

# Maximum Margin Matrix Factorization

Nati Srebro  
University of Toronto

Jason Rennie  
MIT

Tommi Jaakkola  
MIT

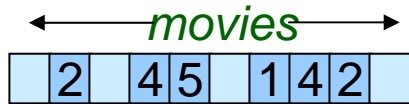
# Collaborative Prediction

Based on partially observed matrix:

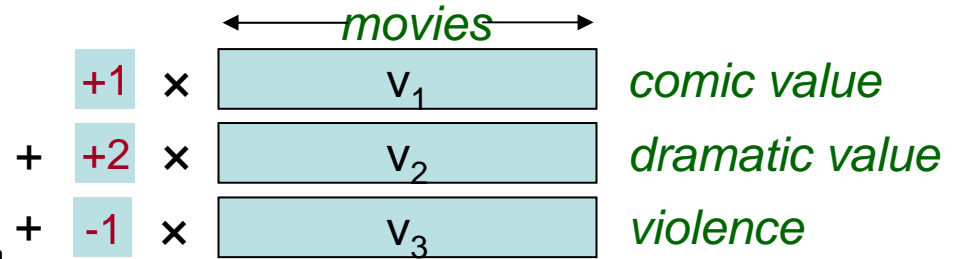
⇒ Predict unobserved entries “Will user  $i$  like movie  $j$ ?”

	movies											
users		2		1			4				5	
		5		4				?		1		3
			3		5			2				
		4		?			5		3		?	
			4		1	3				5		
				2				1	?			4
		1					5		5		4	
			2		?	5		?		4		
		3		3		1		5		2		1
		3				1				2		3
		4			5	1				3		
			3				3	?				5
		2	?		1		1					
			5			2	?		4		4	
		1		3		1	5		4		5	
	1		2			4			5	?		

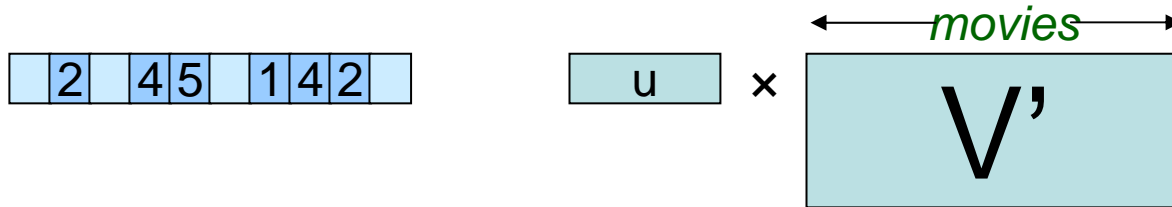
# Linear Factor Model



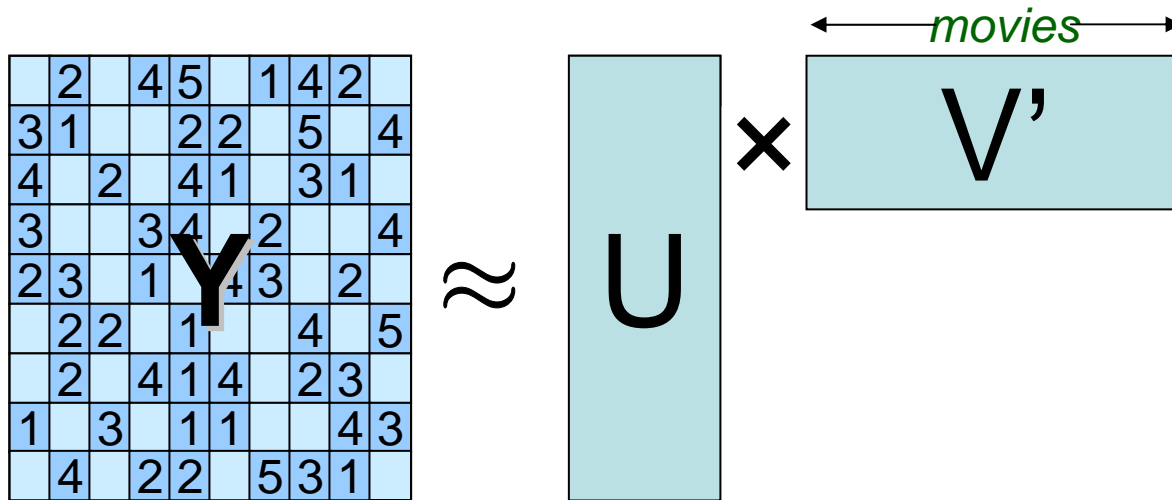
*preferences of a specific user*



# Linear Factor Model

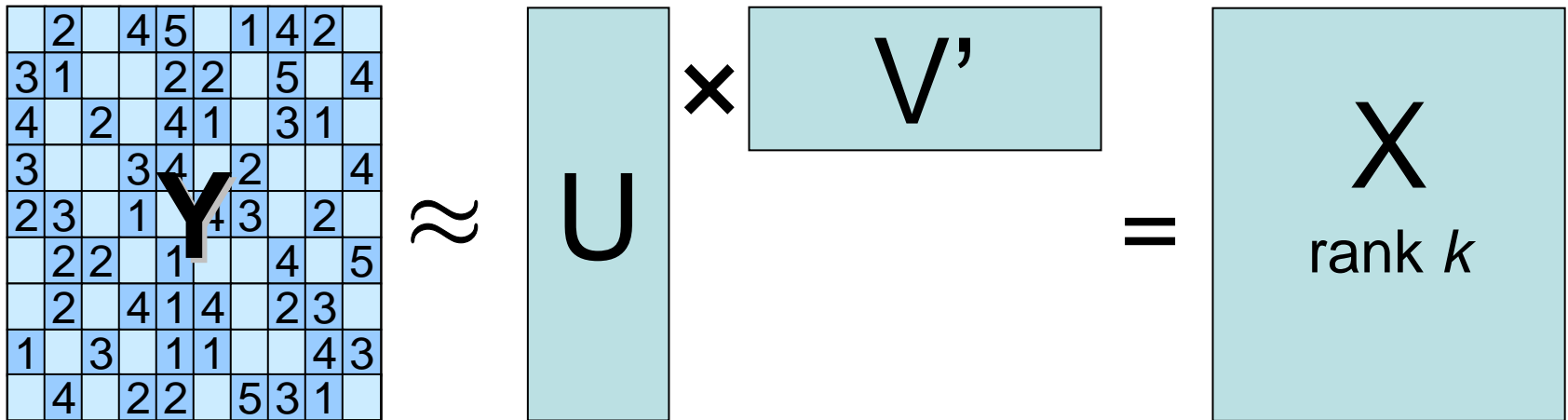


# Linear Factor Model



# Matrix Factorization

## Unconstrained: Low Rank Approximation

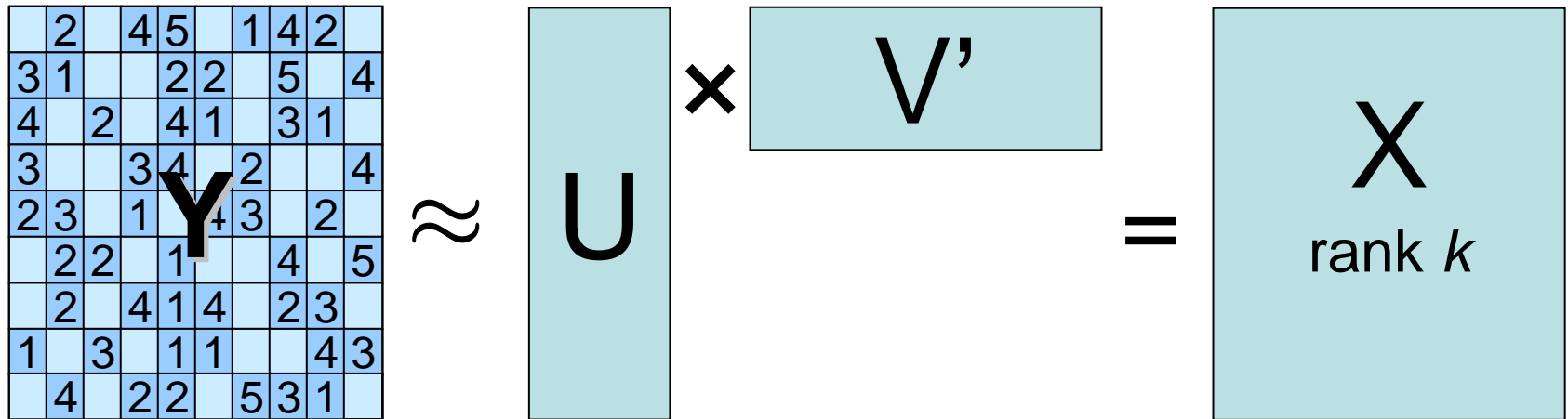


- Additive Gaussian noise: minimize  $\|Y - UV'\|_{\text{Fro}}$
- General additive noise
- General conditional models
  - Multiplicative noise, Exponential-PCA [Collins+01], Multinomial (pLSA [Hofmann01]), etc
- General loss functions
  - Hinge loss, loss functions appropriate for ratings, etc [Gordon03]

Unconstrained  $U, V$ ,  
fully observed  $Y$   
→ use SVD

**non-convex,  
no explicit solution**

# Matrix Factorization



- Non-Negativity [LeeSeung99]
- Stochasticity (convexity) [LeeSeung97] [Hofmann01]
- Sparsity
  - Clustering as an extreme (when rows of  $U$  sparse)

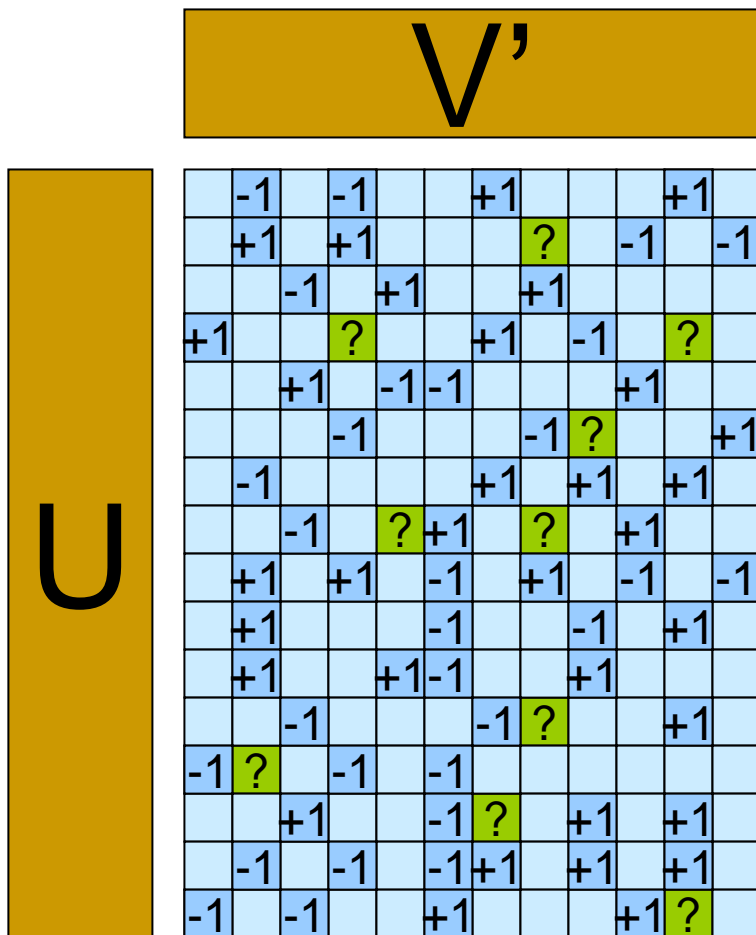
**Overall number of factors still constrained**  
**Non-convex optimization problems**

# Outline

- **Maximum Margin Matrix Factorization**
  - Unbounded number of factors
  - Convex!
- Learning MMMF: Semidefinite Programming
- Generalization Error Bounds

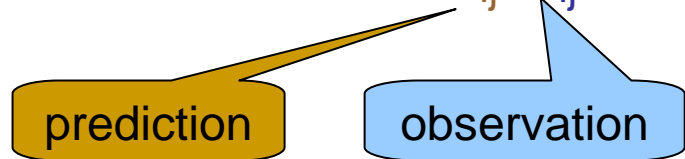


# Collaborative Prediction with Matrix Factorization



Fit factorizable (low-rank) matrix  $X=UV'$  to observed entries.

minimize  $\sum \text{loss}(X_{ij}; Y_{ij})$



Use matrix  $X$  to predict unobserved entries.

# Collaborative Prediction with Matrix Factorization

feature vectors

1.3	0.4	-1.5
8.3	2.5	-4.8
0.7	-0.2	3.4
1.7	-5.2	1.6
-3.7	2.1	0.9
4.3	-0.5	2.7
4.7	0.2	6.4
6.0	0.3	-5.8
-1.5	-3.7	0.4
-4.8	4.3	2.5
3.4	4.7	-0.2
1.6	6.0	-5.2
0.9	1.3	2.1
2.7	8.3	-0.5
6.4	0.7	0.2
-5.8	1.7	0.3

linear classifiers

	v1	v2	v3	v4	v5	v6	v7	v8	v9	v10	v11	v12
		-1		-1			+1				+1	
		+1		+1						-1		-1
			-1		+1			+1				
	+1					+1			-1			
			+1		-1	-1				+1		
				-1				-1				+1
		-1				+1		+1	+1		+1	
		+1		+1		-1	+1		-1	-1		-1
		+1				-1			-1		+1	
		+1			+1	-1			+1			
			-1				-1				+1	
	-1			-1		-1						
			+1			-1			+1		+1	
		-1		-1		-1	+1		+1		+1	
	-1		-1			+1				+1		

When  $U$  is fixed, each row is a linear classification problem:

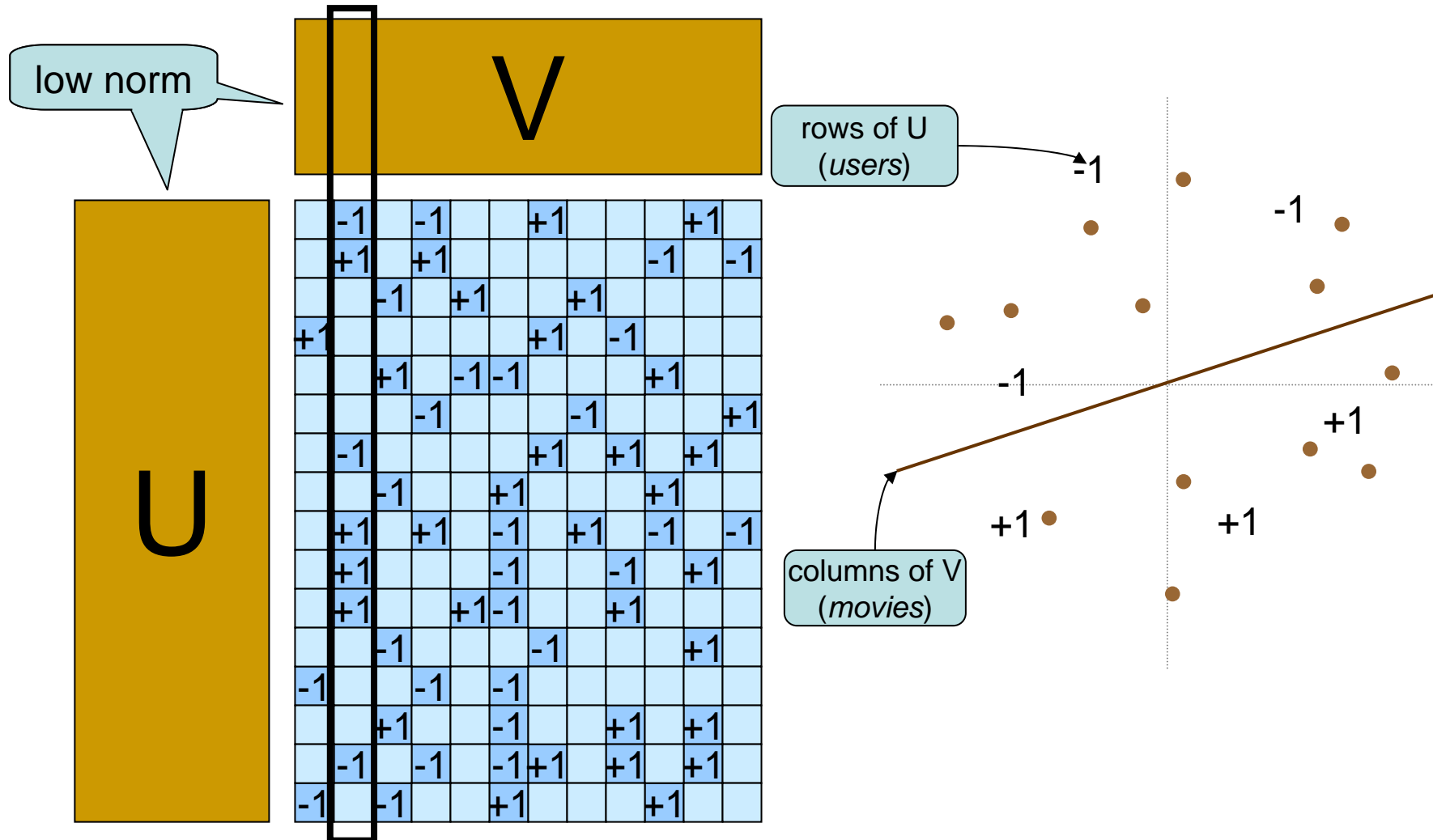
- rows of  $U$  are feature vectors
- columns of  $V$  are linear classifiers

Fitting  $U$  and  $V$ :

Learning features that work well across all classification problems.

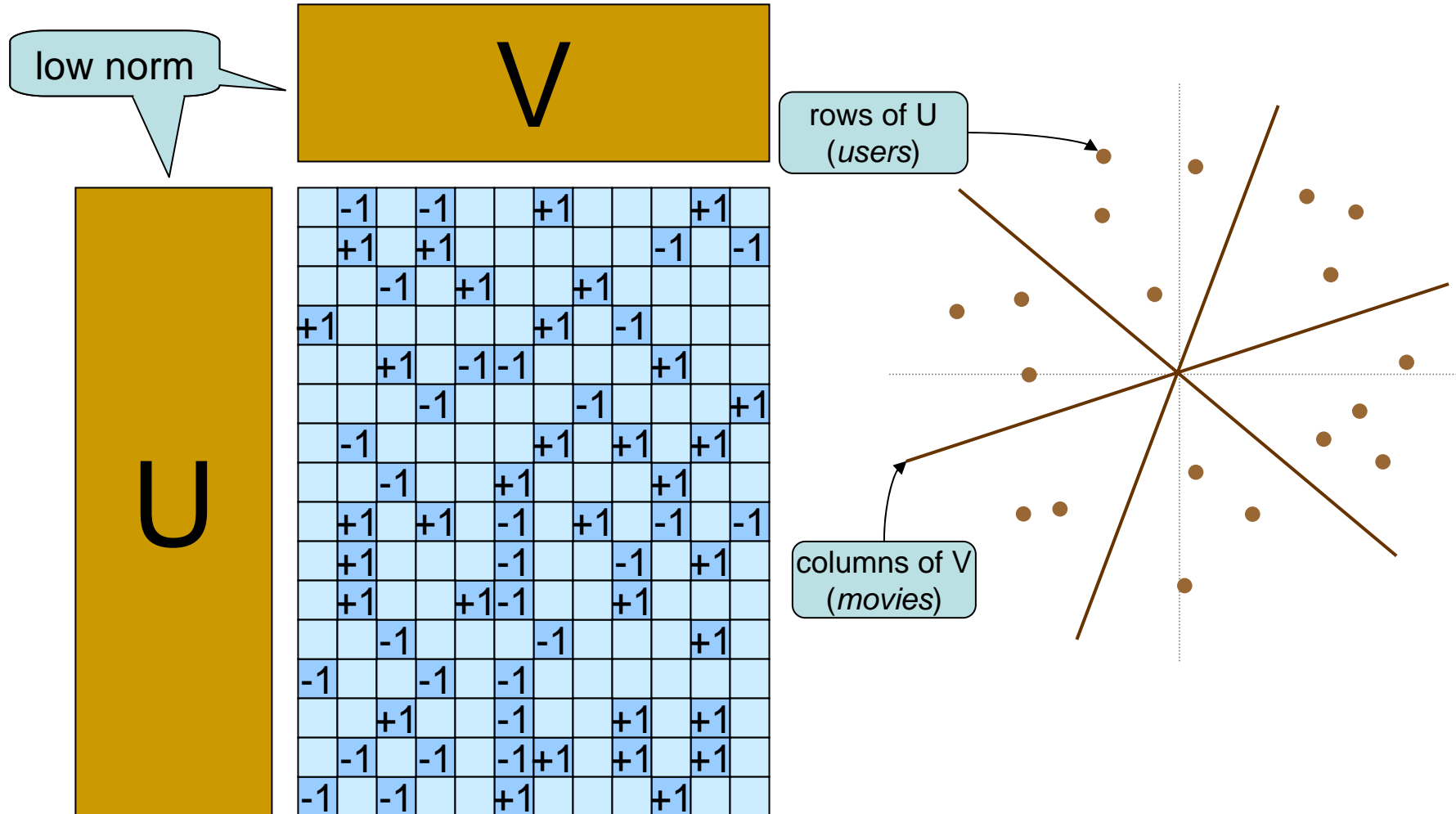
# Geometric Interpretation:

## Co-embedding Points and Separating Hyperplanes



# Geometric Interpretation:

## Co-embedding Points and Separating Hyperplanes



# Max-Margin Matrix Factorization:

Bound norms of  $U, V$  instead of their dimensionality

low norm

$V$

bound norms uniformly:

$$(\max_i |U_i|^2) (\max_j |V_j|^2) \leq 1$$

rows of  $U$   
(users)

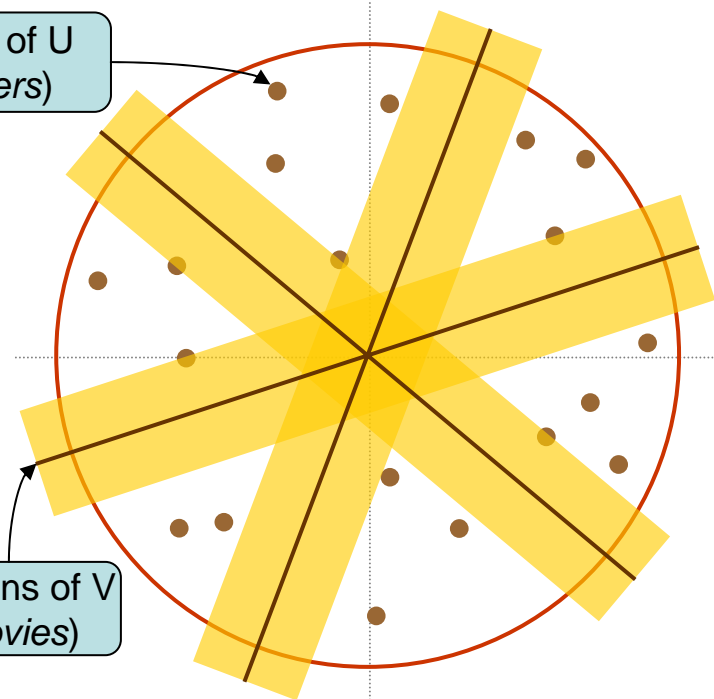
	-1	-1		+1		+1	
	+1	+1				-1	-1
		-1	+1		+1		
+1				+1	-1		
		+1	-1	-1		+1	
		-1			-1		+1
-1				+1	+1	+1	
	-1		+1		+1		
+1	+1		-1	+1	-1	-1	
+1			-1		-1	+1	
+1			+1	-1	+1		
	-1			-1			+1
-1		-1	-1				
		+1		-1	+1	+1	
-1		-1	-1	+1	+1	+1	
-1	-1		+1		+1		

columns of  $V$   
(movies)

For observed  $Y_{ij} \in \pm 1$ :

$$Y_{ij} X_{ij} \geq \text{Margin}$$

$\langle U_i, V_j \rangle$



# Max-Margin Matrix Factorization:

Bound norms of  $U, V$  instead of their dimensionality

low norm

$V$

$U$

	-1	-1		+1		+1	
	+1	+1				-1	-1
		-1	+1		+1		
+1				+1	-1		
		+1	-1	-1		+1	
		-1			-1		+1
-1				+1	+1	+1	
	-1		+1		+1		
	+1	+1	-1	+1	-1	-1	
	+1		-1		-1	+1	
	+1		+1	-1		+1	
		-1		-1			+1
-1		-1	-1				
		+1	-1		+1	+1	
	-1	-1	-1	+1	+1	+1	
-1	-1		+1			+1	

bound norms uniformly:

$$(\max_i |U_i|^2) (\max_j |V_j|^2) \leq 1$$

bound norms on average:

$$(\sum_i |U_i|^2) (\sum_j |V_j|^2) \leq 1$$

$U$  is fixed:

each column of  $V$  is SVM

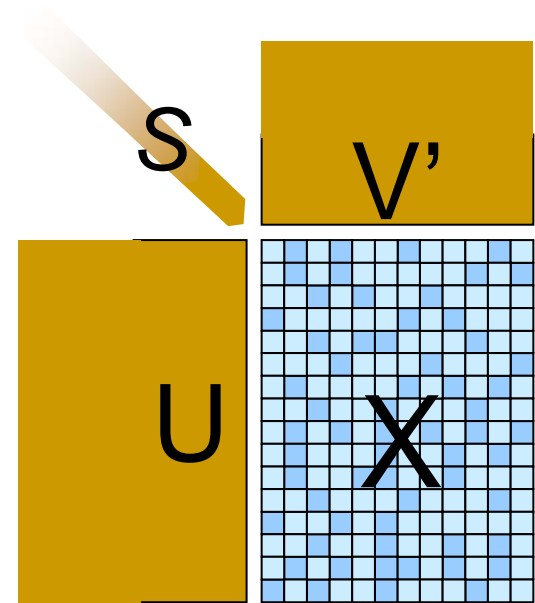
For observed  $Y_{ij} \in \pm 1$ :

$$Y_{ij} X_{ij} \geq \text{Margin}$$

$\langle U_i, V_j \rangle$

# Finding Max-Margin Matrix Factorizations

$$\begin{aligned} &\text{maximize } M \\ &Y_{ij} X_{ij} \geq M \\ &X = UV \\ &\underbrace{(\sum_i |U_i|^2) (\sum_j |V_j|^2)} \leq 1 \\ &|X|_{\text{tr}} = \sum (\text{singular values of } X) \end{aligned}$$



Unlike  $\text{rank}(X) \leq k$ , this a convex constraint!

# Finding Max-Margin Matrix Factorizations

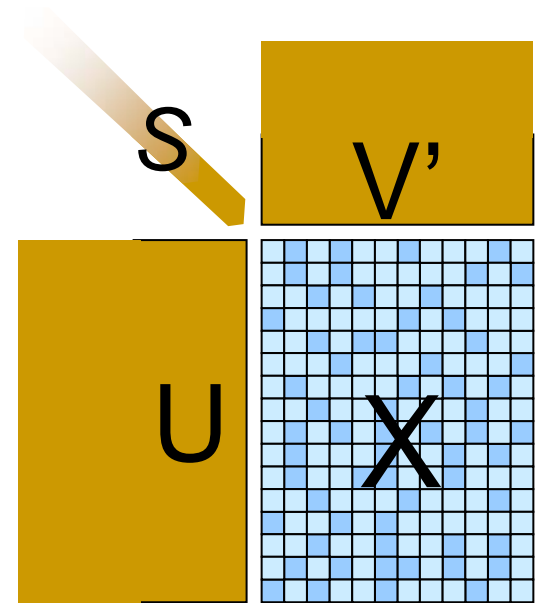
maximize  $M$

$$Y_{ij} X_{ij} \geq M$$

$$X = UV$$

$$\underbrace{(\sum_i |U_i|^2) (\sum_j |V_j|^2)} \leq 1$$

$$\|X\|_{\text{tr}} = \sum (\text{singular values of } X)$$



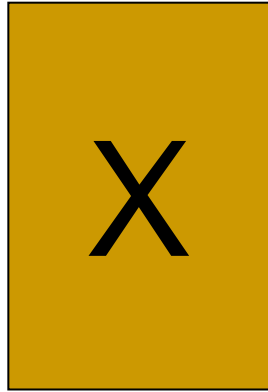
minimize  $\text{tr}(A) + \text{tr}(B)$

$$Y_{ij} X_{ij} \geq 1$$

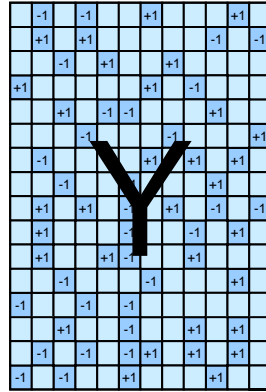
$$\begin{pmatrix} A & X \\ X' & B \end{pmatrix} \text{ p.s.d.}$$



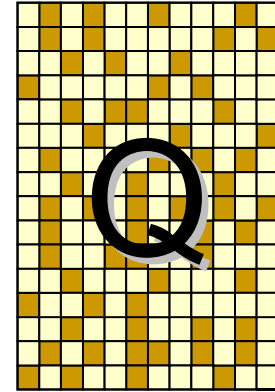
# Finding Max-Margin Matrix Factorizations



dense primal



sparse observations  
(constraints)



sparse dual

minimize  $\text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B})$

$$Y_{ij} X_{ij} \geq 1$$

$$\begin{pmatrix} \mathbf{A} & \mathbf{X} \\ \mathbf{X}' & \mathbf{B} \end{pmatrix} \text{ p.s.d.}$$

Dual variable  $Q_{ij}$  for each observed  $(i,j)$

maximize  $\sum Q_{ij}$

$$0 \leq Q_{ij}$$

$$\|\mathbf{Q} \otimes \mathbf{Y}\|_2 \leq 1$$

sparse elementwise product  
(zero for unobserved entries)

# Experimental Results on MovieLens Subset

	<b>MMMF</b>	K-median	Low Rank	per-user median rating
holdout set	1 <b>43.8</b>	41.5	42.5	41.3
	2 <b>44.7</b>	43.3	46.8	41.6
	3 <b>47.3</b>	43.9	45.7	43.2
	4 <b>45.0</b>	42.8	44.7	41.4
overall	<b>45.2</b>	42.9	44.2	41.9

% rating agreement (ratings are 1,2,3,4,5)

100 users × 100 movies, 7030 ratings

# Generalization Error Bounds

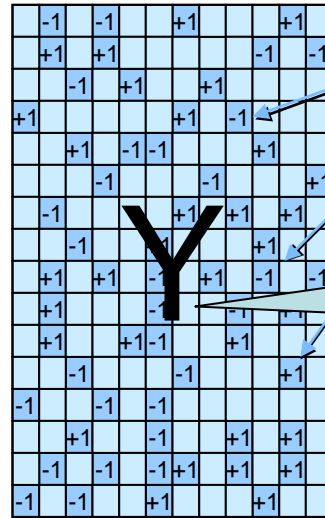
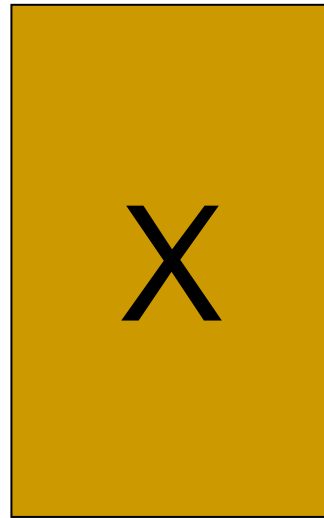
$$D(\mathbf{X}; \mathbf{Y}) = \frac{\#\{ij(\mathbf{X}_{ij} \cdot \mathbf{Y}_{ij} < 0)\}}{nm}$$

generalization error

$$D_S(\mathbf{X}; \mathbf{Y}) = \frac{\#\{ij \in S(\mathbf{X}_{ij} \cdot \mathbf{Y}_{ij} < 1)\}}{|S|}$$

empirical error

$$\forall \mathbf{Y} \Pr_S \left( \forall_{\text{rank-}k \mathbf{X}} D(\mathbf{X}; \mathbf{Y}) < D_S(\mathbf{X}; \mathbf{Y}) + \varepsilon \right) > 1 - \delta$$



random

unknown, assumption-free

$$(\sum |U_i|^2/n)(\sum |V_i|^2/m) \leq R^2:$$

$$\text{rank}(X) \leq k:$$

$$\varepsilon = K^4 \sqrt{\ln m} \sqrt{\frac{R^2(n+m)\log n + \log 1/\delta}{|S|}}$$

$$\varepsilon = \sqrt{\frac{k(n+m)\log \frac{8em}{k} + \log 1/\delta}{2|S|}}$$

# Maximum Margin Matrix Factorization as a Convex Combination of Classifiers

$$\begin{aligned} & \{ UV \mid (\sum |U_i|^2)(\sum |V_j|^2) \leq 1 \} \\ & = \text{convex-hull}( \underbrace{\{ uv' \mid u \in \mathbb{R}^n, v \in \mathbb{R}^m \mid |u|=|v|=1 \}}_{\text{rank-one, unit-norm, matrices}} ) \end{aligned}$$

$$\begin{aligned} & \text{conv}( \{ uv' \mid u \in \{\pm 1\}^n, v \in \{\pm 1\}^m \} ) \\ & \subset \{ UV \mid (\max |U_i|^2)(\max |V_j|^2) \leq 1 \} \\ & \subset 2 \text{conv}( \underbrace{\{ uv' \mid u \in \{\pm 1\}^n, v \in \{\pm 1\}^m \}}_{\text{rank-one sign matrices}} ) \end{aligned}$$

# Maximum Margin Matrix Factorization

Unbounded number of factors  
Learning is a convex problem!

- Correspondence with large margin linear classification
- Generalization error bounds
- Learning: Sparse SDP

poster {  
Making predictions using the dual solution  
Other data types and loss functions  
SDP for max-norm formulation  
Slack

- Applicable in other domains where low-rank approximations are currently used
- Direct optimization of dual would enable large-scale applications