

Generalization Error Bounds for Collaborative Prediction with Low-Rank Matrices

Nathan Srebro
Department of Computer Science, University of Toronto

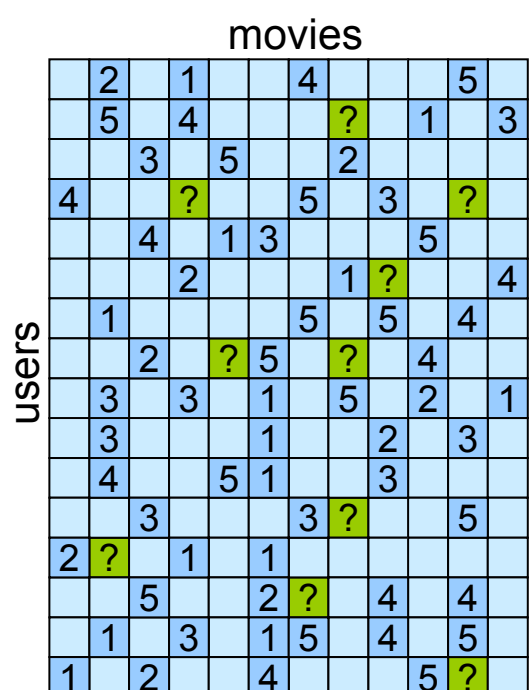
Noga Alon
School of Mathematical Sciences, Tel Aviv University

Tommi Jaakkola
Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology

Results

Collaborative Prediction

Based on partially observed matrix
→ Predict unobserved entries



Generalization Error Bounds

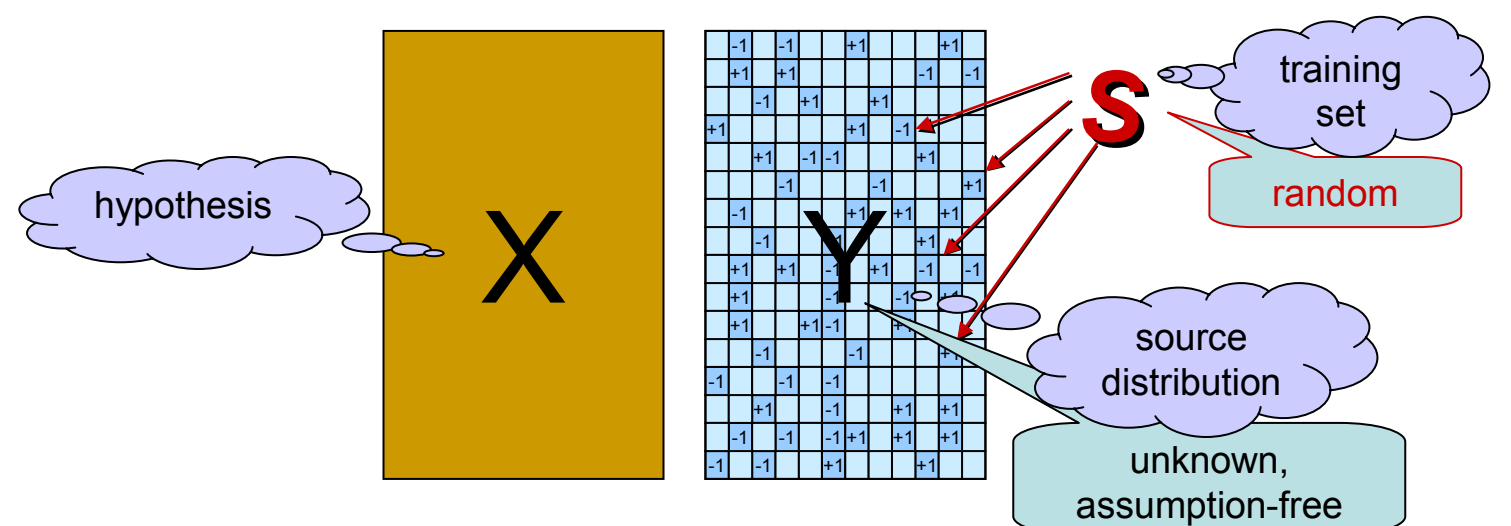
$$D(\mathbf{X}; \mathbf{Y}) = \sum_{ij \in S} \text{loss}(X_{ij}; Y_{ij}) / nm$$

generalization error

$$D_S(\mathbf{X}; \mathbf{Y}) = \sum_{ij \in S} \text{loss}(X_{ij}; Y_{ij}) / |S|$$

empirical error

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



$$\text{loss}(x, y) = \begin{cases} 0 & \text{sign}(x) = y \\ 1 & \text{sign}(x) \neq y \end{cases}$$

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

monotone $\text{loss}(x, y) \leq 1$:

$$\epsilon = 6 \sqrt{\frac{k(n+m) \log \frac{8em}{k} \log \frac{|S|}{k(n+m)} + \log \frac{1}{\delta}}{|S|}}$$

Prior work

- Assuming a low-rank structure (eigengap) in \mathbf{Y} , predict entries:
 - Asymptotic behavior [Azar, Fiat, Karlin, McSherry Saia Spectral analysis of data STOC 2001]
 - Sample complexity, query strategy [Drineas, Kerenidis, Raghavan Competitive recommendation systems STOC 2002]
 - Bounds on residual errors, no assumptions on \mathbf{Y} : [Shaw-Taylor, Cristianini, Kandola On the concentration of spectral properties NIPS 2002]
 - Subset of rows fully observed, bound is on distance of new rows to learned subspace
- In this work: collaborative prediction analysis (entry prediction), no assumptions on \mathbf{Y} .

Major Assumption: Random Observations

Although we did not make any assumptions about the true preferences \mathbf{Y} , we made a very strong assumption about the set S of observed entries: we assumed entries as selected uniformly at random. Although the uniformity requirement can be relaxed:

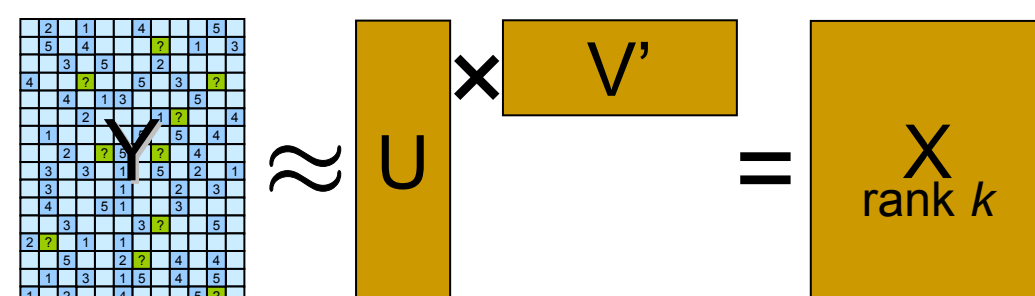
$$D(\mathbf{X}; \mathbf{Y}) = E_{ij} [\text{loss}(X_{ij}; Y_{ij})] \quad D_S(\mathbf{X}; \mathbf{Y}) = \sum_{ij \in S} \text{loss}(X_{ij}; Y_{ij}) / |S|$$

same observation distribution

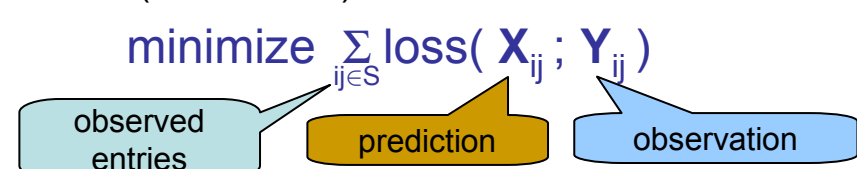
$$\forall \mathbf{Y} \Pr_S (\forall_{\text{rank-}k \mathbf{X}} D(\mathbf{X}; \mathbf{Y}) < D_S(\mathbf{X}; \mathbf{Y}) + \epsilon) > 1 - \delta$$

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

Low-Rank Matrix Factorization



Fit low-rank (factorizable) matrix $\mathbf{X} = \mathbf{U}\mathbf{V}$ to observed entries.



Use matrix \mathbf{X} to predict unobserved entries.

- [Sarwar, Karypis, Konstan, Riedl Applications of dimensionality reduction in recommender systems—a case study WebKDD 2000]
- [Hoffman Latent semantic models for collaborative filtering ACM Trns. Inf. Syst. 2004]
- [Marlin, Zemel Modeling user rating profiles for collaborative filtering NIPS 2003]
- [Cannedy GAP: A Factor Model for Discrete Data SIGIR 2004]
- many others...

Different low-rank methods differ in how they relate real-valued entries in \mathbf{X} to the observations (preferences) \mathbf{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) [Hoffman+99] | $Y_{ij} X_{ij} \sim \text{Bin}(N, X_{ij})$ | $P(Y_{ij}=1) = X_{ij}$ |
| Natural parameterization unconstrained X_{ij} | SDR [Globerson+02] | $Y_{ij} X_{ij} \sim \text{Bin}(N, g(X_{ij}))$ | Logistic Low Rank Approximation [Schein+03] |
| | | Exponential PCA: [Collins+01] $p(Y_{ij} X_{ij}) \propto \exp(Y_{ij} X_{ij} + F(Y_{ij}))$ | hinge loss |

row features most informative about columns

$$g(x) = 1/(1+e^x)$$

Proofs

Binary Labels, Zero-One Loss

For particular \mathbf{X}, \mathbf{Y} :

$$\text{loss}(X_{ij}; Y_{ij}) \sim \text{Bernoulli}(D(\mathbf{X}; \mathbf{Y}))$$

$$\Pr (D_S(\mathbf{X}; \mathbf{Y}) < D(\mathbf{X}; \mathbf{Y}) - \epsilon) < e^{-2|S|\epsilon^2}$$

Union bound over all possible \mathbf{X} s: $\forall \mathbf{X} \Pr_S (D(\mathbf{X}; \mathbf{Y}) < D_S(\mathbf{X}; \mathbf{Y}) + \epsilon) > 1 - \delta$

$$\epsilon = \sqrt{\frac{\log(\# \text{ possible } \mathbf{X}\text{s}) + \log \frac{1}{\delta}}{2|S|}}$$

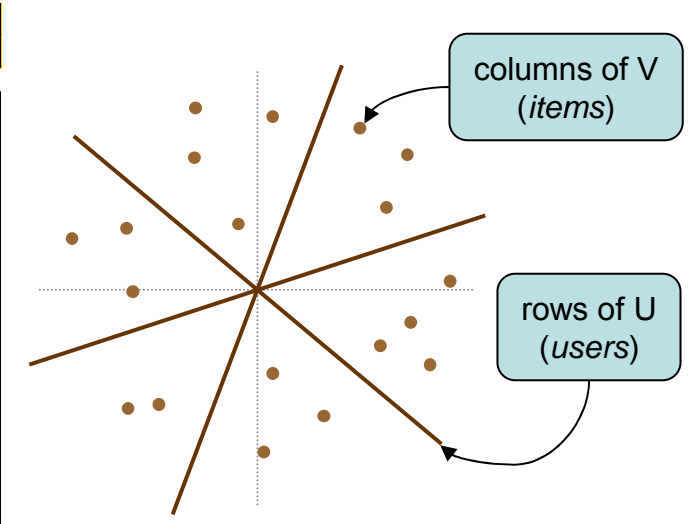
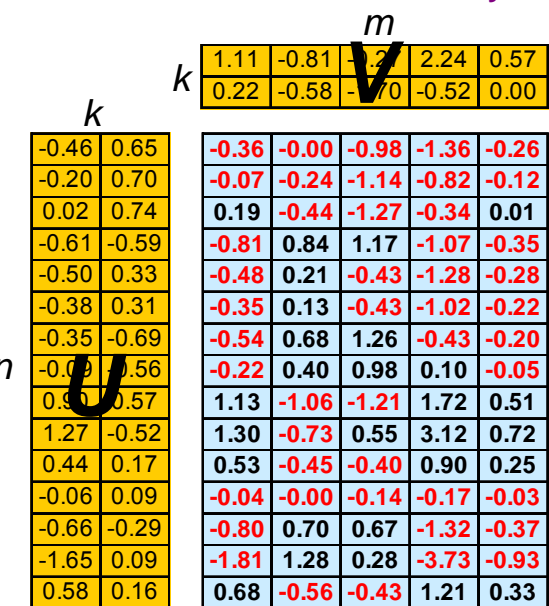
The bound rests on bounding the number of possible \mathbf{X} s. The behavior of $\text{loss}(X_{ij}; Y_{ij})$ only depends on $\text{sign}(\mathbf{X})$, and so it is enough to bound the number of sign configurations:

$$F(n, m, k) = \{ \text{sign}(\mathbf{X}) \in \{-, +\}^{n \times m} \mid \mathbf{X} \in \mathbb{R}^{n \times m}, \text{rank } \mathbf{X} \leq k \}$$

$$f(n, m, k) = \#F(n, m, k)$$

Sign Configurations of Low-Rank Matrices

Following [Alon Tools from higher algebra Handbook of Combinatorics 1995], similar to [Alon, Frankl, Rödl Geometric realization of set systems and probabilistic communication complexity FOCS 1985]



Warren (1968): The number of connected components of $\{ \mathbf{x} \mid \forall_i P_i(\mathbf{x}) \neq 0 \}$

$$\text{is at most } \frac{4e \cdot (\text{degree}) \cdot (\# \text{polys})}{(\# \text{variables})} \cdot (\# \text{variables})$$

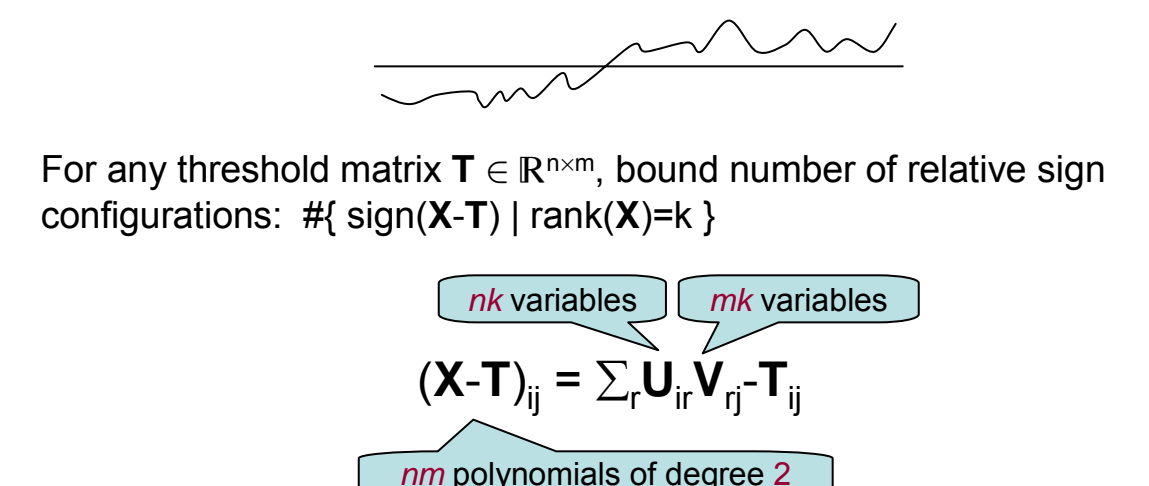
Corollary: The number of sign configurations of the polynomials is at most $\frac{4e \cdot (\text{degree}) \cdot (\# \text{polys})}{(\# \text{variables})} \cdot (\# \text{variables})$

$$\mathbf{X}_{ij} = \sum_r \mathbf{U}_{ir} \mathbf{V}_{rj}$$

$$f(n, m, k) \leq \frac{4e \cdot 2 \cdot nm}{k(n+m)}^{k(n+m)} = 2^{k(n+m) \log(8em/k)}$$

General Bounded Loss Functions

For 0/1 loss, behavior of entries of \mathbf{X} around zero enough. More generally, need to bound complexity of behavior everywhere.



$$\# \{ \text{sign}(\mathbf{X}-\mathbf{T}) \mid \text{rank}(\mathbf{X})=k \} \leq \frac{4e \cdot 2 \cdot nm}{k(n+m)}^{k(n+m)} = 2^{k(n+m) \log(8em/k)}$$

Viewing matrices as a mappings from index pairs to values: $(i, j) \mapsto X_{ij}$, this gives us a bound of $k(n+m) \log(8em/k)$ on the pseudo-dimension of rank- k matrices. We can now invoke standard results bounding the generalization error in terms of the pseudo-dimension.

A class \mathcal{F} of real-valued functions pseudo-shatters the points x_1, \dots, x_n with thresholds t_1, \dots, t_n , if for every binary labeling of the points $(s_1, \dots, s_n) \in \{+, -\}^n$ there exists $f \in \mathcal{F}$ s.t. $f(x_i) \leq t_i$ iff $s_i = -$. The pseudo-dimension of a class \mathcal{F} is the supremum over n for which there exist n points and thresholds that can be pseudo-shattered.

A Weaker Bound Using Realizable Oriented Matroids

For a fixed \mathbf{V} , each row is linear classification of columns in \mathbf{V} , and there are $\leq 2(k+1)m^{k-1}$ such classifications. Overall, for each fixed \mathbf{V} , the number of possible sign matrices is bounded by:

$$\# \{ \text{sign}(\mathbf{U}\mathbf{V}) \mid \mathbf{U} \in \mathbb{R}^{n \times k} \} \leq (2(k+1)m^{k-1})^n$$

This should be multiplied by the number of \mathbf{V} s, or rather the number of \mathbf{V} s yielding different sets of possible classification vectors:

$$M(\mathbf{V}) = \{ \text{sign}(u^T \mathbf{V}) \mid u \in \mathbb{R}^k \}$$

$$\# \{ M(\mathbf{V}) \mid \mathbf{V} \in \mathbb{R}^{k \times m} \} \leq m^{k(k+1)m} \cdot \text{extraneous } k^2 \text{ term!}$$

$$f(n, m, k) \leq (2(k+1)m^{k-1})^n m^{k(k+1)m} \leq 2^{k(n+m) \log(2m) + k^2 m \log(2m)}$$

Why not treat as combined classifiers?

- For MMMF (max-margin/low-norm matrix factorizations), generalization error bounds obtained by viewing MMMF as a "combined" classifier, a convex combination of unit-norm rank-1 matrices. Rank- k matrices can be viewed as "combined", or "voting" classifiers, each combining k rank-1 matrices. Can a similar approach be taken for low-rank matrices?
 - Scale-sensitive complexity (log covering numbers, Rademacher complexity) carries over to convex hull. (scale-invariant complexity certainly not conserved for convex hull)
 - VC-dimension scales gracefully with k for combinations of k classifiers
 - generalization error bounds for linear combinations of signs of low-rank matrices
 - Pseudo-dimension of a linear combinations of k functions from a low-pseudo-dimension class?

Counter Example: A family \mathcal{F} of functions closed under scalar multiplication, with pseudo-dimension 3, such that $\{ f_1, f_2 \mid f_1, f_2 \in \mathcal{F} \}$ has infinite pseudo-dimension:

$$\mathcal{F} = \{ \alpha f_A, \alpha g_A \mid \alpha \in \mathbb{R}, A \in \mathbb{N} \}$$

$$f_A(x) = 2^{A+1} x^A$$

$$g_A(x) = g^{A^2}$$

Consider a 1:1 mapping $\mathbb{N} \rightarrow 2^{\mathbb{N}}$. 'A' denotes both a number, and the corresponding subset.

Related Work

Warren's Theorem and Configuration Counting

Warren's Theorem, and a weaker result of Milnor, have a long history in combinatorics and learning theory:

[Goodman Pollack Upper bounds for configurations and polytopes in \mathbb{R}^d Disc Comp Geom 1986]
[Alon 1986 The number of polytopes, configurations and real matroids Mathematika 1986]
Bound on the number of non-equivalent point configurations (realizable oriented matroids). Can be used to obtain weak bound on number of sign configurations of low-rank matrices (green panel).

[Ben-David, Lindenbaum Localization vs. identification of semi-algebraic sets COLT 1993]
VC-dimension of set of transformations of an image, used to analyze sample complexity of determining location.

[Goldberg, Jerrum Bounding the VC dimension of concept classes parameterized by real numbers COLT 1993]
VC-dimension of any concept classes parameterized by real numbers, where each concept can be written as logical formula over polynomial inequalities

≤ 2 (# of params describing each concept) $\log(8e)$ (degree of polys used) (# of polys used)

Can be applied to collaborative prediction with low-rank matrices, where:

$$X_{(i,j)} = v_{ij} \quad (i=j \Rightarrow v_{ij} = \sum_r U_{ir} V_{rj}, v_{ij} > 0)$$

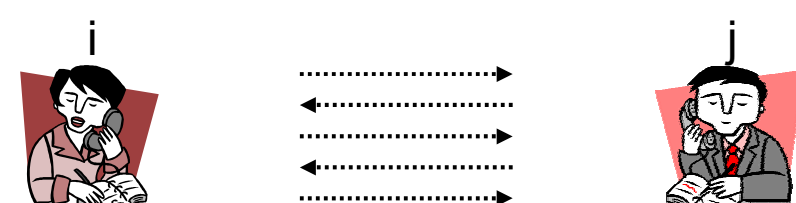
yielding: VC-dim(rank- k matrices) $\leq 2k(n+m) \log(8e \cdot 2 \cdot 3nm) \leq 2k(n+m) \log(48em)$

By directly applying Warren's Theorem we:

- avoided symmetrization (for 0/1 error)
- avoided Sauer's lemma, and a $\log|S|$ term in the generalization error bound
- bounded the pseudo-dimension and obtained generalization error bounds for general loss functions

More on Sign Configurations of Low-Rank Matrices

Unbounded Error Communication Complexity



$$Y(i, j) \in \{+1, -1\}$$

Unbounded error communication complexity C
= Randomized protocol, always $\leq C$ bits, $P(\text{correct answer}) > \frac{1}{2}$

[Paturi, Simon Probabilistic communication complexity FOCS 84]
 $\log \text{rank } \mathbf{X} \leq C \leq \lceil \log \text{rank } \mathbf{X} \rceil, \text{sign}(\mathbf{X}) = \mathbf{Y}$

[Alon, Frankl, Rödl Geometric realization of set systems and probabilistic communication complexity FOCS 1985]
Bound # sign configurations

counting arguments $\Rightarrow \exists \mathbf{Y}$ with $\text{rank}(\mathbf{X}) > n/32 \Rightarrow \exists \mathbf{Y}: \{0, 1\}^n \rightarrow \{0, 1\}$ with $C > \tau \cdot 5$

Embedability as Linear Classification

Can all concept classes be embedded as linear classifications in a low dimensional space?
 $C = \{c_1, \dots, c_n\}$ can be embedded as k -dimensional linear classification $\Leftrightarrow \text{Rank-}k \mathbf{X}$, s.t. $c(i) = \text{sign}(X_{ij})$
counting arguments $\Rightarrow \exists$ small concept class, not embeddable as low dimensional linear classification

Explicit Examples

These counting arguments provide only existence proofs, not explicit constructions of sign configurations with no low-rank realization (i.e. functions with high unbounded error communication complexity, or concept classes that cannot be embedded as low-dimensional linear classification).

[Forester A linear bound on the unbounded error communication complexity CCC 2001]
 $\text{rank}(\mathbf{X}) \geq n / \|\text{sign}(\mathbf{X})\|_2$ (spectral norm of $\text{sign}(\mathbf{X})$)
In particular, the $2^{\times 2}$ Hadamard matrix cannot be realized with $\text{rank}(\mathbf{X}) < 2^{1/2}$

In this example, $\text{rank}(\mathbf{X}) \geq \sqrt{n}$. No known explicit example with $\text{rank}(\mathbf{X}) = \Omega(n)$.