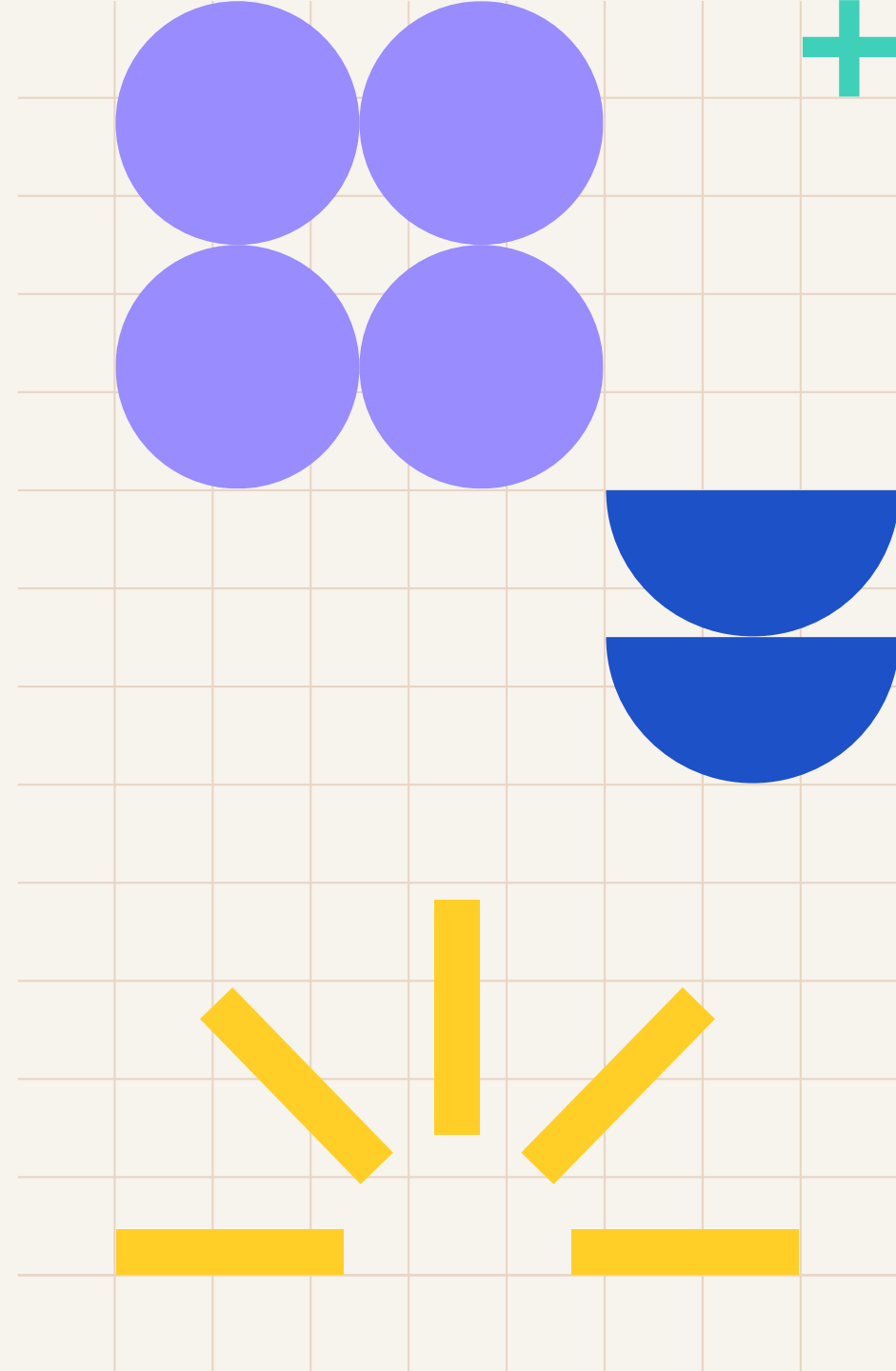




Practical algorithms for Recommender Systems (Part 2)

Evgeny Frolov

Research Scientist, Skoltech



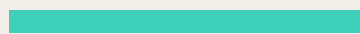
Today

01 More on Matrix Factorization

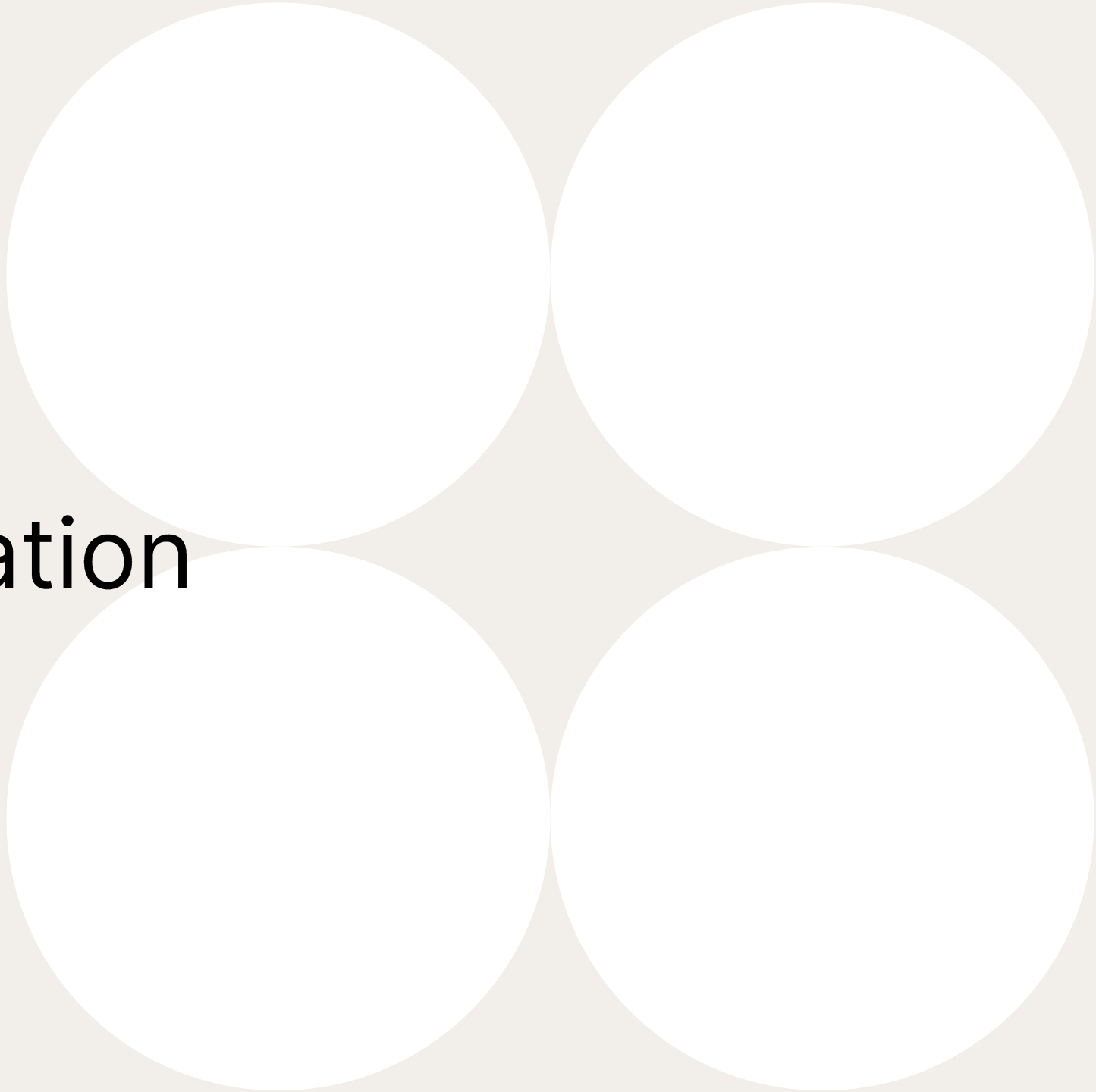
02 Hybrid Recommender Systems

03 Context-awareness

01



Matrix Factorization (continued)



Previous lecture – a general view on latent factors models

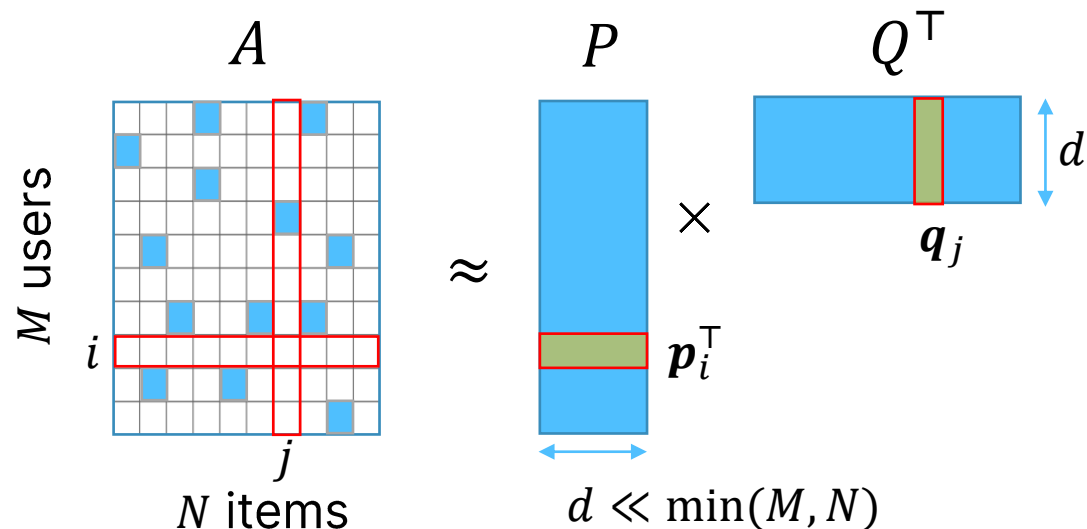
Task:

- find a relevance function

$$f_R: r_{ij} \approx \mathbf{p}_i^\top \mathbf{q}_j = \sum_{k=1}^d p_{ik} q_{jk}$$

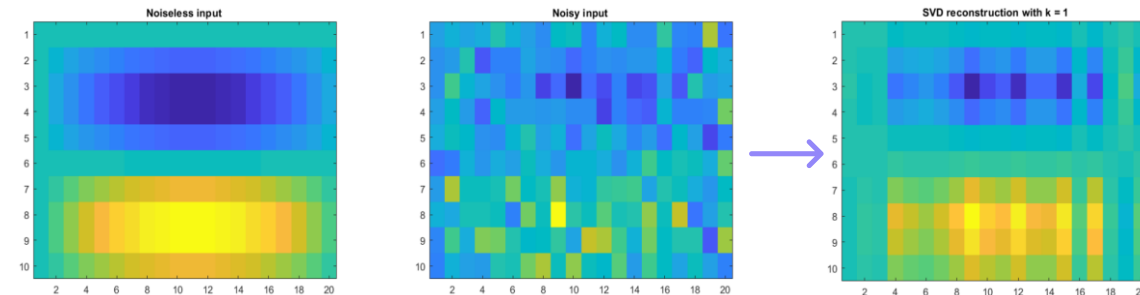
- via an optimization problem:

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



Components of the solution:

- Utility function to generate R
- Optimization objective \mathcal{L}
- Optimization algorithm



Previous lecture - PureSVD model

$$\|A_0 - R\|_F^2 \rightarrow \min, \text{ s.t. } \text{rank}(R) = d$$

$$[A_0]_{ij} = \begin{cases} a_{ij}, & \text{if known} \\ 0, & \text{otherwise} \end{cases}$$

$$f_R: R = A_0 V_d V_d^T$$

Efficient computation with Lanczos algorithm:

- iterative process
- requires *only sparse matrix-vector* (matvec) multiplications (fast with CSR format);
- training complexity $O(\text{nnz} \cdot d) + O((M + N) \cdot d^2)$

Efficient implementations in Python:

- SciPy Sparse svds, Scikit-Learn TruncatedSVD.
- core functionality is also implemented in Spark.

In distributed setups, randomized SVD is used.

```
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)
```

More general MF optimization scheme

Optimization objective:

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

Model parameters: $\Theta = \{P, Q\}$

$\Omega(\Theta)$ - additional constraints, e.g. L_2 regularization

Typical optimization algorithms:

stochastic gradient descent (SGD)

alternating least squares (ALS)

ALS:

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

GD:

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

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
- Coordinate Descent can be a good alternative (e.g., eALS by [He et al. 2016])

SGD

- Sensitive to hyper-parameters
- Requires special treatment of learning rate
- Lower complexity but slower convergence, using adaptive learning rate schedule (ADAM, Adagard, etc.) helps
- Inherently sequential (parallelization is tricky for RecSys)
- Hogwild! algorithm is not directly applicable in CF settings

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

* For both standard and randomized implementations [71].

(ALS, SGD) vs SVD:

- **More involved optimization (no rank truncation).**
- **Allow for custom optimization objectives.**

General rule of thumb: avoid distributed setups.

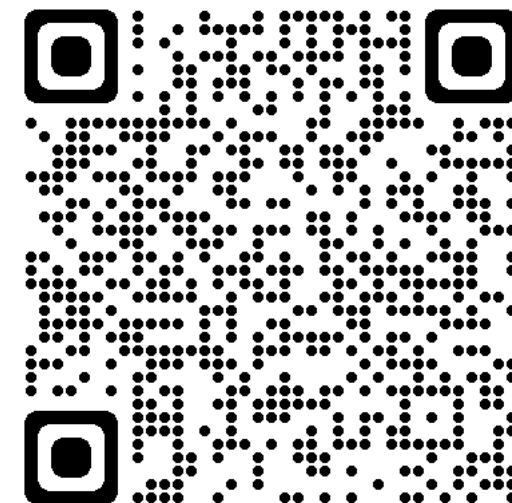
Working with imbalanced data

- SGD:
 - negative sampling
- iALS:
 - confidence weights

$$\mathcal{L} = \sum_{ij} w(a_{ij}) \cdot l(s_{ij} - r_{ij})^2$$

- PureSVD
 - data normalization

$$\tilde{A} = DA_0, \quad [D]_{ii} = \|a_i\|^{f-1}$$



All these methods aim to balance contribution of positive and negative items!

Case study: Yandex Zen

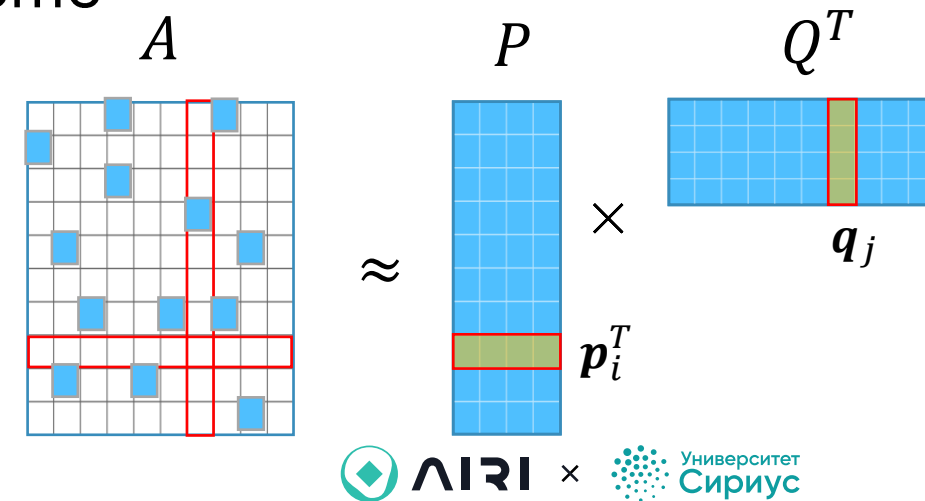
Company manages different types of media content (news, search, etc.).
Goal: have a unified user representation across all domains.

Solution:

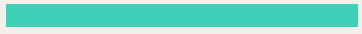
- A DNN embeds unstructured content into a shared latent space
- Users are updated through the “half”-ALS scheme

Algorithm:

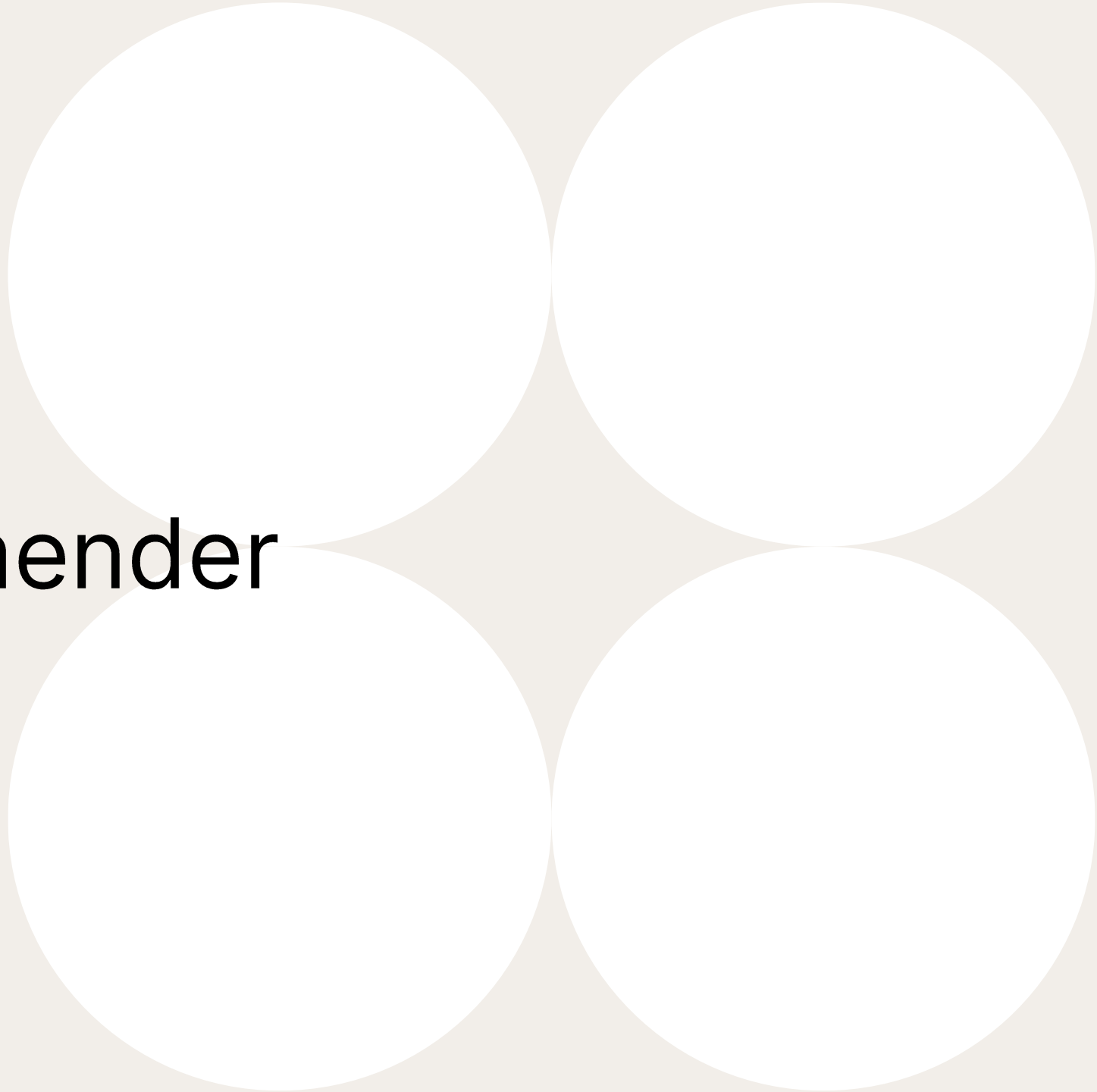
- Get Q from external source (DNN)
- Update P based on most recent Q



02

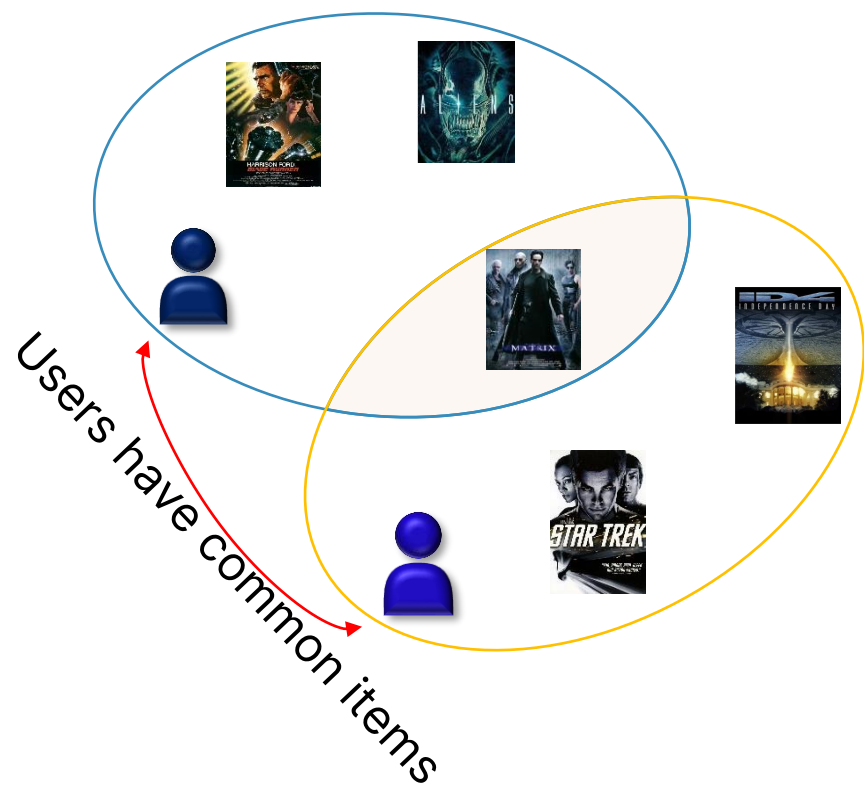


Hybrid Recommender Systems



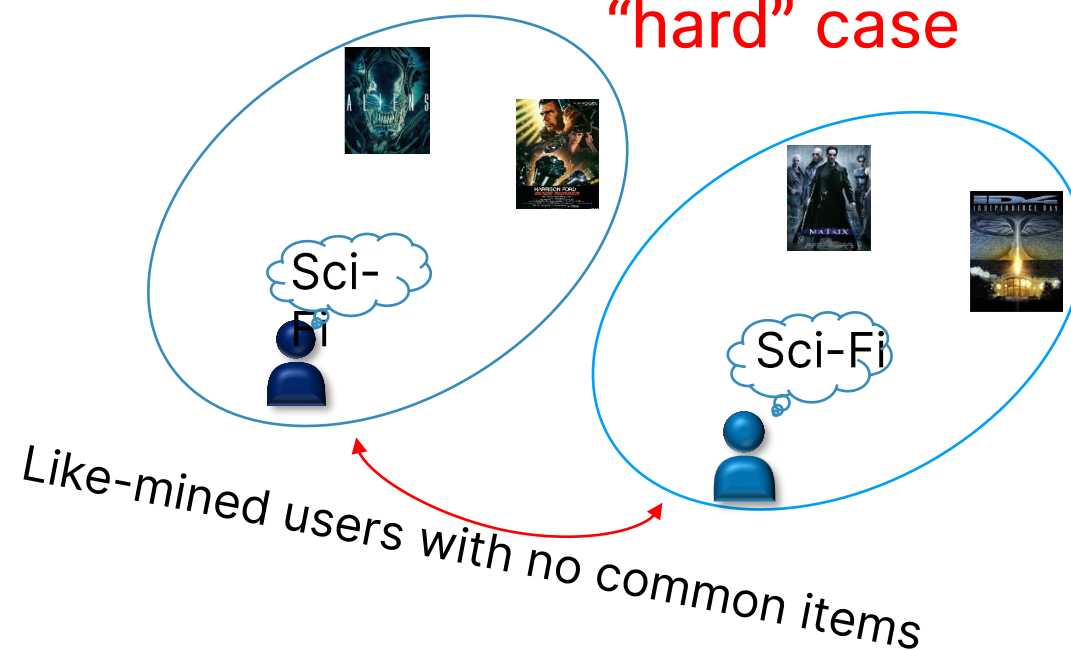
The problem of rare interactions

“easy” case

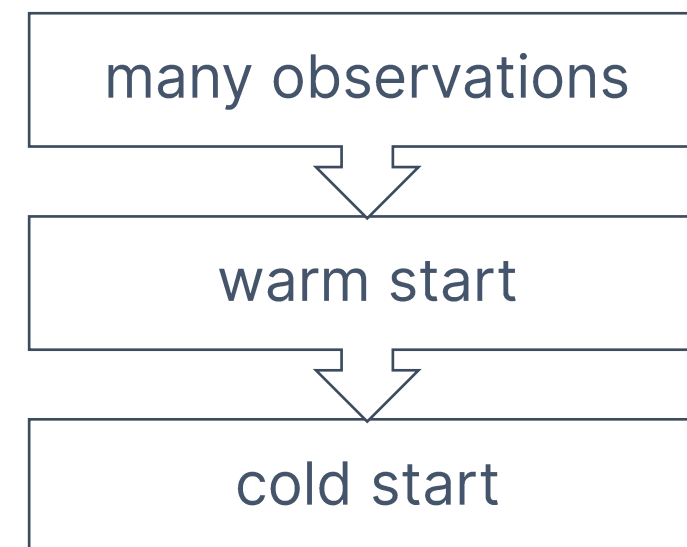
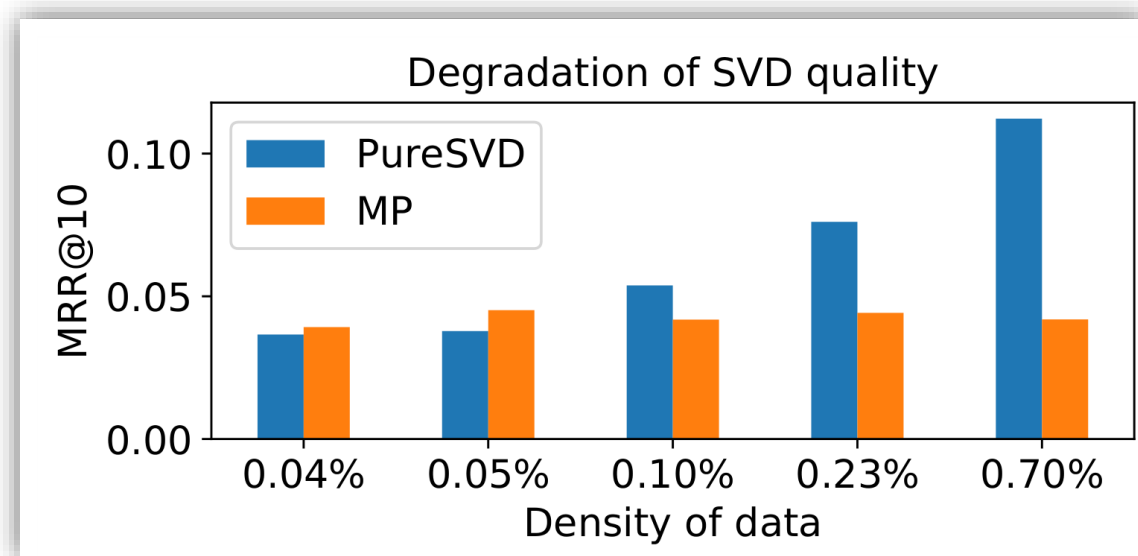


VS.

“hard” case



The problem of rare interactions



How to mitigate that?

Example of content features

Users

The screenshot shows a Goodreads user profile for 'Steve'. The profile includes a circular profile picture of a man reading a book, a 'Follow' button, and an 'Add friend' button. Below the profile picture, it shows '262 ratings (3.77 avg)', '114 reviews', and 'more photos (1)'. A section titled 'STEVE'S FAVORITE BOOKS' displays three book covers. A larger, semi-transparent overlay shows a '2021 reading journey' for Steve, featuring a book cover and the text 'This is my journey in books for 2021!'. It lists '6,938 pages read' and '19 books read'. Below this, it compares 'Shortest Book' (228 pages) and 'Longest Book' (704 pages). At the bottom, it states 'Average book length in 2021' is 365 pages.

More attributes:

- demographics
- location
- occupation
- ...

Items

The screenshot shows a Goodreads book page for 'Mindset: The New Psychology of Success' by Carol S. Dweck. The page features the book cover, a 'Want to Read' button, and a 'Rate this book' section with a 4.07 star rating. The description states: 'After decades of research, world-renowned Stanford University psychologist Carol S. Dweck, Ph.D., discovered a simple but groundbreaking idea: the power of mindset. In this brilliant book, she shows how success in school, work, sports, the arts, and almost every area of human endeavor can be dramatically influenced by how we t ...more'. Below the description are 'Amazon' and 'Stores' buttons. The page also includes a 'FRIEND REVIEWS' section, a 'READER Q&A' section, and a 'GENRES' list with user counts: Psychology (2,453 users), Nonfiction (2,158 users), Self Help (1,469 users), Business (737 users), Self Help > Personal Development (646 users), Education (597 users), Leadership (328 users), Parenting (327 users), Productivity (177 users), and Teaching (173 users). A 'READERS ALSO ENJOYED' section shows covers for 'Atomic Habits' and 'GRIT'.

Other features:

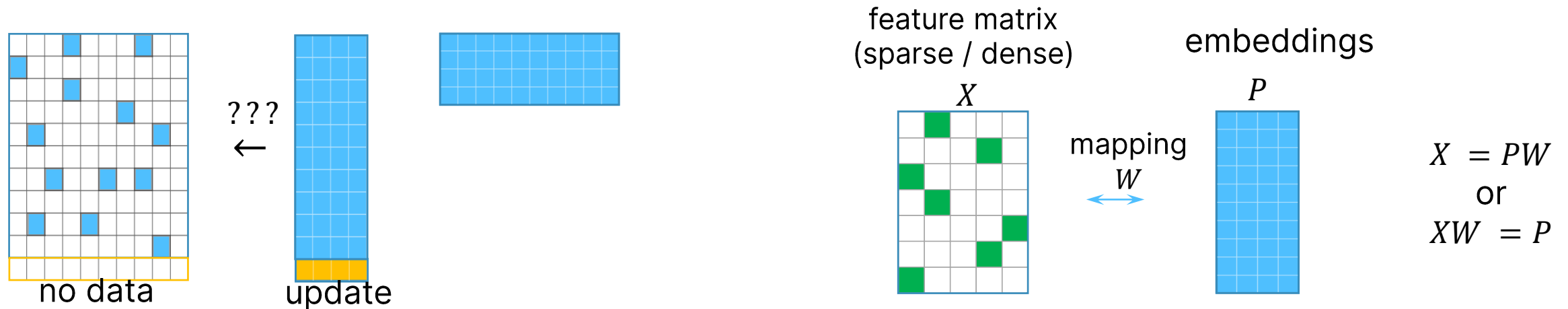
- price
- format/style
- language
- ...

Mitigating cold-start problems with hybrid approach

Pure content-based filtering may not be effective:

- noisy / incomplete side information,
- overspecialization.

Pure collaborative filtering is inapplicable in cold start!



How can we use side information for recovering latent features?

Simple linear regression model

$f(\text{user features, item features}) \rightarrow \text{feedback}$

$$r = \mathbf{w}^T \mathbf{z} + \epsilon$$



Does it provide personalized recommendations?

$$r_{ui} = b_{\text{user}} + b_{\text{item}} + \mathbf{w}_{\text{user}}^T \mathbf{x}_u + \mathbf{w}_{\text{item}}^T \mathbf{y}_i$$

How to add personalization?

We need a way to entangle user and item features!

Example:

- users are described with 2 features based on age group, e.g. [>18 , >65]
- items are described with 3 features based on book genre, e.g., [is action, is romance, is drama]

How to add personalization?

- let's encode all possible combinations of features via Cartesian product (bias terms are omitted for simplicity):

$$Z = \mathbf{x}\mathbf{y}^\top, \quad \mathbf{x} = [x_1, \dots, x_{m_x}]^\top, \quad \mathbf{y} = [y_1, \dots, y_{n_y}]^\top$$

- new model:

$$r = \mathbf{w}^\top \mathbf{z}, \quad \mathbf{z} = \text{vec}(Z)$$

- or equivalently:

$$r = \mathbf{x}^\top W \mathbf{y},$$
$$\text{vec}(W) = \mathbf{w}, \quad W \in \mathbb{R}^{m_x \times n_y}$$

Improved top- n ranking

- the model:

$$r_{xy} = \mathbf{x}^\top W \mathbf{y} = \sum_{i=1}^{m_x} \sum_{j=1}^{n_y} w_{ij} \cdot x_i y_j$$

- ranking now depends on the association strength w_{ij} between user features x_i and item features y_j

$$r_{xy} - r_{xy'} = \sum_{i=1}^{m_x} x_i \sum_{j=1}^{n_y} w_{ij} \cdot (y_j - y'_j)$$

- # learned parameters of the global model: $O(\quad)$

Limited expressiveness of the model

still, there're problems:

- various items/users may have the same features - *overspecialization*
- if ratings are different → ill-posed problem

What additional information will help?

Improved personalized regression model

- personalization issue fix – encode user and item ids:

$$r_{xy} = \mathbf{x}^\top W \mathbf{y}$$

$$\mathbf{x}^\top = [\mathbf{x}_{\text{id}}^\top \quad \mathbf{x}_{\text{feat}}^\top], \quad \mathbf{y}^\top = [\mathbf{y}_{\text{id}}^\top \quad \mathbf{y}_{\text{feat}}^\top]$$



- matrix form:

$$R =$$

- resolves expressiveness problem
- # learned parameters:

$$O(\quad)$$

Combating data sparsity

Structural problem:

- the weights matrix W can become *restrictively large*
- conversely, there's only a small number of known user-item interactions

How can we deal with that?

- Imposing low rank structure:

$$W = PQ^T$$

- yields:

$$R = XP(YQ)^T$$

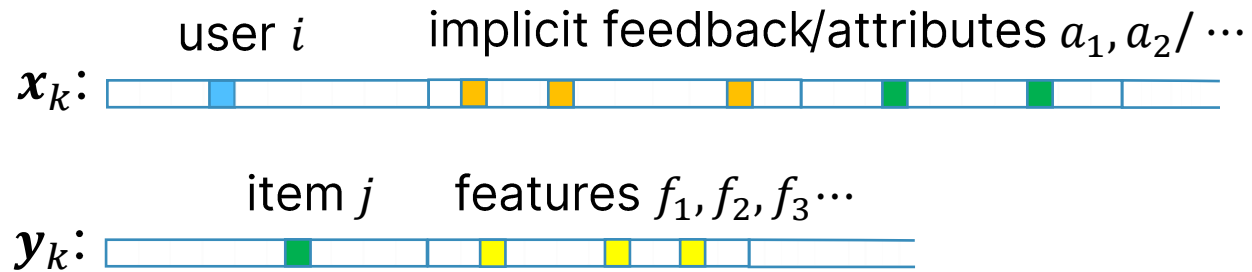


SVDFeature

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

$$R = (XP)(YQ)^T$$

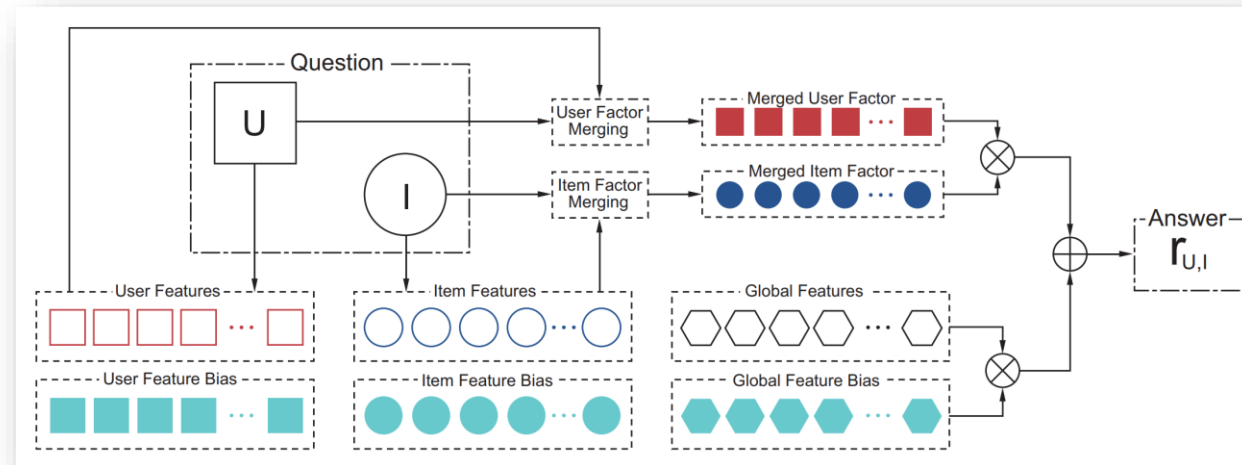
$$X = [X_1 \ X_2 \ \dots \ X_m], \quad Y = [Y_1 \ Y_2 \ \dots \ Y_n]$$



Including bias terms:

$$r = b_0 + \mathbf{g}^T \mathbf{x} + \mathbf{f}^T \mathbf{y} + \mathbf{x}^T P Q^T \mathbf{y}$$

- $b_0 = \sum_{g \in G} \gamma_g \mu_g$ is precomputed.
- model parameters: $\Theta = \{\mathbf{g}, \mathbf{f}, P, Q\}$
- optimized with ALS, SGD, BPR.



LightFM

- Same general approach as in SVDFeature
- Option to pick logistic loss
- SGD with BPR/WARP optimizers

The model is parameterised in terms of d -dimensional user and item feature embeddings \mathbf{e}_f^U and \mathbf{e}_f^I for each feature f . Each feature is also described by a scalar bias term (b_f^U for user and b_f^I for item features).

The latent representation of user u is given by the sum of its features' latent vectors:

$$\mathbf{q}_u = \sum_{j \in f_u} \mathbf{e}_j^U$$

The same holds for item i :

$$\mathbf{p}_i = \sum_{j \in f_i} \mathbf{e}_j^I$$

The bias term for user u is given by the sum of the features' biases:

$$b_u = \sum_{j \in f_u} b_j^U$$

The same holds for item i :

$$b_i = \sum_{j \in f_i} b_j^I$$

The model's prediction for user u and item i is then given by the dot product of user and item representations, adjusted by user and item feature biases:

$$\hat{r}_{ui} = f(\mathbf{q}_u \cdot \mathbf{p}_i + b_u + b_i) \quad (1)$$

There is a number of functions suitable for $f(\cdot)$. An identity function would work well for predicting ratings; in this paper, I am interested in predicting binary data, and so after Rendle *et al.* [16] I choose the sigmoid function

$$f(x) = \frac{1}{1 + \exp(-x)}.$$

The optimisation objective for the model consists in maximising the likelihood of the data conditional on the parameters. The likelihood is given by

$$L(\mathbf{e}^U, \mathbf{e}^I, \mathbf{b}^U, \mathbf{b}^I) = \prod_{(u,i) \in S^+} \hat{r}_{ui} \times \prod_{(u,i) \in S^-} (1 - \hat{r}_{ui}) \quad (2)$$

M. Kula, "Metadata embeddings for user and item cold-start recommendations." 2015.
<https://github.com/lyst/lightfm>

Factorization Machines

Idea: polynomial expansion

$$f(\mathbf{z}) = b_0 + \mathbf{b}^\top \mathbf{z} + \mathbf{z}^\top \mathbf{H} \mathbf{z} + \dots$$

Feature vector \mathbf{x}															Target y							
$\mathbf{x}^{(1)}$	1	0	0	...	1	0	0	0	...	0.3	0.3	0.3	0	...	13	0	0	0	0	...	5	$y^{(1)}$
$\mathbf{x}^{(2)}$	1	0	0	...	0	1	0	0	...	0.3	0.3	0.3	0	...	14	1	0	0	0	...	3	$y^{(2)}$
$\mathbf{x}^{(3)}$	1	0	0	...	0	0	1	0	...	0.3	0.3	0.3	0	...	16	0	1	0	0	...	1	$y^{(2)}$
$\mathbf{x}^{(4)}$	0	1	0	...	0	0	1	0	...	0	0	0.5	0.5	...	5	0	0	0	0	...	4	$y^{(3)}$
$\mathbf{x}^{(5)}$	0	1	0	...	0	0	0	1	...	0	0	0.5	0.5	...	8	0	0	1	0	...	5	$y^{(4)}$
$\mathbf{x}^{(6)}$	0	0	1	...	1	0	0	0	...	0.5	0	0.5	0	...	9	0	0	0	0	...	1	$y^{(5)}$
$\mathbf{x}^{(7)}$	0	0	1	...	0	0	1	0	...	0.5	0	0.5	0	...	12	1	0	0	0	...	5	$y^{(6)}$
	A	B	C	...	TI	NH	SW	ST	...	TI	NH	SW	ST	...	Time	TI	NH	SW	ST	...		
	User				Movie					Other Movies rated					Last Movie rated							

Factorization Machines

$$r = b_0 + \mathbf{g}^\top \mathbf{x} + \mathbf{f}^\top \mathbf{y} + \mathbf{x}^\top \mathbf{P} \mathbf{Q}^\top \mathbf{y}$$

$$\mathbf{b} = \begin{bmatrix} t \\ \mathbf{f} \end{bmatrix}, \quad \mathbf{z} = \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}$$

$$r(\mathbf{z}) = b_0 + \mathbf{b}^\top \mathbf{z} + \mathbf{z}^\top \mathbf{H} \mathbf{z} + \dots$$

characterizes relations between all types of encoded entities

Data is sparse → impose low-rank structure on H

H is symmetric positive semi-definite

$$H = VV^\top \quad V \text{ embeds all users, items and their side information}$$


2nd order FM:

$$r(\mathbf{z}) = b_0 + \sum_{k=1}^K b_k z_k + \sum_{k=1}^K \sum_{k'=k+1}^K \mathbf{v}_k^\top \mathbf{v}_{k'} \cdot z_k z_{k'}$$

Model parameters: $\Theta = \{b_0, \mathbf{z}, V\}$



Factorization Machines computation

$$\mathbf{z} = \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \\ \mathbf{f} \end{bmatrix} \quad V = \begin{bmatrix} V_x \\ V_y \\ V_f \end{bmatrix} \quad V \in \mathbb{R}^{K \times d} \quad \mathbf{z} \in \mathbb{R}^K:$$


$$\mathbf{z}^T V V^T \mathbf{z} = \left(\begin{bmatrix} \mathbf{x}^T & \mathbf{y}^T & \mathbf{f}^T \end{bmatrix} \begin{bmatrix} V_x \\ V_y \\ V_f \end{bmatrix} \right) \left(\begin{bmatrix} V_x^T & V_y^T & V_f^T \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \\ \mathbf{f} \end{bmatrix} \right) = (\mathbf{x}^T V_x + \mathbf{y}^T V_y + \mathbf{f}^T V_f)(\mathbf{x}^T V_x + \mathbf{y}^T V_y + \mathbf{f}^T V_f)^T =$$

“self-interaction” terms

actual 2nd order FM model $r(\mathbf{z})$ (w/o biases)

$$= \mathbf{x}^T V_x V_x^T \mathbf{x} + \mathbf{y}^T V_y V_y^T \mathbf{y} + \mathbf{f}^T V_f V_f^T \mathbf{f} + 2 \left(\underbrace{\mathbf{x}^T V_x V_y^T \mathbf{y}}_{\text{user-item}} + \underbrace{\mathbf{x}^T V_x V_f^T \mathbf{f}}_{\text{user-feature}} + \underbrace{\mathbf{y}^T V_y V_f^T \mathbf{f}}_{\text{item-feature}} \right)$$

interactions


$$\mathbf{v}_x = V_x^T \mathbf{x}, \mathbf{v}_y = \dots$$

$$r(\mathbf{z}) = \frac{1}{2} \left[(\mathbf{v}_x + \mathbf{v}_y + \mathbf{v}_f)^T (\mathbf{v}_x + \mathbf{v}_y + \mathbf{v}_f) - (\|\mathbf{v}_x\|^2 + \|\mathbf{v}_y\|^2 + \|\mathbf{v}_f\|^2) \right] =$$

$$= \frac{1}{2} \sum_{l=1}^d \left(\left(\sum_{k=1}^K v_{kl} z_k \right)^2 - \sum_{k=1}^K (v_{kl} z_k)^2 \right)$$

reduces the number of operations

Matrix form of FM

$$\mathbf{z} = \begin{bmatrix} x \\ y \\ f \end{bmatrix} \quad V = \begin{bmatrix} V_x \\ V_y \\ V_f \end{bmatrix} \quad V \in \mathbb{R}^{K \times d} \quad \mathbf{z} \in \mathbb{R}^K:$$


$$r(\mathbf{z}) = \frac{1}{2} \mathbf{z}^T \left(\begin{bmatrix} V_x \\ V_y \\ V_f \end{bmatrix} \begin{bmatrix} V_x^T & V_y^T & V_f^T \end{bmatrix} - \begin{bmatrix} V_x & 0 \\ 0 & V_y & V_f \end{bmatrix} \begin{bmatrix} V_x & 0 \\ 0 & V_y & V_f \end{bmatrix}^T \right) \mathbf{z} =$$

$$= \frac{1}{2} \mathbf{z}^T \begin{bmatrix} 0 & V_x V_y^T & V_x V_f^T \\ V_y V_x^T & 0 & V_y V_f^T \\ V_f V_x^T & V_f V_y^T & 0 \end{bmatrix} \mathbf{z} = \mathbf{z}^T H \mathbf{z}$$

Connection between FM, LightFM and MF

$$V = \begin{bmatrix} V_x \\ V_y \\ V_{f_x} \\ V_{f_y} \end{bmatrix}, \quad \mathbf{z} = \begin{bmatrix} x \\ y \\ f_x \\ f_y \end{bmatrix}$$



$$V \in \mathbb{R}^{K \times d}, \quad K = M + N + m_x + n_y$$

$$H = \frac{1}{2} \begin{bmatrix} 0 & V_x V_y^T & V_x V_{f_x}^T & V_x V_{f_y}^T \\ V_y V_x^T & 0 & V_y V_{f_x}^T & V_y V_{f_y}^T \\ V_{f_x} V_x^T & V_{f_x} V_y^T & 0 & V_{f_x} V_{f_y}^T \\ V_{f_y} V_x^T & V_{f_y} V_y^T & V_{f_y} V_{f_x}^T & 0 \end{bmatrix}$$

Challenges

Finding linear map is a hard task in general:

- it becomes part of a main optimization routine.

If there're too many different types of real features

- the latent space size may explode.

Incorporating side information into correlations

“similarity” of users i and j depends on co-occurrence of items in their preferences

$$C = AA^T = U\Sigma^2U^T \leftrightarrow c_{ij} = a_i^T a_j$$

Key idea: replace scalar products with a bilinear form.

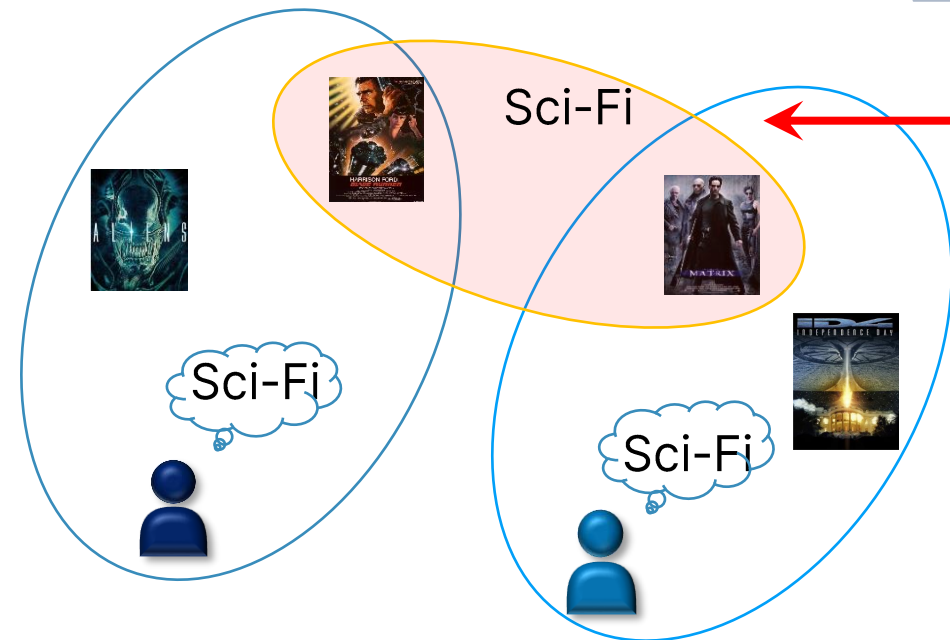
$$\text{sim}(i, j) \sim a_i^T S a_j$$

Similarity matrix S



1			
	1	0.5	
	0.5	1	
			1

“Virtual” connection based on item features.



HybridSVD – formal problem statement

1. Build SPD similarity matrices K, S for users and items based on *side information*.
2. Solve a *new eigen-decomposition* problem:

$$\begin{cases} A S A^T = U \Sigma^2 U^T \\ A^T K A = V \Sigma^2 V^T \end{cases}$$

Σ is a diagonal matrix of singular values.

HybridSVD solution

$$\begin{cases} AA^T = U\Sigma^2U^T \\ A^T A = V\Sigma^2V^T \end{cases} \Rightarrow \begin{cases} ASA^T = U\Sigma^2U^T \\ A^T KA = V\Sigma^2V^T \end{cases}$$

Solution:

via SVD of an auxiliary matrix
[Abdi 2007; Allen et al. 2014]:

$$L_K^T AL_S = \hat{U}\Sigma\hat{V}^T, \quad L_K L_K^T = K, \quad L_S L_S^T = S$$

link to the original latent space

$$L_K^{-T} \hat{U} = U, \quad L_S^{-T} \hat{V} = V$$

Properties:

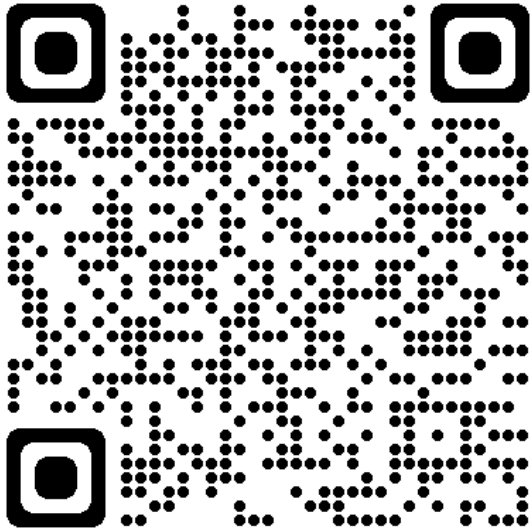
latent space structure:

$$U^T K U = I, \quad V^T S V = I$$

“hybrid” folding-in:

$$\mathbf{p} = L_S^{-T} \hat{V} \hat{V}^T L_S^T \mathbf{a}.$$

Computation example (naïve)



```
# rating matrix
a = np.array([[1, 0, 1, 1, 0],
              [1, 1, 0, 1, 0],
              [1, 0, 0, 1, 0]])
```

```
# vector of a new user's preferences
p = np.array([1, 0, 0, 1, 1])
```

```
# recommendations for the user with standard folding-in approach
recs = p.dot(v).dot(v.T)
recs
# output: [ 0.8904344 ,  0.31234752,  0.31234752,  0.8904344 ,  0. ]
```

```
# ===== Hybrid Model =====

d = 0.5 # off-diagonal similarity factor

# item similarity matrix
# non-zero off-diagonal values denote similarity between items 3 and 5
s = np.array([[1, 0, 0, 0, 0],
              [0, 1, 0, 0, 0],
              [0, 0, 1, 0, d],
              [0, 0, 0, 1, 0],
              [0, 0, d, 0, 1]])

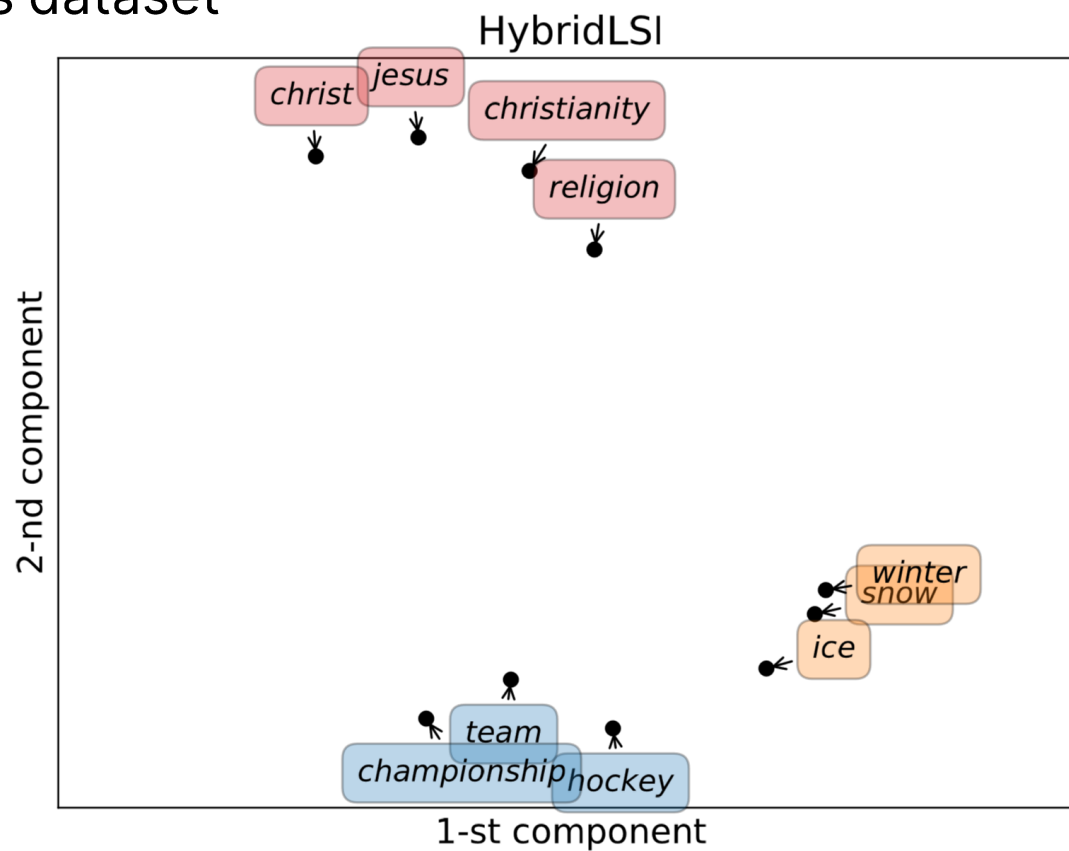
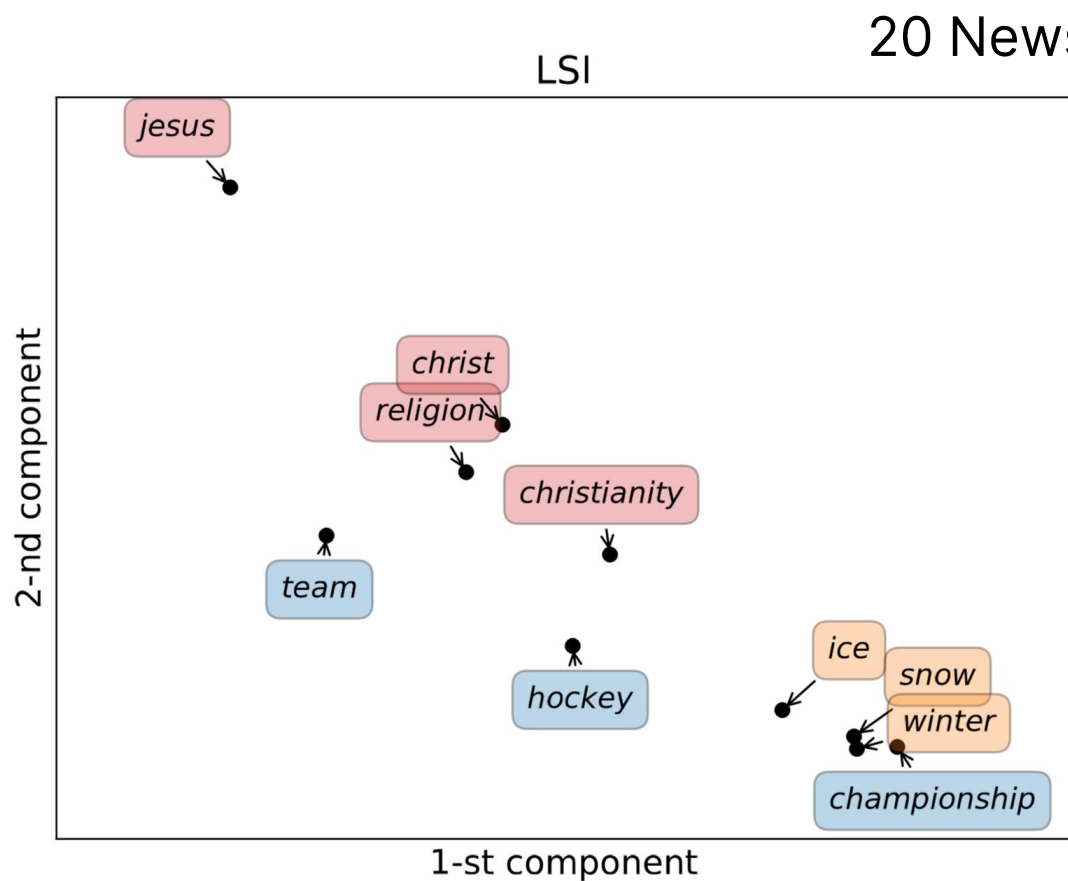
# finding Cholesky factors
L = np.linalg.cholesky(s)
u2, s2, v2 = np.linalg.svd(a.dot(L), full_matrices=False)
v2 = v2.T[:, :rank]

# preparing for hybrid folding-in calculation
lv = L.dot(v2)
rv = spsolve_triangular(csr_matrix(L.T), v2, lower=False)

# recommendations for the user with hybrid model
recs2 = p.dot(lv).dot(rv.T)
recs2
# output: [0.96852129, 0.08973892, 0.58973892, 0.96852129, 0. ]
```

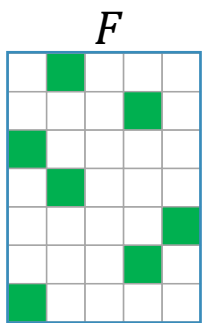
Latent space structure with HybridSVD

- Use general semantic similarity of words based on a global model, e.g. word2vec.

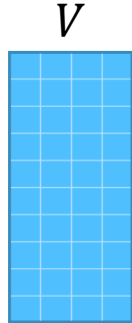


Solving cold start with HybridSVD

Item features
(sparse / dense)



embeddings



mapping
 W

Using the S-orthogonality property:

$$VW = F \rightarrow W = V^T S F$$

← analytic solution

Given any feature vector f , we find the corresponding embedding v from:

$$W^T v = f$$

← quick to solve

Relevance scores prediction:

$$p = U \Sigma v = A V v$$

Works for PureSVD as well by setting $S = I, K = I$.

Notes on HybridSVD scalability

Auxiliary matrix $L_K^T A L_S$ is likely to become dense:

- can be avoided via matvec in the Lanczos procedure.

Building similarity matrices K and S can also be prohibitively expensive:

- use sparse QR / Cholesky decompositions (via [scikit-sparse](#)) or
- compute similarities in the reduced dimension + QR /Cholesky via fast symmetric factorization [Ambikasaran et al. 2014].

03



Context-awareness



Context vs Content

There's no sharp boundary!

Content is typically:

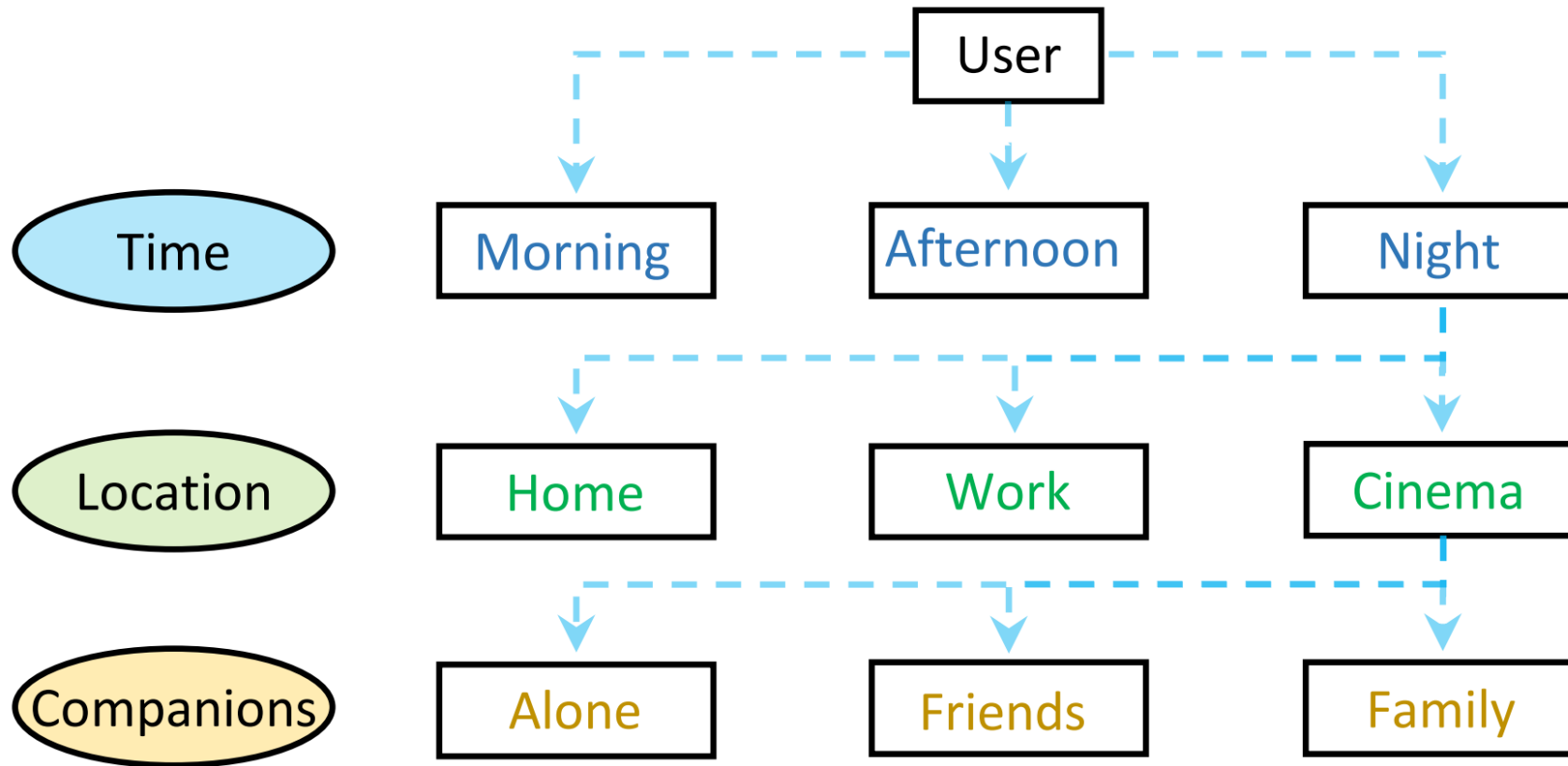
- static,
- fixed to an entity it describes.

Context is typically:

- situational / dynamic,
- characterizes interaction between entities.

Context or content?
(user, movie, *genre*, *tag*)

Examples of context in RecSys



Also: folksonomies, cross-domain RS, temporal models, etc.

Contextual recommendations

$f_U: \text{User} \times \text{Item} \rightarrow \text{Relevance}$

↓

$f_U: \text{User} \times \text{Item} \times \text{Context} \rightarrow \text{Relevance}$

↓

$f_U: \text{User} \times \text{Item} \times \text{Context}_1 \times \cdots \times \text{Context}_f \rightarrow \text{Relevance}$

Suggest an approach for contextual modeling.

Higher order contextual models

FM models pairwise (or 2-way) relations

$$f_U: \text{User} \times \text{Item} \times \text{Context} \rightarrow \text{Relevance}$$

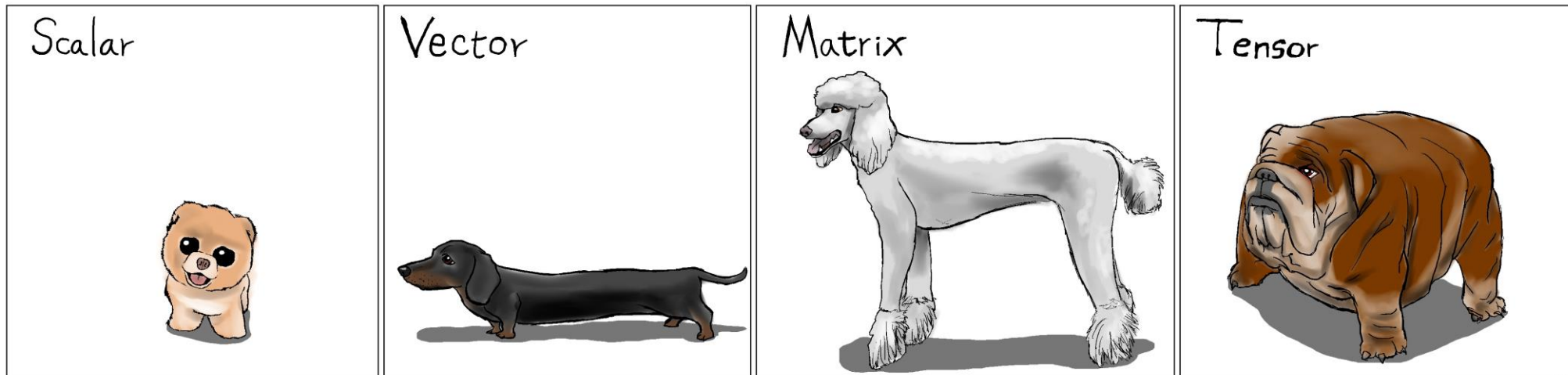
How does FM capture triplet interactions?

Multiway (multi-aspect) learning

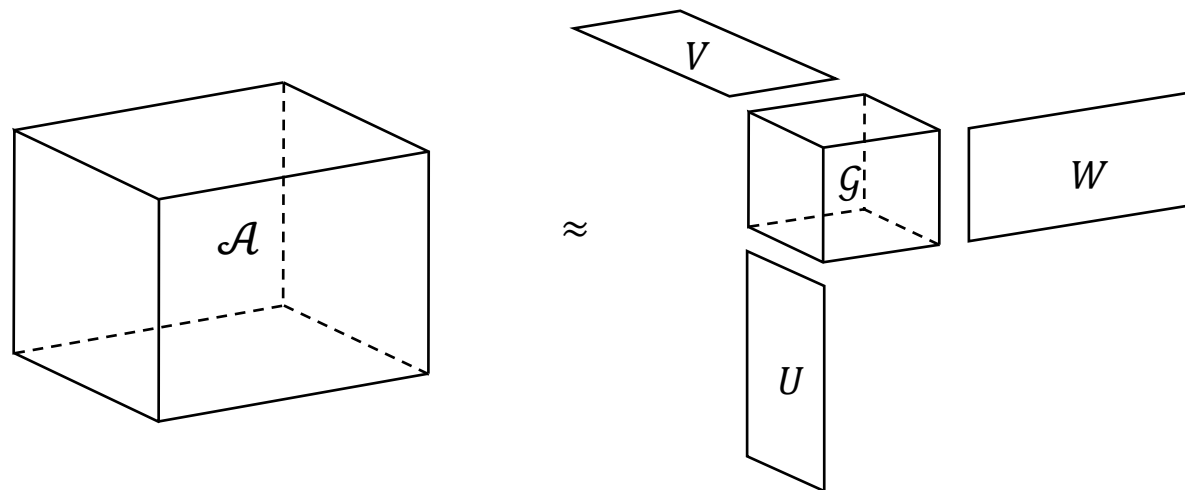
$$f_U: \text{User} \times \text{Item} \times \text{Context}_1 \times \dots \times \text{Context}_f \rightarrow \text{Relevance}$$

In the paradigm of contextual modeling:

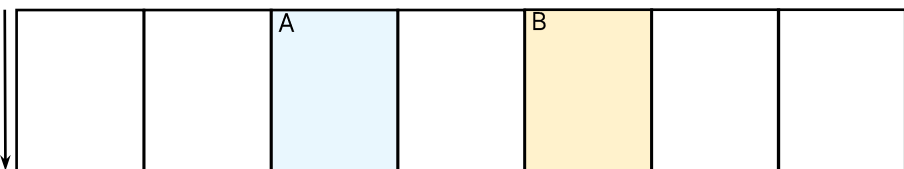
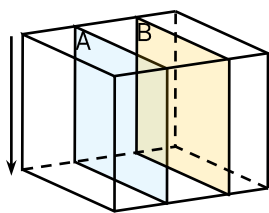
- data seem to be better described via multiway relations
- multiway relations can be naturally encoded via tensor formats



Higher Order SVD (HOSVD)



Tensor matricization
(unfolding)



\times_n denotes a tensor product, e.g., for tensor \mathcal{A} and matrix B :

$$[\mathcal{A} \times_n B]_{i_1 \dots i_{n-1} j i_{n+1} \dots i_m} = \sum_i a_{i_1 i_2 \dots i_m} b_{ji_n}$$

$$\underset{G, U, V, W}{\text{minimize}} \|\mathcal{A}_0 - \mathcal{R}\|_F^2$$

$$\mathcal{R} = \mathcal{G} \times_1 U \times_2 V \times_3 W$$

$$U^\top U = I, \quad V^\top V = I, \quad W^\top W = I$$

HOSVD procedure

compute from tensor unfoldings:

$$U \leftarrow d_1 \text{ left singular vectors of } \mathcal{A}_0^{(1)}$$

$$V \leftarrow d_2 \text{ left singular vectors of } \mathcal{A}_0^{(2)}$$

$$W \leftarrow d_3 \text{ left singular vectors of } \mathcal{A}_0^{(3)}$$

$$\mathcal{G} \leftarrow \mathcal{A}_0 \times_1 U \times_2 V \times_3 W$$

return U, V, W, G

Contextual top- n recommendations scenarios

recommend the best items within a selected context

- e.g., best restaurant based on location

$$\text{toprec}(u, c, n) := \arg \max_i^n r_{uic}$$

recommend the best context for a target item

- e.g., find best distribution channel

$$\text{toprec}(u, i, n) := \arg \max_c^n r_{uic}$$

User feedback peculiarities

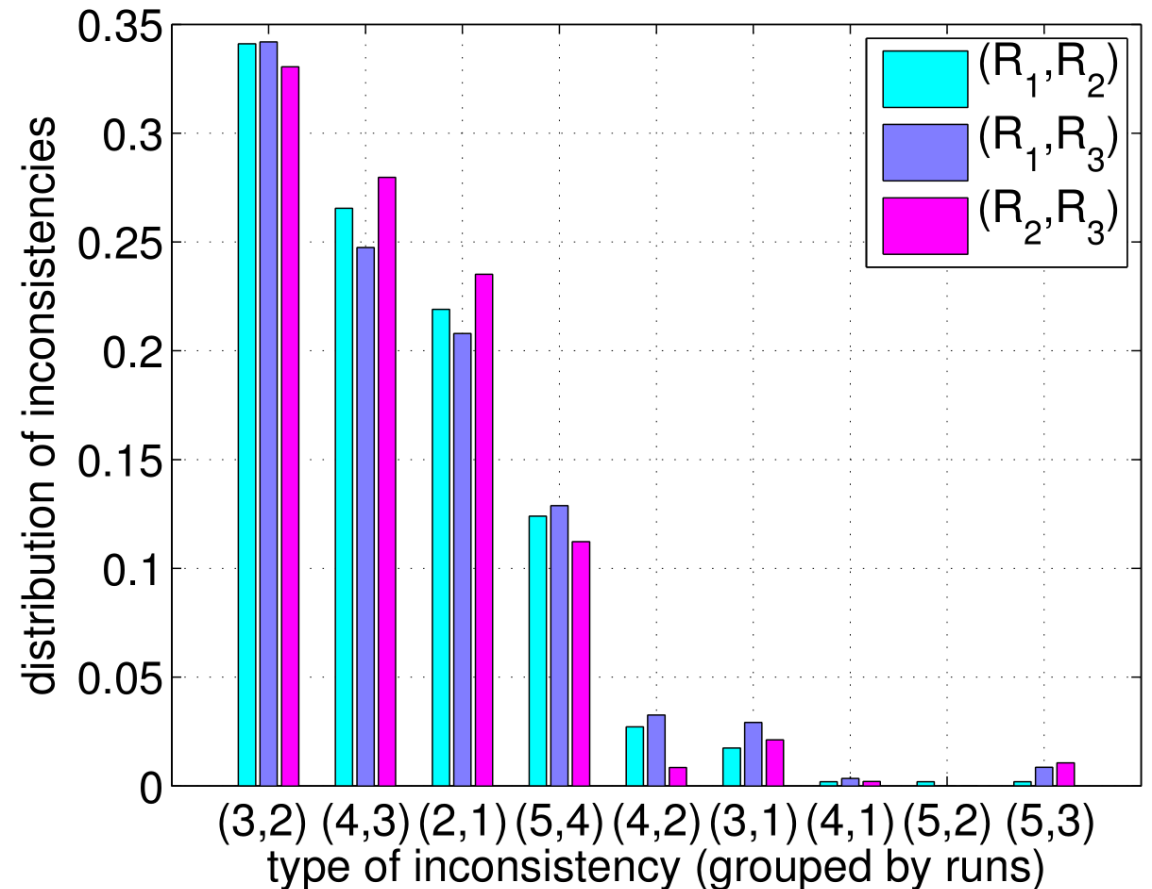


2.5x better?



Traditional recommender models treat ratings as cardinal numbers.

From neoclassical economics: utility is an ordinal concept.

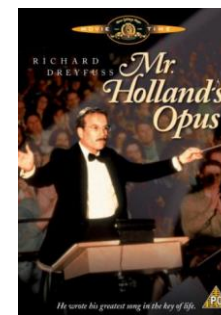


Negative feedback problem

What is likely to be recommended in this case?

User feedback is negative!

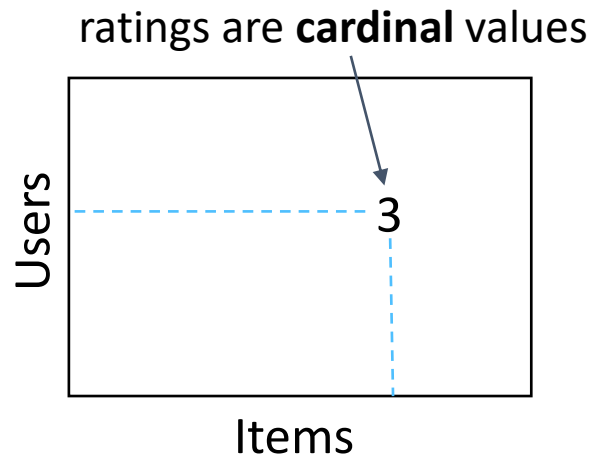
Probably the user doesn't like criminal movies.



Redefining the utility function

Standard MF model

$$f_U: \text{User} \times \text{Item} \rightarrow \text{Rating}$$

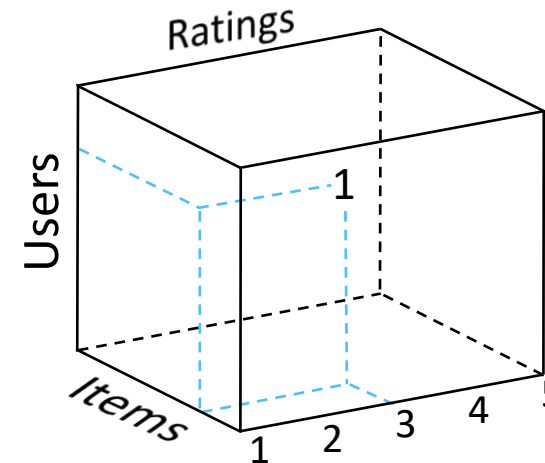


$$\|A_0 - R\|_F^2 \rightarrow \min$$

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

Collaborative Full Feedback model – CoFFee*

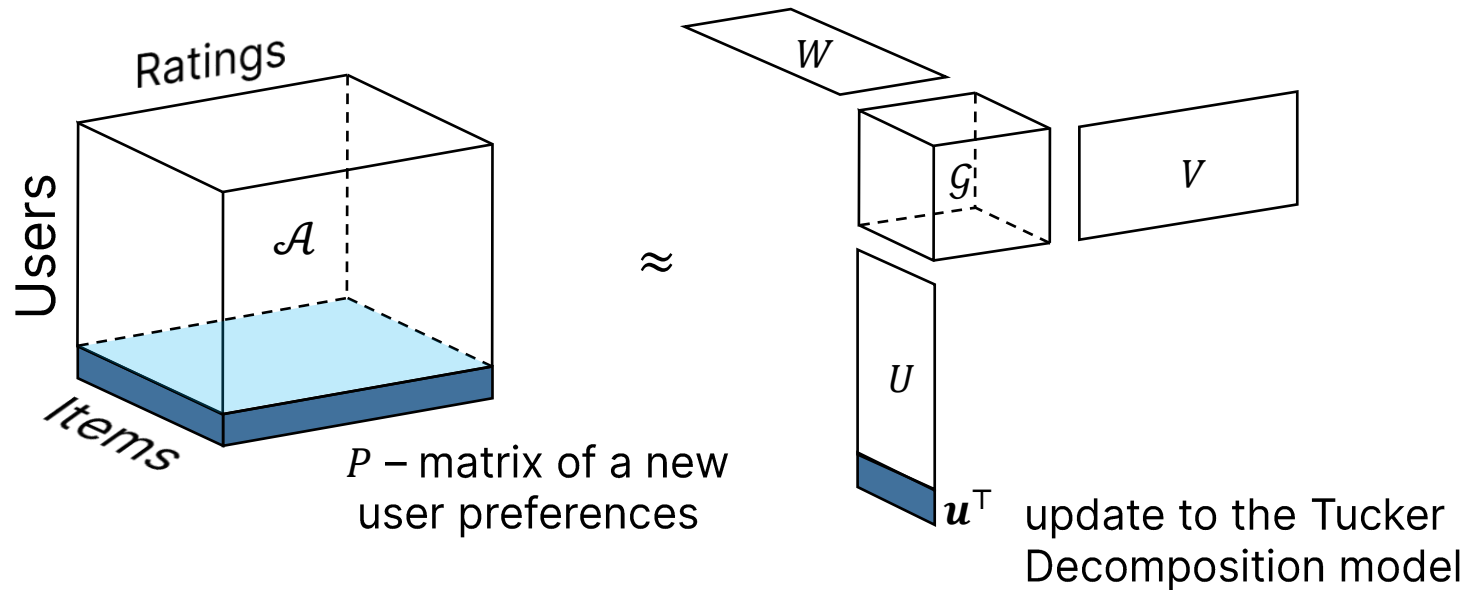
$$f_U: \text{User} \times \text{Item} \times \text{Rating} \rightarrow \text{Relevance Score}$$



$$\|A_0 - R\|_F^2 \rightarrow \min$$

$$R = G \times_1 U \times_2 V \times_3 W$$

Higher order folding-in



$$R \approx VV^T P W W^T \quad \text{predictions matrix}$$

Compare to SVD:

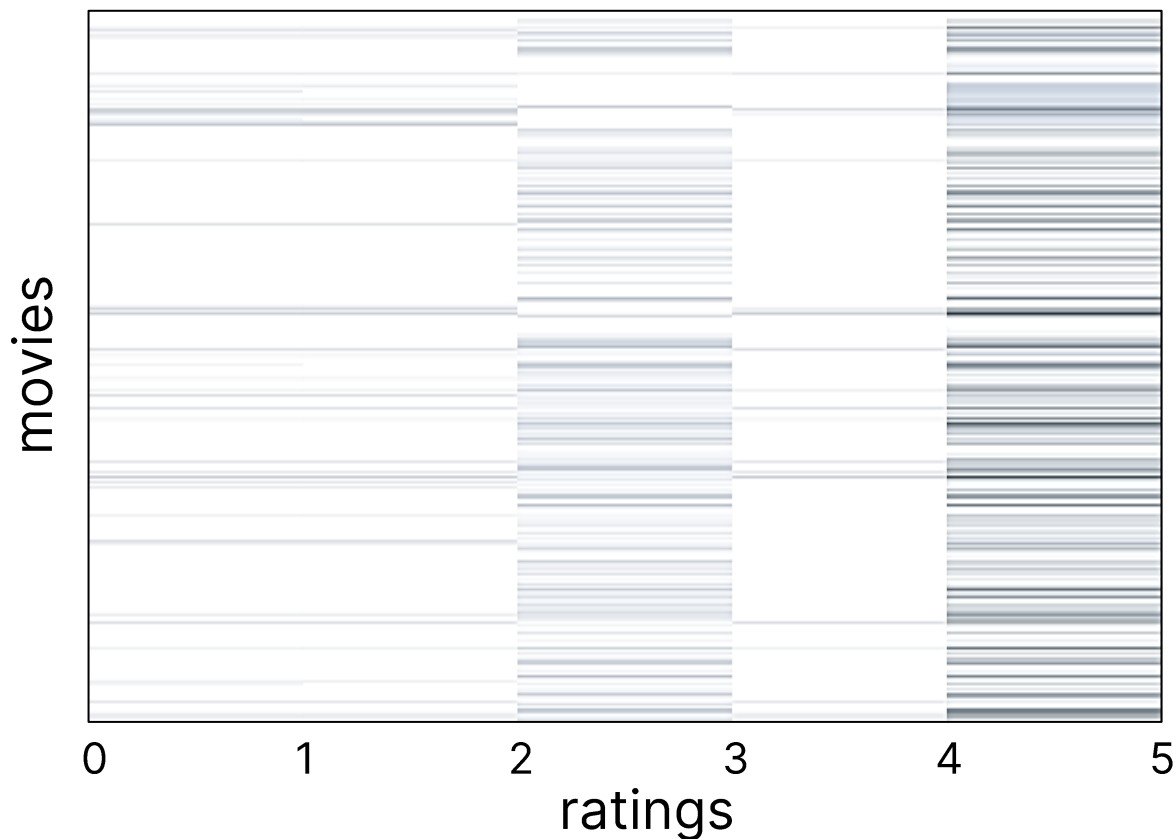
$$r = VV^T p \quad \text{predictions vector}$$

“Shades” of ratings

More dense colors correspond to higher relevance score.

$$R = VV^T PWW^T$$

matrix of known
user preferences



- Solves both tasks:
- ranking
 - rating prediction

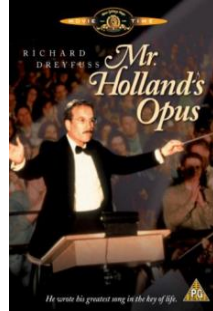
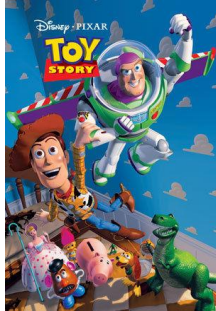


Granular view of user preferences,
concerning all possible ratings.



Model is equally sensitive
to any kind of feedback.

Warm-start with CoFFee



Uncovers new recommendation modes:

“users who like this also like...”



“users who **dislike** this, do like...”

	Scarface ★☆☆☆☆	LOTR: The Two Towers ★☆☆☆☆	Star Wars: Episode VII - The Force Awakens ★★★★★
CoFFee	Toy Story Mr. Holland's Opus Independence Day	Net, The Cliffhanger Batman Forever	Dark Knight, The Batman Begins Star Wars: Episode IV - A New Hope
SVD	Reservoir Dogs Goodfellas Godfather: Part II, The	LOTR: The Fellowship of the Ring Shrek LOTR: The Return of the King	Dark Knight, The Inception Iron Man



Artificial Intelligence Research Institute

airi.net



[airi_research_institute](https://t.me/airi_research_institute)



[AIRI_Institute](https://vk.com/AIRI_Institute)



[AIRI_Institute](https://www.youtube.com/AIRI_Institute)



[AIRI_inst](https://twitter.com/AIRI_inst)



[artificial-intelligence-research-institute](https://www.linkedin.com/company/artificial-intelligence-research-institute)