

Maximum Margin Matrix Factorization

Nathan Srebro
Department of Computer Science
University of Toronto

Jason Rennie
Computer Science and Artificial Intelligence Laboratory
Massachusetts Institute of Technology

Tommi Jaakkola

Matrix Factorization

Unconstrained: Low Rank Approximation

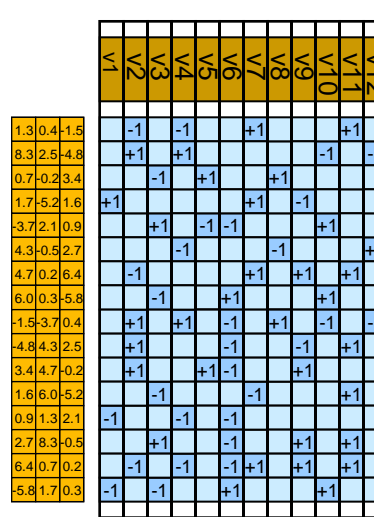
$$Y \approx UV^T = X_{\text{rank } k}$$

- Additive Gaussian noise: minimize $\|Y - UV\|_{Fro}$
- General additive noise
- General conditional models
- Multiplicative noise, Exponential-PCA (Collins+01), Multinomial (pLSA (Hofman01), etc)
- General loss functions
- Hinge loss, loss functions appropriate for ratings, etc

Constrained U, V: recover more factors

- Non-Negativity [LeeSeung99]
 - Stochasticity (convexity) [LeeSeung97] [Hofman01]
 - Sparsity
 - Clustering as an extreme (when rows of U sparse)
- Overall number of factors still constrained
Non-convex optimization problems

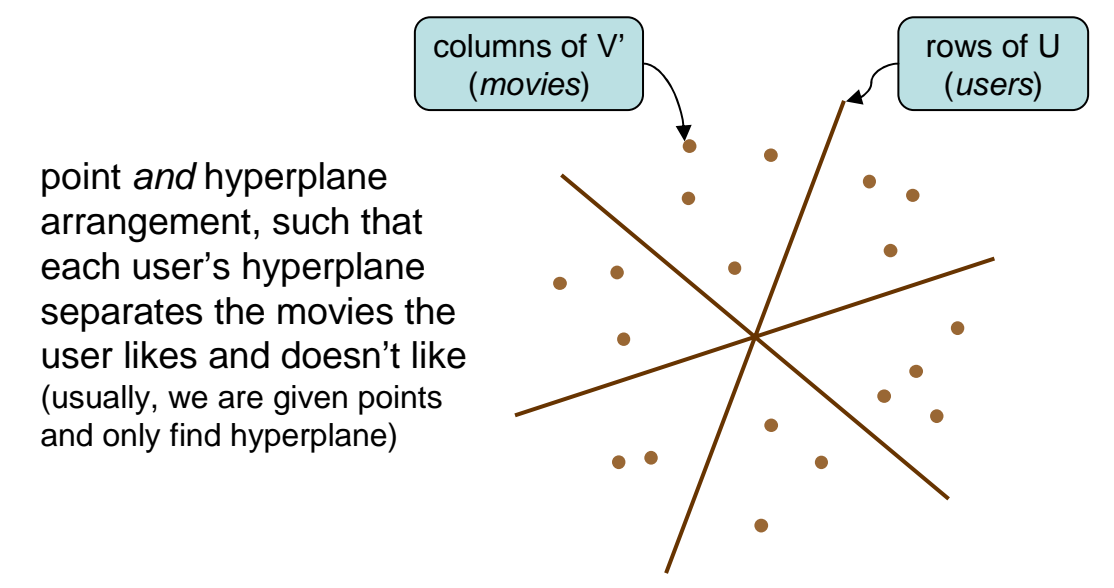
Matrix Factorization as Feature Learning



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.



Max-Margin Matrix Factorization

Instead of bounding dimensionality of U, V, bound norms of U, V

low norm V'

bound norms on average: $(\sum_i |U_i|^2) (\sum_j |V_j|^2) \leq 1$

bound norms uniformly: $(\max_i |U_i|^2) (\max_j |V_j|^2) \leq 1$

For observed $Y_{ij} \in \pm 1$: $Y_{ij} X_{ij} \geq \text{Margin}$

U is fixed: each column of V is SVM

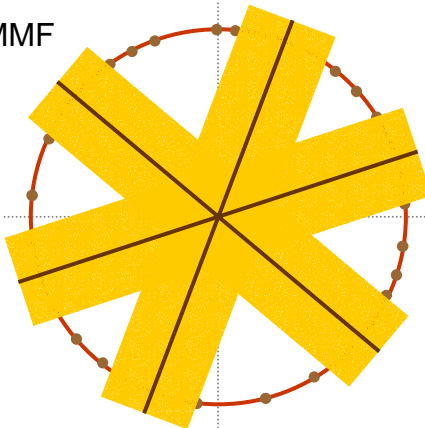
Unlike $\text{rank}(X) \leq k$, these are convex constraints!

$$\alpha U_1 V_1 + (1-\alpha) U_2 V_2 = (\sqrt{\alpha} U_1, \sqrt{1-\alpha} U_2) \cdot \left(\frac{\sqrt{\alpha} V_1}{\sqrt{1-\alpha} V_2} \right)$$

Geometric Interpretation

for max-norm (uniform) MMMF

point and hyperplane arrangement in infinite dimensional unit sphere, such that hyperplane separate according to Y with large margin.

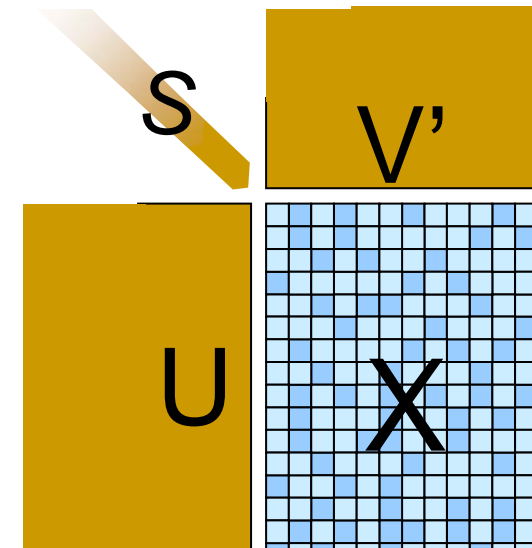


MMMF, Rank and the SVD

$\text{rank}(X) = |\text{singular values}|_0$

trace-norm (average) MMMF minimizes: $\|X\|_{tr} = |\text{singular values}|_1$

[Fazel Hindi Boyd 2001] suggest $\|X\|_{tr}$ as convex surrogate to $\text{rank}(X)$. Here, we justify it directly through the connection with max-margin classification, and by providing generalization error bounds.



Finding Max-Margin Matrix Factorizations

trace-norm (average) MMMF

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

$$\begin{aligned} &\text{primal SDP} \\ &\text{minimize } \text{tr}(A) + \text{tr}(B) + c \sum \xi_{ij} \\ &Y_{ij} X_{ij} \geq 1 - \xi_{ij} \\ &\begin{pmatrix} A & X \\ X^T & B \end{pmatrix} \text{ p.s.d.} \end{aligned}$$

$$\begin{aligned} \|X\|_{tr} &= \sum (\text{singular values of } X) \\ &= \min_{X=UV} \sqrt{(\sum_i |U_i|^2) (\sum_j |V_j|^2)} \\ &= \min_{X=UV} \frac{1}{2} (\text{tr}(A) + \text{tr}(B)) \\ &= \min_{A,B} \frac{1}{2} (\text{tr}(A) + \text{tr}(B)) \\ &\text{s.t. } \begin{pmatrix} U^T A U & X \\ X^T & V^T B V \end{pmatrix} \text{ p.s.d.} \end{aligned}$$

[Fazel Hindi Boyd 2001]

$$\begin{aligned} &\text{dual SDP} \\ &\text{maximize } \sum Q_{ij} \\ &0 \leq Q_{ij} \leq c \\ &\|Q \otimes Y\|_2 \leq 1 \\ &\text{sparse elementwise product (zero for unobserved entries)} \end{aligned}$$

max-norm (uniform) MMMF

$$\begin{aligned} &\text{maximize } M \\ &Y_{ij} X_{ij} \geq M \\ &X = UV^T \\ &(\max_i |U_i|) (\max_j |V_j|) \leq 1 \\ &\|X\|_{\max} \end{aligned}$$

$$\begin{aligned} &\text{primal SDP} \\ &\text{minimize } t + c \sum \xi_{ij} \\ &Y_{ij} X_{ij} \geq 1 - \xi_{ij} \\ &\begin{pmatrix} A & X \\ X^T & B \end{pmatrix} \text{ p.s.d.} \\ &A_{ii} \leq t, B_{jj} \leq t \end{aligned}$$

Reconstructing Primal X^* from Dual Q^*

X^* spanned by $Q^* \otimes Y$ SVD components of singular value 1

For trace-norm problems without slack, the primal optimal X^* can be extracted from dual optimal Q^* :

- Compute the SVD: $Q^* \otimes Y = U \Lambda V^T$
- Let U^*, V^* be components of U, V with value 1
- Primal optimal is of the form $X^* = U^* R R^T V^{*T}$
- Solve linear equations in RR^T , with $Q^*_{ij} > 0 \Rightarrow X^*_{ij} = Y_{ij}$

Querying Primal X^*_{ij} from Dual Q^*

Add constraint $X_{ij} > 0$ to primal \Rightarrow Add variable Q_{ij} to dual Q^* still feasible.

BUT: No optimal solution with $X^*_{ij} > 0 \Rightarrow Q^*$ not optimal Q^* still optimal $\Rightarrow X^*_{ij} > 0$

To query if $\text{sign}(X^*_{ij})$:

Add $Q^*_{ij} = 0$ to Q^* with $Y_{ij} = 1$ and reoptimize

Add $Q^*_{ij} = 0$ to Q^* with $Y_{ij} = -1$ and reoptimize

MMMF as a Convex Combination

$$\begin{aligned} \{X = UV^T \mid (\sum |U_i|^2) (\sum |V_j|^2) \leq 1\} \\ = \text{conv-hull}(\{u v^T \mid u \in \mathbb{R}^n, v \in \mathbb{R}^m, |u|=|v|=1\}) \\ \text{conv}(\{u v^T \mid u \in \pm 1^n, v \in \pm 1^m\}) \\ \subset \{X = UV^T \mid (\max |U_i|^2) (\max |V_j|^2) \leq 1\} \\ \subset 2 \text{conv}(\{u v^T \mid u \in \pm 1^n, v \in \pm 1^m\}) \end{aligned}$$

Grothendieck's Inequality

Major Assumption: Random Observations

Although we did not make any assumptions about the true preferences Y, we made a very strong assumption about the set S of observed entries: we assumed entries as selected uniformly at random.

For $(\sum |U_i|^2/n) (\sum |V_j|^2/m) \leq R^2$, uniformity crucial.

For $(\max |U_i|^2) (\max |V_j|^2) \leq R^2$ and $\text{rank}(X) \leq k$, S need not be uniform:

$$\begin{aligned} D(X; Y) &= E_{ij} [\text{loss}(X_{ij}; Y_{ij})] & D_S(X; Y) &= \sum_{ij \in S} \text{loss}(X_{ij}; Y_{ij}) / |S| \\ \forall Y \Pr_S (\forall X D(X; Y) < D_S(X; Y) + \epsilon) &> 1 - \delta \end{aligned}$$

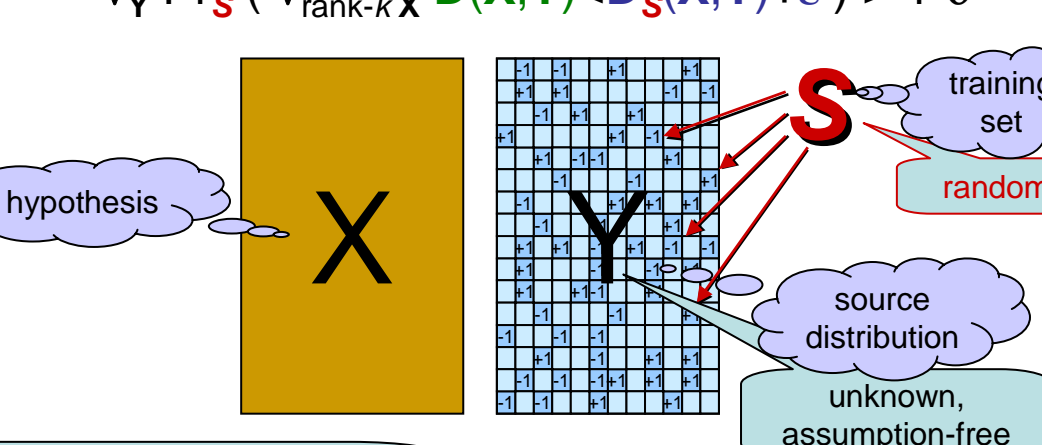
same observation distribution

Not very satisfying: we are guaranteed good generalization only on items the user is likely to observe on its own—not on items we might recommend.

Generalization Error Bounds

$$\begin{aligned} D(X; Y) &= \#_{ij \in S} (X_{ij} \neq Y_{ij}) / |S| & \text{generalization error} \\ D_S(X; Y) &= \#_{ij \in S} (X_{ij} \neq Y_{ij}) / |S| & \text{empirical error} \end{aligned}$$

$$\forall Y \Pr_S (\forall \text{rank-}k \text{ } X D(X; Y) < D_S(X; Y) + \epsilon) > 1 - \delta$$



Universal constant from bound on spectral norm of random matrix [Singer09]

$$(\sum |U_i|^2/n) (\sum |V_j|^2/m) \leq R^2: \quad \epsilon = K \sqrt{\ln m} \frac{R^2 (n+m) \log n + \log \frac{1}{\delta}}{|S|}$$

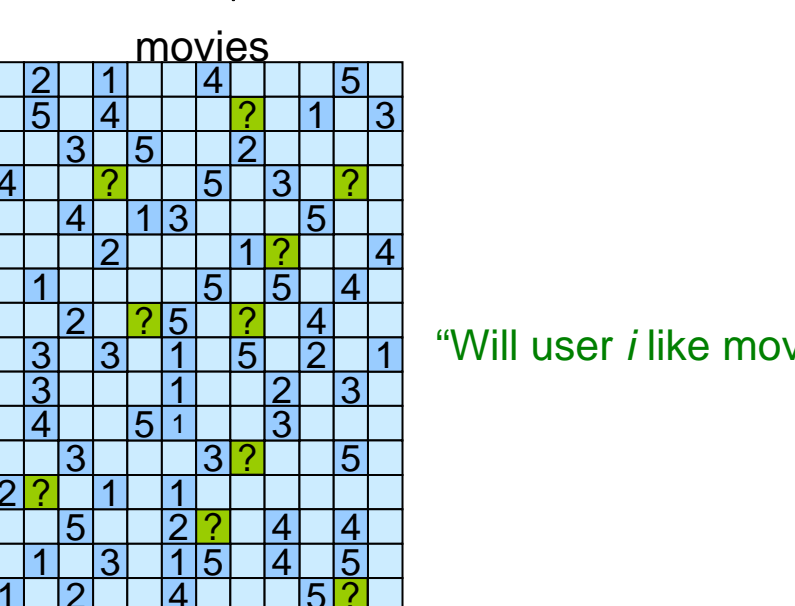
$$(\max |U_i|^2) (\max |V_j|^2) \leq R^2: \quad \epsilon = 12 \sqrt{\frac{R^2 (n+m) + \log \frac{1}{\delta}}{|S|}}$$

Compare with the low-rank bound: [Poster tomorrow!]

$$\text{rank}(X) \leq k: \quad \epsilon = \sqrt{\frac{k(n+m) \log \frac{nm}{k} + \log \frac{1}{\delta}}{2|S|}}$$

Collaborative Prediction

Based on partially observed matrix \Rightarrow Predict unobserved entries



Fit low-rank matrix $X = UV^T$ to observed entries.

$$\text{minimize } \sum_{j \in S} \text{loss}(X_{ij}; Y_{ij})$$

Use matrix X to predict unobserved entries.

Maximum Margin Matrix Factorization is as an alternative to Low Rank methods:

- Allows an unbounded number of factors
- Convex optimization problem: sparse SDP
- Correspondence with large margin linear classification
- Generalization error bounds
- Applicable in other applications where low-rank approximations are currently used

Direct optimization of dual would enable large-scale applications

Experiments

Preliminary experiments on 100 user \times 100 movie subset of MovieLens

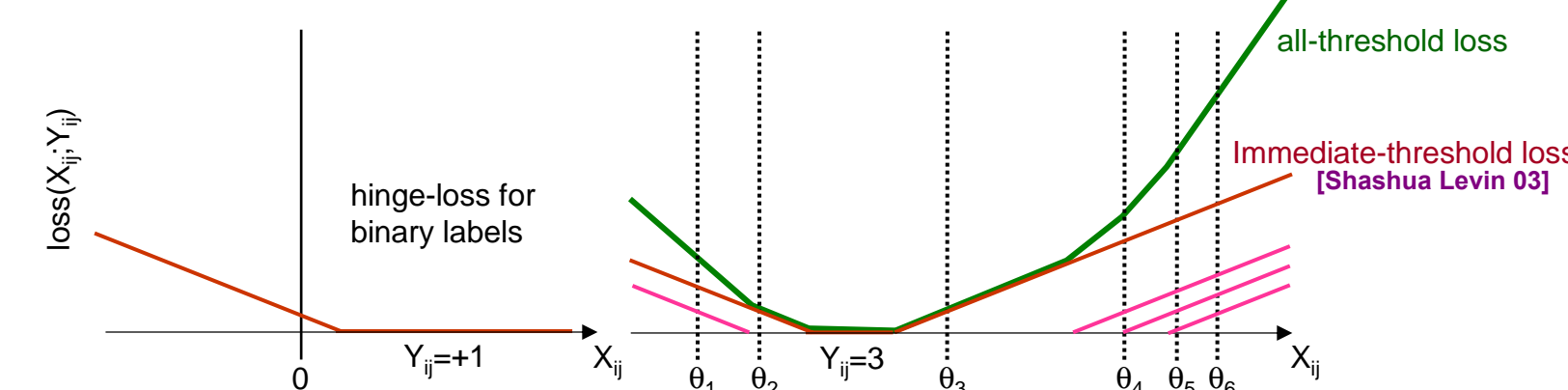
Two level cross validations:

- Train all variants, with various regularization parameters, on 50% of ratings
- Validate on 25% of data to select best variant and parameters (3-fold CV on 75% of data)
- Evaluate single variant and parameters on held out 25% of data

Compare trace-norm and max-norm MMMF to low-rank approximation minimizing sum-squared error and to K-medians clustering of users.

holdout set	rank agreement error			mean rank difference		
	rank-2 aprox			2-medians		
1	0.575	0.562	max-norm, c=0.12	0.691	0.677	max-norm, c=0.12
2	0.562	0.552	trace-norm, c=0.24	0.683	0.681	max-norm, c=0.11
3	0.543	0.527	max-norm, c=0.12	0.681	0.646	max-norm, c=0.12
4	0.558	0.550	max-norm, c=0.12	0.696	0.686	max-norm, c=0.12
average	0.558	0.548		0.687	0.673	

Generalizations of hinge-loss for ordinal labels



- All-threshold loss is a bound on the absolute rank-difference
- We experimented with both: "all-thresholds" consistently outperformed "immediate-threshold"
- For both loss functions: learn per-user θ 's (no extra cost to SDP)

MATLAB code available @ <http://www.cs.toronto.edu/~nati/mmmf>

Different matrix factorization methods differ in how they relate real-valued entries in X to the observations (preferences) Y, possibly through a probabilistic model, and in the associated contrast (loss) functions.

Low-rank models of co-occurrence or frequency data

	Multinomial	Independent Binomials	Independent Bernoulli
Mean parameterization $0 \leq X_{ij} \leq 1$ $E[Y_{ij} X_{ij}] = X_{ij}$	Aspect Model (pLSA) [Hofman+99]	$Y_{ij} X_{ij} \sim \text{Bin}(N, X_{ij})$	$P(Y_{ij}=1) = X_{ij}$
Natural parameterization unconstrained X_{ij}	SDR [Giberson+02]	$Y_{ij} X_{ij} \sim \text{Bin}(N, g(X_{ij}))$	Logistic Low Rank Approximation [Schein+03]

row features most informative about columns

$g(x) = 1/(1+e^{-x})$