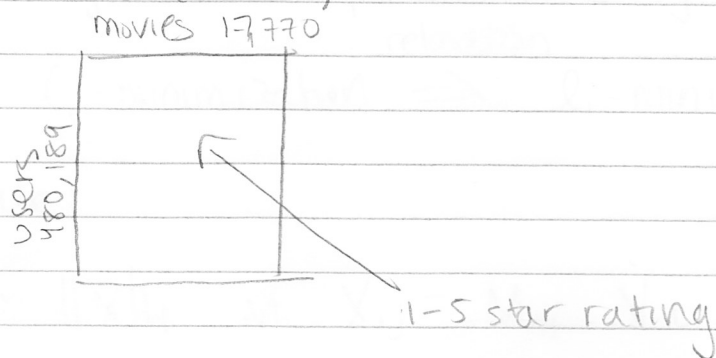


## Lecture #21: Matrix Completion

Netflix Prize (2006)



Can you improve their predictor's RMSE by 10%?

Eventually won by Bellkor's Pragmatic Chaos team in 2009, many lessons about aggregating predictors

e.g. Memento vs. Forest Gump

Popular abstraction:

- ① unknown matrix  $M$  (low-rank, incoherent)
- ② observe  $m$  entries of  $M$  chosen u.a.r.

Can we recover  $M$  exactly or approximately?

Natural, nonconvex approach

$$\min \text{rank}(X) \quad \text{s.t.} \quad X_{ij} = M_{ij} \quad \forall (i,j) \in \Omega$$

observed entries  
↓

Given a hypothesis  $X$  we can write its SVD as

$$X = U \Sigma V^T$$

where  $\Sigma = \text{diag}(\vec{\sigma})$ . Then  $\text{rank}(X) = \|\vec{\sigma}\|_0$

The same way as in compressed sensing  
relaxation

$\ell_0$ -minimization  $\Rightarrow$   $\ell_1$ -minimization

we can write:

$$\min \|X\|_* \quad \text{s.t.} \quad X_{ij} = M_{ij} \quad \forall (i,j) \in \Omega$$

where  $\|X\|_* =$  sum of singular values of  $X = \|\vec{\sigma}\|_1$ ,  
 $\uparrow$   
nuclear norm

This is now a convex relaxation, but does it work?

some intuition:

Lemma: The unit nuclear norm ball - i.e.

$$P = \{X \mid \|X\|_* \leq 1\} = \text{conv} \{ab^T \mid \|a\|_1, \|b\|_1 \leq 1\}$$

Proof: For any  $X \in P$  we have

$$X = \sum_{i=1}^r \sigma_i u_i v_i^T$$

$\uparrow$   
coefficients in convex combination

Conversely  $\|\cdot\|_*$  is a convex function so

$$\left\| \sum_i \lambda_i a_i b_i^T \right\|_* \leq \sum_i \lambda_i \|a_i b_i^T\|_* \leq 1 \quad \square$$

Today, we will take an empirical processes approach:

Suppose the true  $M \in^{n_1 \times n_2}$  satisfies  $\|M\|_* \leq r$  and  
 $\|M\|_\infty \leq \sqrt{\frac{r}{n_1 n_2}}$

many choices in literature  
↓

Here  $M$  is some type of incoherence — how aligned are the singular vectors of  $M$  with standard basis

Let  $X$  be the minimizer of (P), we have

$$\|X\|_* \leq \|M\|_* \leq r$$

Some terminology:

training error:  $\frac{1}{|\Omega|} \sum_{(i,j) \in \Omega} |X_{ij} - M_{ij}| = 0$  ↓ by assumption

test error:  $\frac{1}{n_1 n_2} \sum_{ij} |X_{ij} - M_{ij}|$

Q: If the training error is small, how can we bound the test error?

This must depend on the nuclear norm, following [Srebro, Shraibman]

$$\|X\|_* \leq r \approx \text{"simple" explanation for the observations}$$

$$\text{generalization error} = \text{test error} - \text{training error}$$

Consider the following hypothesis class:

$$\mathcal{H} = \left\{ X: [n_1] \times [n_2] \rightarrow \mathbb{R} \mid \|X\|_* \leq r, \|X\|_{\infty} \leq \sqrt{\frac{r}{n_1 n_2}} \right\}$$

↑  
interpreting  $X$  as a function from domain  $[n_1] \times [n_2]$  to a prediction

We will bound the following:

$$\sup_{x \in \mathcal{Z}} \left| \frac{1}{n_1 n_2} \sum_{ij} |X_{ij} - M_{ij}| - \frac{1}{|\mathcal{Z}|} \sum_{(ij) \in \mathcal{Z}} |X_{ij} - M_{ij}| \right|$$

[Koltchinskii], [Bartlett, Mendelson], ...

Theorem: Consider a set of classifiers

$$G: \mathcal{Z} \rightarrow [a, b]$$

with some distribution  $D$  on  $\mathcal{Z}$  with probability  $\geq 1 - \delta$

$$\sup_{g \in G} \left| \underbrace{\mathbb{E}[g(Y)]}_{\substack{\uparrow \\ \text{rv from } D}} - \frac{1}{m} \sum_{i=1}^m g(Y_i) \right| \quad (*)$$

$$\leq 2 R_m(G) + (b-a) \sqrt{\frac{\log 2/\delta}{2m}}$$

where  $R_m(G) = \mathbb{E}_{S_m} \left[ \underbrace{\mathbb{E}_\sigma}_{\substack{\uparrow \\ \text{random } \pm 1 \text{ signs}}} \left[ \sup_{g \in G} \frac{1}{m} \sum_{i=1}^m \sigma_i g(Y_i) \right] \right]$

$\uparrow$   
 rademacher complexity  
 set of  $m$  samples

Thought experiment: If  $g: \mathcal{Z} \rightarrow \{\pm 1\}$  then

rademacher complexity = how well does best hypothesis in  $G$  agree with random function on  $S$

Intuitively: If we can match random function  $\Rightarrow$  hypothesis class too expressive to generalize

How would we map this to matrix completion?

$$Z = \text{domain} = [n_1] \times [n_2]$$

$$g(i, j) = |X_{ij} - M_{ij}| = \text{loss on prediction for } (i, j)$$

Proof: (of most interesting part)

$$\text{Let } \hat{\mathbb{E}}_S[g] = \frac{1}{m} \sum_{i=1}^m g(Y_i) \text{ then}$$

$$\mathbb{E}_S \left[ \sup_{g \in G} (\mathbb{E}[g] - \hat{\mathbb{E}}[g]) \right] =$$

$$\mathbb{E}_S \left[ \sup_{g \in G} (\mathbb{E}_{S'} [\hat{\mathbb{E}}_{S'}[g]] - \hat{\mathbb{E}}_S[g]) \right]$$

where  $S' = Y'_1, \dots, Y'_m \stackrel{i.i.d.}{\sim} D$  "ghost samples"

convexity

$$\leq \mathbb{E}_{S, S'} \left[ \sup_{g \in G} (\hat{\mathbb{E}}_{S'}[g] - \hat{\mathbb{E}}_S[g]) \right]$$

$$= \mathbb{E}_{S, S'} \left[ \sup_{g \in G} \frac{1}{m} \left( \sum_i g(Y'_i) - g(Y_i) \right) \right]$$

$$= \mathbb{E}_{\sigma, S, S'} \left[ \sup_{g \in G} \frac{1}{m} \left( \sum_i \sigma_i (g(Y'_i) - g(Y_i)) \right) \right]$$

$$\leq \mathbb{E}_{\sigma, S, S'} \left[ \sup_{g \in G} \frac{1}{m} \sum_i \sigma_i g(Y'_i) + \sup_{g \in G} \sum_i (-\sigma_i) g(Y_i) \right]$$

$$= 2R_m(G)$$

This bands expected generalization error, can prove deviation bands using McDiarmid's inequality.

(can also prove other side of band.  $\square$ )

Aside: If  $G$  were finite, can bound generalization error through union bound + Chernoff bound

But  $G$  is infinite in our setting (and whenever you use a convex relaxation to get a predictor)

Can also bound the generalization error using

$$(*) \leq 2R_S(G) + 3(b-a) \sqrt{\frac{\log 4/\epsilon}{2m}}$$

$$\text{where } R_S(G) = \mathbb{E} \left[ \sup_{g \in G} \frac{1}{m} \sum_{i=1}^m \sigma_i g(Y_i) \right]$$

empirical rademacher complexity

This way we're not making any assumptions about the distribution, except that  $S$  comes from it

But how do we bound the rademacher complexity?

Let's consider a simpler problem:

$$R_m(H) = \mathbb{E} \left[ \mathbb{E} \left[ \sup_{x \in H} \frac{1}{m} \sum_{(i,j) \in \Omega} \sigma_{ij} X_{ij} \right] \right]$$

Let's drop the  $\|\cdot\|_2$  constraints

$$\mathbb{E} \left[ \frac{1}{\sigma} \mathbb{E} \left[ \sup_{\|X\|_2 \leq r} \sum_{(i,j) \in \Omega} \sigma_{ij} X_{ij} \right] \right] \leq r \mathbb{E} \left[ \frac{1}{\sigma} \mathbb{E} \left[ \sup_{\|X\|_2 \leq 1} \sum_{(i,j) \in \Omega} \sigma_{ij} X_{ij} \right] \right]$$

$$= \mathbb{E} \left[ \frac{1}{\sigma} \mathbb{E} \left[ \sup_{\|a\|_1, \|b\|_1 \leq 1} \sum_{(i,j) \in \Omega} \sigma_{ij} a_i b_j \right] \right]$$

$$= \mathbb{E} \left[ \frac{1}{\sigma} \mathbb{E} \left[ \sup_{\|a\|_1, \|b\|_1 \leq 1} \langle \Sigma, ab^T \rangle \right] \right]$$

$$= \mathbb{E} \left[ \|\Sigma\| \right]$$

$\uparrow$   
 $m$  nonzero entries chosen u.a.r,  
 random  $\pm 1$ s

The important point is:

$$\mathbb{E} \left[ \|\Sigma\| \right] \stackrel{\text{matrix concentration}}{\ll} \|\Sigma\|_F \|ab^T\|_F = \|\Sigma\|_F = \sqrt{m}$$

$\Rightarrow$  good generalization

$\uparrow$   
 trivial bound