

# Collaborative Filtering and Matrix Factorization

# About me

- Evgeny Frolov  
[evgeny.frolov@skoltech.ru](mailto:evgeny.frolov@skoltech.ru)
- Received PhD in SkolTech last week; you can find materials at  
<https://www.skoltech.ru/en/2018/09/phd-thesis-defense-evgeny-frolov/>
- Working in the group of prof. Ivan Oseledets + Sberbank AI Laboratory
- Previously graduated from MSU, Physics Department
- Approx. 5 years experience in IT industry

# What is a recommender system?



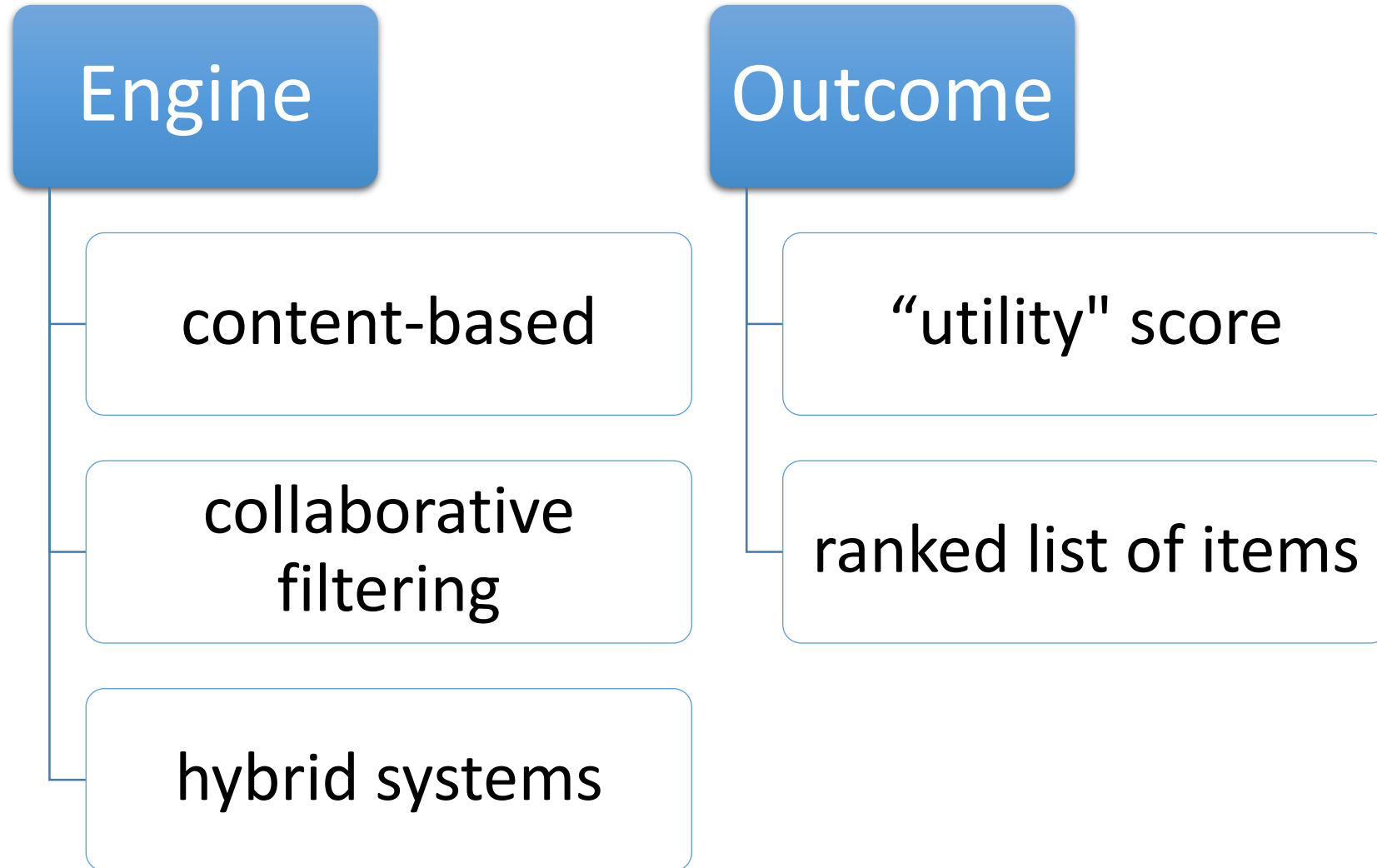
## Examples:

- Amazon
- Netflix
- Pandora
- Spotify
- etc.

**Many different areas:** e-commerce, news, tourism, entertainment, education...

Goal: predict user preferences given some prior information on user behavior.

# Recommender systems internals



# General workflow

Goal: predict user preferences based on prior user feedback and collective user behavior.

collect data

			
	?	?	3
	5	5	?
	4.5	?	4

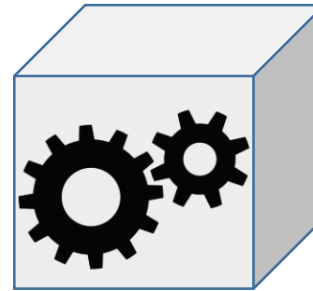
user-movie matrix  $R^{M \times N}$

$r_{ij}$  is a rating of  $i^{th}$  user for  $j^{th}$  movie

? - missing (unknown) values

build model

$$f_U: User \times Item \rightarrow Rating$$



$f_U$  - utility function

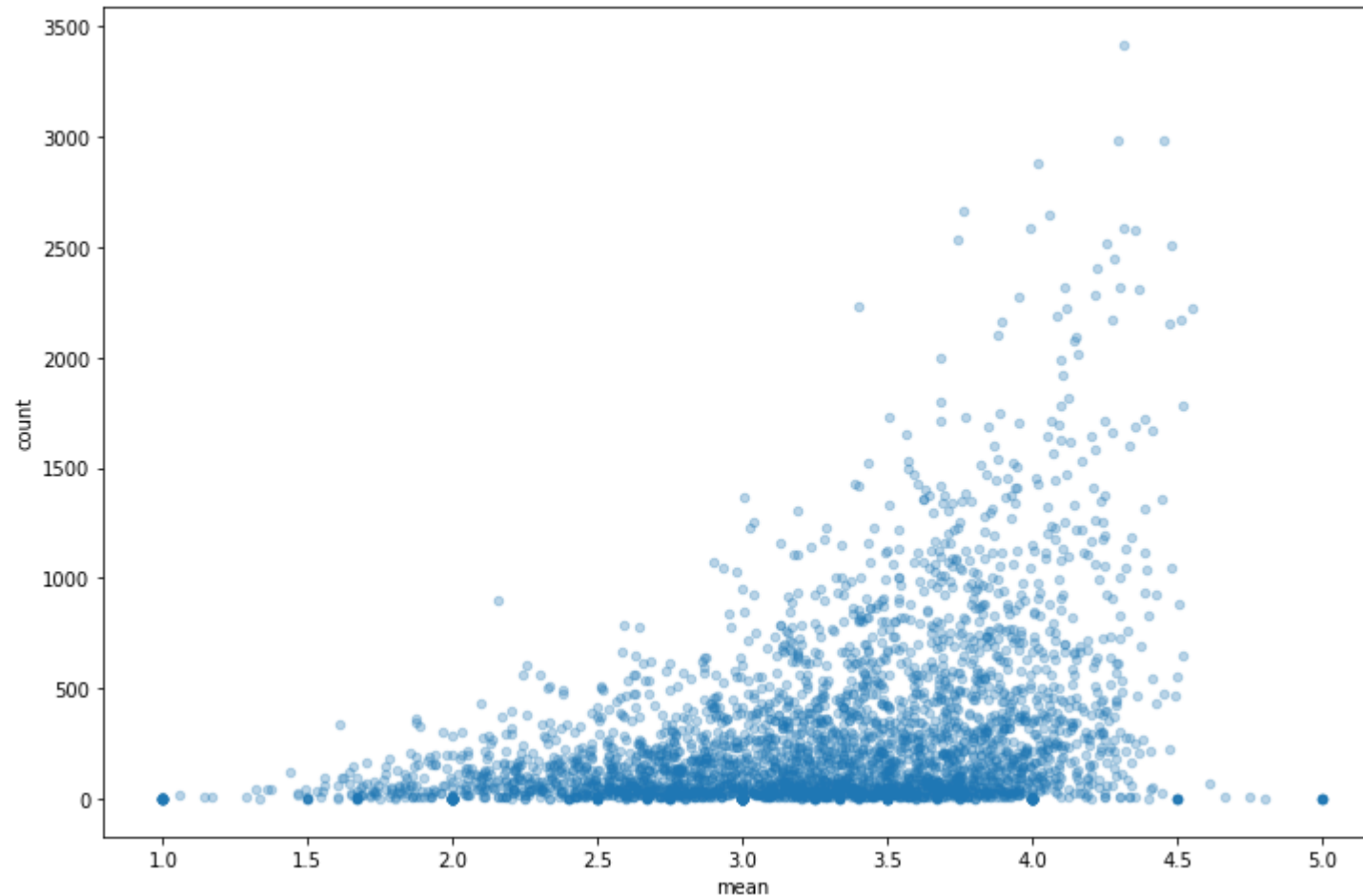
generate recommendations



unknown user

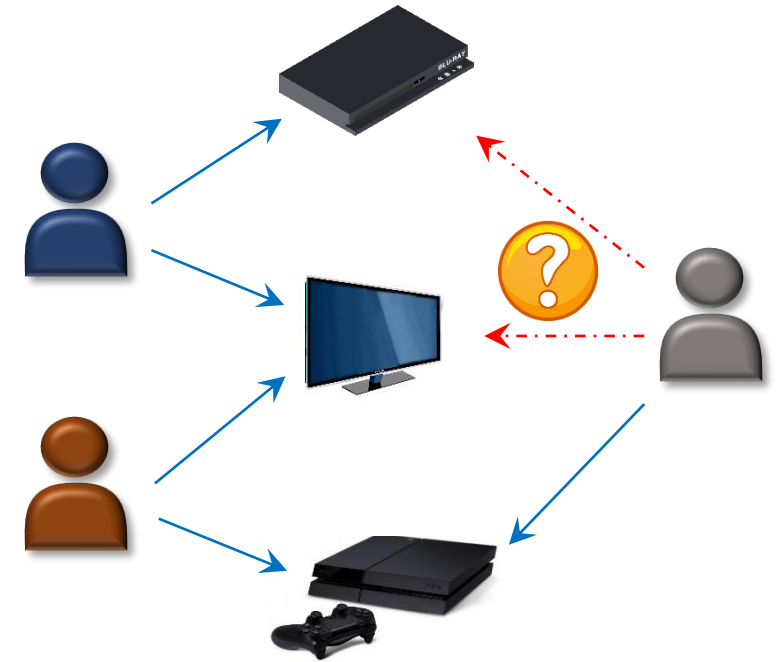
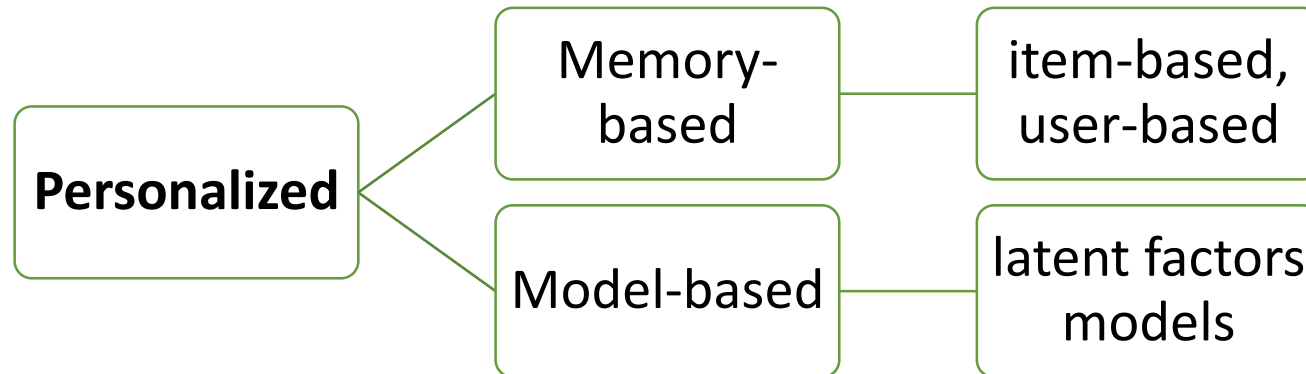
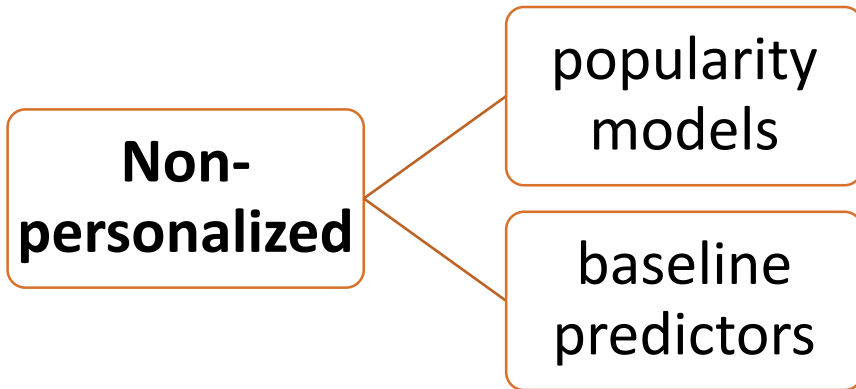
# Popularity-based model w/ averaged ratings

What would be a necessary modification for this rating pattern?

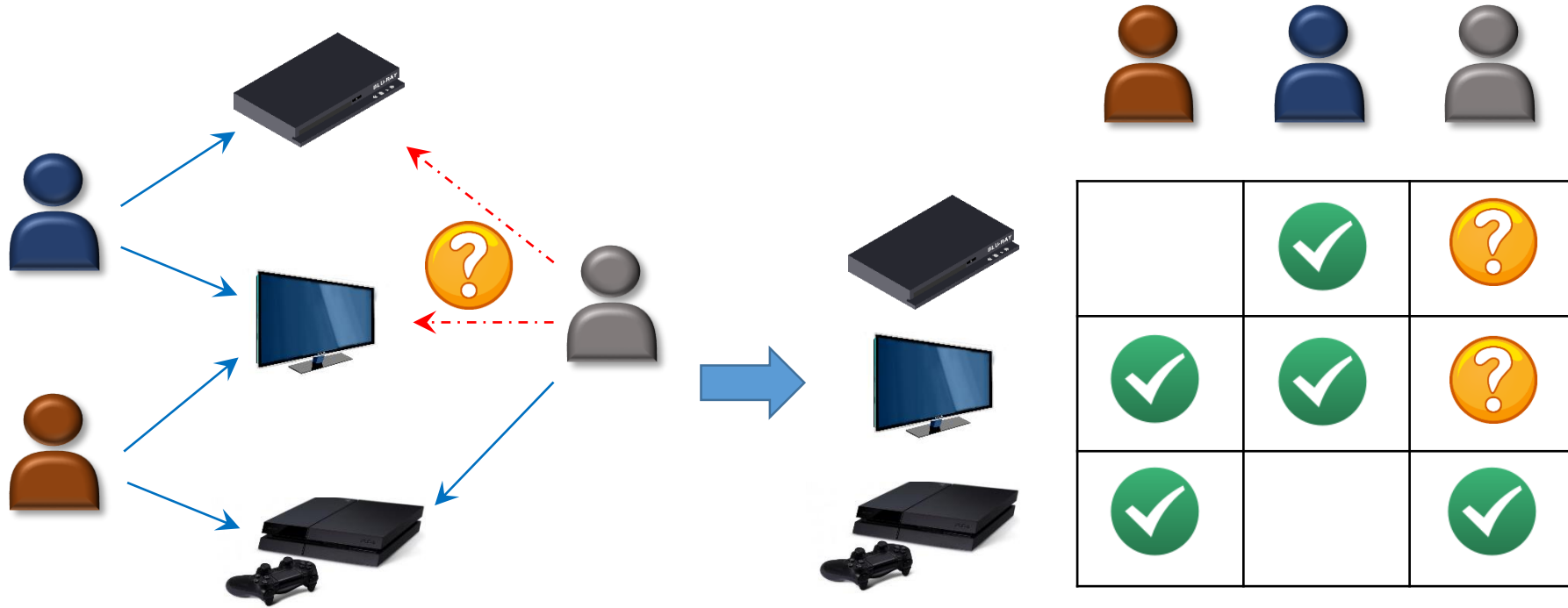


$$\text{Popularity} = (\text{accumulated\_rating}) / (\text{damping\_factor} + \text{total\_number\_of\_ratings})$$

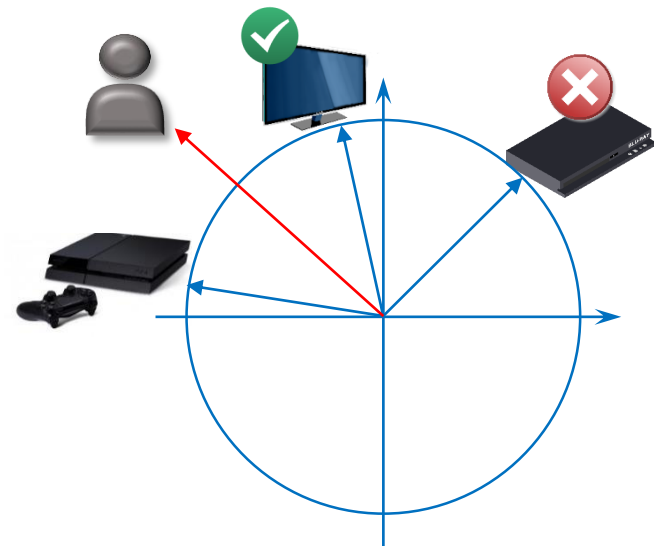
# Collaborative Filtering



# Collaborative filtering – basic idea



$$\text{Sim}(\text{Blu-ray}, \text{PS3}) = \text{Cos}([0, 1, 0], [1, 0, 1]) = 0$$





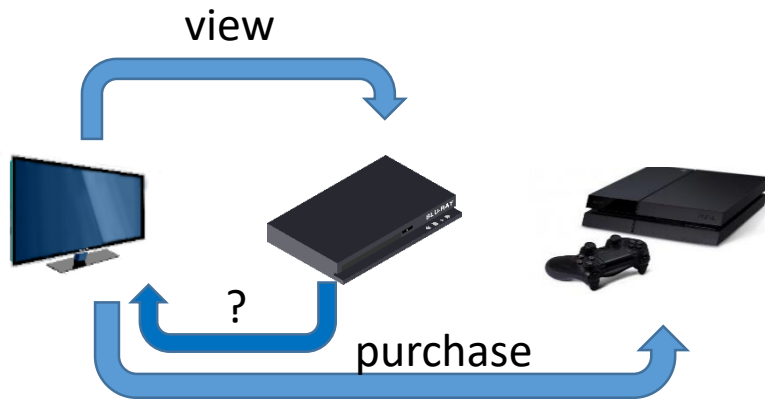
# Memory-based techniques

a.k.a neighborhood models or similarity-based models

# Pure Item-to-item

- Item-to-item is still a favorite choice in many cases, especially when data is sparse

Count co-occurrence of items:



Pair count:

- Symmetric
- Asymmetric

userid	itemid	transact.
0	575	view
0	1881	view
0	846	basket
1	1878	purchase
1	576	view
...	...	...

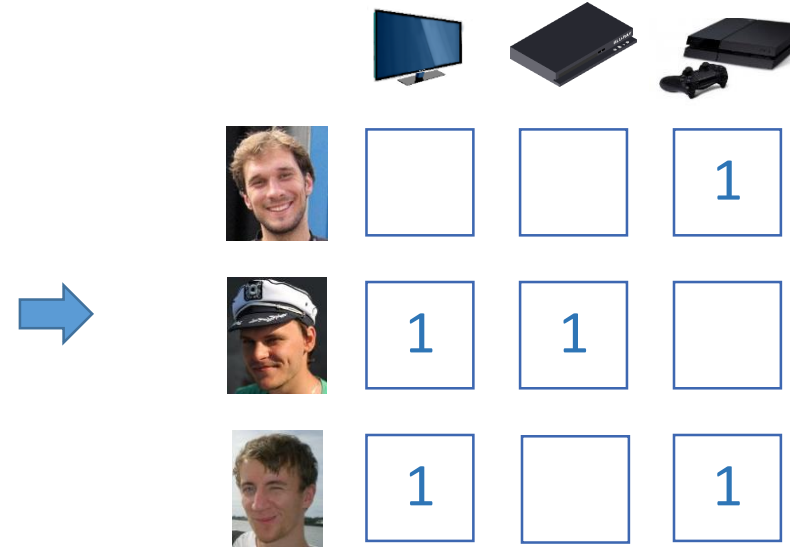
Scaling of values (e.g. count data):

- thresholding
- weighting
- sigmoid
- log*
- TFIDF

# Pure Item-to-item

- Convenient representation of the data – sparse matrix

userid	itemid	transact.
0	575	view
0	1881	view
0	846	basket
1	1878	purchase
1	576	view



Can be efficiently stored in CSR or CSC formats.

Also allows for efficient computations  
(especially useful for experiments).



How to compute item-to-item co-occurrence matrix  
in symmetric case?

How to compute similarity scores in that case?

# Computing scores

$$A = R^T R - \text{diag}(R^T R)$$

“self-links”

if  $p$  is a vector of known user preferences, then vector of predicted scores is computed as:

$$\text{scores} = Ap$$

Recommendations list:

$$\text{toprec}(i, n) := \arg \max_j^n \hat{r}_{ij}$$

$\hat{r}_{ij}$  - is the predicted score of a  $j$ -th item for  $i$ -th user

# Amazon item-to-item

2017 [\[PDF\] Amazon.com recommendations item-to-item collaborative filtering](https://www.cs.umd.edu/~samir/498/Amazon-Recommendations.pdf)  
<https://www.cs.umd.edu/~samir/498/Amazon-Recommendations.pdf> ▼  
Cited by 4095 - Related articles

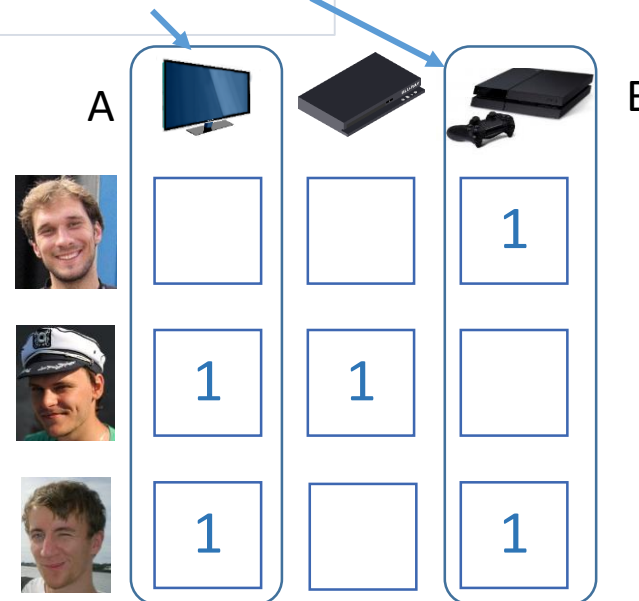
2018 [\[PDF\] Amazon.com recommendations item-to-item collaborative filtering](https://www.cs.umd.edu/~samir/498/Amazon-Recommendations.pdf)  
<https://www.cs.umd.edu/~samir/498/Amazon-Recommendations.pdf> ▼  
Cited by 5365 - Related articles

```
For each item in product catalog,  $I_1$ 
  For each customer  $C$  who purchased  $I_1$ 
    For each item  $I_2$  purchased by
      customer  $C$ 
      Record that a customer purchased  $I_1$ 
        and  $I_2$ 
  For each item  $I_2$ 
    Compute the similarity between  $I_1$  and  $I_2$ 
```

Iterative algorithm

$$\text{similarity}(\vec{A}, \vec{B}) = \cos(\vec{A}, \vec{B}) = \frac{\vec{A} \cdot \vec{B}}{\|\vec{A}\| * \|\vec{B}\|}$$

Computes similarity of items based on user purchases.



# Scalability trick

$$\text{sim}(i_p, i_q) = \frac{\vec{i}_p \cdot \vec{i}_q}{|\vec{i}_p| \cdot |\vec{i}_q|} = \frac{\sum_{u \in U} r_{u,p} r_{u,q}}{\sqrt{\sum r_{u,p}^2} \sqrt{\sum r_{u,q}^2}}$$

$$\text{sim}(i_p, i_q) = \frac{\text{pairCount}(i_p, i_q)}{\sqrt{\text{itemCount}(i_p)} \sqrt{\text{itemCount}(i_q)}}$$

$$\text{itemCount}(i_p) = \sum r_{u,p} \quad \text{pairCount}(i_p, i_q) = \sum_{u \in U} \text{co-rating}(i_p, i_q)$$

$$\begin{aligned} \text{sim}'(i_p, i_q) &= \frac{\text{pairCount}'(i_p, i_q)}{\sqrt{\text{itemCount}'(i_p)} \sqrt{\text{itemCount}'(i_q)}} = \\ &= \frac{\text{pairCount}(i_p, i_q) + \Delta \text{co-rating}(i_p, i_q)}{\sqrt{\text{itemCount}(i_p) + \Delta r_{u_p}} \sqrt{\text{itemCount}(i_q) + \Delta r_{u_q}}} \end{aligned}$$

# Item-to-item problems

- somewhat obvious recommendations
- low generalization on very sparse data
- no hidden relations modelling

## **Workaround:**

compute cosine-similarity on top of i2i matrix.

But!

matrix becomes dense – increased storage requirements.

Idea:

*Maybe limit the number of stored similar items - remember only nearest neighbours?*

*Remark:* i2i matrix can also become dense if there are too many interactions per user.

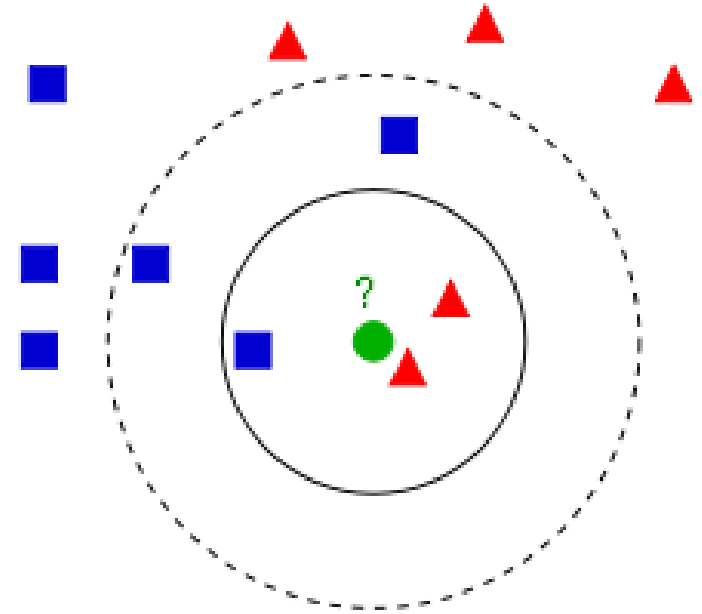
# Similarity-based models

a.k.a neighborhood models

Memory-based approach!

Types of models:

- User-based
- Item-based



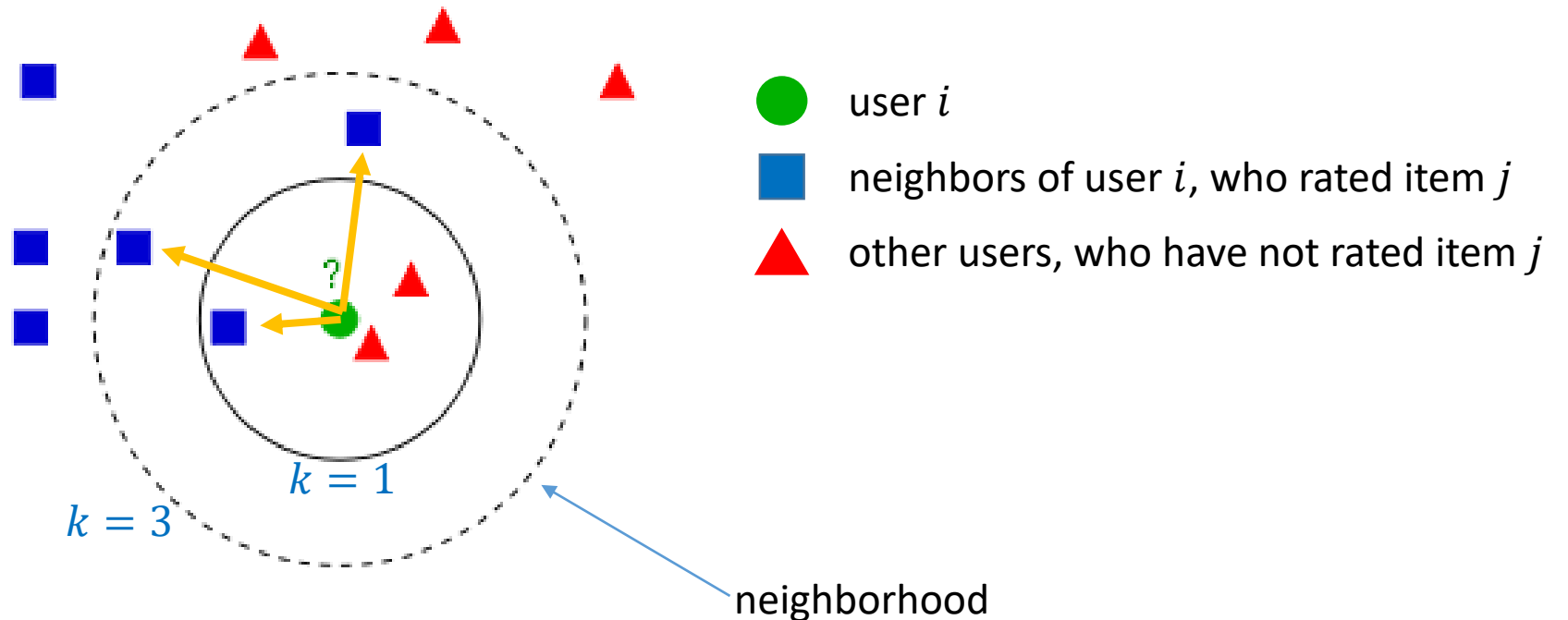


# User-based models

$$\hat{r}_{ij} = \text{agg } r_{uj} \\ u \in \mathcal{N}_j(i)$$

$\mathcal{N}_j(i)$  denotes *k-nearest-neighbors* (*k*-NN) of user  $i$ , who have also rated item  $j$

takes into account opinion  
of like-minded users



# Neighborhood formation

$$\hat{r}_{ij} = \frac{1}{|U|} \sum_{u \in \mathcal{N}_j(i)} r_{uj}$$

$U$  – set of all users

$$\hat{r}_{ij} = k \sum_{u \in \mathcal{N}_j(i)} \text{sim}(u, i) r_{uj}$$

$\text{sim}(u, i)$  – similarity between users  $u$  and  $i$

$k$  – scaling factor  $k = \frac{1}{\sum_{u \in \mathcal{N}(i)} |\text{sim}(u, i)|}$   $|\mathcal{N}_i| \approx 20$

Most general form: 
$$\hat{r}_{ij} = \bar{r}_i + k \sum_{u \in \mathcal{N}_j(i)} \text{sim}(u, i) (r_{uj} - \bar{r}_u)$$
  $\bar{r}_i$  - average rating of user  $i$

Item-based approach is similar: 
$$\hat{r}_{ij} = \bar{r}_j + k \sum_{v \in \mathcal{N}_i(j)} \text{sim}(j, v) (r_{iv} - \bar{r}_v)$$
  $k = \frac{1}{\sum_{v \in \mathcal{N}_i(j)} |\text{sim}(j, v)|}$

**Similarity functions:** Cosine Similarity, Pearson Correlation, Spearman's Rank Correlation, etc...

# When to use user-based or item-based?

- Depends on #users and #items
- Depends on dynamics
- Item-based recommendations are easier to explain.
- User-based recommendations increase serendipity.

# Recap: Item-based or user-based approach

## Key advantages:

- Easy to implement
- Intuitive explanations
- Good baseline

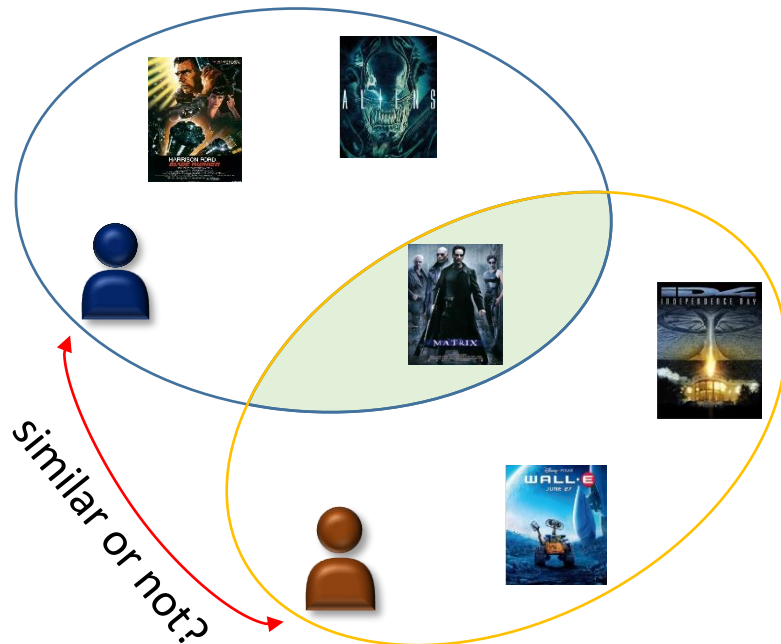
## Scalability:

$O(n^2)$  or  $O(m^2)$  complexity in the worst case

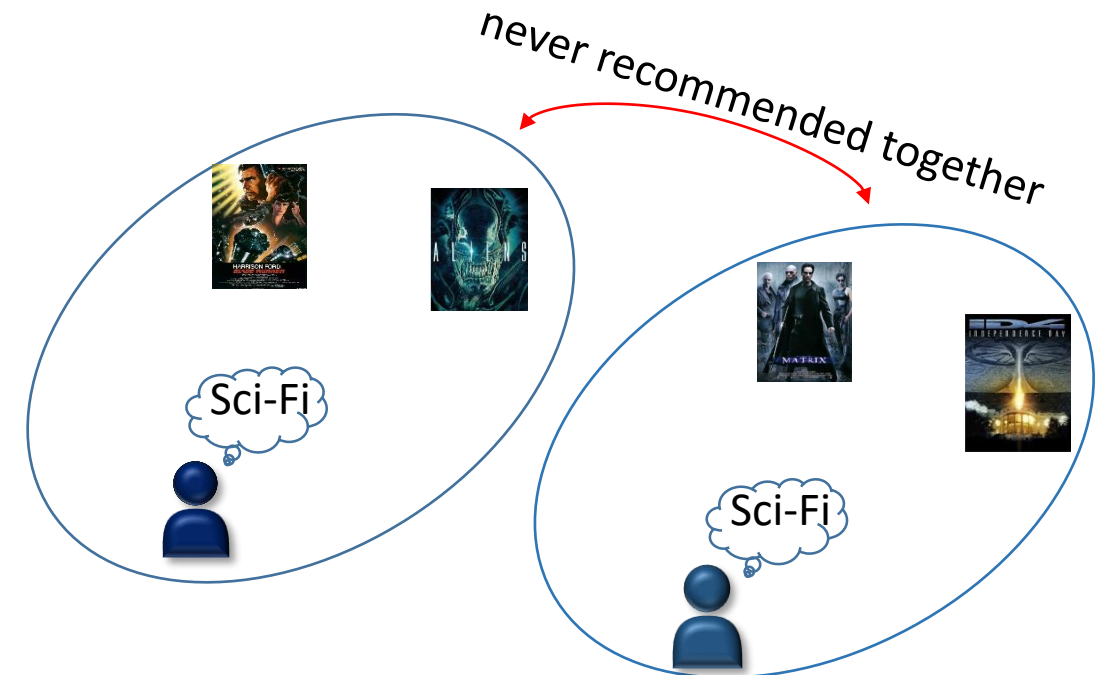
- due to sparsity real complexity is close to *linear*
- could store only limited number of neighbors
- make incremental updates

## Limited coverage problems

### Unreliable correlations



### Weak generalization

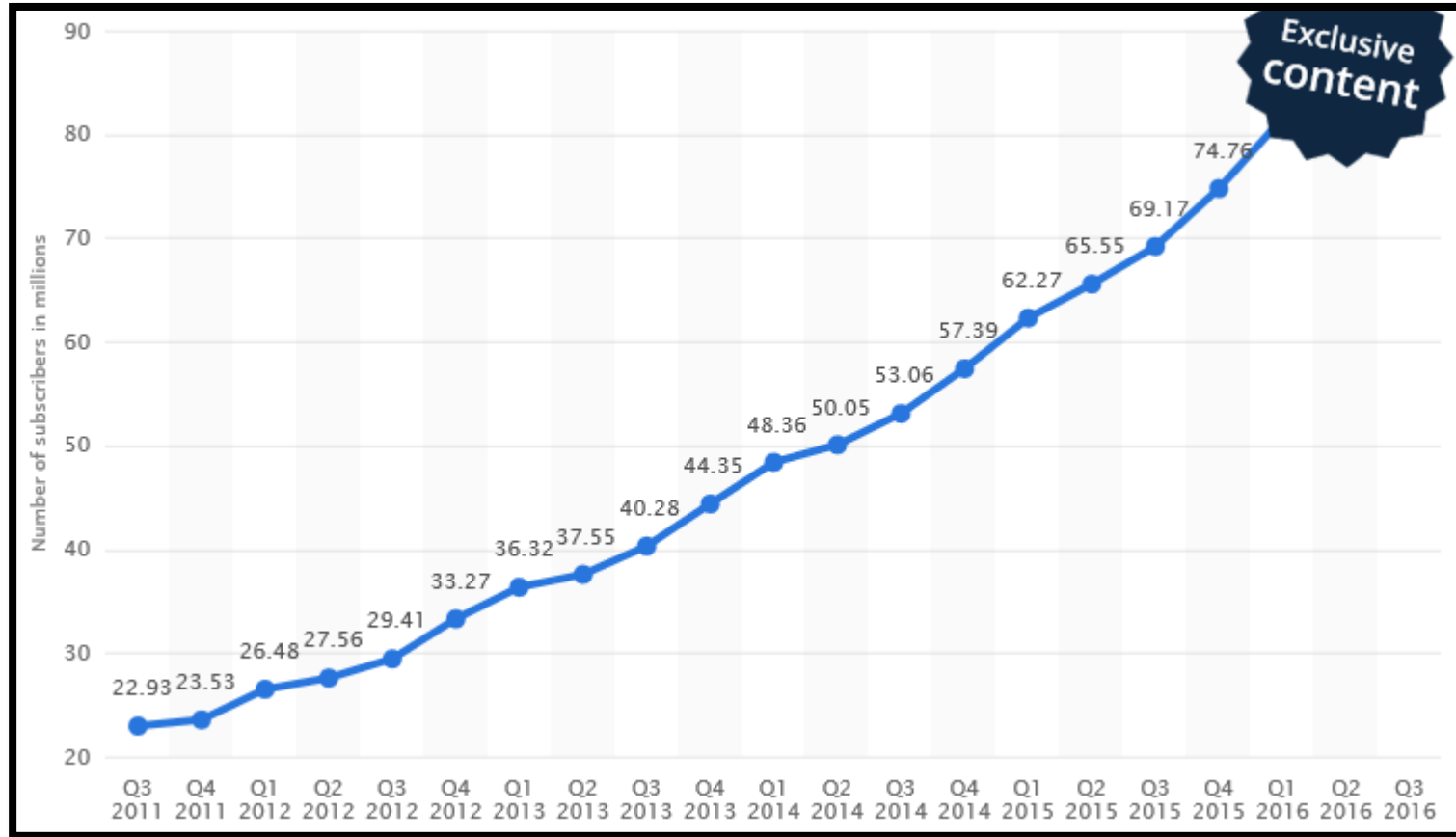


# Model-based approach

Latent factor models

# Netflix

## Netflix's audience



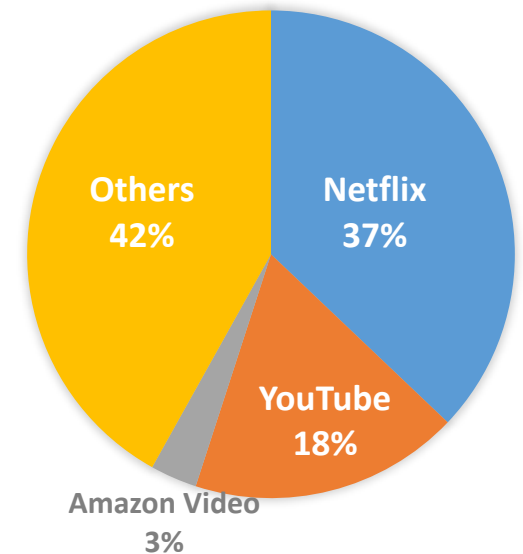
Data sources:

<http://www.internetphenomena.com/tag/amazon-video/>

<https://www.statista.com/statistics/250934/quarterly-number-of-netflix-streaming-subscribers-worldwide/>



## Internet video traffic share



# Netflix prize story

October 2, 2006 - June 26, 2009



**Contest:** Given a database of movies rated by users, beat Netflix's recsys by at least 10%

**Award:** \$1,000,000



**Key to success:** ensemble of models.

Actual solution was never implemented!

<https://www.techdirt.com/blog/innovation/articles/20120409/03412518422/why-netflix-never-implemented-algorithm-that-won-netflix-1-million-challenge.shtml>

But latent factors models based on **matrix factorization** gained popularity.

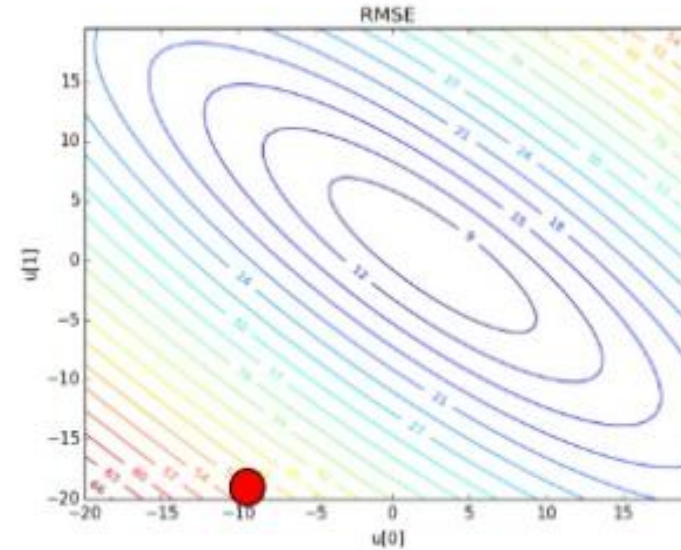
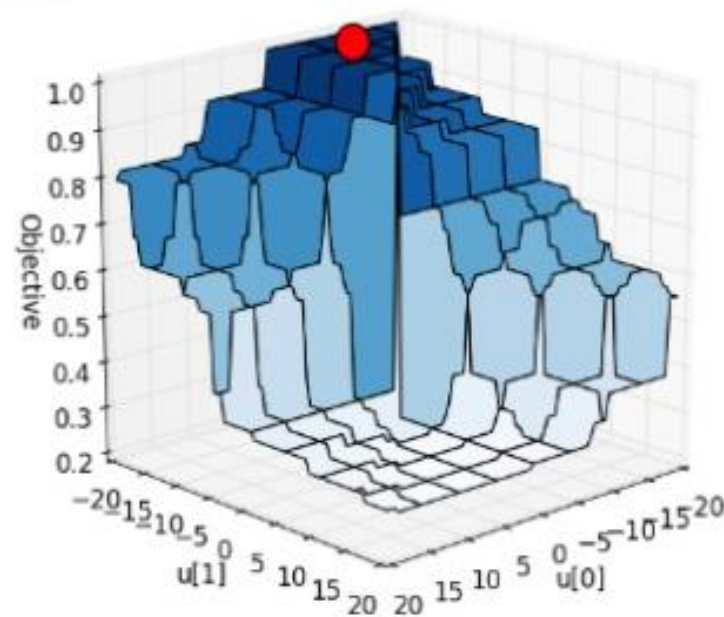
## Issues:

- treated as a pure matrix completion problem
- evaluation metric: RMSE (Precision, Recall, MAP, NDCG, etc. are more adequate)
- proposed matrix factorization models are not better than pure SVD in many practical cases

# Error minimization issue

## AP vs RMSE

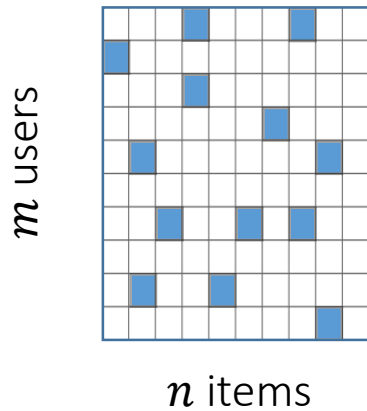
Clip slide





# A general view on matrix factorization

utility matrix  $A$



Incomplete data:

- known entries
- unknown entries

**Task:** find utility (or relevance) function  $f_U$  such that:

$$f_U: \text{Users} \times \text{Items} \rightarrow \text{Relevance score}$$

As optimization problem with some *loss function*  $\mathcal{L}$ :

$$\mathcal{L}(A, R) \rightarrow \min$$

Any factorization model consists of:

- Utility function to generate  $R$
- Optimization objective defined by  $\mathcal{L}$
- Optimization method (algorithm)

Can you guess any form of  $R$  and  $\mathcal{L}$ ?

# A general view on matrix factorization

Let's make the following assumption about observations:

there is a *small* number of common patterns in human behavior + *individual variations*

$$A_{full} = R + E$$

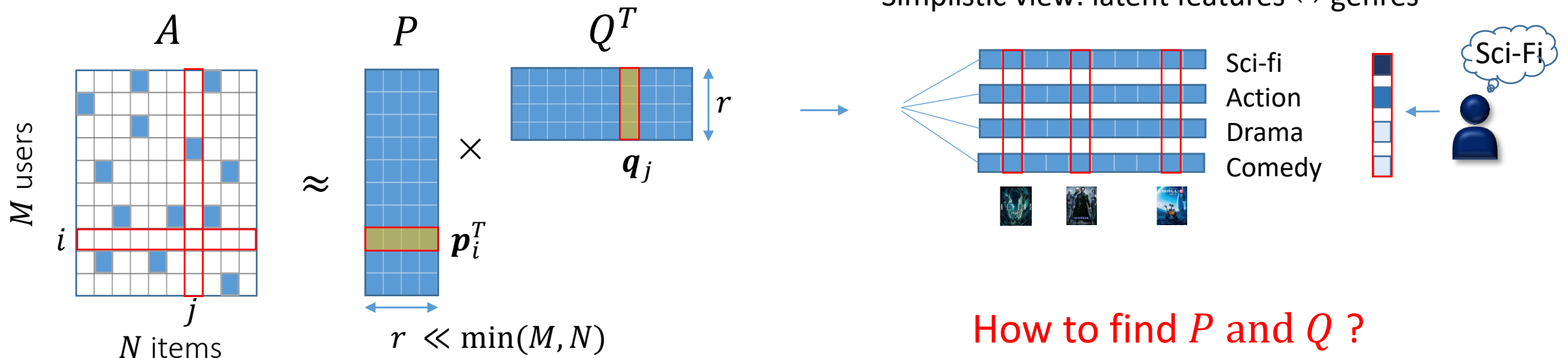
$$R = PQ^T$$

rows of  $P$  and  $Q$  give *embeddings* of users and items onto a latent feature space

predicted utility of item  $j$  for user  $i$

$$r_{ij} \approx \mathbf{p}_i^T \mathbf{q}_j = \sum_{k=1}^r p_{ik} q_{jk}$$

$\mathbf{p}_i$  - latent feature vector for user  $i$   
 $\mathbf{q}_j$  - latent feature vector for item  $j$



# Singular Value Decomposition

Used in LSA/LSI, PCA...

$$\|A - R\|_F^2 \rightarrow \min$$

$$\|X\|_F^2 = \sum_{ij} x_{ij}^2$$

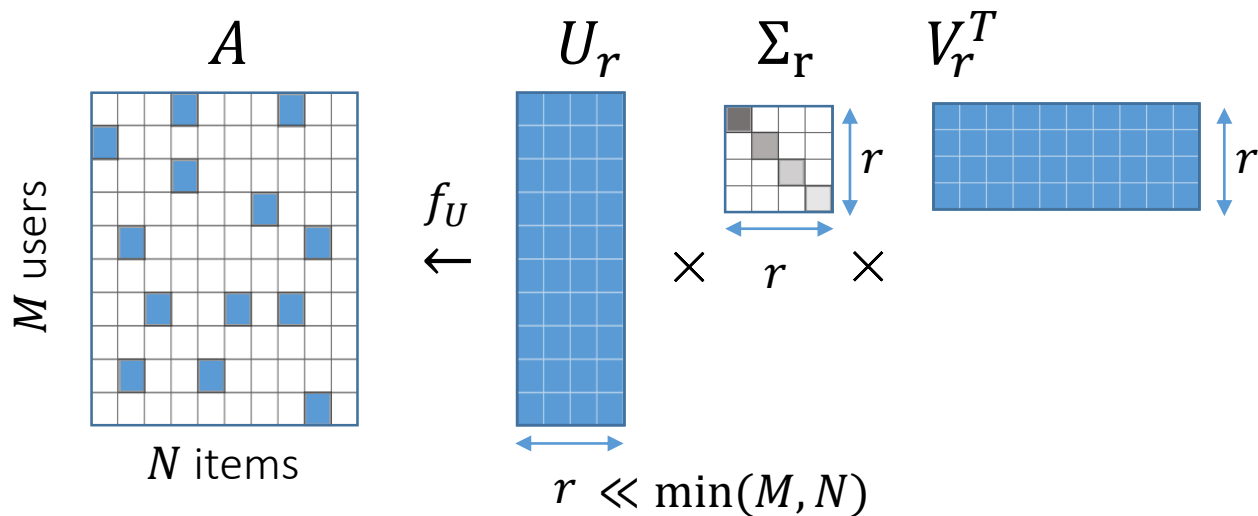
### Analytical solution: SVD (Eckart-Young theorem)

$$A = U \Sigma V^T$$

$$U \in R^{M \times M}, V \in R^{N \times N} \text{--orthogonal}$$

$$\Sigma \in \mathbb{R}^{M \times N} \text{--diagonal}$$

## Truncated SVD of rank $r$



Undefined for incomplete matrix!

## Let's impute zeros - PureSVD model.

$$A_0 = U\Sigma V^T$$

$$R = U_r \Sigma_r V_r^T$$

$$A_0 V_r V_r^T = U \Sigma V^T V_r V_r^T = U_r \Sigma_r V_r^T = R$$

## Important notes:

values are highly biased towards 0

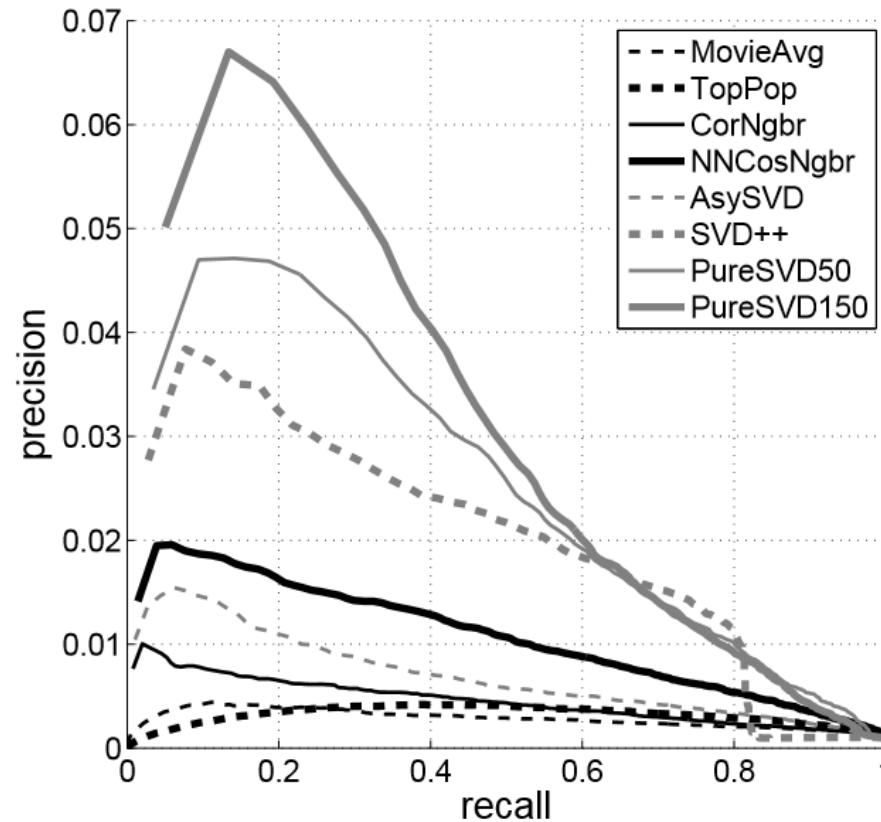
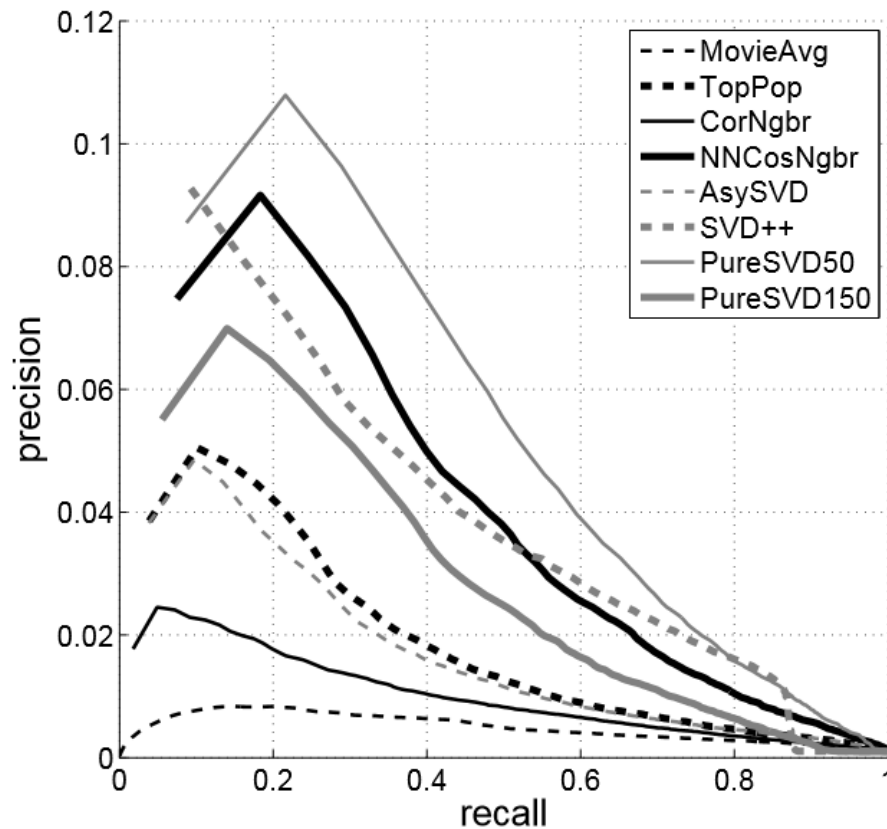
not good for rating prediction

its not a big problem for **top- $n$**  recommendations

## top- $n$ recommendations task:

$$\text{toprec}(i, n) := \arg \max_j^n r_{ij}$$

# PureSVD – quality of recommendations

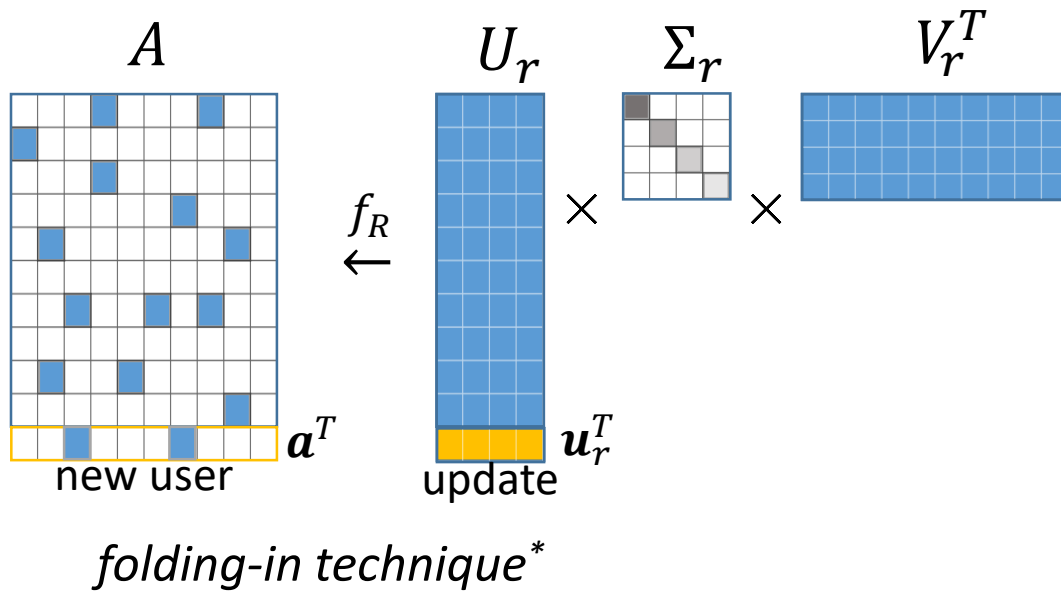


Netflix data: complete dataset (left) and "long-tail" (right).

P. Cremonesi, Y.Koren, R.Turrin, "Performance of Recommender Algorithms on Top-N Recommendation Tasks",  
Proceedings of the 4th ACM conference on Recommender systems, 2011.

Note: Funk SVD, SVD++, TimeSVD++, Asymmetric SVD ... **are not** the SVD!

# PureSVD – recommending online



$$\|a_0^T - u^T \Sigma V^T\|_2^2 \rightarrow \min$$

$$u^T = a_0^T V \Sigma^{-1}$$

new user embedding

$$r^T = u_r^T \Sigma_r V_r^T = a_0^T V_r \Sigma_r^{-1} \Sigma_r V_r^T = a_0^T V_r V_r^T$$

allows for real-time  
recommendations

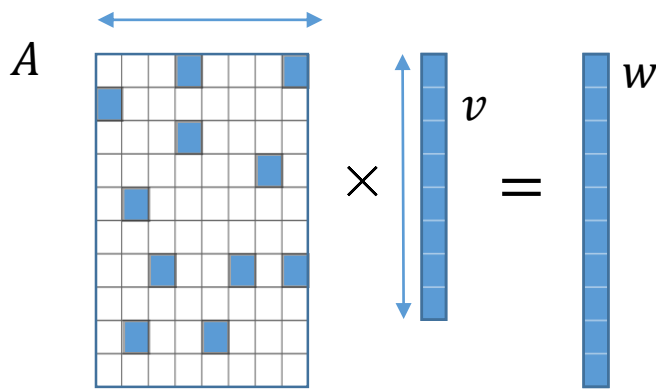
**vector of predicted item scores**

$$r \approx V_r V_r^T a_0$$

$O(Nr)$  complexity

# PureSVD computation

- Efficient computation with Lanczos algorithm
  - iterative process
  - requires only *sparse* matrix-vector (matvec) multiplications (fast with CSR format)
  - complexity  $\sim O(\text{nnz} \cdot r) + O((m + n) \cdot r^2)$
- Implemented in many languages, e.g. MATLAB, Python (SciPy).
- Only core functionality is implemented in Spark.



$O(\text{nnz})$  operations for *sparse* matvec  
 $\text{nnz}$  – number of non-zeros of  $A$

```
In [1]: import numpy as np
        from scipy.sparse import csr_matrix
        from scipy.sparse.linalg import svds
```

```
In [2]: # convert sparse matrix into efficient CSR format
        A = csr_matrix([[0, 1, 1, 0, 0, 0],
                        [0, 1, 0, 1, 0, 0],
                        [0, 0, 1, 0, 1, 1],
                        [0, 1, 0, 1, 0, 1]], dtype=np.float64)

        A
```

```
Out[2]: <4x6 sparse matrix of type '<type 'numpy.float64'>'
        with 10 stored elements in Compressed Sparse Row format>
```

```
In [3]: # compute sparse SVD of rank 2
        rank = 2
        U, S, Vt = svds(A, k=rank)
```

```
In [4]: # check orthogonality
        np.testing.assert_almost_equal(U.T.dot(U), np.eye(rank), decimal=15)
        np.testing.assert_almost_equal(Vt.dot(Vt.T), np.eye(rank), decimal=15)
```

# Demo

Recommender system in 3 lines of code

# Customizing PureSVD

PureSVD is equivalent to an *eigenproblem for the scaled cosine similarity matrix*\*

$$A_0 = U\Sigma V^T \rightarrow A_0^T A_0 = DCD = V\Sigma^2 V^T \quad (\text{similarly for rows})$$

$$C = \left[ \frac{a_i^T a_j}{d_i d_j} \right]_{i,j=1}^N, \quad D = \text{diag}\{d_i\}, \quad d_i = \|a_i\|_2$$

## Options for customization:

- replace cosine similarity with another similarity measure  $S$ ,
- vary scaling factors  $p$ .

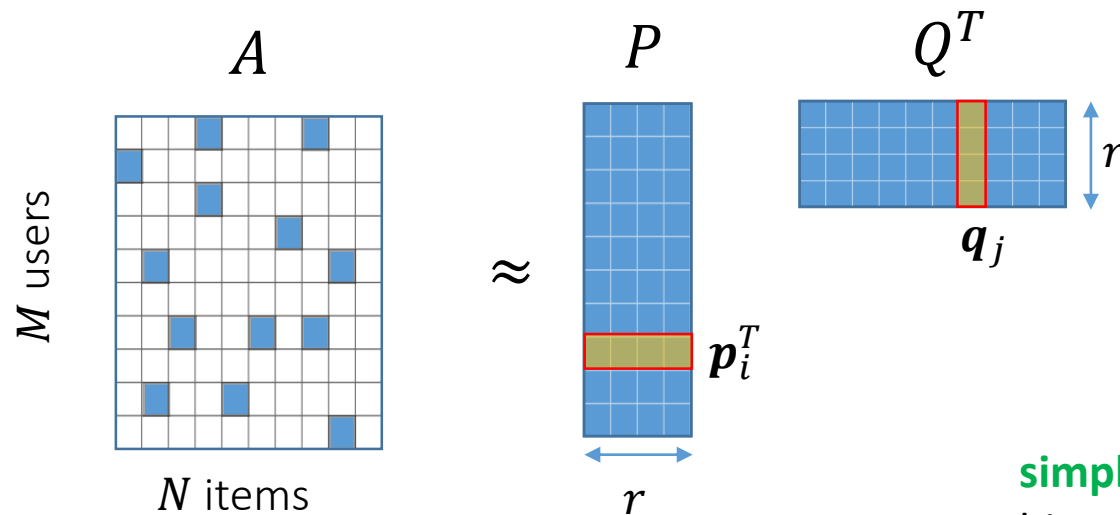
$$DCD \rightarrow D^p S D^p$$

What is the effect of  $p$   
on the model behavior?

\* Nikolakopoulos A. N., Kalantzis V. G., Garofalakis J. D., “EIGENREC: An Efficient and Scalable Latent Factor Family for Top-N Recommendation”, 2015



# Weighted matrix factorization



SVD:  $\mathcal{L} = \|A_0 - R\|_F^2$   $R = U\Sigma V^T$

MF:  $\mathcal{L} = \|W \odot (A - R)\|_F^2$   $R = PQ^T$

Hadamard (element-wise) product

$$\begin{cases} w_{ij} = 1, & \text{if } a_{ij} \text{ is known,} \\ w_{ij} = 0, & \text{otherwise.} \end{cases}$$

simplest case -  
binary weights

elementwise form:

$$\mathcal{L}(A, \Theta) = \frac{1}{2} \sum_{i,j \in S} (a_{ij} - \mathbf{p}_i^T \mathbf{q}_j)^2$$

**Uses only the observed data!**

$$S = \{(i, j): w_{ij} \neq 0\}$$

$$\mathcal{J}(\Theta) = \mathcal{L}(A, \Theta) + \Omega(\Theta)$$

$$\Theta = \{P, Q\}$$

additional constraints on factors

Typical optimization algorithms:

stochastic gradient descent (SGD)

alternating least squares (ALS)

GD:

$$\begin{cases} \mathbf{p}_i \leftarrow \mathbf{p}_i - \eta \nabla_{\mathbf{p}_i} \mathcal{J} \\ \mathbf{q}_j \leftarrow \mathbf{q}_j - \eta \nabla_{\mathbf{q}_j} \mathcal{J} \end{cases}$$

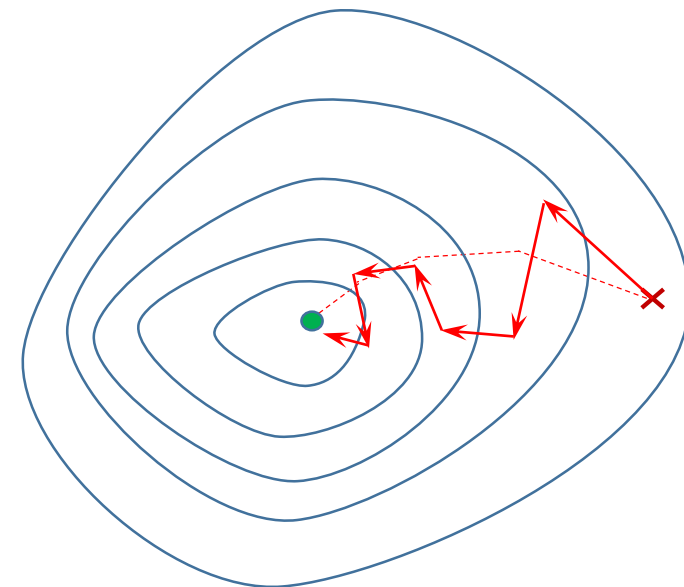
ALS:

$$\begin{cases} P = \arg \min_Q \mathcal{J}(\Theta) \\ Q = \arg \min_P \mathcal{J}(\Theta) \end{cases}$$

# Optimization with SGD

$$J(P, Q) = \frac{1}{2} \sum_{i,j \in S} \overbrace{(a_{ij} - \mathbf{p}_i^T \mathbf{q}_j)^2}^{l_{ij}} + \lambda (\|\mathbf{p}_i\|^2 + \|\mathbf{q}_j\|^2)$$

determined by cross-validation



“true” gradient:

$$\frac{\partial J}{\partial \mathbf{p}_i} = - \sum_{j \in S(i, \cdot)} (a_{ij} - \mathbf{p}_i^T \mathbf{q}_j) \mathbf{q}_j + \lambda \mathbf{p}_i$$

ratings of user  $i$

can be inefficient with large data

$$\begin{cases} \mathbf{p}_i \leftarrow \mathbf{p}_i - \eta \cancel{\frac{\partial J}{\partial \mathbf{p}_i}} \frac{\partial l_{ij}}{\partial \mathbf{p}_i} \\ \mathbf{q}_j \leftarrow \mathbf{q}_j - \eta \cancel{\frac{\partial J}{\partial \mathbf{q}_j}} \frac{\partial l_{ij}}{\partial \mathbf{q}_j} \end{cases}$$

approximate gradients

$$\frac{\partial l_{ij}}{\partial \mathbf{p}_i} = - \underbrace{(a_{ij} - \mathbf{p}_i^T \mathbf{q}_j)}_{e_{ij}} \mathbf{q}_j + \lambda \mathbf{p}_i$$

$$\frac{\partial l_{ij}}{\partial \mathbf{q}_j} = - \underbrace{(a_{ij} - \mathbf{p}_i^T \mathbf{q}_j)}_{e_{ij}} \mathbf{p}_i + \lambda \mathbf{q}_j$$

Algorithm

Initialize  $P$  and  $Q$ .

Iterate until stopping criteria met:

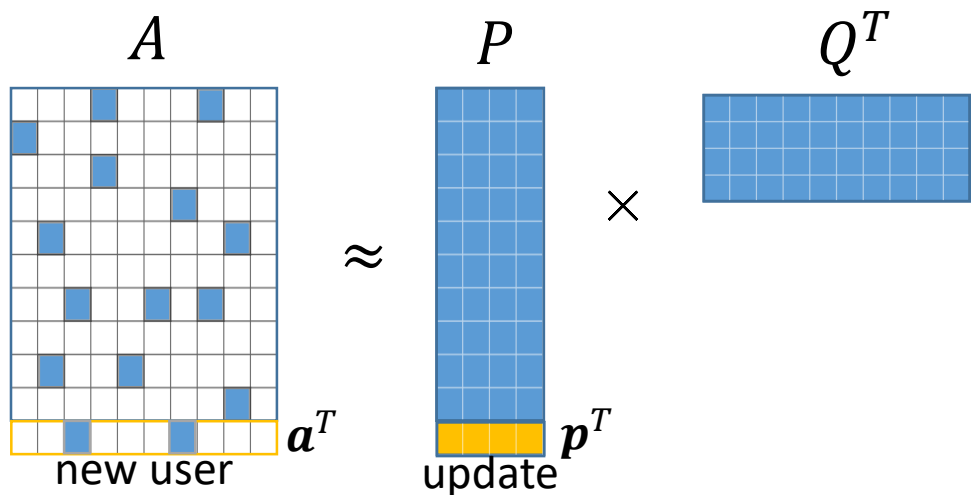
for each pair  $i, j \in S$  (shuffled):

compute  $e_{ij}$

$$\begin{cases} \mathbf{p}_i \leftarrow \mathbf{p}_i + \eta (e_{ij} \mathbf{q}_j - \lambda \mathbf{p}_i) \\ \mathbf{q}_j \leftarrow \mathbf{q}_j + \eta (e_{ij} \mathbf{p}_i - \lambda \mathbf{q}_j) \end{cases}$$

Complexity:  
 $O(\text{nnz} \cdot r)$

# Incremental updates



What are the key differences between SGD-based and SVD-based folding-in?

Folding-in

in SVD:  $\mathbf{u} = \Sigma^{-1} V^T \mathbf{a}_0$

via SGD:

Initialize  $\mathbf{p}$

Iterate until stopping criteria met:

For all ratings in  $\mathbf{a}$ :

$$e_{aj} = a_j - \mathbf{p}^T \mathbf{q}_j$$

$$\mathbf{p} \leftarrow \mathbf{p} + \eta (e_{aj} \mathbf{q}_j - \lambda \mathbf{p})$$

$$O(\text{nnz}_a \cdot r)$$

# of non-zero elements of  $\mathbf{a}$

$$O(\text{nnz}_a \cdot r)$$

# Including bias terms



popularized by Simon Funk  
during the Netflix Prize competition

# NETFLIX

- critical users tend to rate movies lower than average user
- popular movies on average receive higher ratings

pre-calculated global average

tendency of a user to rate movies higher (or lower)

$$r_{ij} = \mu + t_i + f_j + \mathbf{p}_i^T \mathbf{q}_j$$

how favorable is an item in general

$$\underset{\mathbf{p}_i, \mathbf{q}_j, b_i, b_j}{\text{minimize}} \sum_{i,j \in S} e_{ij}^2 + \lambda (\|\mathbf{p}_i\|^2 + \|\mathbf{q}_j\|^2 + t_i^2 + f_j^2) \quad e_{ij} = ?$$

$$\begin{cases} \mathbf{p}_i \leftarrow \mathbf{p}_i + \eta(e_{ij}\mathbf{q}_j - \lambda\mathbf{p}_i) \\ \mathbf{q}_j \leftarrow \mathbf{q}_j + \eta(e_{ij}\mathbf{p}_i - \lambda\mathbf{q}_j) \\ t_i \leftarrow t_i + \eta(e_{ij} - \lambda t_i) \\ f_j \leftarrow f_j + \eta(e_{ij} - \lambda f_j) \end{cases}$$

# Matrix form

Can you incorporate bias terms into matrix?

**Hint:** resulting rank is  $r+2$ .

$$[P \ e \ t][Q \ f \ e]^T = PQ^T + ef^T + te^T$$

# Alternative SGD optimization scheme

- Instead of learning “by entity”, one could use “by component” scheme.

Think, what are pros and cons of such approach?

# Optimization with ALS

$$\mathcal{J}(\Theta) = \mathcal{L}(A, \Theta) + \Omega(\Theta) \quad \mathcal{L} = \frac{1}{2} \|W \odot (A - PQ^T)\|_F^2 \quad \Omega(\Theta) = \frac{1}{2} \lambda (\|P\|_F^2 + \|Q\|_F^2)$$

“user-oriented” form:

$$\mathcal{J}(\Theta) = \frac{1}{2} \sum_i \|\mathbf{a}_i - Q \mathbf{p}_i\|_{W^{(i)}}^2 + \frac{1}{2} \lambda \sum_i \|\mathbf{p}_i\|_2^2 + \frac{1}{2} \lambda \|Q\|_F^2$$

$$\|\mathbf{x}\|_W^2 = \mathbf{x}^T W \mathbf{x}$$

$$W^{(i)} = \text{diag}\{w_{i1}, w_{i2}, \dots, w_{iN}\}$$

Block-coordinate descent:

$$\begin{aligned} \frac{\partial \mathcal{J}(\Theta)}{\partial P} &= 0 \\ \frac{\partial \mathcal{J}(\Theta)}{\partial Q} &= 0 \end{aligned} \quad \begin{array}{c} \text{entity-wise} \\ \text{updates} \end{array}$$

$$(Q^T W^{(i)} Q + \lambda I) \mathbf{p}_i = Q^T W^{(i)} \mathbf{a}_i$$

← row of A

$$(P^T W^{(i)} P + \lambda I) \mathbf{q}_j = P^T W^{(i)} \bar{\mathbf{a}}_j$$

← column of A

system of linear equations (SLA):  
 $A\mathbf{x} = \mathbf{b}$

“embarrassingly parallel”

[https://en.wikipedia.org/wiki/Embarrassingly\\_parallel](https://en.wikipedia.org/wiki/Embarrassingly_parallel)

Algorithm

Initialize  $P$  and  $Q$ .

Iterate until stopping criteria met:

$$P = \arg \min_Q \mathcal{J}(\Theta)$$

$$Q = \arg \min_P \mathcal{J}(\Theta)$$

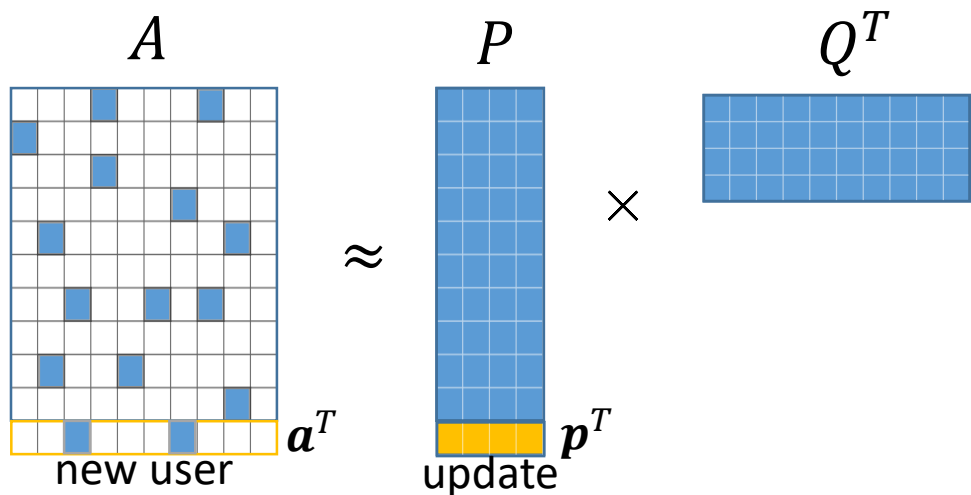
**Complexity:**

$$O(\text{nnz}_A \cdot r^2) + O((M + N)r^3)$$

can be improved with  
approximate SLA solvers (e.g. *Conjugate Gradient*)

or switch to *Coordinate Descent*

# Incremental updates



Folding-in

in SVD:  $\mathbf{u} = \Sigma^{-1} V^T \mathbf{a} \quad O(\text{nnz}_a \cdot r)$

in MF:  $\|\mathbf{a} - Q\mathbf{p}\|_{W^{(a)}}^2 \rightarrow \min$

$W^{(a)}$  follows nnz pattern of  $\mathbf{a}$

least squares problem

$$\mathbf{p} = (Q^T W^{(a)} Q + \lambda I)^{-1} Q^T W^{(a)} \mathbf{a} \quad O(\text{nnz}_a \cdot r^2 + r^3)$$



# Task

- You have a binary utility matrix (with “true” zeros) resulted from some implicit feedback information.
  - What will be the complexity of SGD?
  - What will be the SGD-based solution if you omit zero values?
  - Is it reasonable to use bias terms?

# ALS vs SGD vs SVD

## ALS

More stable

Fewer hyper-parameters to tune

Higher complexity, however requires fewer iterations

Embarrassingly parallel

Higher communication cost in distributed environment

## SGD

Sensitive to hyper-parameters

Requires special treatment of learning rate

Lower complexity, but slower convergence

Inherently sequential (parallelization is tricky for RecSys)

For binary feedback complexity changes:  $nnz \rightarrow MN$

Unlike SVD:

**More involved model selection (no rank truncation).**

**No global convergence guarantees!**

**Asynchronous SGD is non-deterministic.**

**Allow for custom optimization objectives.**  $\mathcal{L}(A, R) \rightarrow \mathcal{L}(f(A, R))$

For explicit feedback:

Algorithm	Overall complexity	Update complexity	Sensitivity	Optimality
SVD*	$O(nnz_A \cdot r + (M + N)r^2)$	$O(nnz_a \cdot r)$	Stable	Global
ALS	$O(nnz_A \cdot r^2 + (M + N)r^3)$	$O(nnz_a \cdot r + r^3)$	Stable	Local
CD	$O(nnz_A \cdot r)$	$O(nnz_a \cdot r)$	Stable	Local
SGD	$O(nnz_A \cdot r)$	$O(nnz_a \cdot r)$	Sensitive	Local

\* For both standard and randomized implementations [71].

# Remark on user feedback

User feedback is a source for constructing a utility function.

## Implicit feedback

Easy to collect

Lots of data

Intrinsic

Hard to interpret

## Explicit feedback

Hard to collect

Less data

Subjective

“Easy” to interpret

# Explicit feedback peculiarities

- horror movies ratings are typically lower, even if user actually likes it



“Ghostbusters” Is A Perfect Example Of How Internet Movie Ratings Are Broken



- IMDb **average user rating**: 4.1 out of 10, of 12,921 reviewers
- IMDb **average user rating among men**: 3.6 out of 10, of 7,547 reviewers
- IMDb **average user rating among women**: 7.7 out of 10, of 1,564 reviewers

Source: <http://fivethirtyeight.com/features/ghostbusters-is-a-perfect-example-of-how-internet-ratings-are-broken/>

# Confidence-based model (a.k.a iALS, WRMF)

$$\mathcal{L} = \frac{1}{2} \|W \odot (S - PQ^T)\|_F^2$$

$$S: \begin{cases} s_{ij} = 1, & \text{if } a_{ij} \text{ is known,} \\ s_{ij} = 0, & \text{otherwise.} \end{cases}$$

↑  
preference  
confidence

**Weights**  $W = [w_{ij}^{1/2}]$  **are not binary**

$$\begin{cases} w_{ij} = 1 + \alpha f(a_{ij}), & \text{if } a_{ij} \text{ is known,} \\ w_{ij} = 1, & \text{otherwise.} \end{cases}$$

constant

$$f(x) \begin{cases} x \\ \log\left(x + \frac{1}{\epsilon}\right) \end{cases}$$

$$\mathbf{p} = (Q^T W^{(a)} Q + \lambda I)^{-1} Q^T W^{(a)} \mathbf{a} \quad W^{(a)} = \text{diag}\{\mathbf{1} + \alpha f(\mathbf{a})\} - \text{dense!}$$

Naïve approach:  $O(\textcolor{red}{MN} \cdot r^2) + O((M + N)r^3)$

**Trick:** pre-calculate and store  $r \times r$  matrix  $Q^T Q$

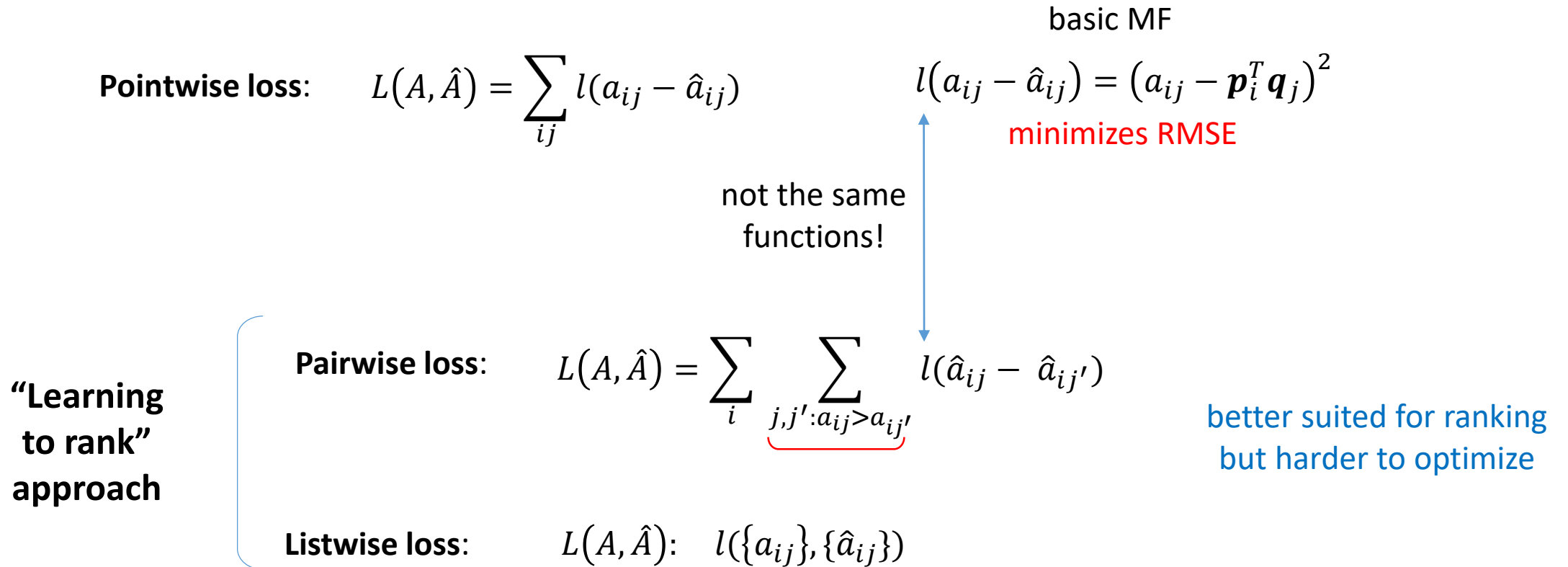
$$Q^T W^{(a)} Q = \textcolor{green}{Q^T Q} + Q^T C^{(a)} Q$$

$$C^{(a)} = W^{(a)} - I \quad \text{is sparse!}$$

← computed only once per epoch

Overall complexity reduces to:  $O(\text{nnz}_A \cdot r^2) + \textcolor{green}{O((M + N)r^2)} + O((M + N)r^3)$

# Remark on optimization objectives



# Other factorization methods

- NMF
- MMF
- PMF
- CAMF
- SLIM
- FISM
- CliMF
- OrdRec
- CofaFi
- ...
- **Hybrid methods** (e.g. *Factorization Machines*, *HybridSVD*)
- **Context-aware** methods (e.g., *Tensor factorization*)

# Helpful resources

## *Books*

- **Recommender Systems Handbook, 2015**, 2<sup>nd</sup> edition; F. Ricci, L. Rokach, B. Shapira
- **Recommender Systems. The Textbook, 2016**; Charu C. Aggarwal
- **Recommender Systems: An Introduction, 2010**; D.Jannach, M.Zanker, A.Felfernig, G.Friedrich
- **Statistical Methods for Recommender Systems, 2016**; Deepak K. Agarwal, Bee-Chung Chen
- **Collaborative Recommendations: Algorithms, Practical Challenges and Applications**; S. Berkovsky, I. Cantador and D. Tikk; *expected May 2019*.

## *Conferences*

**ACM RecSys**

UMAP

WWW

KDD

WSDM

## *Competitions*

RecSys Challenge

CIKM challenge

Kaggle

Tip: Want to find a job? Attend RecSys conferences!



# Helpful resources

## Courses

- <https://www.coursera.org/learn/recommender-systems/home/welcome>  
From pioneers of RecSys field
- Week 3 of “Big Data Applications: Machine Learning at Scale” course in Big Data for Data Engineers Specialization, <https://www.coursera.org/specializations/big-data-engineering>
- Machine Learning: Recommender Systems & Dimensionality Reduction  
<https://www.coursera.org/learn/ml-recommenders> (Amazon Professors)
- Mining Massive Datasets, Chapter on Recommender Systems, Stanford University  
<http://www.mmds.org>

## Video tutorials

- **Machine Learning Summer School 2014**  
[https://www.youtube.com/playlist?list=PLZSO\\_6-bSqHQCIYxE3ycGLXHMjK3XV7Iz](https://www.youtube.com/playlist?list=PLZSO_6-bSqHQCIYxE3ycGLXHMjK3XV7Iz)  
Lectures from Xavier Amatriain and Deepak Agarwal
- **Introduction to Machine Learning 10-701 CMU 2015**, Alex Smola  
<https://www.youtube.com/watch?v=gCaOa3W9kM0>

## Other resources

- RecSys wiki: <http://recsyswiki.com> (currently down)
- Blog: **A Practical Guide to Building Recommender Systems**  
<https://buildingrecommenders.wordpress.com>
- OpenDataScience Slack, #recommender\_systems channel (mostly Russian language)

# Libraries and Frameworks

- Frameworks

- Polara (*Disclaimer: I'm the author*)  
<https://github.com/evfro/polara>
- MyMediaLite  
<http://www.mymedialite.net>
- Collaborative Filtering – Apache Spark  
<http://spark.apache.org/docs/latest/mllib-collaborative-filtering.html>  
(Neighborhood models and MF)
- Surprise  
<https://github.com/NicolasHug/Surprise>
- Turi Create (ex GraphLab Create)  
<https://apple.github.io/turicreate/docs/userguide/recommender/>

- Useful libraries

- Collaborative Filtering for Implicit Feedback Datasets  
<https://github.com/benfred/implicit> (the fastest)  
<https://github.com/quora/qmf> (by Quora)  
<https://github.com/MrChrisJohnson/logistic-mf> (as in Spotify)
- Factorization Machines  
<https://github.com/srendle/libfm>

- Other libraries

- Neural Networks  
<https://github.com/Netflix/vectorflow> (by Netflix)  
<https://github.com/amzn/amazon-dsstne> (Amazon)  
<https://github.com/maciejkula/spotlight>  
<https://github.com/MrChrisJohnson/deep-mf>  
<https://github.com/songgc/TF-recomm>
- Bilinear models  
<https://github.com/lyst/lightfm/>  
<http://www.recsyswiki.com/wiki/SVDFeature>
- Many latent factor models  
<https://github.com/zhangsi/CisRec>
- Simple content-based recommendation engine  
<https://github.com/groveco/content-engine>
- Logistic Matrix Factorization  
<https://github.com/MrChrisJohnson/implicit-mf>
- Hermes (Supports Spark)  
<https://github.com/Lab41/hermes>