

SLIM: Sparse Linear Methods for Top-N Recommender Systems

Group 3: Kunj, Aditi, Jyoti, Garvika

based on the research paper by
Xia Ning and George Karypis

Top-N recommender systems

Memory-based Collaborative Filtering

Item-based k-Nearest-Neighbors: fast but low quality

Model-based Collaborative Filtering

Singular Value Decomposition: high quality but slow

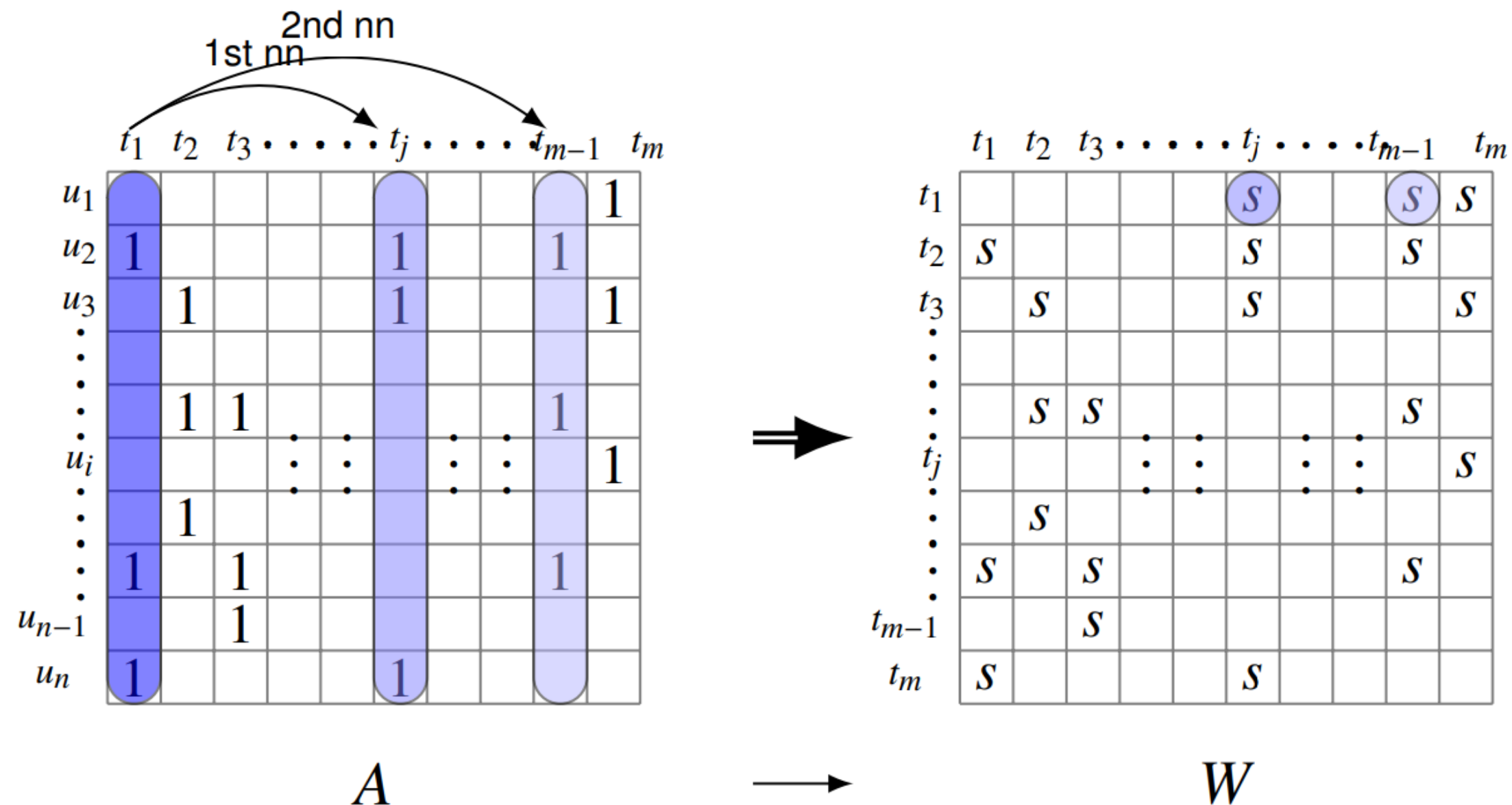
SLIM: Sparse Linear Methods

fast and high quality

Memory-based CF: itemkNN

Identify a set of similar items

Item-item similarity calculated from A using cosine

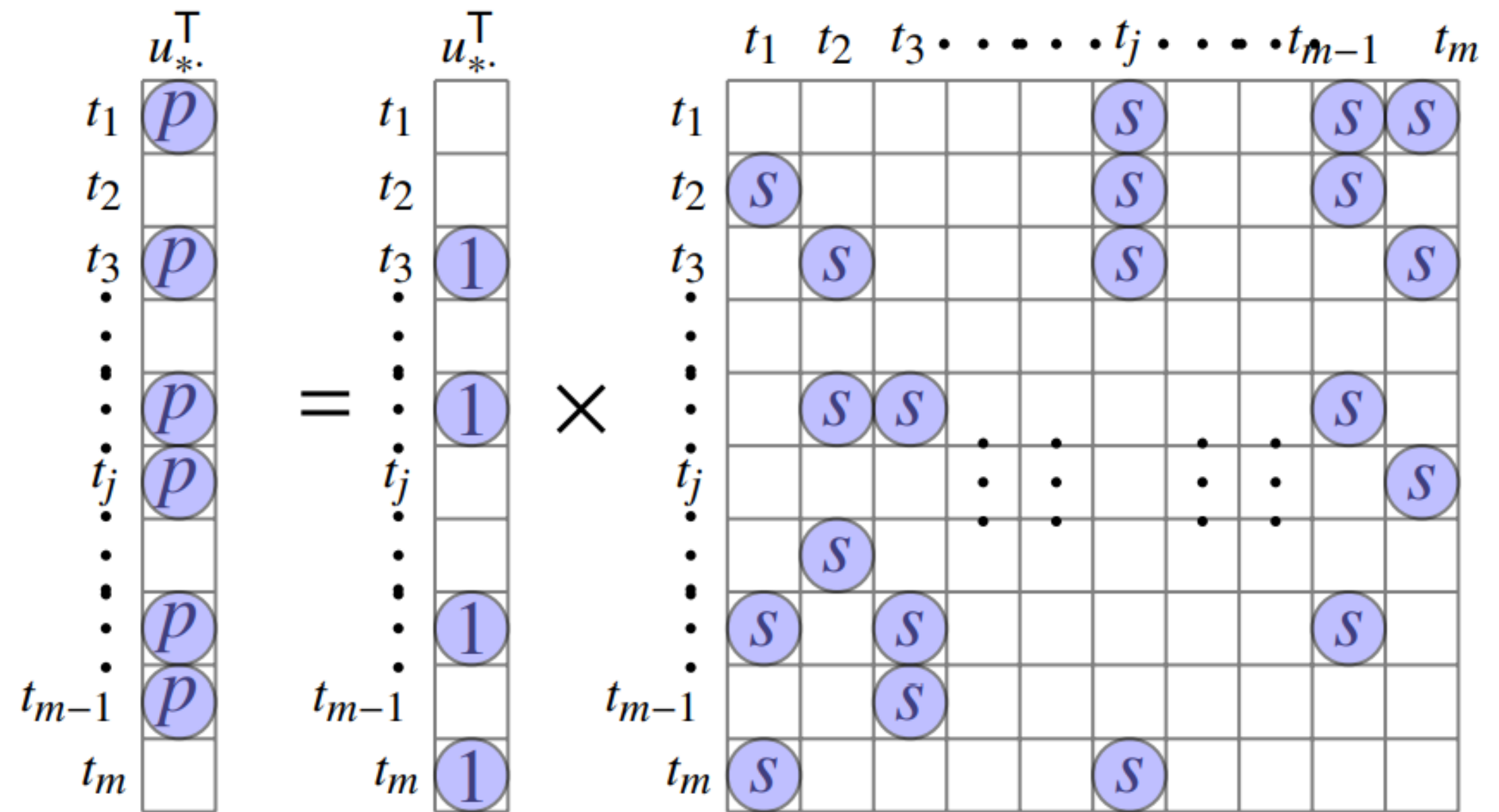


Recommend similar items to what the user has purchased

Fast: sparse item neighborhood

Low quality: no knowledge is learned

$$\tilde{a}_i^T = a_i^T + W$$



$$R(a, u_i) = \frac{\sum_j S(a, b) R(b, u_i)}{\sum_j S(a, b)}$$

$R(a, u_i)$: rating for movie **a** by user **u_i**

$S(a, b)$: similarity b/w movie **a** and **b**

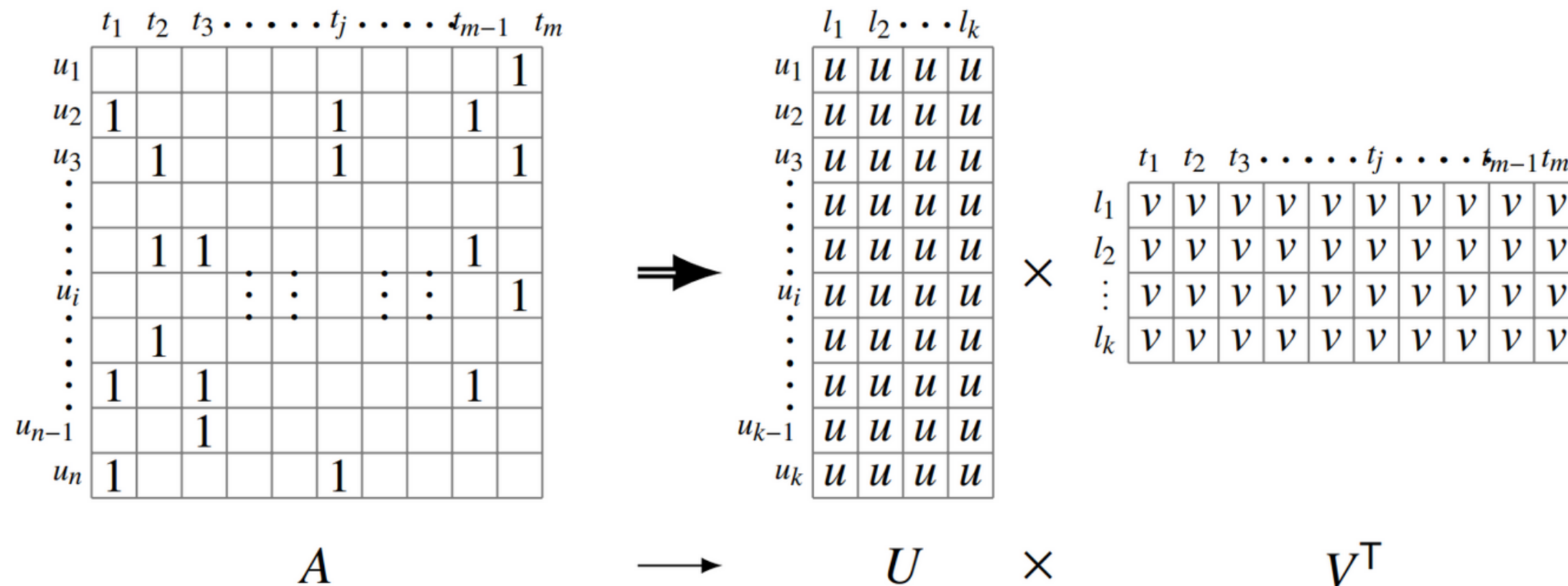
$j \in J$, where **J** is the set of similar movies to **a**

Model-based CF: SVD

Factorize A into low-rank user and item factors that represent user and item characteristics in a common latent space

Formulated as an optimization problem

$$\underset{U, V^T}{\text{minimize}} \quad \frac{1}{2} \|A - UV^T\|_F^2 + \frac{\beta}{2} \|U\|_F^2 + \frac{\lambda}{2} \|V^T\|_F^2$$



Prediction: dot product in the latent space

Slow: dense U and V^T

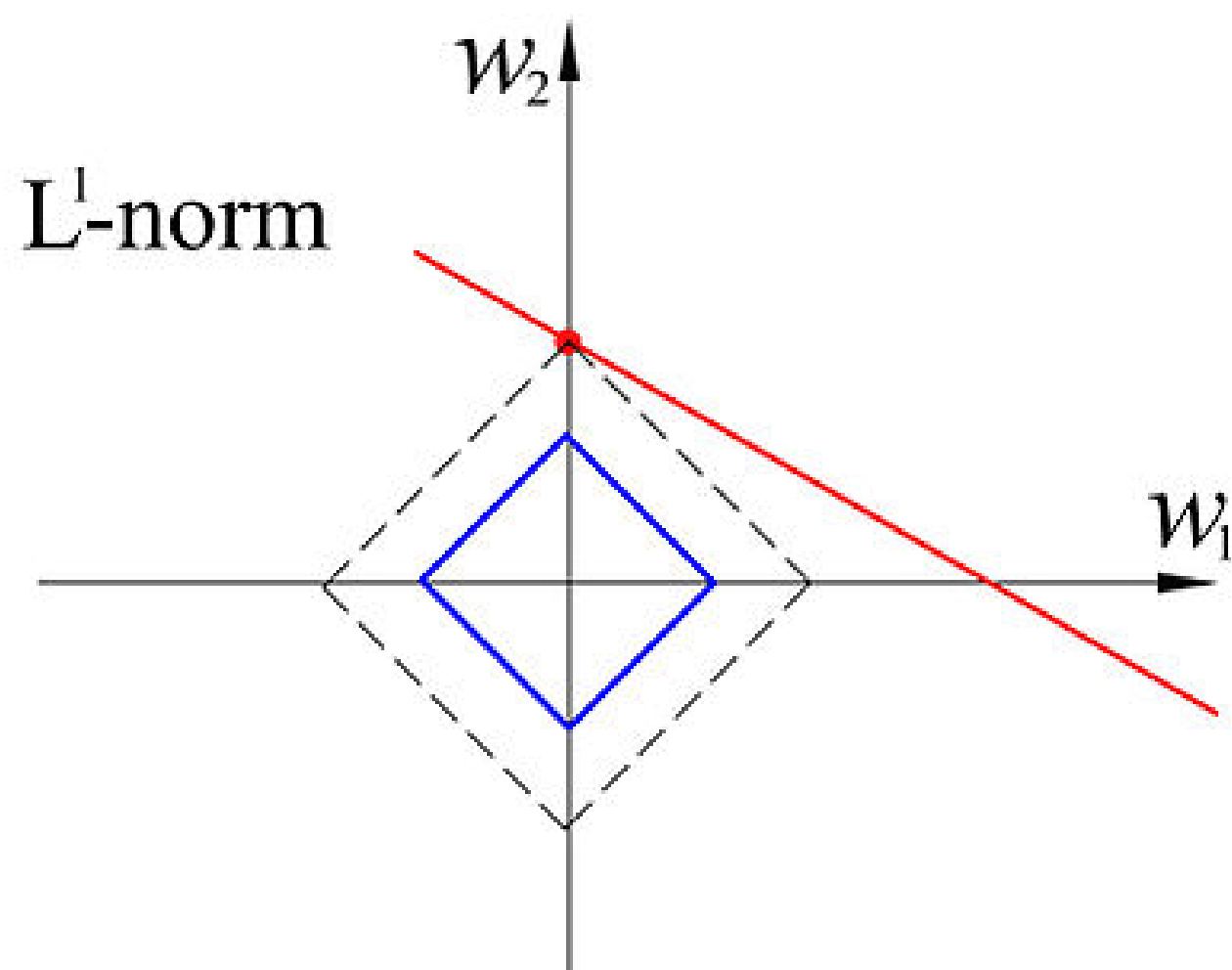
High quality: user tastes and item properties are learned

$$\widetilde{a_{ij}} = U_i^T \cdot V_j$$

The diagram illustrates the dot product calculation for a prediction. It shows a vertical vector U_i^T (labeled u_{*}^T) and a horizontal vector V_j (labeled u_{*}) being multiplied to produce a scalar prediction. The vertical vector U_i^T has elements p and is indexed by $t_1, t_2, t_3, \dots, t_j, \dots, t_{m-1}, t_m$. The horizontal vector V_j has elements u and is indexed by l_1, l_2, \dots, l_k . The resulting prediction is shown as a scalar value.

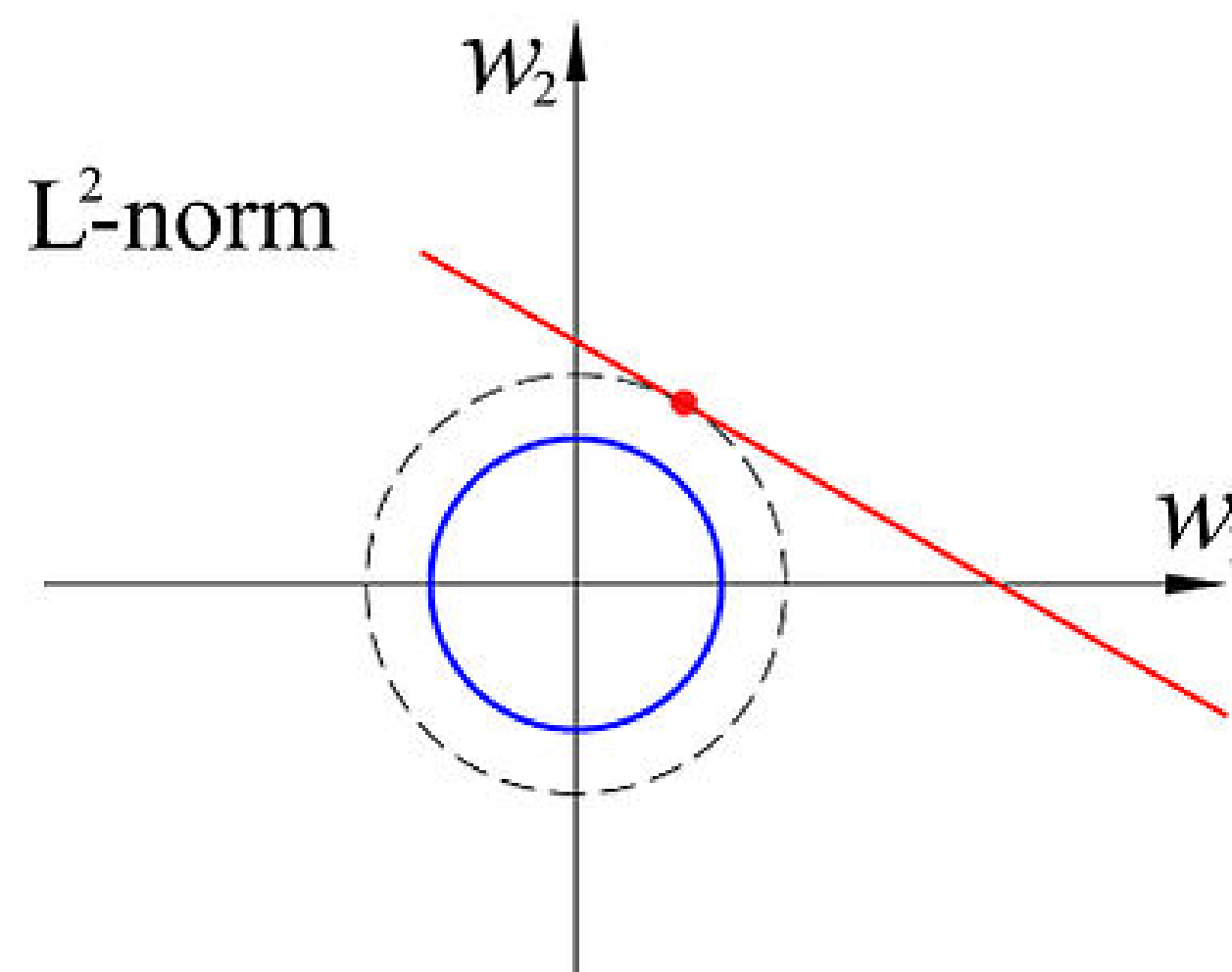
L1 Norm

$$\|\mathbf{x}\|_1 = |x_1| + |x_2| + \cdots + |x_n|$$



L2 Norm

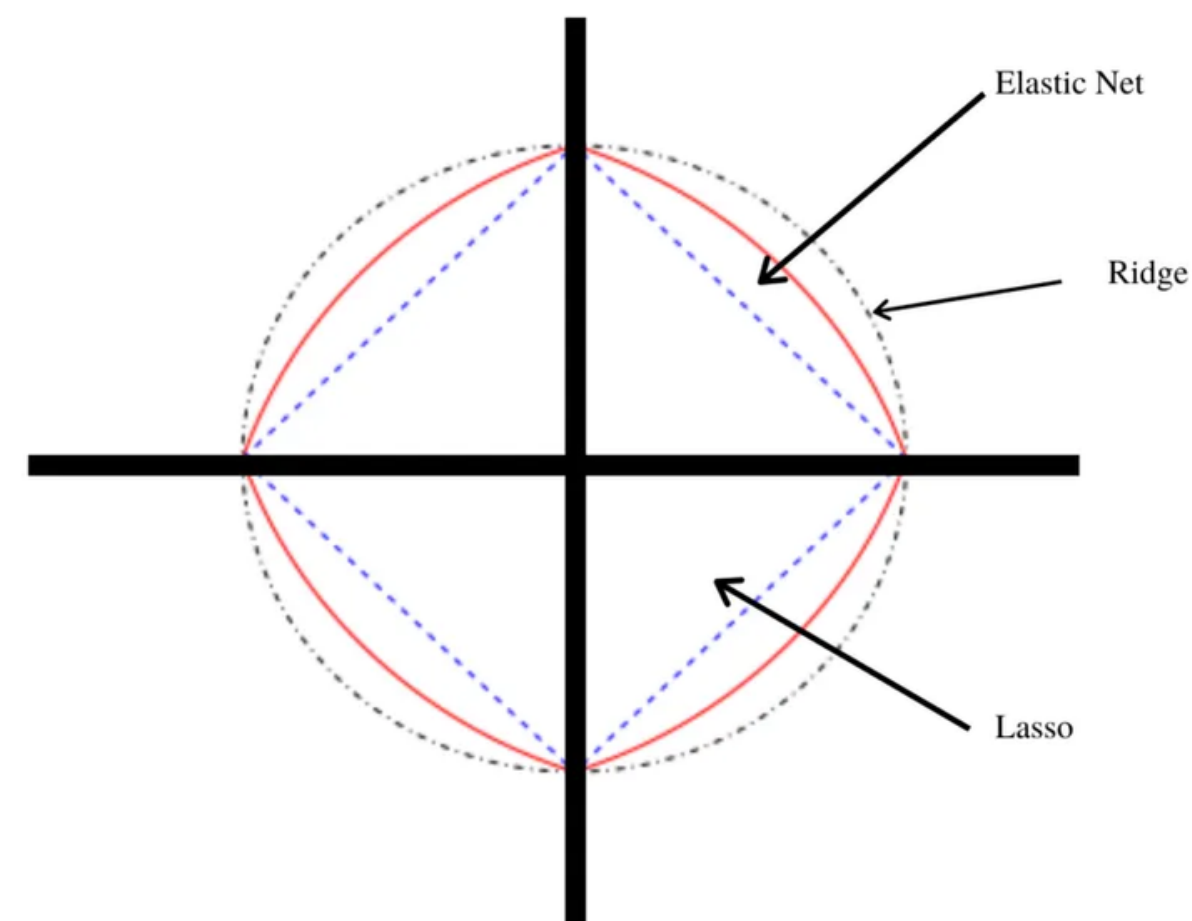
$$\|\mathbf{x}\|_2 = \sqrt{x_1^2 + x_2^2 + \cdots + x_n^2}$$



$$\hat{\beta}_{\text{lasso}} = \arg \min_{\beta} \{ \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \}$$

$$\hat{\beta}_{\text{ridge}} = \arg \min_{\beta} \{ \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \}$$

$$\hat{\beta} = \arg \min_{\beta} \{ \|y - X\beta\|_2^2 + \lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1 \}$$



Minimize p_u and q_i over the loss function

Lasso

$$\min_{(p_u, q_i)} \sum_{(u, i) \in K} \left(r_{ui} - p_u^T q_i \right)^2 + \lambda_1 \|p_u\|_1 + \lambda_2 \|q_i\|_1$$

Ridge

$$\min_{(p_u, q_i)} \sum_{(u, i) \in K} \left(r_{ui} - p_u^T q_i \right)^2 + \rho_1 \|p_u\|_2^2 + \rho_2 \|q_i\|_2^2$$

Elastic Net

$$\begin{aligned} \min_{(p_u, q_i)} \sum_{(u, i) \in K} \left(r_{ui} - p_u^T q_i \right)^2 + \rho_1 \|p_u\|_2^2 + \rho_2 \|q_i\|_2^2 \\ + \lambda_1 \|p_u\|_1 + \lambda_2 \|q_i\|_1 \end{aligned}$$

Motivations:

Recommendations generated fast

High quality recommendations

Key ideas:

retain the nature of itemkNN: sparse W

optimize the recommendation performance: learn W from A

sparsity structures

coefficient values

Def	Descriptions
u_i	user
t_j	item
\mathcal{U}	all users ($ \mathcal{U} = n$)
\mathcal{T}	all items ($ \mathcal{T} = m$)
A	user-item purchase/rating matrix, size $n \times m$
W	item-item similarity matrix/coefficient matrix
\mathbf{a}_i^\top	The i -th row of A , the purchase/rating history of u_i on \mathcal{T}
\mathbf{a}_j	The j -th column of A , the purchase/rating history of \mathcal{U} on t_j

Row vectors are represented by having the transpose T superscript, otherwise by default they are column vectors.

$$\begin{array}{ll}
\underset{W}{\text{minimize}} & \frac{1}{2} \|A - AW\|_F^2 + \frac{\beta}{2} \|W\|_F^2 + \lambda \|W\|_1 \\
\text{subject to} & W \geq 0 \\
& \text{diag}(W) = 0,
\end{array}$$

$$\begin{aligned}
& \underset{\mathbf{w}_j}{\text{minimize}} && \frac{1}{2} \|\mathbf{a}_j - A\mathbf{w}_j\|_2^2 + \frac{\beta}{2} \|\mathbf{w}_j\|_2^2 + \lambda \|\mathbf{w}_j\|_1 \\
& \text{subject to} && \mathbf{w}_j \geq \mathbf{0} \\
& && w_{j,j} = 0,
\end{aligned}$$

Columns of W are independent, hence are easy to parallelize

Learning W using coordinate descent

- The objective function is separable.
- Input features have sparse representations.
- Objective function is non-differentiable
- If the dimensionality of the function is high, in computing the gradient for all parameters can be very expensive.
- If the objective function is not smooth or has discontinuity, then Gradient descent may get stuck in the local minima and oscillate around the optimal solution. In such cases Coordinate descent can be more robust as it updates the parameters in a sequential manner and can avoid getting stuck in a local minima.
- When input features are highly co-related, in such cases GD may suffer slow convergence due to the so-called zig-zagging effect

The SLIM method uses specific regularizers, namely the L1-norm and L2-norm regularizations, for the following reasons:

- **Sparsity Introduction:** The L1-norm regularization is known to introduce sparsity into the solutions. By using the L1-norm of W as a regularizer, SLIM ensures that the learned aggregation coefficient matrix W is sparse. This sparsity allows SLIM to generate recommendations very quickly.
- **Model Complexity Control:** The L2-norm regularization, also known as ridge regression, helps in controlling the complexity of the model and prevents overfitting. By including the L2-norm of W as another regularizer, SLIM ensures that the model does not become too complex, leading to better generalization performance.
- **Implicit Grouping:** The combination of L1-norm and L2-norm regularizations together implicitly groups correlated items in the solutions. This means that items that are related or have similar characteristics are grouped together in the learned aggregation coefficient matrix W .

- **Noise Reduction:** The L1-norm regularization enforces sparsity in W , which helps in capturing the most informative signals while discarding noise. This leads to high-quality recommendations as only relevant information is considered in the aggregation process.

Overall, the choice of L1-norm and L2-norm regularizations in SLIM is aimed at promoting sparsity, controlling model complexity, grouping correlated items, and reducing noise in the recommendations, ultimately leading to efficient and high-quality top-N recommendations.

Coordinate descent

This suggests that for the problem

$$\min_x f(x)$$

where $f(x) = g(x) + \sum_{i=1}^n h_i(x_i)$, with g convex and differentiable and each h_i convex, we can use **coordinate descent**: let $x^{(0)} \in \mathbb{R}^n$, and repeat

$$x_i^{(k)} = \operatorname{argmin}_{x_i} f(x_1^{(k)}, \dots, x_{i-1}^{(k)}, x_i, x_{i+1}^{(k-1)}, \dots, x_n^{(k-1)}),$$
$$i = 1, \dots, n$$

for $k = 1, 2, 3, \dots$

Important note: we always use **most recent information** possible

Consider the **lasso** problem:

$$\min_{\beta} \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1$$

Note that the nonsmooth part is separable: $\|\beta\|_1 = \sum_{i=1}^p |\beta_i|$

Minimizing over β_i , with β_j , $j \neq i$ fixed:

$$0 = X_i^T X_i \beta_i + X_i^T (X_{-i} \beta_{-i} - y) + \lambda s_i$$

where $s_i \in \partial |\beta_i|$. Solution is simply given by soft-thresholding

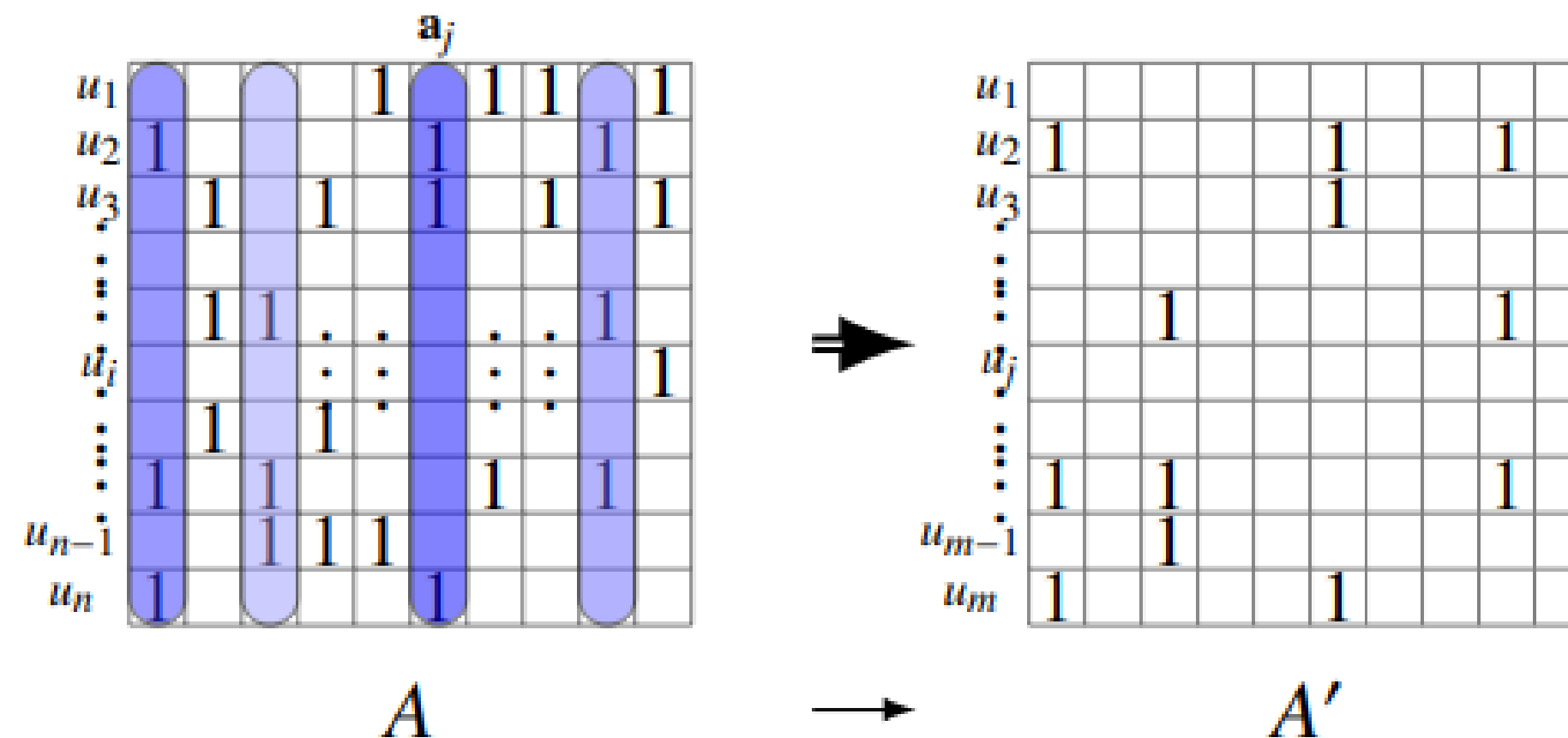
$$\beta_i = S_{\lambda / \|X_i\|_2^2} \left(\frac{X_i^T (y - X_{-i} \beta_{-i})}{X_i^T X_i} \right)$$

Repeat this for $i = 1, 2, \dots, p, 1, 2, \dots$

$$\underset{\mathbf{w}_j}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{a}_j - A\mathbf{w}_j\|_2^2 + \frac{\beta}{2} \|\mathbf{w}_j\|_2^2 + \lambda \|\mathbf{w}_j\|_1$$

□ fsSLIM: SLIM with *feature selection*

- Prescribe the potential non-zero structure of \mathbf{w}_j
- Select a subset of columns from A
 - itemkNN item-item similarity matrix



$$\underset{\mathbf{w}_j}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{a}_j - A'\mathbf{w}_j\|_2^2 + \frac{\beta}{2} \|\mathbf{w}_j\|_2^2 + \lambda \|\mathbf{w}_j\|_1$$

#users	#items	#ratings	rsize	csize	density
1292	313	23806	18.13	74.83	0.0589
2352	620	60976	25.55	96.92	0.0418
4203	1830	181442	42.65	97.97	0.0236
6046	3666	382038	62.52	103.11	0.0172
10213	12971	1276344	123.71	97.40	0.0096

Beer dataset with ~1.6 Million datapoints
Ratings from 1 to 5 with 0.5 granularity

- ❑ Evaluation methodology: Leave-One-Out cross validation
- ❑ Evaluation metrics

- ❑ Hit Rate:
$$HR = \frac{\text{\#hits}}{\text{\#users}}$$

- ❑ Average Reciprocal Hit-Rank (ARHR) [2]:

$$ARHR = \frac{1}{\text{\#users}} \sum_{i=1}^{\text{\#hits}} \frac{1}{p_i}$$

Conclusion

- Recommendation score for a new item is computed by aggregating scores of other items.
- A sparse aggregation coefficient matrix W is optimized for SLIM to enhance speed of aggregation.
- W is determined by solving an optimization problem with l_1 -norm and l_2 -norm regularization to induce sparsity in W .
- The process is designed to be fast and efficient.
- Applying stronger l_1 -norm regularization (larger λ) leads to a sparser W and reduced recommendation time.
- Optimal recommendation quality is obtained when both regularization parameters β and λ are non-zero.
- Recommendation quality varies smoothly with changes in regularization parameters β and λ .

Thank You