

# Kernelized Matrix Factorization for Collaborative Filtering

Xinyue Liu\*, Charu Aggarwal<sup>†</sup>, Yu-Feng Li<sup>#</sup>

Xiangnan Kong\* , Xinyuan Sun\*, Saket Sath<sup>†</sup>

\* Worcester Polytechnic Institute    ‡ Nanjing University

† IBM T.J. Watson Research Center

# Collaborative Filtering Problem

	Movie A	Movie B	Movie C	Movie D
User 1	5	3	?	?
User 2	?	4	2	4
User 3	3	5	?	2
User 4	?	?	?	5
User 5	?	1	?	?

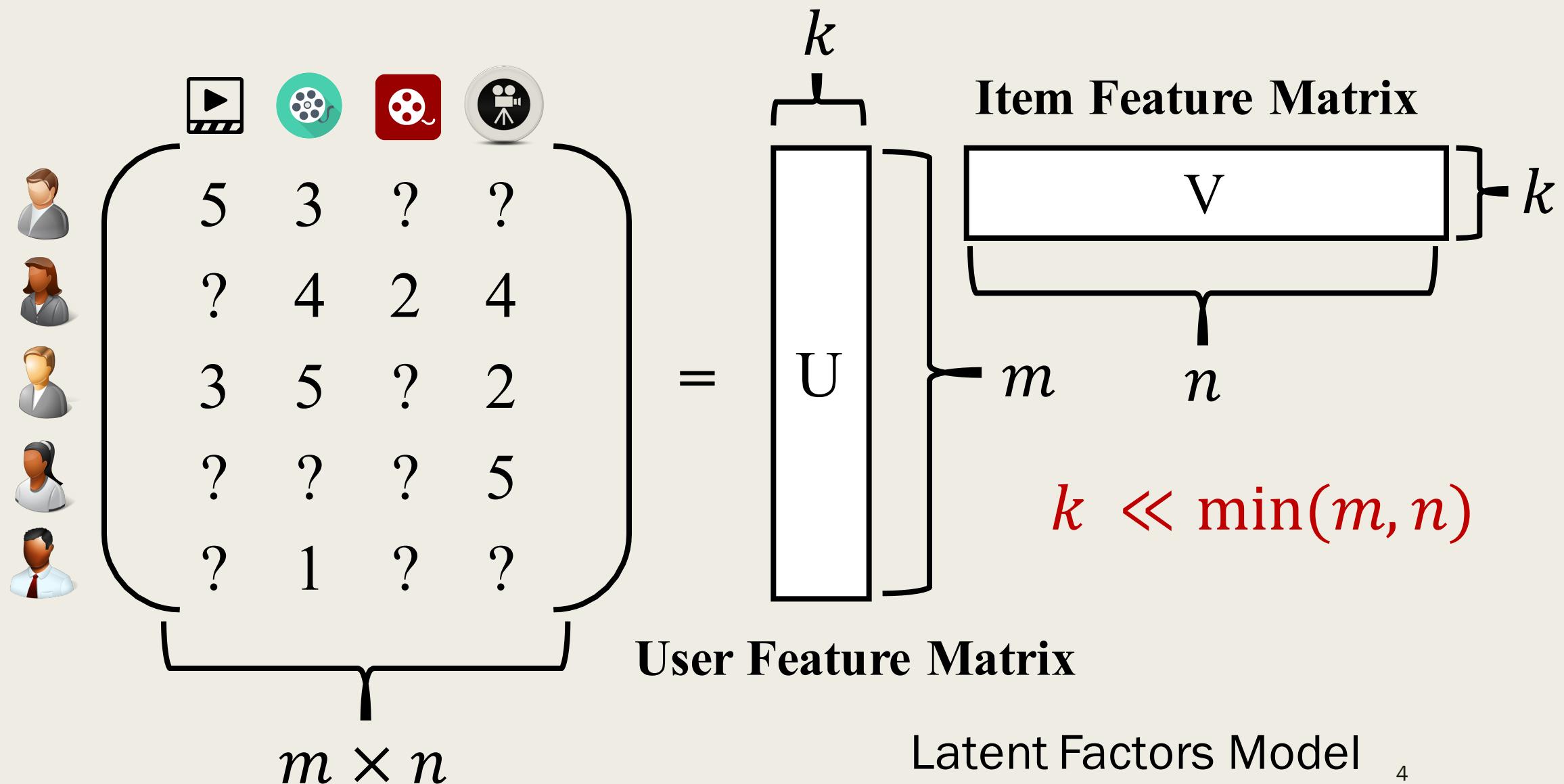
- Given *users* and *items* in a recommender system.
- Assume each user can *rate* each item only *once*.
- Then we can cast the user-item feedbacks into a *matrix*.
- The matrix is always very *sparse*.
  - Netflix Prize dataset (1.18%)
  - MovieLens dataset (6.3%)
- The task is to estimate the unknown ratings (?) based on the observed ones.

# Applications

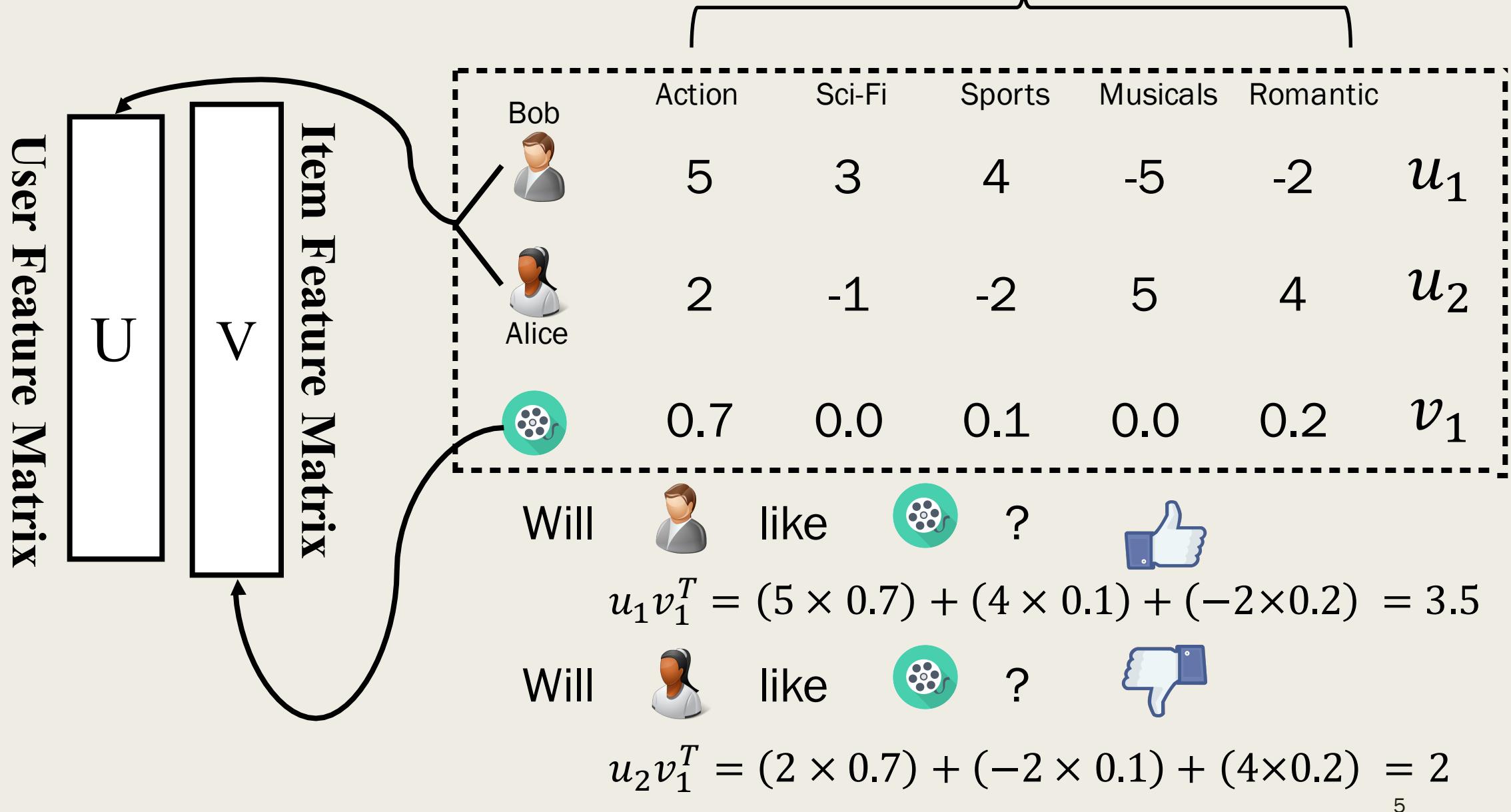
- ❖ Netflix: Movie Recommendations
- ❖ Amazon: Merchandise Recommendations
- ❖ Pandora: Music Recommendations
- ❖ Yelp: Restaurant Recommendations
- ❖ Educational system, dating website, everything where you have choices...



# Matrix Factorization



# The Latent Factors



# Find the Latent Factors

- The goal is to find  $U$  and  $V$  such that:

$$\min_{U,V} \sum_{(i,j) \in R} (r_{ij} - u_i v_j^T)^2$$

Can be solved via  
gradient descent



users

	items									
1		3		5			5		4	
		5	4		4			2	1	3
2	4			1	2		4	3	5	
	2	4	5			4		2		
		4	4	3	4			2		5
1		3		3		2		4		

factors

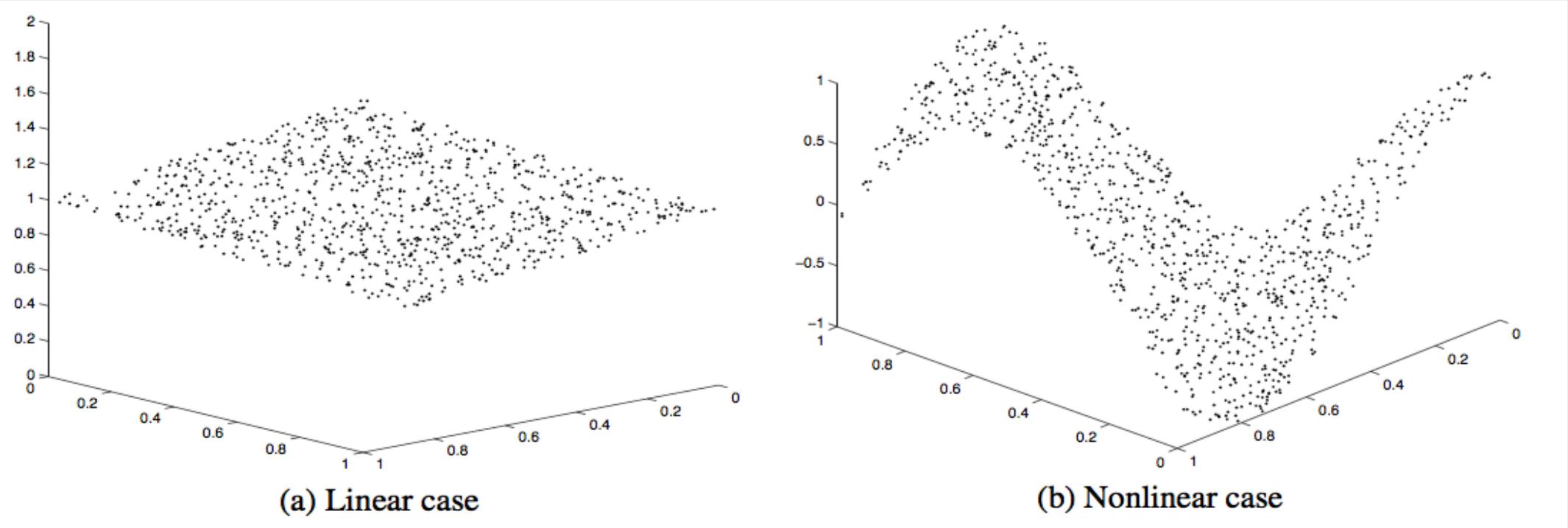
.1	.4	.2
-.5	.6	.5
-.2	.3	.5
1.1	2.1	.3
-.7	2.1	-2
.1	.7	.3

≈ users

items

1.1	-.2	.3	.5	-2	-.5	.8	-.4	.3	1.4	2.4	-9
-.8	.7	.5	1.4	.3	-1	1.4	2.9	-.7	1.2	-.1	1.3
2.1	-.4	.6	1.7	2.4	.9	-.3	.4	.8	.7	-.6	.1

# Limitation of Matrix Factorization



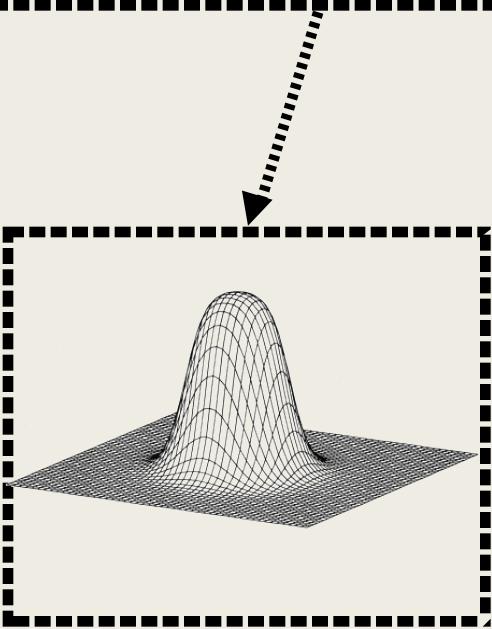
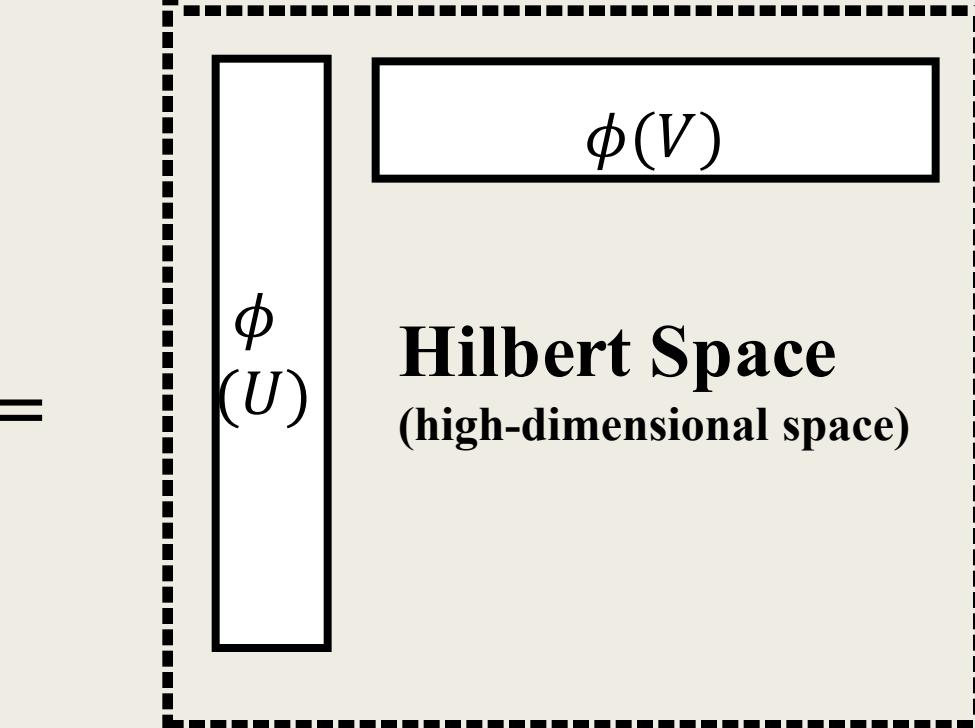
$$\min_{U,V} \sum_{(i,j) \in R} (r_{ij} - \boxed{u_i v_j^T}_{\text{linear}})^2$$

Data can distributed on a non-linear hyper-plane.  
MF is unable to capture the nonlinearity of the data!

	5	3	?	?
	?	4	2	4
	3	5	?	2
	?	?	?	5
	?	1	?	?

Original Space

Kernelized MF



Implicitly define  
→  $\phi(\cdot)$

Mapping function

Kernel

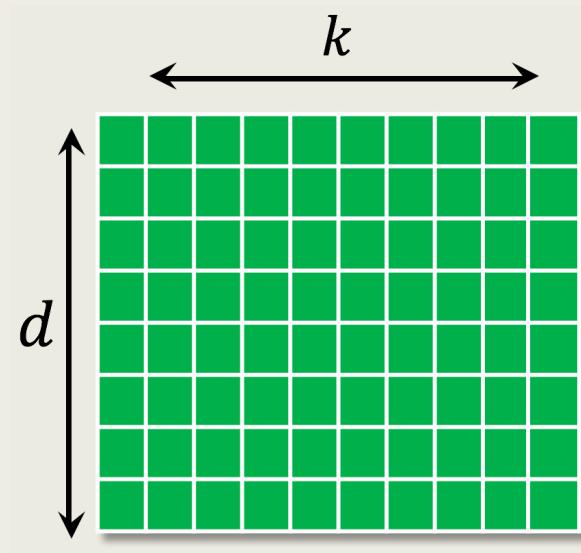
# Challenges

$$\min_{\phi(U), \phi(V)} \sum_{(i,j) \in R} (r_{ij} - \phi(u_i)^T \phi(v_j))^2$$

1. *It is difficult to minimize the objective function over implicitly defined  $\phi(U)$  and  $\phi(V)$*
2. *It is difficult to choose a good kernel function (and kernel parameters) for the data*
3. *It is difficult to use only a single kernel to capture the complex non-linearity embedded in the data*

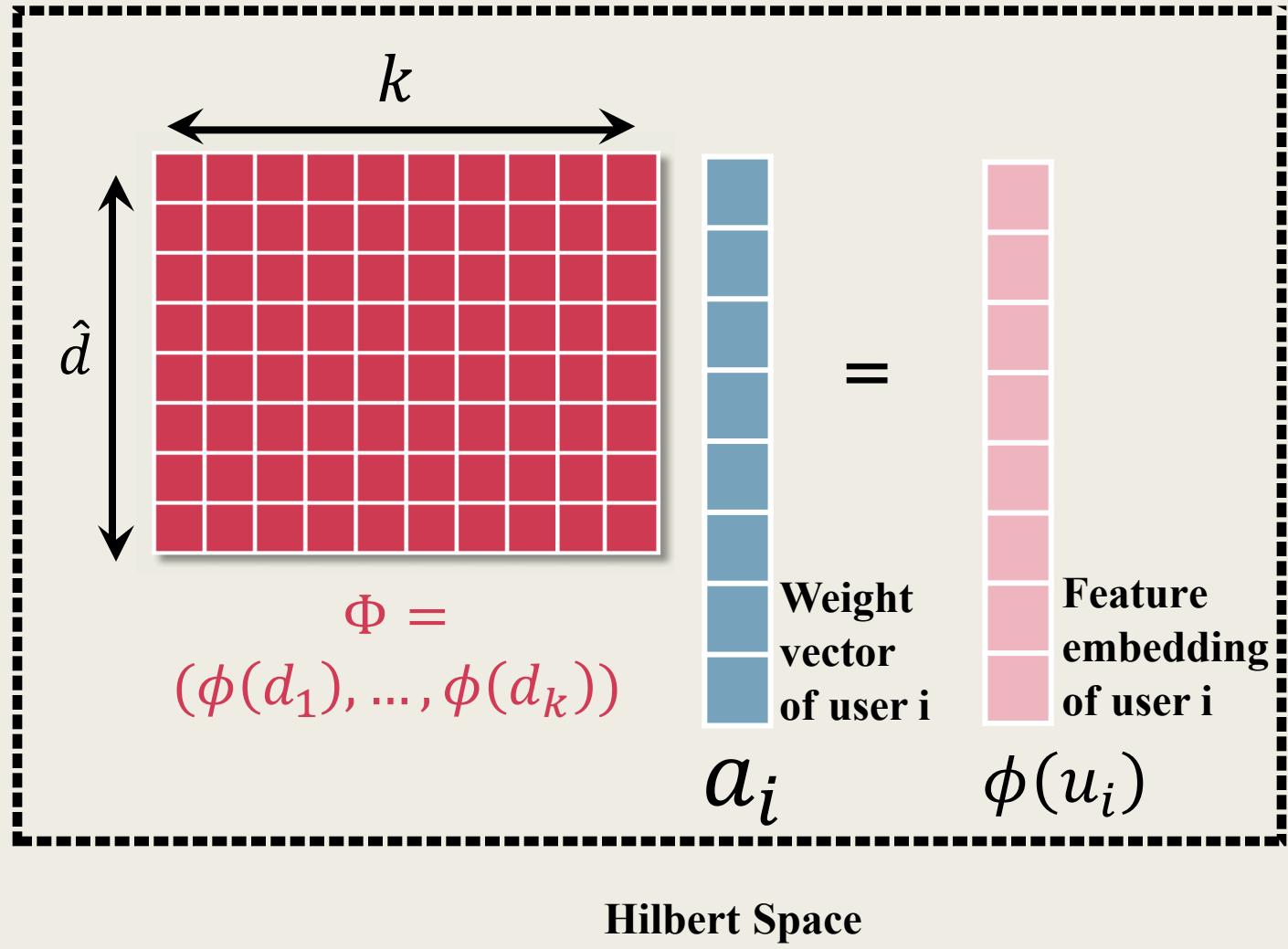
# Dictionary Vectors

$$(d_1, \dots, d_k)$$



$f(\cdot)$

$$d_i \xrightarrow{f(\cdot)} \phi(d_i)$$



# Dictionary Vectors Cont.

$$\min_{\phi(U), \phi(V)} \sum_{(i,j) \in R} (r_{ij} - \phi(u_i)^T \phi(v_j))^2$$

$\downarrow$

$$\begin{aligned}\phi(u_i) &= \Phi \mathbf{a}_i \\ \phi(v_j) &= \Phi \mathbf{b}_j\end{aligned}$$

$$\min_{\mathbf{A}, \mathbf{B}} \sum_{(i,j) \in R} (r_{ij} - \mathbf{a}_i \Phi^T \Phi \mathbf{b}_j)^2$$

$\downarrow$

$$\mathbf{K} = \Phi^T \Phi$$

$$\min_{\mathbf{A}, \mathbf{B}} \sum_{(i,j) \in R} (r_{ij} - \mathbf{a}_i \mathbf{K} \mathbf{b}_j)^2$$

$$\mathbf{K} = \begin{pmatrix} \phi(d_1)^T \phi(d_1) & \phi(d_1)^T \phi(d_k) \\ \vdots & \vdots \\ \phi(d_k)^T \phi(d_1) & \phi(d_k)^T \phi(d_k) \end{pmatrix}$$

Gram Matrix

$$K_{ij} = \phi(d_i)^T \phi(d_j) = f(d_i, d_j)$$

Kernel function

$${}^* \mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_m), \mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_n)$$

$$f(x_1, x_2) = \phi(x_1)^T \phi(x_2) = (x_1^T x_2 + c)^d$$

Polynomial Kernel

# Optimization

L2-norm regularizations

$$\min_{\mathbf{A}, \mathbf{B}} \sum_{(i,j) \in R} (r_{ij} - \mathbf{a}_i^T \mathbf{K} \mathbf{b}_j)^2 + \lambda (\mathbf{a}_i^T \mathbf{K} \mathbf{a}_i + \mathbf{b}_j^T \mathbf{K} \mathbf{b}_j)$$

$$\boxed{\min_B \sum_{(i,j) \in R} (r_{ij} - \mathbf{a}_i^T \mathbf{K} \mathbf{b}_j)^2 + \lambda \mathbf{b}_j^T \mathbf{K} \mathbf{b}_j}$$

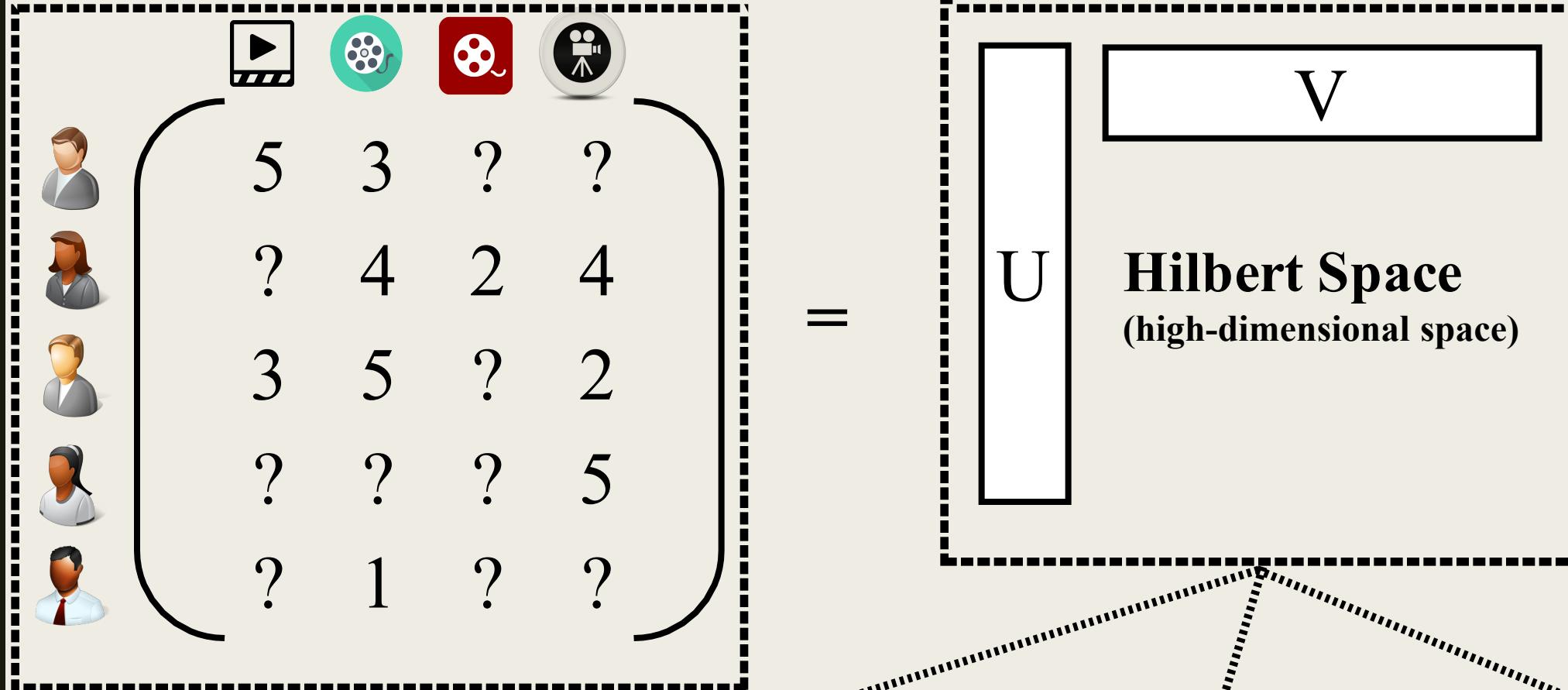


Close Form Solution:  $\hat{\mathbf{b}}_j = (\mathbf{K}^T \mathbf{A} \mathbf{A}^T \mathbf{K} + \lambda \mathbf{K})^{-1} \mathbf{K}^T \mathbf{A} r_j$

# Challenges

$$\min_{\phi(U), \phi(V)} \sum_{(i,j) \in R} (r_{ij} - \phi(u_i)^T \phi(v_j))^2$$

1. *It is difficult to minimize the objective function over implicitly defined  $\phi(U)$  and  $\phi(V)$*
2. *It is difficult to choose a good kernel function (and kernel parameters) for the data*
3. *It is difficult to use only a single kernel to capture all the non-linearity embedded in the data*



Original Space

Multi-Kernel MF

Kernel 1

Kernel 2

Kernel 3

V

U

Hilbert Space  
(high-dimensional space)

# Combine Multiple Kernels

- Given  $p$  different kernels, we can compute  $p$  Gram matrices:

$$\{\mathbf{K}_1, \dots, \mathbf{K}_p\}$$

- We want to learn the kernel weight vector for  $p$  kernels:

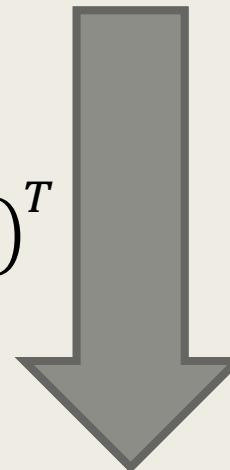
$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^T$$

$$\min_{\boldsymbol{\mu}} \sum_{(i,j) \in R} (r_{ij} - \sum \textcolor{red}{\boldsymbol{\mu}_u} \mathbf{a}_i^T \mathbf{K}_u \mathbf{b}_j)^2 + \lambda (\mathbf{b}_j^T \sum \textcolor{red}{\boldsymbol{\mu}_u} \mathbf{K}_u \mathbf{b}_j + \mathbf{a}_i^T \sum \textcolor{red}{\boldsymbol{\mu}_u} \mathbf{K}_u \mathbf{a}_i)$$

# Combine Multiple Kernels Cont.

$$\min_{\mu} \sum_{(i,j) \in R} (r_{ij} - \sum \mu_u a_i^T K_u b_j)^2 + \lambda (b_j^T \sum \mu_u K_u b_j + a_i^T \sum \mu_u K_u a_i)$$

$$v_{ij} = (a_i^T K_1 b_j, \dots, a_i^T K_p b_j)^T$$
$$y_{ij} = (a_i^T K_1 a_i + b_j^T K_1 b_j, \dots, a_i^T K_1 a_i + b_j^T K_1 b_j)^T$$



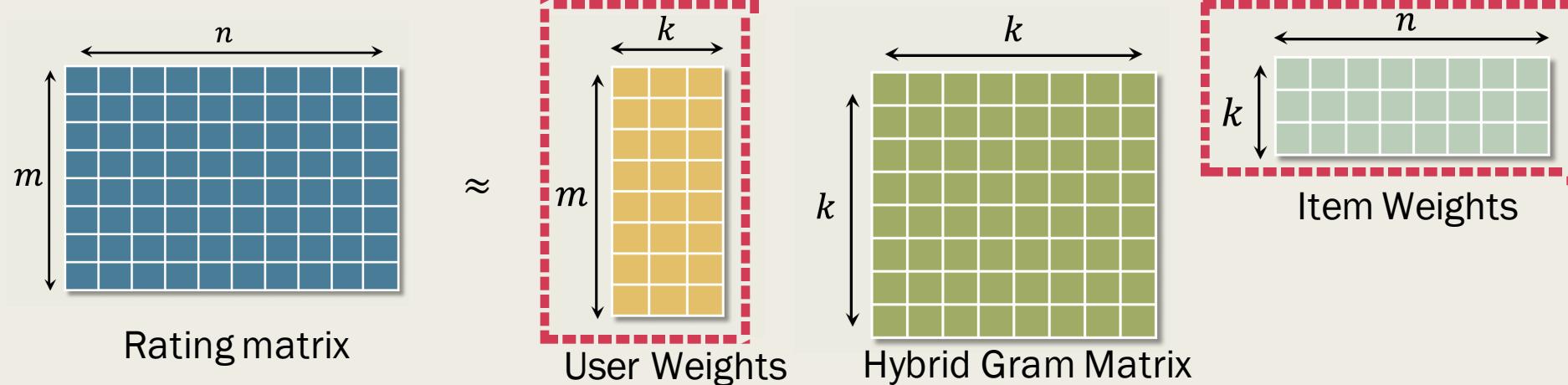
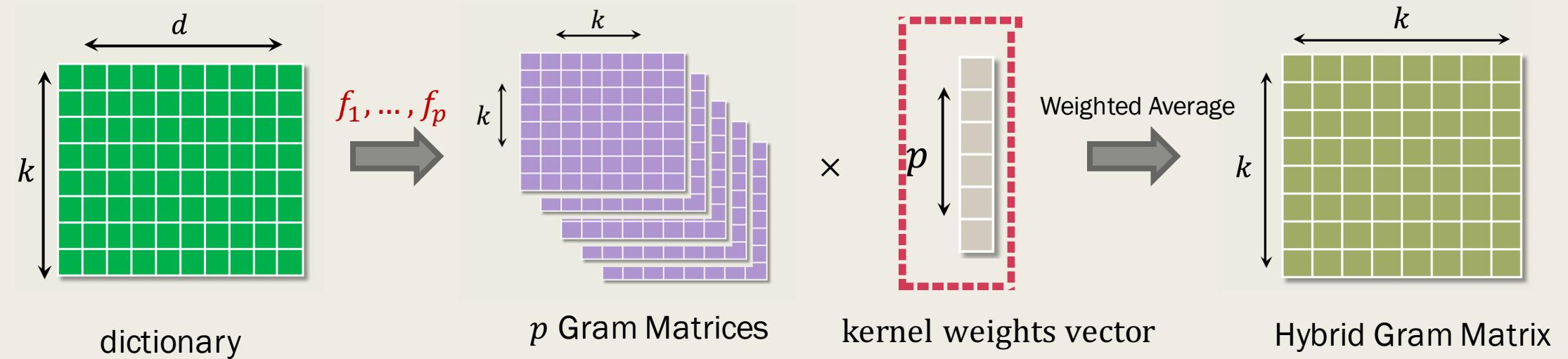
$$Y = \sum_{u,i} v_{ij} v_{ij}^T$$

$$Z = \sum_{u,i} (\lambda - 2r_{ij}) y_{ij}$$

Quadratic Programming

$$\begin{aligned} & \text{minimize } \mu^T Y \mu + Z \mu \\ & \text{subject to } \mu \geq 0 \\ & \quad \mathbf{1}_p^T \mu = 1 \end{aligned}$$

# Multiple Kernel Matrix Factorization



# Datasets

Table 2: Summary of Datasets

Dataset	# of users	# of items	Density (%)
MovieLens	6,040	3,900	6.3
Jester	73,421	100	55.8
Flixster	147,612	48,784	0.11
Dating Agency	135,359	168,791	0.76
Yahoo Music	1,948,882	98,211	0.006
ASSISTments	46,627	179,084	0.073

# Experiments

1. Sample a  $1000 \times 1000$  matrix for each dataset
2. Hold 1 rating for each user for **test**, the remaining as **training** data
3. Repeat this partition **5 times** for each dataset
4. Average of **RMSE** on test data are reported

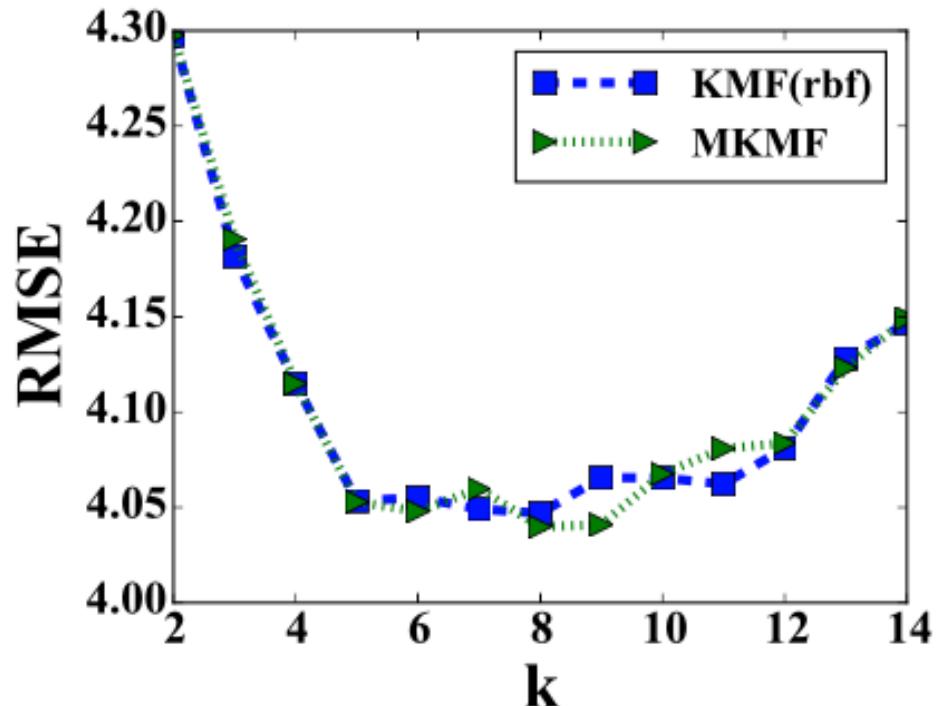
$$\text{RMSE} = \sqrt{\frac{\sum_{(u,i) \in \Omega} (r_{ui} - \hat{r}_{ui})^2}{|\Omega|}}$$

# Results

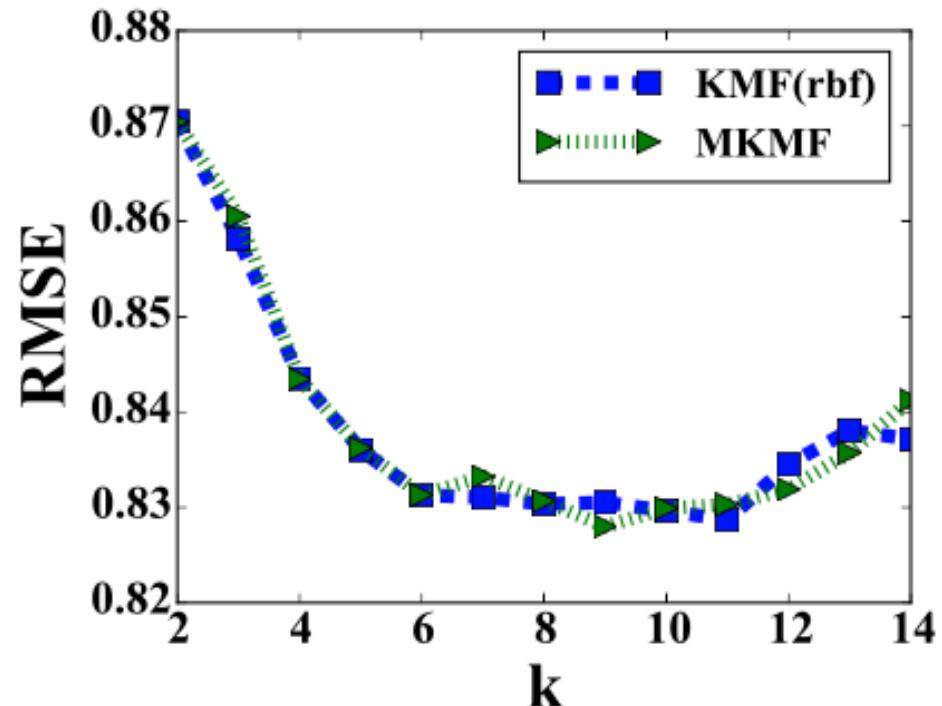
Methods	Dataset										Ave.
	Flixster	Yahoo Music	MovieLens	Jester	Dating	ASSISTments	Rank				
MKMF	0.8270 (1)	18.3036 (1)	0.8223 (1)	4.0398 (1)	1.6697 (1)	0.7920 (2)	1.1667				
KMF (LINEAR)	0.8290 (4)	18.3734 (3)	0.8241 (3.5)	4.0471 (2.5)	1.6804 (4)	0.7933 (4)	3.5				
KMF (POLY)	0.8289 (3)	18.3766 (4)	0.8241 (3.5)	4.0472 (4)	1.6797 (3)	0.7934 (5)	3.75				
KMF (RBF)	0.8292 (5)	18.3883 (5)	0.8246 (5)	4.0471 (2.5)	1.6701 (2)	0.7927 (3)	3.75				
MF	0.8286 (2)	18.3214 (2)	0.8235 (2)	4.0549 (5)	1.6849 (5)	0.8001 (6)	3.6667				
SVD	0.9441 (9)	26.2255 (9)	0.9406 (7)	4.2794 (6)	1.7920 (7)	0.8919 (9)	7.8333				
Ivc-Cos	0.9223 (7)	22.9477 (7)	1.0016 (8)	4.6137 (8)	2.7052 (8)	0.8106 (7)	7.5				
Ivc-PEARSON	0.9226 (8)	22.9486 (8)	1.0020 (9)	4.6142 (9)	2.8209 (9)	0.8446 (8)	8.5				
Avg	0.9006 (6)	22.4159 (6)	0.8887 (6)	4.3867 (7)	1.7349 (6)	0.7658 (1)	5.3333				

- Our methods are in **red box** (KMF & MKMF)
- Kernelized MF is better than other baselines (**green box**)
- Multi-Kernel MF is better than single kernel (**blue box**)

# Parameter Study ( $k$ , # of latent factors)



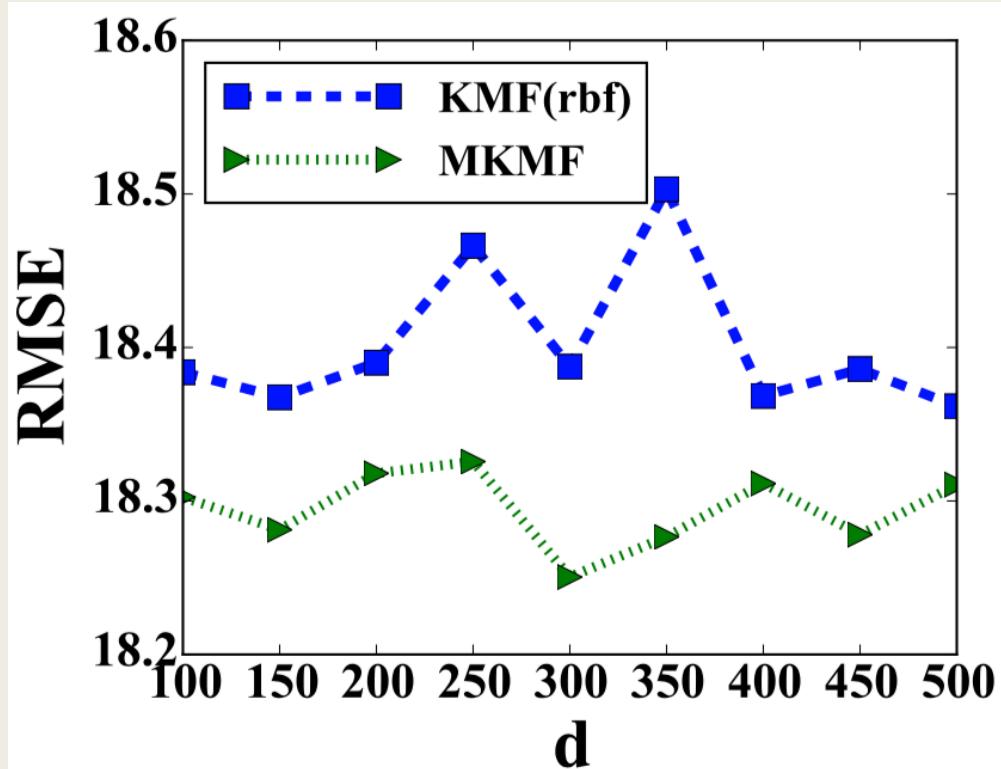
(a) Jester Dataset



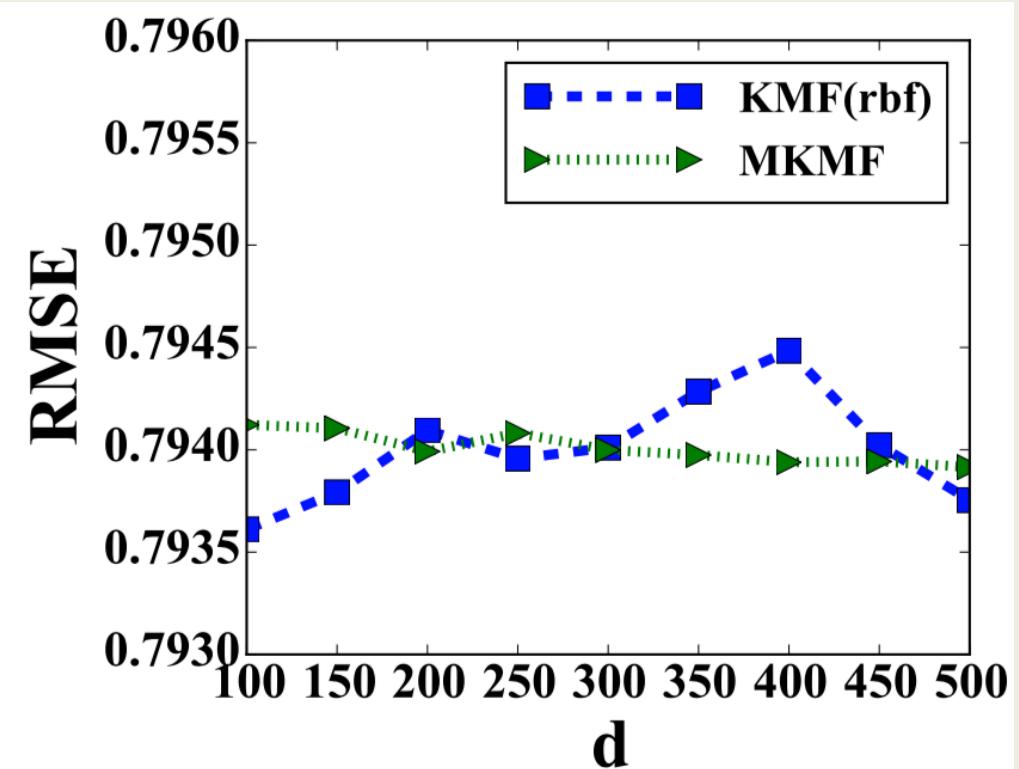
(b) Flixster Dataset

- Small  $k$  ( $\leq 5$ ) could not fit the data very well
- Large  $k$  ( $> 12$ ) could be over-fitting
- Cross-validate on each dataset for the best choice of  $k$

# Parameter Study ( $d$ , the dimension of dictionary vectors)



(a) Yahoo Music Dataset



(b) ASSISTments Dataset

- Not sensitive on  $d$
- Any  $d > 100$  works just fine

# Conclusion & Future Works

- Conclusion:
  - Our proposed KMF and MKMF methods both exploit **the underlying nonlinear correlations** among rows (users) and columns (items) in the rating matrix.
  - KMF/MKMF **improves the overall performance** of predicting the unobserved ratings compared to state-of-art baselines.
- Future Works:
  - Alternative regularizations ( e.g. **L2 norm**) can be applied on the weight learning process.
  - The kernel (Gram) matrix might be **learned directly** from the data without the help of dictionary.
  - **Scale** to big data (e.g. Netflix dataset)