

Matrix Factorizations: A Tale of Two Norms

Nati Srebro

Toyota Technological Institute—Chicago

Maximum Margin Matrix Factorization

S, Jason Rennie, Tommi Jaakkola (MIT), NIPS 2004

Rank, Trace-Norm and Max-Norm

S, Adi Shraibman (Weizmann), COLT 2005

Collaborative Filtering in a Non-Uniform World: Learning with the Weighted Trace Norm

Ruslan Salakhutdinov (MIT), S, NIPS 2010

Practical Large-Scale Optimization for Max-Norm Regularization

Jason Lee (Stanford), Ben Recht (Wisconsin), Ruslan Salakhutdinov (MIT), S,
Joel Tropp (CalTech), NIPS 2010

Concentration-Based Guarantees for Low-Rank Matrix Reconstruction

Rina Foygel (UChicago), S, 2011

Matrix Reconstruction (Collaborative Prediction)

- Observe Y_S , estimate other entries of Y

movies

	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	?	

users

Matrix Reconstruction (Collaborative Prediction)

- Observe Y_S , estimate other entries of Y

$$\hat{X}(S) = \arg \min_{\text{rank}(X) \leq k} |X_S - Y_S|$$

$|X_S - Y_S|_1$ (abs err)
or
 $|X_S - Y_S|_2$ (sq err)

- Convex alternative:

$$\hat{X}(S) = \arg \min_{\|X\|_{\text{tr}} \leq k} |X_S - Y_S|$$

$$\hat{X}(S) = \arg \min_{\|X\|_{\text{max}} \leq k} |X_S - Y_S|$$

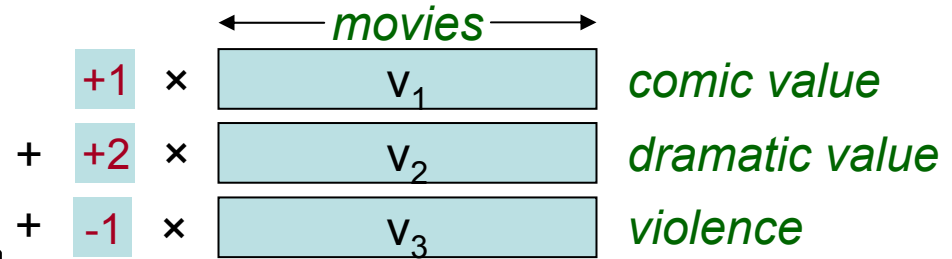
Outline

- Regularized Factorizations: $\|X\|_{\text{tr}}$ and $\|X\|_{\text{max}}$
- Optimizing with $\|X\|_{\text{tr}}$ and $\|X\|_{\text{max}}$
- Low-Rank Reconstruction Guarantees
- *Non-Uniform Sampling*
- Empirical Results

Linear Factor Model

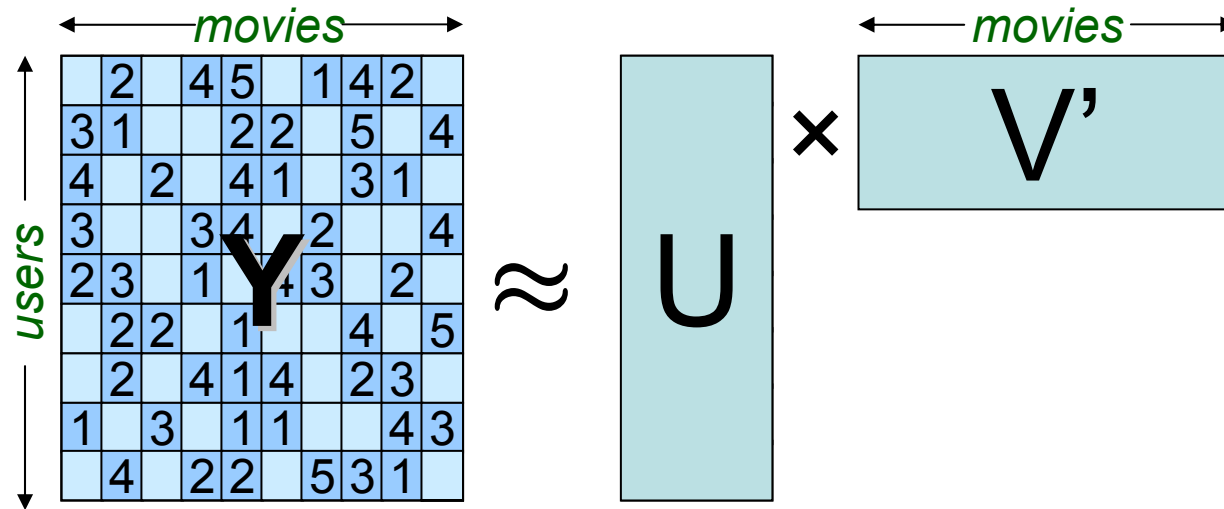


preferences of a specific user

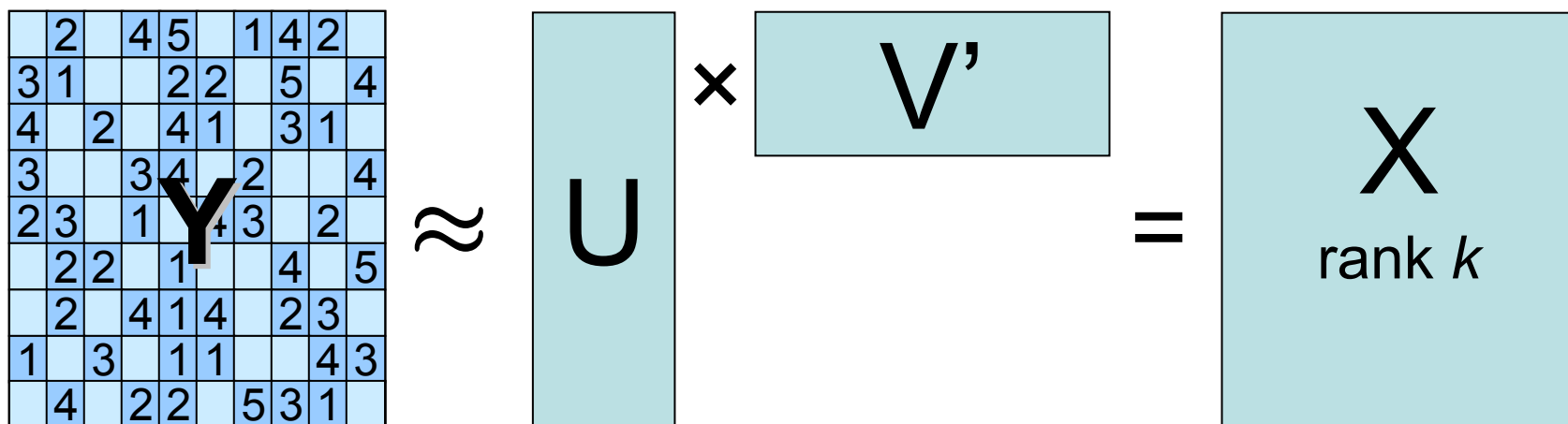


characteristics of the user

Linear Factor Model



Unconstrained Matrix Factorization: Low Rank Approximation



- **Fitting low-rank model is hard**
 - rank is non-convex
 - NP-hard to find low-rank matrix completion, or minimize error s.t. rank constraint
 - Only if Y is **fully observed**, and we want to find low-rank X minimizing **squared-error**, can use SVD
- **Richer models allowing more factors?**

Matrix Factorization Models

Bound dimensionality of factorization:

$$\begin{aligned} \text{rank}(X) &= \min_{X=UV'} \dim(U, V) \\ &= |\text{sing values}|_0 \end{aligned}$$

V														
U	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	

Matrix Factorization Models

low norm



	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	

U

Bound dimensionality of factorization:
 $\text{rank}(X) = \min_{X=UV'} \dim(U, V)$
 $= |\text{sing values}|_0$

Bound avg norm of factorization:

$$\|U\|_F^2 = \sum_i |U_i|^2$$

$$\|X\|_{\text{tr}} = \min_{X=UV'} \|U\|_F \cdot \|V\|_F$$

$$= |\text{sing values}|_1$$

Bound norm of fact. uniformly:

$$\|U\|_{2,\infty} = \max_i |U_i|$$

$$\|X\|_{\text{max}} = \min_{X=UV'} \|U\|_{2,\infty} \cdot \|V\|_{2,\infty}$$

aka $\gamma_2: \ell_1 \rightarrow \ell_\infty$ norm

Matrix Factorization Models: Convexity

Bound dimensionality of factorization:

$$\begin{aligned} \text{rank}(X) &= \min_{X=UV'} \dim(U, V) \\ &= |\text{sing values}|_0 \end{aligned}$$

Bound avg norm of factorization:

$$\|U\|_F^2 = \sum_i |U_i|^2$$

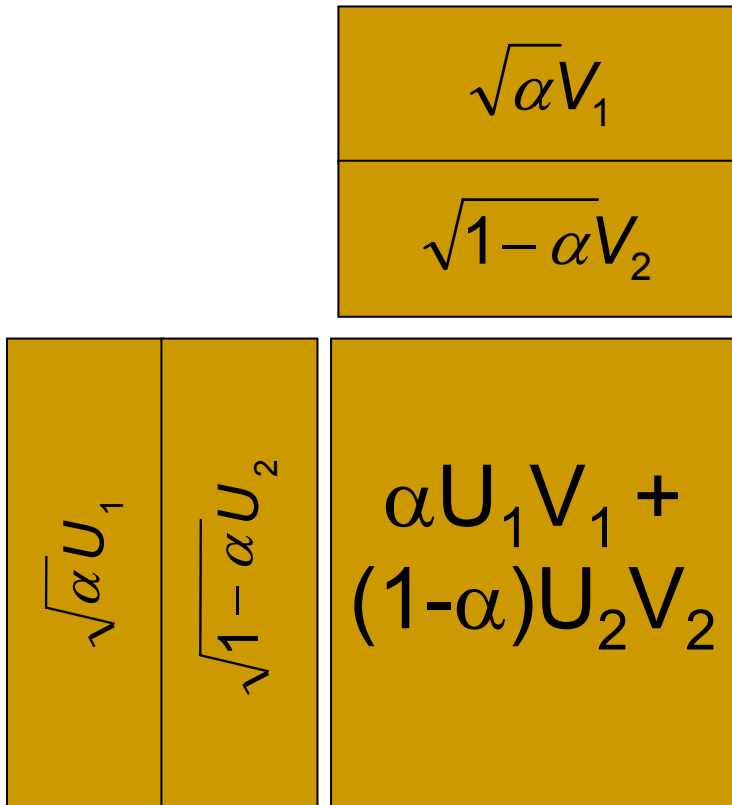
$$\begin{aligned} \|X\|_{\text{tr}} &= \min_{X=UV'} \|U\|_F \cdot \|V\|_F \\ &= |\text{sing values}|_1 \end{aligned}$$

Bound norm of fact. uniformly:

$$\|U\|_{2,\infty} = \max_i |U_i|$$

$$\|X\|_{\max} = \min_{X=UV'} \|U\|_{2,\infty} \cdot \|V\|_{2,\infty}$$

aka $\gamma_2: \ell_1 \rightarrow \ell_\infty$ norm



Trace-Norm and Max-Norm

Trace-Norm (On-Average Constraint):

Regularized factorization: $\|X\|_{\text{tr}} = \min_{X=U'V} \sqrt{\left(\sum_i |U_i|^2\right) \left(\sum_j |V_j|^2\right)}$

Infinite comb. of factors: $\{X \mid \|X\|_{\text{tr}} \leq 1\} = \text{conv} \left(\{uv' \mid |u| = |v| = 1\} \right)$

Lower bound on rank: $\text{tc}(X) \doteq \frac{\|X\|_{\text{tr}}^2}{nm} \leq \text{rank}(X) \cdot \frac{\sum_{ij} |X_{ij}|^2}{nm}$

$\|X\|_{\text{tr}} \leq \text{rank}^{1/2}(X) \|X\|_{\text{Fro}}$

Max-Norm (Uniform Constraint):

Regularized factorization: $\|X\|_{\text{max}} = \min_{X=U'V} (\max_i |U_i|) (\max_j |V_j|)$

Infinite comb. of factors: $\text{conv} \left(\{uv' \mid u \in \pm 1^n, v \in \pm 1^m\} \right) \subseteq$
 $\{X \mid \|X\|_{\text{max}} \leq 1\} \subseteq 2 \text{conv} \left(\{uv' \mid u \in \pm 1^n, v \in \pm 1^m\} \right)$

Grothendiek's Inequality

Lower bound on rank: $\text{mc}(X) \doteq \|X\|_{\text{max}}^2 \leq \text{rank}(X) \cdot \max_{ij} |X_{ij}|^2$

Outline

- Regularized Factorizations: $\|X\|_{\text{tr}}$ and $\|X\|_{\text{max}}$
- **Optimizing with $\|X\|_{\text{tr}}$ and $\|X\|_{\text{max}}$**
- Low-Rank Reconstruction Guarantees
- *Non-Uniform Sampling*
- Empirical Results

Optimization with the Trace-Norm

Covex

$$\begin{aligned} \text{minimize } \text{loss}(\mathbf{X}) + \lambda \cdot \underbrace{\|\mathbf{X}\|_{\text{tr}}}_{\|\mathbf{X}\|_{\text{tr}}} &= \sum (\text{singular values of } \mathbf{X}) = \min_{\mathbf{X}=\mathbf{U}\mathbf{V}} \frac{1}{2}(\sum_i |\mathbf{U}_i|^2 + \sum_j |\mathbf{V}_j|^2) \\ &= \min_{\mathbf{A}, \mathbf{B}} \frac{1}{2}(\text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B})) \end{aligned}$$

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

[Fazel Hindi Boyd 2001]

SDP

$$\text{minimize } \text{loss}(\mathbf{X}) + \lambda/2(\text{tr}(\mathbf{A})+\text{tr}(\mathbf{B}))$$

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

Optimization with the Trace-Norm

Not Convex!

But no local (non-global) minima!

follows from [Burer Monteiro 03]

$$\text{minimize } \text{loss}(\mathbf{U}'\mathbf{V}) + \lambda/2(\sum_i |\mathbf{U}_i|^2 + \sum_j |\mathbf{V}_j|^2)$$
$$\mathbf{X}=\mathbf{UV}$$

- Unconstrained opt on \mathbf{U}, \mathbf{V}
[Rennie S 05]
- Stochastic Grad Desc on \mathbf{U}, \mathbf{V}
[Salakhutdinov Minh 08]

Convex

$$|\mathbf{X}|_{\text{tr}} = \sum (\text{singular values of } \mathbf{X}) = \min_{\mathbf{X}=\mathbf{U}'\mathbf{V}} \frac{1}{2}(\sum_i |\mathbf{U}_i|^2 + \sum_j |\mathbf{V}_j|^2)$$
$$\text{minimize } \text{loss}(\mathbf{X}) + \lambda \cdot \overbrace{|\mathbf{X}|_{\text{tr}}}$$

- Subgradient descent
- Smoothed Subgradient descent
- Singular Value Thresholding
[Cai Candes Shen 08]
- . . .

All require SVD at each iteration

SDP

$$\text{minimize } \text{loss}(\mathbf{X}) + \lambda/2(\text{tr}(\mathbf{A})+\text{tr}(\mathbf{B}))$$

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

- Can use standard SDP code
But only on small problems...

Optimization with the Max-Norm

Not Convex!

But no local (non-global) minima! (if U, V of high enough dim)
follows from [Burer Monteiro 03]

$$\text{minimize } \text{loss}(\mathbf{U}'\mathbf{V}) + \lambda/2(\max_i |\mathbf{U}_i|^2 + \max_j |\mathbf{V}_j|^2)$$

$$\mathbf{X} = \mathbf{U}\mathbf{V}$$

Convex

$$\text{minimize } \text{loss}(\mathbf{X}) + \lambda \cdot \overbrace{|\mathbf{X}|_{\max}} = \cancel{\sum(\text{sing val of } \mathbf{X})} = \min_{\mathbf{X}=\mathbf{U}'\mathbf{V}} \frac{1}{2}(\max_i |\mathbf{U}_i|^2 + \max_j |\mathbf{V}_j|^2)$$

$$= \min_{\mathbf{A}, \mathbf{B}} \frac{1}{2}(\max \text{diag}(\mathbf{A}) + \max \text{diag}(\mathbf{B}))$$

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

SDP

$$\text{minimize } \text{loss}(\mathbf{X}) + \lambda/2(\max \text{diag}(\mathbf{A}) + \max \text{diag}(\mathbf{B}))$$

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

Trace-Norm vs Max-Norm Optimization

	$ X _{\text{tr}}$	$ X _{\text{max}}$
SDP Represent able?	Yes	Yes
Calculating $ X $	sum(SVD)	Requires solving SDP (as hard as training)
1 st Order methods on X	Yes, SVD at each iteration (and also SVT)	No
Unconstrained Opt. of U,V (including Stoch Grad Descent)	Yes	Yes

[Lee Recht Salakhutdinov
S Tropp NIPS'10]

Trace-Norm vs Max-Norm: Empirical Matrix Completion Results

- 2004 experiments on 100x100 subsets of MovieLens and EachMovie using standard SDP code:

MaxNorm better 90% of the time

- 2010 experiments on full NetFlix using stochastic gradient descent on U,V:

	RMSE	%improvement
NetFlix Cinematch:	0.9525	0 (baseline)
TraceNorm:	0.9235	3.04
MaxNorm:	0.9138	4.06
Winning team:	0.8553	10.2

Reconstruction Guarantees

$$Y = M + Z$$

- Using the **trace-norm**:

$$\text{\#samples to get } \frac{1}{nm} \|\hat{X} - Y\|_1 \leq \frac{1}{nm} \|Z\|_1 + \epsilon$$

$$\propto \text{tc}(M) \cdot (n+m) \log(n) / \epsilon^2 \leq \text{rank}(M) \cdot (n+m) \log(n) / \epsilon^2$$

$$\text{tc}(X) = \|X\|_{\text{tr}}^2 / nm \leq \text{rank}(X) \|X\|_F^2 / nm$$

$$|M_{ij}|^2 = O(1) \text{ on average}$$

- Using the **max-norm**:

$$\text{\#samples to get } \frac{1}{nm} \|\hat{X} - Y\|_1 \leq \frac{1}{nm} \|Z\|_1 + \epsilon$$

$$\propto \text{mc}(M) \cdot (n+m) / \epsilon^2 \leq \text{rank}(M) \cdot (n+m) / \epsilon^2$$

$$\text{mc}(X) = \|X\|_{\text{max}}^2 \leq \text{rank}(X) \|X\|_{\infty}^2$$

$$\sigma^2 = \frac{1}{nm} \|Z\|_2^2$$

$$|M_{ij}|^2 = O(1) \text{ for all } i, j$$

$$\text{\#samples to get } \frac{1}{nm} \|\hat{X} - Y\|_2^2 \leq \sigma^2 + \epsilon$$

$$\propto \text{rank}(M) \cdot (n+m) \cdot \frac{\log^3 1/\epsilon}{\epsilon} \cdot \frac{\sigma^2 + \epsilon}{\epsilon}$$

Reconstruction Guarantees

$$Y=M+Z$$

- Using the **trace-norm**:

$$\text{\#samples to get } \frac{1}{nm} \|\hat{X} - Y\|_1 \leq \frac{1}{nm} \|Z\|_1 + \epsilon$$

$$\propto \text{tc}(M) \cdot (n+m) \log(n) / \epsilon^2 \leq \text{rank}(M) \cdot (n+m) \log(n) / \epsilon^2$$

$$\text{tc}(X) = \|X\|_{\text{tr}}^2 / nm \leq \text{rank}(X) (\|X\|_F^2 / nm)$$

Only under uniform sampling.

$$|M_{ij}|^2 = O(1) \text{ on average}$$

- Using the **max-norm**:

$$\text{\#samples to get } \|\hat{X} - Y\|_{L_1(\mathcal{D})} \leq \|Z\|_{L_1(\mathcal{D})} + \epsilon$$

$$\propto \text{mc}(M) \cdot (n+m) / \epsilon^2 \leq \text{rank}(M) \cdot (n+m) / \epsilon^2$$

$$\text{mc}(X) = \|X\|_{\text{max}}^2 \leq \text{rank}(X) \|X\|_{\infty}^2$$

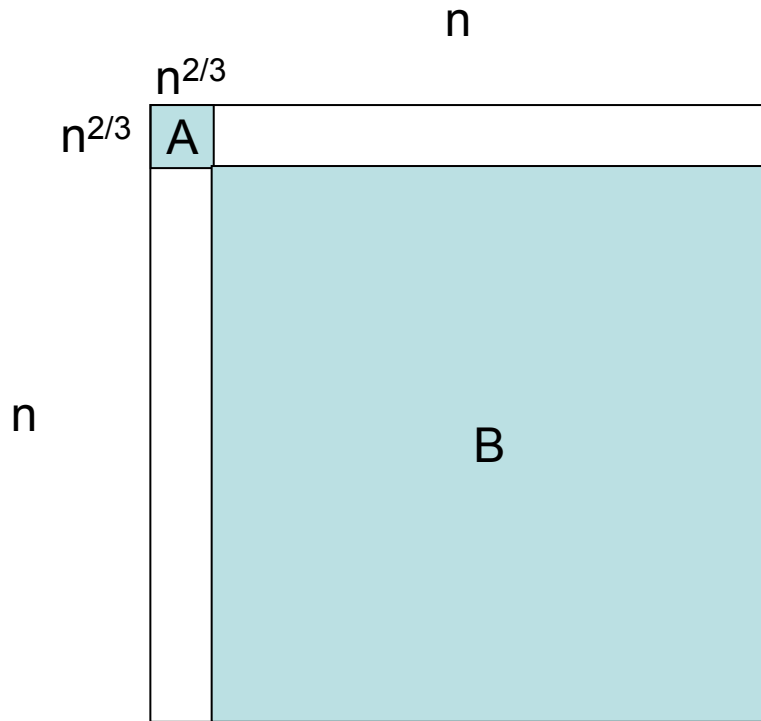
$$|M_{ij}|^2 = O(1) \text{ for all } i, j$$

$$\text{\#samples to get } \|\hat{X} - Y\|_{L_2(\mathcal{D})}^2 \leq \sigma^2 + \epsilon$$

$$\propto \text{rank}(M) \cdot (n+m) \cdot \frac{\log^3 1/\epsilon}{\epsilon} \cdot \frac{\sigma^2 + \epsilon}{\epsilon}$$

Under arbitrary sampling!

The Trace-Norm with Non-Uniform Sampling



- Both A, B of rank 2
- Sampling:
 - uniform in A w.p. $\frac{1}{2}$
 - uniform in B w.p. $\frac{1}{2}$

- Regularizing with the **rank** or with the **max-norm**:
sample complexity $\propto n$
- Regularizing with the **trace-norm**:
sample complexity $\propto n^{4/3}$!!!

Correcting for Non-Uniform Sampling: The Weighted Trace-Norm

[Salakhutdinov S NIPS'10]

- Instead of the Trace-Norm:

$$|X|_{tr} = \min_{X=UV} \frac{1}{2} (\sum_i |U_i|^2 + \sum_j |V_j|^2)$$

- Regularize with the weighted Trace-Norm:

$$|X|_{tr(p)} = \min_{X=UV} \frac{1}{2} (\sum_i \mathbf{p(i)} |U_i|^2 + \sum_j \mathbf{p(j)} |V_j|^2)$$

frequency of choosing row i

frequency of choosing column j

- On NetFlix:

	RMSE	%improvement
NetFlix Cinematch:	0.9525	0 (baseline)
TraceNorm:	0.9235	3.04
MaxNorm:	0.9138	4.06
Weighted TraceNorm:	0.9105	4.41
Winning team:	0.8553	10.20

Max-Norm and Trace-Norm Regularization

- Both lower bound on rank
- Both SDP representable, with practical large-scale opt method
- Reconstruction guarantees for both (better for max-norm)
- Good empirical performance (better for max-norm)
- **Better model than the rank: infinite factor model**
- Not only collaborative filtering!
 - Multi-task and multi-class learning
[Argyriou et al 07] [Abernethy et al 08] [Amit et al 07]
 - Semi-supervised learning following eg [Rie Zhang 05,07]

Maximum Margin Matrix Factorization, S, Rennie, Jaakkola, NIPS 2004

Rank, Trace-Norm and Max-Norm S, Shraibman, COLT 2005

Collaborative Filtering in a Non-Uniform World: Learning with the Weighted Trace Norm
Salakhutdinov, S, NIPS 2010

Practical Large-Scale Optimization for Max-Norm Regularization

Lee, Recht, Salakhutdinov, S, Tropp, NIPS 2010

Concentration-Based Guarantees for Low-Rank Matrix Reconstruction, Foygel, S, 2011