$ Model Based

→ if come new user, we have to recompute similarities

→ can make recommendation for new users

Item based CF - MODEL BASED

doesn't change significantly II
Model Based

Item Based

$$\vec{r}_{ui} = \frac{\sum_{j \in KNN(i)} r_{uj} \cdot s_{ji}}{\sum_{j \in KNN(i)} s_{ji}}$$

S_{II} = similarity matrix

\vec{r}_u = profile of user for which we are making recommendations

- $S_{II} = f(URM)$
- estimated ratings = $g(S_{II}, \vec{r}_u)$

adding a new user doesn't change the model

item	item		
	-	0.3	0.15
0.3	-	0.6	0.43
0.15	0.6	-	0.98
0.2	0.43	0.98	-

$S_{UU} = f(URM)$
 $\text{estimated ratings} = g(S_{UU}, \vec{r}_u)$

user	-	0.3	0.15	0.2
-	0.3	-	0.6	0.43
0.3	-	0.6	-	0.98
0.15	0.6	-	0.98	-
0.2	0.43	0.98	-	-



Memory Based 

User Based

$$\tilde{r}_{ui} = \frac{\sum_{v \in KNN(u)} r_{vi} \cdot s_{vu}}{\sum_{v \in KNN(u)} s_{vu}}$$

S_{UU} = similarity matrix
 \vec{r}_u = profile of user for which we are making recommendations

 **Digital
MASTER SCHOOL** Basic Recom
ender Systems – prof. Maurizio Ferrari Dacrema  **POLITECNICO
MILANO 1863**

More significantly change model,
because we will have new column
and row \Rightarrow [Memory based]

But in Item based, if we add new item \Rightarrow we also have to recompute S_{II} , Because there will be new column.

Association Rules

$$P(\text{item } i \mid \text{history of ratings}) = P(i \mid j_1, j_2, j_3)$$

implicit rating case: $P(i \mid j) \approx \frac{\# \langle i, j \rangle}{\# \langle j \rangle} = s_{ij}$

$$P(i \mid j) \approx \frac{\# \langle i, j \rangle}{\# \langle j \rangle + z} = s_{ij}$$

By using frequency based approach, we measure how frequently happen that a user likes both item i and item j . by knowing that user liked item j .

The conditional probability, that a user will like item i , by knowing that he /she previously liked item j . is equivalent to count, how many times a user liked both item i and item j . divided by the number

of times, he or she liked item j.

this is new definition of s_{ij} . similarity.

ITEM BASED collaborative

Filtering as an Optimization

problem (SITM)

$$\hat{r}_{ui} = \sum_j r_{uj} s_{ji}$$

$$\hat{R} = RS$$

$$S \in R^{I \times J}$$

$$\begin{matrix} \hat{R} \\ \left[\begin{matrix} \hat{r}_{ui} \end{matrix} \right] \end{matrix} = h \begin{matrix} R \\ \left[\begin{matrix} \dots & r_u & \dots \end{matrix} \right] \end{matrix} \times \begin{matrix} S \\ \left[\begin{matrix} \dots & s_i & \dots \end{matrix} \right] \end{matrix}$$

Instead heuristic to compute

s_{ij} , we can learn it from

data.

$$\hat{R} = RS \quad \text{ML model} \quad s_{ij} - \text{parameter.}$$

that we will

As loss function we

will use MSE, because other aren't differentiable

SLIM - sparse Linear method for Top-N Recommender systems

$$E(S) = \sum (r_{ui} - \tilde{r}_{ui})^2$$

$u_i \in R^+$ — only existing ratings

$$\tilde{r}_{ui} = \sum_j r_{uj} s_{ij}$$

Not normalised

$$E(S) = \|R - \underbrace{RS}_R\|_2$$

$$S^* = \underset{S}{\operatorname{argmin}} \|R - RS\|_2 = \underset{S}{\operatorname{argmin}} E(S)$$

But if $S = I \Rightarrow \underline{\text{degenerate solution}}$

⇒ add regularization term weight

$$E(S) = \|R - RS\|_2 + \lambda \|S\|_2$$

$$\|S\|_2 = \sum_{ij} s_{ij}^2$$

$$\lambda \rightarrow 0 \Rightarrow S = I$$

$$\lambda \rightarrow \infty \Rightarrow S = 0$$

as sparse
as possible
↑

$$\|S\|_1 = \sum_{ij} |s_{ij}| - \lambda \text{ regularization}$$

$$E(s) = \|R - Rs\|_2 + d_2 \|S\|_2 + d_1 \|S\|_1$$

$$d_1, d_2 > 0$$

SLIM

Lasso (L1)	Ridge (L2)
Absolute value of the parameter	Squared value of the parameter
Can force the parameter to be zero	Shrinks the values of the less important parameters towards zero .
Can lead to a sparse model	Does not lead to a sparse model
Not differentiable , can be solved numerically (gradient descent)	Differentiable , can be solved analytically

$$\frac{\partial E(s)}{\partial s} = -2(R - Rs)R^T + 2d_2 s + d_1 \frac{\partial L_1}{\partial s}$$

$$\frac{\partial L_1}{\partial s_{ij}} = \begin{cases} -1, & s_{ij} < 0 \\ 0, & s_{ij} = 0 \\ 1, & s_{ij} > 0 \end{cases} = \text{sign}(s_{ij})$$

Steps of S updating

- Sample $u, i, r_{ui} \in R^+$ (existing ratings)
- Compute Gradient

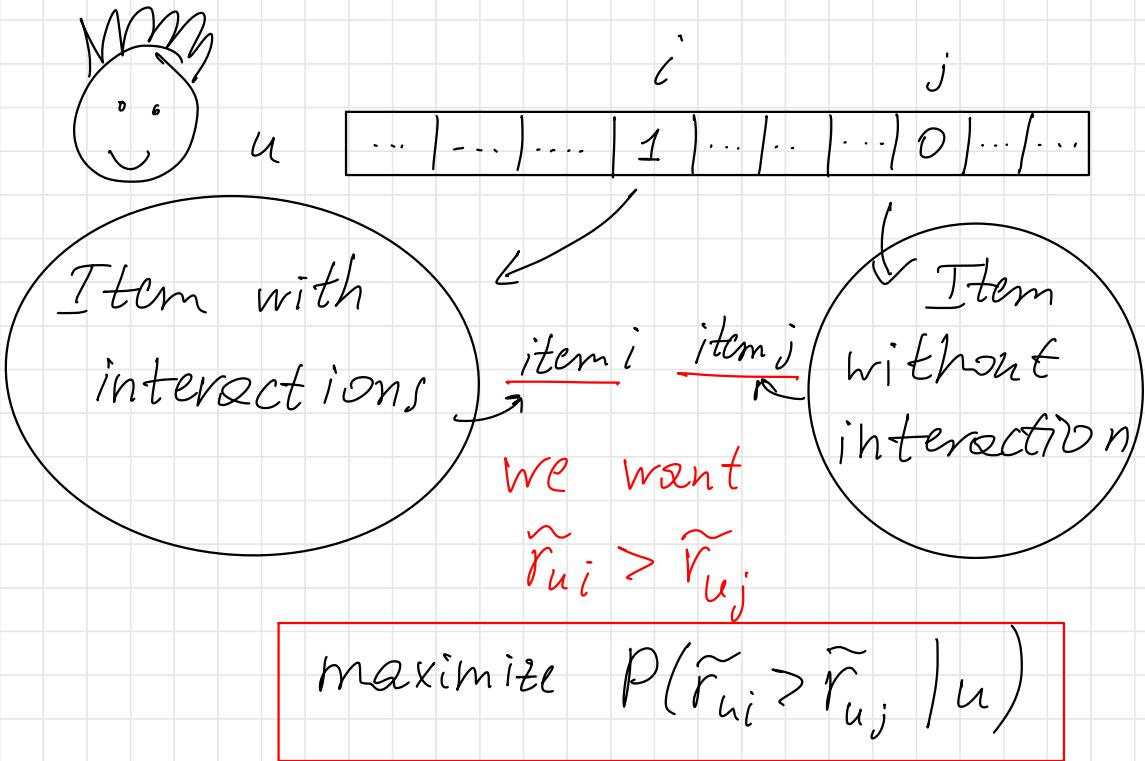
$$\frac{\partial E(s)}{\partial s_{i*}} = -2(r_{ui} - r_u s_{i*})r_u^T + 2\alpha_2 s_{i*} + \alpha_i \text{sign}(s_{i*})$$

i-th column

- Update parameters $s_{i*} = s_{i*} - \eta \cdot \frac{\partial E(s)}{\partial s_{i*}}$
learning rate
- Repeat

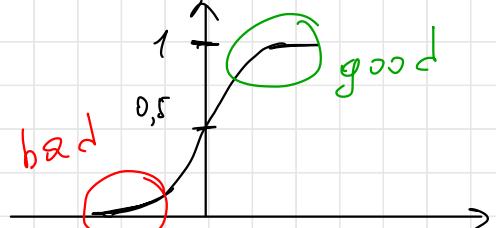
BPR Bayesian Personalized Ranking

Optimize pairwise ranking



$$x_{uij} = \tilde{r}_{ui} - \tilde{r}_{uj} \quad P(\tilde{r}_{ui} > \tilde{r}_{uj} | u) \equiv \delta(x_{uij})$$

$$\delta(x_{uij}) = \frac{1}{1 + e^{-x_{uij}}}$$



θ - model parameters

i - item relevant for user u

j - item non relevant for user u

$$P(\tilde{r}_u(\theta) > \tilde{r}_{uj}(\theta) | u) \equiv \delta(x_{uj}(\theta))$$

Assume that probability for each pair of items are independent (but they are not in real)

$$\underset{\theta}{\operatorname{argmax}} \prod_{u,i,j} P(\tilde{r}_{ui}(\theta) > \tilde{r}_{uj}(\theta) | u)$$

maximise log likelihood

$$\underset{\theta}{\operatorname{argmax}} \log \left[\prod_{u,i,j} P(\tilde{r}_{ui}(\theta) > \tilde{r}_{uj}(\theta) | u) \right] =$$

$$= \underset{u,i,j}{\operatorname{argmax}} \sum \log (P(\tilde{r}_{ui}(\theta) > \tilde{r}_{uj}(\theta) | u)) =$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_{u_{ij}} \log \delta(x_{u_{ij}}(\theta)) \iff$$

$$\underset{\theta}{\operatorname{argmin}} - \sum_{u_{ij}} \log \delta(x_{u_{ij}}(\theta)) + \lambda \|\theta\|_2$$

regularization

$$\frac{\partial}{\partial \theta} \left\{ - \sum_{u_{ij}} \log \delta(x_{u_{ij}}(\theta)) \right\}$$

$$\ln n = \frac{1}{x}$$

$$= - \sum \frac{\partial}{\partial \theta} \log (\delta(x_{u_{ij}}(\theta))) =$$

$$= - \sum \frac{1}{\delta(x_{u_{ij}}(\theta))} \delta(x_{u_{ij}}(\theta)) (1 - \delta(x_{u_{ij}}(\theta))) \frac{\partial x_{u_{ij}}}{\partial \theta}$$

$$\frac{\partial}{\partial x} \frac{1}{1 + e^{-x}} = \frac{(1 + e^{-x})^2}{(1 + e^{-x})^2} =$$

$$= \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{e^{-x}}{1 + e^{-x}} \quad \frac{1}{1 + e^{-x}} = \frac{1}{e^{-x} + 1} \cdot \delta(x) =$$

$$= \tilde{G}(-x) \tilde{G}(x) = \underbrace{(1 - G(x))}_{(1 - G(x))} G(x)$$

$$\frac{e^{-x} + 1 - 1}{e^{-x} + 1} = 1 - \frac{e^{-x}}{e^{-x} + 1} = (1 - G(x))$$

$$X_{uij} = r_u \cdot s_{i*} - r_u s_{j*} = r_u (s_{i*} - s_{j*})$$

$$\frac{\partial X_{uij}}{\partial s_{i*}} = r_u$$

$$\frac{\partial X_{uij}}{\partial s_{j*}} = -r_u$$

OPTIMIZE BPR

(Missing es
negative)



MAXIMIZE AUC
in the ROC curve

• Sample data point $u_{i,j}$

$$u_{i,j} \in R^+$$

$$u_{i,j} \in R^-$$

• Compute gradients

$$\frac{\partial E(s)}{\partial s_{i*}} = \frac{e^{-x_{uij}}}{1+e^{-x_{uij}}} r_u^T$$

• Update parameters $s_{i*} = s_{i*} - l \frac{\partial E(s)}{\partial s_i}$

• Repeat

WARP (Weighted Approximate-Rank Pairwise)

loss.

Idea: When sampling $u_{i,j}$ check if model predictions would already produce the right ranking, i.e. $\tilde{r}_{u_i} > \tilde{r}_{u_j}$. If so, draw another j . The goal is to look for training samples that violate the correct ranking

BPR exhibits a very strong popularity bias

- Popular items will be sampled a lot, their predictions will tend to be high
- Unpopular items will be sampled rarely, their predictions will tend to be low

Offline evaluation often uses test data splits that exhibit popularity bias (the same as the training data).

BPR may look good offline but perform poorly online

Matrix factorization

In CB we don't have content attributes.

Ideas: each user and item will be associated to set of latent features which are learned from the collaborative data (also called latent factors, or embeddings)

User preferences

$$\xrightarrow{\quad} \vec{x}_u = \begin{array}{ccccccccc} & & & & x_{uk} & & & & \\ | & | & | & \diagup & | & | & | & | & | \end{array}$$

k features

x_{uk} - how much user u

likes feature k.

Item description



$$\xrightarrow{\quad} \vec{y}_i = \begin{array}{ccccccccc} & & & & y_{ik} & & & & \end{array}$$

k features

y_{ik} : how much feature k is important for item

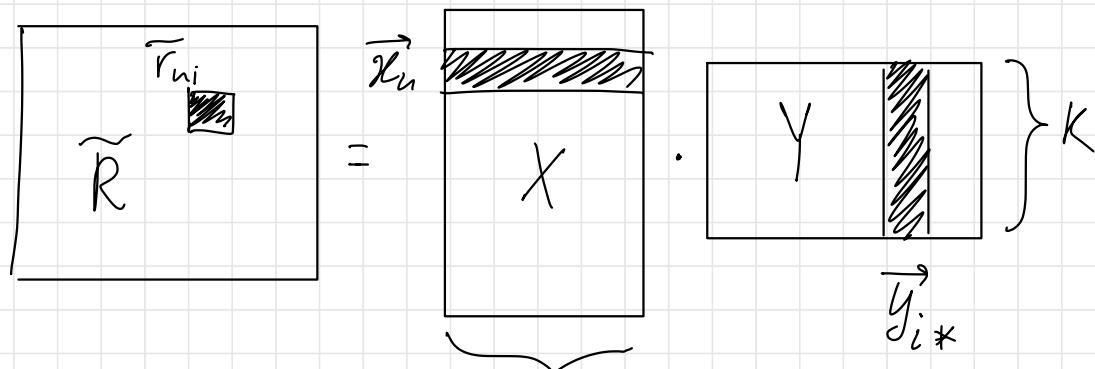
$$\tilde{r}_{ui} = \sum_k x_{uk} \cdot y_{ik} = \vec{x}_u \cdot \vec{y}_i$$

estimated rating for user u on item i

$X \in \mathbb{R}^{|\mathcal{U}| \times K}$: user-feature matrix

$Y \in \mathbb{R}^{K \times |\mathcal{I}|}$: feature-item matrix

$$\tilde{R} = X \cdot Y$$



$$X^*, Y^* = \underset{X, Y}{\operatorname{argmin}} \| R - \tilde{R} \|_2$$

(one of the possible loss function)

$$X - |\mathcal{U}| \times |K|, Y - |K| \times |\mathcal{I}|$$

in SLIM and BPR we had $|\mathcal{U}| \times |\mathcal{U}|$ params

now we have $|k| \cdot (|u| + |I|)$ -
which is much less than $|u| \cdot |I|$

K - hyper parameter

Number of Latent Factors

How many latent factors? Common values are between 10 and 200

- Large number of latent factors
- Few ratings

- Risk of overfitting
- Poor scalability of the model to large data sets

- Small number of latent factors
- Many ratings

- Personalization capabilities of the system are reduced
- Biased towards popular items
- For $K = 1$ we have Top Popular

Regularization: $X^*, Y^* = \underset{X, Y}{\operatorname{argmin}} \|R - XY\|_2 + \lambda_x \|X\|_2 + \lambda_y \|Y\|_2$

gradient: $\frac{\partial}{\partial X_u} \left((r_{ui} - X_u Y_{i*})^2 + \lambda_x \|X\|_2 \right) =$

$$= -2(r_{ui} - X_u Y_{i*}) Y_{i*}^T + 2\lambda_x X_u$$

- Sample data point (x_i, y_i, r_{hi})

- Compute gradient on X :

$$\frac{\partial E(x, y)}{\partial x_n} = -2(r_{hi} - x_n y_{i*}) y_{i*}^T + 2d_x x_n$$

- Compute gradient on Y :

$$\frac{\partial E(x, y)}{\partial y_{i*}} = -2(r_{hi} - x_n y_{i*}) x_n^T + 2d_y y_{i*}$$

- Update parameters X

$$x_n = x_n - \ell \frac{\partial E(x, y)}{\partial x_n}$$

- Update parameters Y

$$y_{i*} = y_{i*} - \ell \frac{\partial E(x, y)}{\partial y_{i*}}$$

- Repeat

Since most of the URM values are zeros, it's important to decide how to consider these values, when choosing the best feature matrices X and Y

- Missing as random. We assume that it's equally likely that a user will like or dislike an item, which hasn't been rated ($50/50 - 0/1$)

Funk SVD: minimize error function only on the non-zero elements

Missing
as random

- Missing-as-negative. It's more likely for a user to dislike an item that has not been rated yet. (optimise Precis, Recall)

Pure SVD: minimize the error function on all the matrix

Missing-as-negative

Alternating Least Squares (ALS)

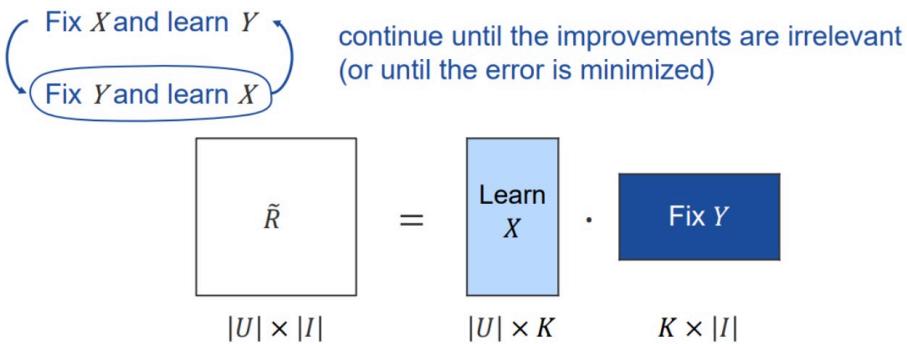
$$X^*, Y^* = \underset{X, Y}{\operatorname{argmin}} \|R - XY\|_2 + \rho_x \|X\| + \rho_y \|Y\|$$

loss - degree 4, gradient degree 3
difficult to find zeros.

ALS - reduces the loss to quadratic
and gradient to linear.

Alternating Least Squares (ALS)

At the beginning, the matrices X and Y are filled with random values

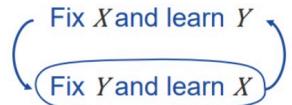


Iteratively update only one of
matrices and each iteration update
new one.

Funk SVD

The latent factors are not computed all at once.

- At the beginning $K = 1$
 X and Y are tuned until the error is minimized



$$K = 1 \quad \tilde{R} = \begin{matrix} X \\ |U| \times K \end{matrix} \cdot \begin{matrix} Y \\ K \times |I| \end{matrix}$$

first feature

Funk SVD

The latent factors are not computed all at once.

- At the beginning $K = 1$, X and Y are tuned until the error is minimized
- K is increased and the method tunes again X and Y
- K is increased until the desired value

$$K = 3 \quad \tilde{R} = \begin{matrix} X \\ |U| \times K \end{matrix} \cdot \begin{matrix} Y \\ K \times |I| \end{matrix}$$

✓ Learning on latent factor at time
is efficient and reduces overfitting

- ✗ Missing es Random Assumption
- ✗ Well for rating prediction, but poorly for Top-N (some zero ratings)
- ✗ No latent factors for new users or new items, ^{retrain} model

SVD++

In Funk SVD

$$\tilde{r}_{ui} = \sum_k x_{uk} \cdot y_{ik}$$

SVD++

$$\tilde{r}_{ui} = \underbrace{\mu + b_u + b_i}_{\text{global effects, learned}} + \sum_k x_{uk} \cdot y_{ik}$$

global effects, learned

$\mu \in \mathbb{R}$ - mean rating, overall average rating

$b_u \in \mathbb{R}^{|U|}$ - user bias, average rating &

user gave to all the rated items

$b_i \in \mathbb{R}^{|I|}$ - item bias, average rating on item received

SVD++ can be better with explicit ratings than Funk SVD, because it consider global effects

$$\begin{aligned} \mu^*, b_u^*, b_i^*, x^*, y^* &= \operatorname{argmin} \|R - \tilde{R}\|_2 + \lambda_x \|x\|_2 + \lambda_y \|y\|_2 \\ \mu^*, b_u^*, b_i^*, x^*, y^* &+ \dots \end{aligned}$$

this model has user bias ↓ item bias ↓
 $K(|U| + |I|) + |U| + |I| + 1$ parameters ↑
↑ μ

Not practical to optimize with ALS,
so it is solved with stochastic Gradient
Descent.

Not Model Based

- ✗ New user → need to recompute X
- ✗ New item → need to recompute Y

$$\tilde{R} = X \cdot Y$$

new user →

Adapting SVD++.

How to solve problem of new users and items - define the user factors as function of the items they interacted with.

$|I| \times |K|$

$$x_{uk} = \sum_j r_{uj} \cdot z_{jk} \quad z \in \mathbb{R}$$

the user preference for a feature is estimated based on:

r_{uj} : how much user u likes item j
(ask the user for some initial preferences)

z_{jk} : how much the feature k is present in the item j

approximate x_{uk} for new users.

$$x_{uk} = \sum_j r_{uj} z_{jk} = \vec{r}_u \cdot \vec{z}_k$$

$$\vec{r}_u = \boxed{\quad | 3 | 5 | 2 | \quad}$$

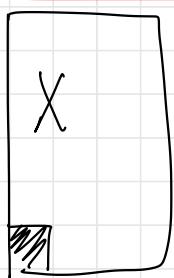
We ask the user some preferences



Movie	Titanic	Iron Man	Shrek	Finding Nemo
Liking	3	5	2	4

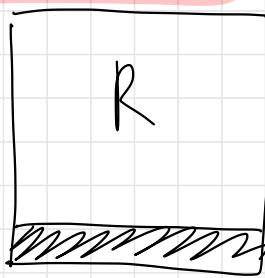
$$X = R \cdot Z$$

$$\tilde{R} = X \cdot Y = R \cdot Z \cdot Y$$

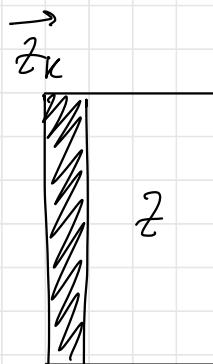


=

$$\vec{r}_u$$



.



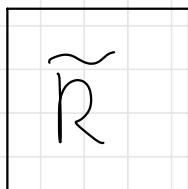
$$x_{uk}$$

$$|U| \times K$$

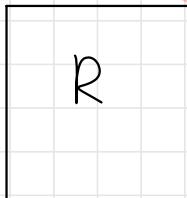
$$|U| \times |I|$$

$$|I| \times K$$

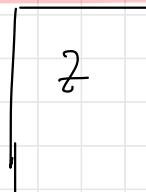
$$\tilde{r}_{ui} = \sum_k x_{uk} \cdot y_{ki} \rightarrow \tilde{r}_{ui} = \sum_k \left(\sum_j r_{uj} z_{jk} \right) y_{ki}$$



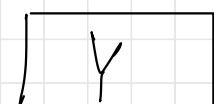
=



.



.



$$|U| \times |I|$$

$$|U| \times |I|$$

$$|I| \times K$$

$$K \times |I|$$

r_{ui} : recorded user preference for item j

z_{jk} : latent feature preference factor

y_{ki} : how much latent feature k is present in item J

$2K|I|$ parameters

Asymmetric SVD with global effects

$$\tilde{r}_{ui} = \mu + b_i + b_u + \sum_k \left(\sum_j r_{uj} z_{jk} \right) y_{ki}$$

$2K|I| + |I| + |U| + 1$ parameters

$$\boxed{\tilde{R}} = \boxed{R} \cdot \boxed{Z} \cdot \boxed{Y} =$$

$$= \boxed{R} \cdot \boxed{S} \quad \text{not similar matrix}$$

$(|I| \times |I|)$

Pare SVD

Singular Value Decomposition

$$\begin{matrix} r_{ui} \\ \hline \hline \bar{r}_u \end{matrix}$$

$$= \begin{vmatrix} | & | & | \\ \vec{u}_1 & \vec{u}_2 & \dots & \vec{u}_n \\ | & | & | \end{vmatrix} \cdot$$

$$\begin{matrix} \sigma_1 & \sigma_2 & \dots \\ \hline \hline \underbrace{\sigma_{\text{RANK}}}_{0} \end{matrix}$$

$$\begin{matrix} -\sigma_1^T- \\ -\sigma_2^T- \\ \dots \\ -\sigma_n^T- \end{matrix}$$

$m \times n$

$n \times m$

$n \times n$

$n \times m$

$$R = U \sum V^T$$

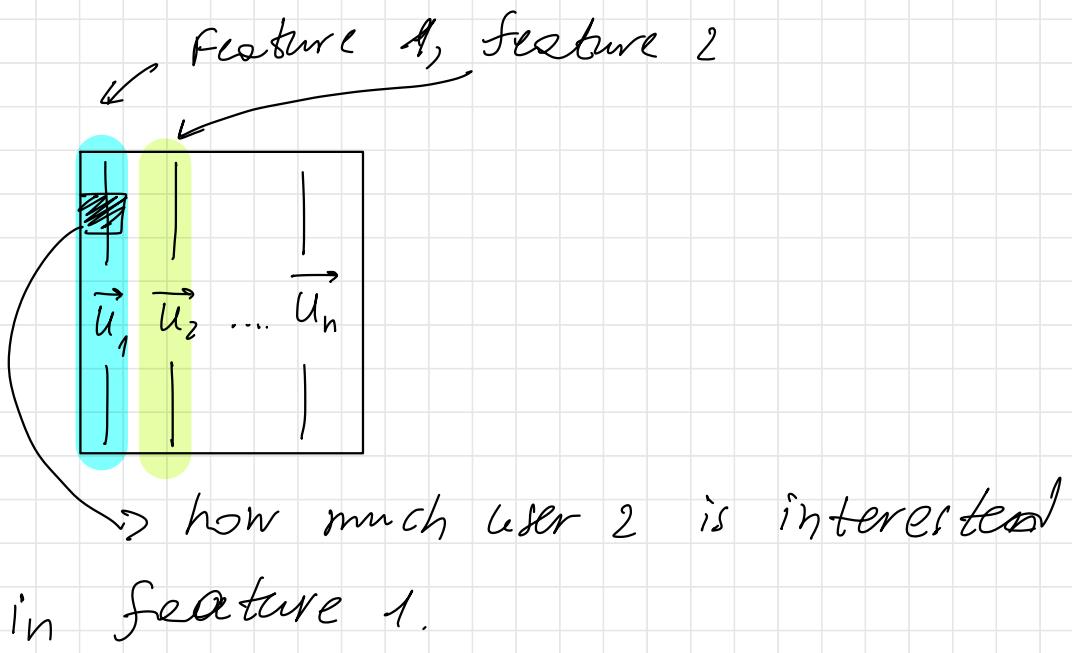
n users, m items

U : left singular vectors

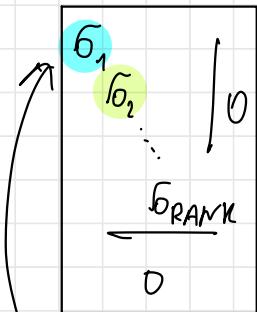
$$U^T U = I \quad U U^T = I$$

Columns

- Eigenvectors of R (left)
- Hierarchically arranged by their ability to describe the variance in columns of R .



Σ : Singular values, diagonal matrix, rank (R) elements.



Singular values of R .

- non negative
- hierarchically arranged

$$\sigma_1 \geq \sigma_2 \dots \geq \sigma_{\text{Rank}(R)} \geq 0$$

Strength of each feature in R

$$\begin{bmatrix} - \vec{v}_1^T - \\ \hline \text{feature 1} \\ \vec{v}_2^T = \\ \hline \cdots \\ \hline - \vec{v}_n^T - \end{bmatrix}$$

$m \times m$

V : Right singular vectors
Orthogonal $V^T V = V V^T = I$
Rows of V^T

- Eigen vectors of R (right)
- Hierarchically arranged by their ability to describe the variance of rows of R .

how much feature 2 describe item 1

Truncated SVD

$$R \cdot U \cdot \Sigma \cdot V^T \approx U_K \Sigma_K V_K^T$$

first K highest singular values

$$U_K \Sigma_K V_K^T =$$

$$\begin{array}{c|c|c|c} & & & \\ \vec{u}_1 & \vec{u}_2 & & \diagup \\ \hline & & & \end{array}$$

$n \times K$

$$\cdot \quad K \quad \left\{ \begin{array}{c|c} \sigma_1 & \\ \hline & \sigma_2 \\ \hline & \end{array} \right. \quad \left. \begin{array}{c|c} & \\ \hline & \end{array} \right\} \quad \cdot$$

$K \times K$

$$\begin{array}{c|c} \vec{v}_1^T & \\ \hline & \vec{v}_2^T \\ \hline & \diagup \\ \hline & \end{array}$$

$K \times m$

Frobenius norm:

$$\|A\|_F = \sqrt{\sum_{i,j} \sigma_{ij}^2}$$

all values

$$\|A\|_F^2 = \sum_k \sigma_k^2(A)$$

sum
of
squares

of singular
values

The best low-dimensional approximation of matrix A by a rank K under the Frobenius norm is given by truncated SVD

The approximation error is given by the $K + 1$ singular value

$$\|A - A_K\|_F = \|A - U_K \Sigma_K V_K^T\|_F = \sigma_{K+1}(A)$$

if we apply SVD decomposition on R
we obtain user and item factors

But, this will not work for recommendation

The full rank SVD decomposition is exact,
we will perfectly reconstruct R, missing
interactions included \Rightarrow Overfitting.

Pure SVD uses a truncated SVD decomposition with K components:

- Approximate error of R allows the model to generalize
- Contains the main features and reduces overfitting

$U_K \in \mathbb{R}^{|U| \times K}$: user-feature matrix

$\Sigma_K \in \mathbb{R}^K$: singular values

$V_K^T \in \mathbb{R}^{K \times |I|}$: feature-item matrix

$$R \approx \tilde{R} = U_K \Sigma_K V_K^T$$

$$\text{if } P = U \Sigma$$

$$R = U \Sigma V^T \quad | \cdot V$$

$$VV^T = I \quad \downarrow$$
$$V^T = V^{-1}$$

$$RV = U\Sigma = P$$

$$R = U \Sigma V^T = P V^T = RVV^T = R \circledcirc$$

similarity

Pure SVD equivalent to an item based CF with $S = V \cdot V^T$

By using truncated SVD, PureSVD optimizes the Frobenius norm

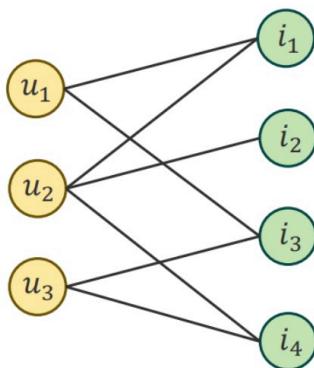
The Frobenius norm accounts for all values of R , including the missing ratings which are treated as zeros

PureSVD relies on the Missing-as-Negative Assumption

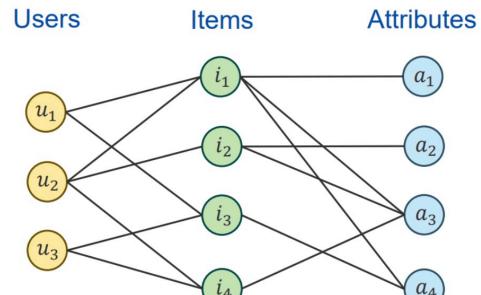
Graph Based

Bipartite User-Item Graph

Users Items



Tripartite User-Item-Attributes Graph



A graph is defined as $G = (N, E)$

N the set of nodes

E the set of edges

Collaborative Model

$$N = U \cup I$$

$$E = \{u, i \mid u, i \in R^+\}$$

Hybrid Content-Collaborative Model

$$N = U \cup I \cup A$$

$$E = \{u, i \mid u, i \in R^+\} \cup \{i, a \mid i, a \in C^+\}$$

Adjacency Matrix Notation

Adjacency Matrix $G \in \mathbb{R}^{|N| \times |N|}$, $g_{xy} = \begin{cases} 1, & \text{node } x \text{ is connected to node } y \\ 0, & \text{otherwise} \end{cases}$

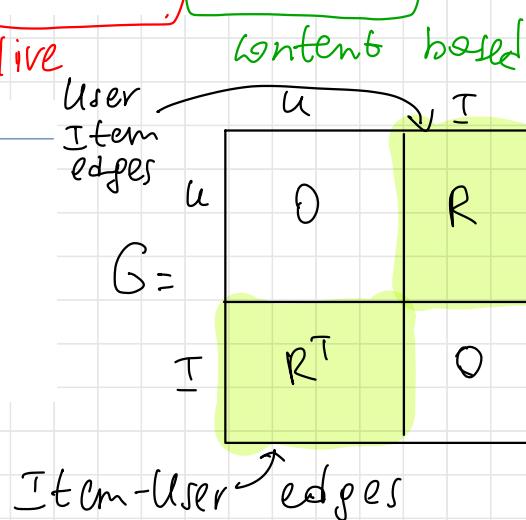
- **Undirected Graph**
 G is symmetric, $g_{xy} = g_{yx}$



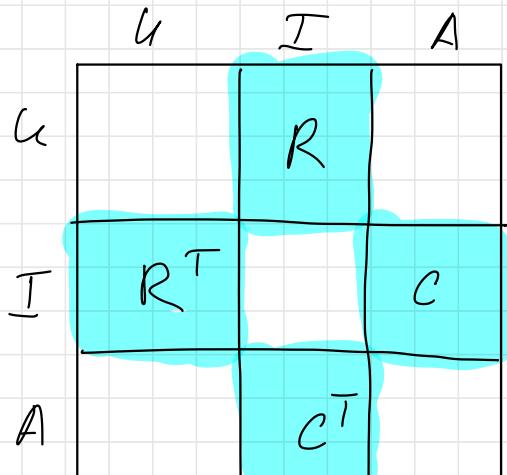
- **Directed Graph**
 G is asymmetric, $g_{xy} = 1$, $g_{yx} = 0$



$$G = \begin{bmatrix} 0 & R \\ R^T & 0 \end{bmatrix}$$



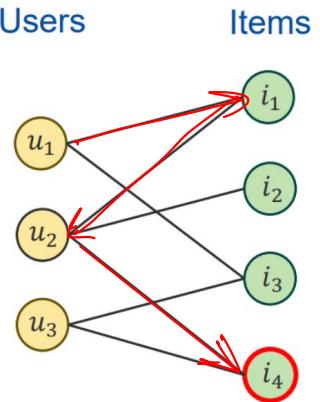
$$G = \begin{pmatrix} 0 & R & 0 \\ R^T & 0 & C \\ 0 & C^T & 0 \end{pmatrix}$$



Random walk

$u_1 \rightarrow i_1 \rightarrow u_2 \rightarrow i_4$

Recommendation
(at 3 hops)
(steps)



We will compute probability to transition between two nodes and a probability distribution on being in each node

We need two things:

- A **transition probability** matrix, $P \in \mathbb{R}^{|N| \times |N|}$, $\sum_y p_{xy} = 1, \forall x$
- A **state probability** vector (row), $\Pi \in \mathbb{R}^{|N|}$, $\sum_i \Pi_i = 1$

Given $D \in \mathbb{R}^{|N|}$, d_x the number of outgoing edges from node x (its **degree**):

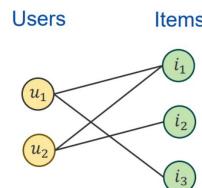
$$d_x = \sum_{y \in N} g_{xy}$$

$$p_{xy} = \frac{g_{xy}}{d_x}$$

G is symmetric!

$$G = \begin{bmatrix} 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$D = \begin{bmatrix} 2 & 2 & 2 & 1 & 1 \end{bmatrix}$$



$$P = \begin{bmatrix} 0 & 0 & 1/2 & 0 & 1/2 \\ 0 & 0 & 1/2 & 1/2 & 0 \\ 1/2 & 1/2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Not symmetric

State Probability

The probability of being in state i can be computed, given:

- The state we were in at the previous hop (0)
- The probability to jump from it to i

probability of being on state 1 and go to state i

$$\Pi_i^{(1)} = \Pi_1^{(0)} p_{1i} + \Pi_2^{(0)} p_{2i} + \dots = \sum_{j \in N} \Pi_j^{(0)} p_{ji}$$

$$\Pi^{(h)} = \boxed{\Pi^{(h-1)} \cdot P}$$

probability being in state 2 and go to state i

Repeat as needed

loop

$$\Pi^{(0)} = [1 \ 0 \ 0 \ 0 \ 0]$$

$$\Pi^{(1)} = \Pi^{(0)} \cdot P = \begin{bmatrix} 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \end{bmatrix}$$

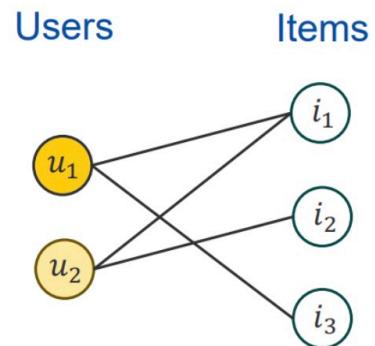
$$\Pi^{(2)} = \Pi^{(1)} \cdot P = \begin{bmatrix} \frac{3}{4} & \frac{1}{4} & 0 & 0 & 0 \end{bmatrix}$$

$$\Pi^{(3)} = \Pi^{(2)} \cdot P = \begin{bmatrix} 0 & 0 & \frac{4}{8} & \frac{1}{8} & \frac{3}{8} \end{bmatrix}$$

$$\Pi^{(4)} = \Pi^{(3)} \cdot P = \begin{bmatrix} \frac{10}{16} & \frac{6}{16} & 0 & 0 & 0 \end{bmatrix}$$

After few steps we can use this scores and rank them

$$\tilde{r}_{ui} = \Pi_{l=1+i}^{(3)}$$



$$P = \begin{bmatrix} 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 0 & \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Recommendations i_1, i_2, i_3

Number of hops h usually 3, 5, 7

- if h is high I travel too far from target user
- Over smoothing - tend to be popularity-based

State Probability: Steady State

Every graph has a state Π such that:

$$\Pi = \Pi \cdot P$$

The state Π is called a **steady state** and is an **eigenvector of P**

If the random walk continues for $t \rightarrow \infty$

- We will reach the steady state regardless for the starting state
- For n -partite graphs we will not, we will find a *periodic* sequence of states

Page-Rank: Random Walk with Restart

Infinite random walk with probability $1 - \gamma$ to restart from the initial state (user)

$$\Pi = \boxed{\gamma \Pi \cdot P} + \boxed{(1 - \gamma) \Pi^{(0)}}$$

Random walk ↗ ↙ Restart

The choice of γ allows to control how **personalized** it is

Π is found by solving a linear system of equations, but is expensive

$\gamma \rightarrow 1 \Rightarrow$ steady state

$\gamma \rightarrow 0 \Rightarrow$ we don't moving and don't predict

Computing P^3 is expensive (and dense) but:

- We only need a specific meta-path: $U \rightarrow I \rightarrow U \rightarrow I$
- P is a block matrix, we can compute only the portion we need

$$P = \begin{bmatrix} 0 & P_{UI} \\ P_{IU} & 0 \end{bmatrix} \quad P_{UI} = \text{diag} \left[\frac{1}{D_U} \right] \cdot R$$
$$P_{IU} = \text{diag} \left[\frac{1}{D_I} \right] \cdot R^T$$

$$P = \begin{bmatrix} 0 & P_{UI} \\ P_{IU} & 0 \end{bmatrix}$$

$U \rightarrow I \rightarrow U, \epsilon [0, 1]^{l_U \times l_U}$

$$P^2 = \begin{bmatrix} P_{hI} \cdot P_{IU} & 0 \\ 0 & P_{IU} \cdot P_{hI} \end{bmatrix}$$

$I \rightarrow U \rightarrow I \in [0, 1]^{l_I \times l_I}$

$$P^3 = \begin{bmatrix} 0 & P_{hI} & P_{IU} \\ P_{IU} & P_{hI} & P_{IU} \\ 0 & 0 & 0 \end{bmatrix}$$

$I \rightarrow U \rightarrow I \rightarrow U$

$U \rightarrow I \rightarrow U \rightarrow I$
 we need only this
 to predict

The meta-path we are interested in is $U \rightarrow I \rightarrow U \rightarrow I$

In the following Π_I is the state probability for item nodes, Π_U for user nodes

$$\Pi^{(3)} = \Pi^{(0)} \cdot P^3 = \Pi^{(0)} \cdot \begin{bmatrix} 0 & P_{UI} \cdot P_{IU} \cdot P_{UI} \\ P_{IU} \cdot P_{UT} \cdot P_{IU} & 0 \end{bmatrix}$$

$$\Pi_I^{(3)} = \boxed{\Pi_U^{(0)} \cdot P_{UI}} \cdot \boxed{P_{IU} \cdot P_{UI}}$$

This is an item-item similarity

This is the (starting) user profile,
 $\Pi_U^{(0)}$ acts as a one-hot encoding

The meta-path we are interested in is $U \rightarrow I \rightarrow U \rightarrow I$

In the following Π_I is the state probability for item nodes, Π_U for user nodes

$$\Pi^{(3)} = \Pi^{(0)} \cdot P^3 = \Pi^{(0)} \cdot \begin{bmatrix} 0 & P_{UI} \cdot P_{IU} \cdot P_{UI} \\ P_{IU} \cdot P_{UI} \cdot P_{IU} & 0 \end{bmatrix}$$

$$\Pi_I^{(3)} = \Pi_U^{(0)} \cdot P_{UI} \cdot P_{IU} \cdot P_{UI} = \Pi_U^{(0)} \cdot P_{UI} \cdot S$$

$$S_{ij} = [P_{IU} \cdot P_{UI}]_{ij} = \sum_{u \in U} \frac{r_{ui} r_{uj}}{d_i d_j}$$

similarity matrix
Item-Item



For implicit interactions P^3 is the Cosine similarity without shrinkage

In P_α^3 the transition probabilities are elevated to a power α

$$\left. \begin{array}{l} P_{UI} = \left(\text{diag} \left[\frac{1}{D_U} \right] \cdot R \right)^\alpha \\ P_{IU} = \left(\text{diag} \left[\frac{1}{D_I} \right] \cdot R^T \right)^\alpha \end{array} \right\} S_{ij} = \sum_{u \in U} \left(\frac{r_{ui} r_{uj}}{d_i d_j} \right)^\alpha$$

If $\alpha > 0$ small transition probabilities and similarities are attenuated

- P_α^3 tends to exhibit a strong popularity bias

$$P_{UI} = \left(\text{diag} \left[\frac{1}{D_U} \right] \cdot R \right)$$

$$P_{IU} = \left(\text{diag} \left[\frac{1}{D_I} \right] R^T \right)$$

RP_β^3 tries to reduce the popularity bias

- **Problem:** it is (too) easy for a random walk to end up in items with high degree
- **Idea:** penalize the similarity of items with a high degree. Divide the similarity by the popularity of the destination item elevated to a power β

$$S_{ij} = \frac{1}{d_j^\beta} \sum_{u \in U} \left(\frac{r_{ui} r_{uj}}{d_i d_j} \right)^\alpha$$

Adding Side Information

However, if we add both item attributes A_I and user attributes A_U we encounter problematic meta-paths:

- $U \rightarrow I \rightarrow U \rightarrow I$
- $U \rightarrow I \rightarrow A_I \rightarrow I$
- $U \rightarrow A_U \rightarrow U \rightarrow I$

New (desired) meta-path

- $U \rightarrow I \rightarrow U \rightarrow A_U$
- $U \rightarrow A_U \rightarrow U \rightarrow A_U$

These meta-paths do not end in the right node!

Deep Learning recommender system

Basic Layer

In a typical fully-connected layer, one computes a weighted average of the input, adds a bias and then applies a non-linear activation function f

$$x_o = f_o(x_i \cdot W_o + b_o)$$

n : Input size

$x_i \in \mathbb{R}^n$: Input data

$W \in \mathbb{R}^{n \times m}$: Weight matrix

m : Output size

$x_o \in \mathbb{R}^m$: Output data

$b \in \mathbb{R}^m$: Bias vector

The parameters



Binary Cross Entropy

If we use implicit ratings and the model computes probabilities:

$$\operatorname{argmin}_{\theta} -\frac{1}{N} \sum_{i=1}^N [r_{ui} \log \tilde{p}_{ui}(\theta) + (1 - r_{ui}) \log(1 - \tilde{p}_{ui}(\theta))]$$

r_{ui} : Ground truth label (1/0)

$\tilde{p}_{ui}(\theta)$: Estimated probability that the label is 1, or $\tilde{p}_{ui}(\theta) = \sigma(\tilde{r}_{ui}(\theta))$

θ : Parameters of the model

N : Number of samples the loss is computed on

This loss function can be applied to most of the recommender algorithms we discussed in the previous lectures (SLIM, MF...)

How do we sample the data?

- If we use only positive interactions, $u, i \in R^+$, we do not learn to rank
- If we use all ground truth, there are too many negatives! **Need to subsample**, for example draw a sample among positives or negatives with 50% probability

Autoencoders

Autoencoders for Recommendation

- The **input data** is the full user profile r_u
- The **embedding** will be a compressed representation of the user interests, (hopefully) sufficient to reconstruct their interactions
- The **reconstructed input** contains the predictions of the model for all items, which we rank and use to decide what to recommend

r_u — input

$e_u = g_e(r_u)$ — embedding after encoder

$\tilde{r}_u = g_d(e_u)$ — reconstructed input from decoder.

Encoder, Decoder will both reconstruct ratings for all users \Rightarrow collaborative filtering (Model Based)

Autoencoders as Item-Item Similarity Models

Consider a shallow autoencoder with no hidden layers and embedding size K

If $f = I$, $b_e, b_d = \bar{0}$ and share the parameters of encoder and decoder $W_d = W_e^T$:

$$e_u = f_e(r_u \cdot W_e + b_e) = r_u \cdot W_e$$

$$\tilde{r}_u = f_d(e_u \cdot W_d + b_d) = e_u \cdot W_d = r_u \cdot W_e \cdot W_d = r_u \cdot W_e \cdot W_e^T$$

$$W_e \in \mathbb{R}^{|I| \times K}, W_d \in \mathbb{R}^{K \times |I|}$$

Factorized **symmetric**
item-item similarity

EASE^R

EASE^R is an item-based similarity collaborative filtering model

$$S^* = \underset{S}{\operatorname{argmin}} \|R - RS\|_F + \lambda \|S\|_F$$

$$\text{s.t. } \operatorname{diag}(S) = 0$$

1. To find the optimal S we set the gradient of the loss to zero (it is a function of $\vec{\gamma}$)

$$S^*(\vec{\gamma}) = (R^T \cdot R + \lambda I_{|I|})^{-1} \cdot (R^T \cdot R - \operatorname{diagMat}(\vec{\gamma}))$$

2. To find the multipliers that satisfy the constraints, we solve $\operatorname{diag}(S(\vec{\gamma})) = 0$

$$P = (R^T \cdot R + \lambda I_{|I|})^{-1}$$

$$S^*(\vec{\gamma}) = P \cdot (P^{-1} - \lambda I_{|I|} - \operatorname{diagMat}(\vec{\gamma})) = I_{|I|} - P \cdot (\lambda I_{|I|} + \operatorname{diagMat}(\vec{\gamma}))$$

$$\vec{\gamma} = \vec{1} \oslash \operatorname{diag}(P)$$

3. The optimal S that will also satisfy the constraints is:

$$S^* = I_{|I|} - P \cdot \operatorname{diagMat}(\vec{1} \oslash \operatorname{diag}(P))$$

Element-wise division

- EASE^R is fast and highly effective
- Computing $R^T \cdot R$ and the inverse required for P is very memory intensive

Why does EASE^R refer to autoencoders?

The loss function aims to minimize the Frobenius norm of $R - RS$, hence a perfect solution would be an S such that $R = RS$

To quote the paper: "*This is done by reproducing the input as its output, as is typical for autoencoders.*" despite EASE^R not having their typical architecture

Denoising Autoencoders

One of the issues when training autoencoders is that the input space (user profiles) is very sparse

- **Risk:** the encoder might create poor embeddings for new user profiles
- **Risk:** the decoder may not know how to reconstruct correctly portions of the embedding space

Denoising Autoencoders

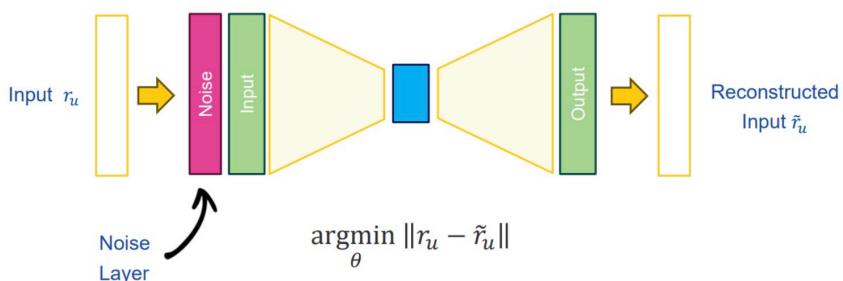
Idea: train the autoencoder to reconstruct the correct user profile from a noisy one

Types of noise Salt & Pepper:

- Randomly remove (dropout) a certain portion of the positive interactions
- Randomly add a (very) small number of positive interactions

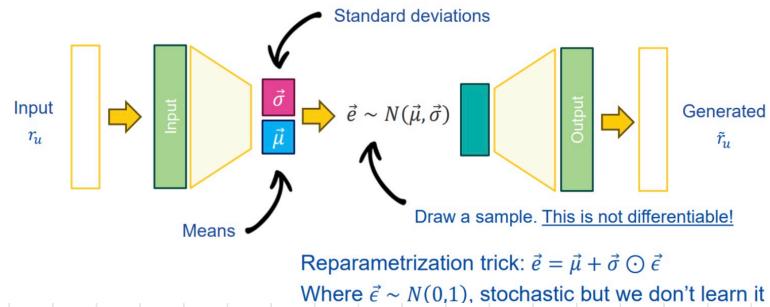
Both encoder and decoder will see many more user profiles

Denoising Autoencoders



VAE

Variational Autoencoders



Variational Autoencoders

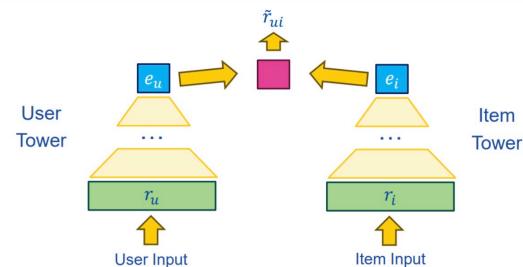
$$P(\theta) = \int P(\theta|z)P(z) dz$$

Annotations explain the equation:

- "Ground Truth" points to $P(\theta|z)$.
- "We choose it to be Normal distribution $\sim N(0, I)$ " points to $P(z)$.
- "We learn it!" and "Approximated with a deep learning model" point to the integral symbol.

TWO - TOWER MODEL

Two-Tower Model



Consider a two-tower model with no hidden layers, embedding size K and both user and item input one-hot encoded x_u, x_i .

$$e_u = f^U(x_u \cdot W^U + b^U)$$

$$e_i = f^I(x_i \cdot W^I + b^I)$$

$$\left. \right\} \quad \tilde{r}_{ui} = e_u \cdot e_i$$

$$x_u \in \mathbb{R}^{|U|}, x_i \in \mathbb{R}^{|I|}$$

$$W^U \in \mathbb{R}^{|U| \times K}, W^I \in \mathbb{R}^{|I| \times K}$$

Two-Tower Model as... Matrix Factorization

Consider a two-tower model with no hidden layers, embedding size K and both user and item input one-hot encoded x_u, x_i . If $f = I$ and $b^U, b^I = \bar{0}$:

$$\begin{aligned} e_u &= f^U(x_u \cdot W^U + b^U) = x_u \cdot W^U = W_u^U \\ e_i &= f^I(x_i \cdot W^I + b^I) = x_i \cdot W^I = W_i^I \end{aligned}$$

$$\left. \begin{aligned} e_u &= f^U(x_u \cdot W^U + b^U) = x_u \cdot W^U = W_u^U \\ e_i &= f^I(x_i \cdot W^I + b^I) = x_i \cdot W^I = W_i^I \end{aligned} \right\} \tilde{r}_{ui} = e_u \cdot e_i = W_u^U \cdot W_i^I$$

User and Item
latent factors

$$x_u \in \mathbb{R}^{|U|}, x_i \in \mathbb{R}^{|I|}$$

$$W^U \in \mathbb{R}^{|U| \times K}, W^I \in \mathbb{R}^{|I| \times K}$$

Two-Tower Model as... Matrix Factorization

A two-tower model with one-hot encoded input is **memory based**

We cannot integrate new users or new interactions without retraining it

The issue can be overcome if one uses other encodings, for example the entire user profile. The model becomes equivalent to Asymmetric SVD

Factorization machines

Input Data Representation

Features									Target
Users			Items			R			Rating
1	0	0	0	1	0	0	0	5	Rating
0	1	0	0	0	1	0	0	3	
0	0	1	0	0	0	0	0	1	2
0	1	0	0	0	0	0	1	0	4
...									

Input Data Representation: Example

Different format, same information as the URM:

Users			Items			R
0	0	1	0	1	0	0
0	1	0	0	0	1	0
1	0	0	0	0	0	1
...						
User 1			Item 4			r_{ui}

$i = 4$

$u = 1$ $i = 4$ $r_{ui} = 2$

Diagram illustrating the mapping from the matrix representation to the row representation. A 4x9 matrix is shown with a highlighted entry at row 1, column 4. This corresponds to the row representation where the first row contains values 0, 0, 1, 0, 1, 0, 0, 0, 5. The index $i = 4$ is shown above the matrix, and the rating $r_{ui} = 2$ is shown next to the highlighted entry.

The Model to Estimate Ratings

$$\tilde{r}^{(k)} = \omega_0 + \sum_{i=1}^n \omega_i x_i^{(k)} + \sum_{i=1}^n \sum_{j=i+1}^n \omega_{ij} x_i^{(k)} x_j^{(k)}$$

index of row (user-item interaction)

i, j : Index of the columns

Users			Items			R
1	0	0	0	1	0	0
0	1	0	0	0	1	0
0	0	1	0	0	0	1
...						

w_0, w_i, w_{ij} - to learn

$\omega_0 \in \mathbb{R}$ is the intercept of the model function

- it represents the global bias
- if we estimate only ω_0 , we obtain the average of all the ratings

$$\tilde{r}^{(k)} = \omega_0 + \sum_{i=1}^n \omega_i x_i^{(k)} + \sum_{i=1}^n \sum_{j=i+1}^n \omega_{ij} x_i^{(k)} x_j^{(k)}$$

n : number of columns/features ($|U| + |I|$)

$\omega_i \in \mathbb{R}$: models the strength of the i -th input feature

- This component is the linear combination of the variables

Given a row k , the equation can be:

$$\tilde{r}^{(k)} = \omega_0 + \omega_u + \omega_i$$

$$\tilde{r}_{u,i} = \mu + b_u + b_i \quad \leftarrow \text{equation to Global Effects}$$

Quadratic Term

$$\tilde{r}^{(k)} = \omega_0 + \sum_{i=1}^n \omega_i x_i^{(k)} + \sum_{i=1}^n \sum_{j=i+1}^n \omega_{ij} x_i^{(k)} x_j^{(k)}$$

Notice this

ω_{ij} models the interaction between i -th and j -th features

- Assumption: matrix is symmetric $\omega_{ij} = \omega_{ji}$
- In theory one could add third degree terms, forth and so on...

$$\underset{\omega}{\operatorname{argmin}} E(\omega) = \underset{\omega}{\operatorname{argmin}} \|r^{(k)}, \tilde{r}^{(k)}\|$$

$$\tilde{r}^{(k)} = \omega_0 + \boldsymbol{\omega} \cdot \mathbf{x}^{(k)} + \mathbf{x}^{(k)T} \cdot \mathbf{W} \cdot \mathbf{x}^{(k)}$$

$$\omega_0 \in \mathbb{R}$$

$$\boldsymbol{\omega} \in \mathbb{R}^n$$

$$\mathbf{W} \in \mathbb{R}^{n \times n}$$

Number of parameters: $1 + n + \frac{n^2 - n}{2}$

A lot!!!

Factorizing W

$$\boxed{W} = \boxed{v}^f \cdot \boxed{v^T}^r$$

$$\omega_{ij} = \mathbf{v}_i \cdot \mathbf{v}_j = \sum_{h=0}^f v_{ih} v_{jh}$$

$$W \in \mathbb{R}^{n \times n}$$

$$V \in \mathbb{R}^{n \times f}$$

where $f \ll n$

f = #latent factors

Final Factorization Machine Model

$$\tilde{r}^{(k)} = \omega_0 + \sum_{i=1}^n \omega_i x_i^{(k)} + \sum_{i=1}^n \sum_{j=i+1}^n \omega_{ij} x_i^{(k)} x_j^{(k)}$$

$$\omega_{ij} = \mathbf{v}_i \cdot \mathbf{v}_j = \sum_{h=0}^f v_{ih} v_{jh}$$

$$\tilde{r}^{(k)} = \omega_0 + \sum_{i=1}^n \omega_i x_i^{(k)} + \sum_{i=1}^n \sum_{j=i+1}^n \sum_{h=0}^f v_{ih} v_{jh} x_i^{(k)} x_j^{(k)}$$

Final Number of Parameters

$$\tilde{r}^{(k)} = \omega_0 + \sum_{i=1}^n \omega_i x_i^{(k)} + \sum_{i=1}^n \sum_{j=i+1}^n \sum_{h=0}^f v_{ih} v_{jh} x_i^{(k)} x_j^{(k)}$$

$$\tilde{r}^{(k)} = \omega_0 + \omega \cdot x^{(k)} + x^{(k)T} \cdot V \cdot V^T \cdot x^{(k)}$$

Number of parameters: $1 + n + nf$

Much more manageable!

Equivalence with SVD++

$$\tilde{r}^{(k)} = \omega_0 + \omega_i + \omega_u + V_i \cdot V_u^T$$

If we use Factorization Machines with only collaborative data, we obtain a model which is equivalent to SVD++

We need to do more!

