

Matrix Completion

Stanley Chan

Purdue University

ECE/STAT 695 Spring 2018

The Netflix Challenge

Netflix Prize **COMPLETED**

Home Rules Leaderboard Update

Leaderboard

Showing Test Score. [Click here to show quiz score](#)

Display top 20 leaders.

Rank	Team Name	Best Test Score	% Improvement	Best Submit Time
Grand Prize - RMSE = 0.8567 - Winning Team: BellKor's Pragmatic Chaos				
1	BellKor's Pragmatic Chaos	0.8567	10.06	2009-07-26 18:18:28
2	The Ensemble	0.8567	10.06	2009-07-26 18:38:22
3	Grand Prize Team	0.8582	9.90	2009-07-10 21:24:40
4	Opera Solutions and Vandelay United	0.8588	9.84	2009-07-10 01:12:31
5	Vandelay Industries I	0.8591	9.81	2009-07-10 00:32:20
6	PragmaticTheory	0.8594	9.77	2009-06-24 12:06:56
7	BellKor in BigChaos	0.8601	9.70	2009-05-13 08:14:09
8	Dace	0.8612	9.59	2009-07-24 17:18:43
9	Feeds2	0.8622	9.48	2009-07-12 13:11:51
10	BigChaos	0.8623	9.47	2009-04-07 12:33:59
11	Opera Solutions	0.8623	9.47	2009-07-24 00:34:07
12	BellKor	0.8624	9.46	2009-07-26 17:19:11

Progress Prize 2008 - RMSE = 0.8627 - Winning Team: BellKor in BigChaos

- Predict user rating
- Missing data

Missing Data Problem

	Dirty Dancing	Meet the Parents	Top Gun	The Sixth Sense	Catch Me If You Can	The Royal Tenenbaums	Con Air	Big Fish	The Matrix	A Few Good Men
Customer 1	•	•	•	•	4	•	•	•	•	•
Customer 2	•	•	3	•	•	•	3	•	•	3
Customer 3	•	2	•	4	•	•	•	•	2	•
Customer 4	3	•	•	•	•	•	•	•	•	•
Customer 5	5	5	•	•	4	•	•	•	•	•
Customer 6	•	•	•	•	•	2	4	•	•	•
Customer 7	•	•	5	•	•	•	•	3	•	•
Customer 8	•	•	•	•	•	2	•	•	•	3
Customer 9	3	•	•	•	5	•	•	5	•	•
Customer 10	•	•	•	•	•	•	•	•	•	•

Optimization

The general problem of matrix completion is

$$\hat{\mathbf{Z}} = \arg \min_{\mathbf{Z}} \|\mathbf{Z} - \mathbf{M}\|_F^2 \quad \text{subject to} \quad \Phi(\mathbf{Z}) \leq c. \quad (1)$$

Examples of $\Phi(\mathbf{Z})$:

- $\|\mathbf{Z}\|_1$: Sparse matrix
- $\text{rank}(\mathbf{Z})$: Low rank
- $\|\mathbf{Z}\|_*$: Nuclear norm
- $\mathbf{Z} = \mathbf{L} + \mathbf{S}$, \mathbf{S} sparse: Decomposition

Rank Minimization

Candes and Recht:

$$\begin{array}{ll} \underset{\mathbf{X}}{\text{minimize}} & \text{rank}(\mathbf{X}) \\ \text{subject to} & X_{ij} = M_{ij}, \quad (i, j) \in \Omega. \end{array} \quad (2)$$

- $\Omega = \{(i, j) \mid M_{ij} \text{ is available.}\}$
- Finds the matrix with the minimum rank
- NP-hard

Nuclear Norm Minimization

The nuclear norm minimization problem is defined as

$$\begin{aligned} & \underset{\mathbf{X}}{\text{minimize}} && \|\mathbf{X}\|_* \\ & \text{subject to} && X_{ij} = M_{ij}, \quad (i, j) \in \Omega. \end{aligned} \tag{3}$$

Comparison between Nuclear Norm and Frobenius Norm:

- $\|\mathbf{X}\|_* = \text{Tr} \left(\sqrt{\mathbf{X}^T \mathbf{X}} \right) = \sum_{i=1}^n \sigma_i(\mathbf{X})$
- $\|\mathbf{X}\|_F = \sqrt{\text{Tr}(\mathbf{X}^T \mathbf{X})} = \sqrt{\sum_{i=1}^n \sigma_i^2(\mathbf{X})}$
- $\|\mathbf{X}\|_*$ promotes sparsity
- Nuclear norm minimization is the tightest convex relaxation of the rank minimization.
- $\{\mathbf{X} \mid \|\mathbf{X}\|_* \leq 1\}$ is the convex hull of set of rank-one matrices with spectral norm bounded by 1.

Random Orthogonal Model

Candes and Recht 2008:

$$\mathbf{M} = \sum_{k=1}^r \sigma_k \mathbf{u}_k \mathbf{v}_k^T, \quad (4)$$

- $\{\mathbf{u}_k\}_{k=1}^r$ is selected *uniformly* at random among all families of orthonormal vectors.
- So is $\{\mathbf{v}_k\}$
- That is: \mathbf{M} has an intrinsic low rank

Main Theoretical Result

Candes and Recht 2008:

Theorem

Let $\mathbf{M} \in \mathbb{R}^{n_1 \times n_2}$ be a rank r matrix sampled from the random orthogonal model. Let $n = \max(n_1, n_2)$. Suppose we observe m entries of \mathbf{M} uniformly at random. Then, there are constants c and C such that if

$$m \geq Cn^{5/4}r \log n,$$

then the minimizer to the Nuclear norm minimization is unique and equal to \mathbf{M} with probability $1 - cn^{-3}$.

Implications:

- Under the hypothesis of the Theorem, there is a unique low-rank matrix which is consistent with the observed entries
- This matrix can be recovered by a convex optimization algorithm.
- So the nuclear norm relaxation is formally equivalent to the NP hard problem.

Solving the Nuclear Norm Problem

- Can be translated into the semi-definite programming
- Standard packages exist, e.g., CVX, SDPT3
- Typically solves the problem using interior point method
- Polynomial runtime

Projection Operator

Define a projection operator \mathcal{P}_Ω :

$$[\mathcal{P}_\Omega(\mathbf{X})]_{ij} = \begin{cases} X_{ij}, & (i,j) \in \Omega \\ 0, & (i,j) \notin \Omega \end{cases} \quad (5)$$

Then, the nuclear norm minimization can be written as

$$\begin{array}{ll} \underset{\mathbf{X}}{\text{minimize}} & \|\mathbf{X}\|_* \\ \text{subject to} & \mathcal{P}_\Omega(\mathbf{X}) = \mathcal{P}_\Omega(\mathbf{M}). \end{array} \quad (6)$$

We can relax the problem by

$$\begin{array}{ll} \underset{\mathbf{X}}{\text{minimize}} & \tau \|\mathbf{X}\|_* + \frac{1}{2} \|\mathbf{X}\|_F^2 \\ \text{subject to} & \mathcal{P}_\Omega(\mathbf{X}) = \mathcal{P}_\Omega(\mathbf{M}). \end{array} \quad (7)$$

As $\tau \rightarrow \infty$, the second problem becomes the first problem.

Solving the Problem

Recall:

$$\begin{aligned} & \underset{\mathbf{X}}{\text{minimize}} && \tau \|\mathbf{X}\|_* + \frac{1}{2} \|\mathbf{X}\|_F^2 \\ & \text{subject to} && \mathcal{P}_\Omega(\mathbf{X}) = \mathcal{P}_\Omega(\mathbf{M}). \end{aligned} \tag{8}$$

Define the Lagrangian:

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}) = \tau \|\mathbf{X}\|_* + \frac{1}{2} \|\mathbf{X}\|_F^2 + \langle \mathbf{Y}, \mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}) \rangle \tag{9}$$

Uzawa's Algorithm:

$$\begin{aligned} \mathbf{X}^k &= \arg \min_{\mathbf{X}} \mathcal{L}(\mathbf{X}, \mathbf{Y}^{k-1}) \\ \mathbf{Y}^k &= \mathbf{Y}^{k-1} + \delta_k \mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}^k) \end{aligned} \tag{10}$$

Subproblem Solution

$$\begin{aligned}\mathbf{X}^k &= \arg \min_{\mathbf{X}} \mathcal{L}(\mathbf{X}, \mathbf{Y}^{k-1}) \\ &= \arg \min_{\mathbf{X}} \tau \|\mathbf{X}\|_* + \frac{1}{2} \|\mathbf{X}\|_F^2 + \langle \mathbf{Y}^{k-1}, \mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}) \rangle \\ &= \arg \min_{\mathbf{X}} \tau \|\mathbf{X}\|_* + \frac{1}{2} \|\mathbf{X} - \mathcal{P}_\Omega \mathbf{Y}^{k-1}\|_F^2.\end{aligned}$$

Theorem (Theorem 2.1, Cai-Candes-Shen 2008)

For every $\tau \geq 0$ and $\mathbf{Y} \in \mathbb{R}^{n_1 \times n_2}$, it holds that

$$\arg \min_{\mathbf{X}} \tau \|\mathbf{X}\|_* + \frac{1}{2} \|\mathbf{X} - \mathcal{P}_\Omega \mathbf{Y}\|_F^2 = \mathcal{D}_\tau(\mathbf{Y}), \quad (11)$$

where $\mathcal{D}_\tau(\mathbf{Y})$ is the singular value shrinkage operator:

$$\mathcal{D}_\tau(\mathbf{Y}) = \mathbf{U} \mathcal{D}_\tau(\mathbf{\Sigma}) \mathbf{V}^T, \quad \mathcal{D}_\tau(\mathbf{\Sigma}) = \text{diag}((\sigma_i - \tau)_+). \quad (12)$$

The SVT Algorithm

We now have the Singular Value Thresholding algorithm

$$\begin{aligned}\mathbf{X}^k &= \mathcal{D}_\tau(\mathbf{Y}^{k-1}) \\ \mathbf{Y}^k &= \mathbf{Y}^{k-1} + \delta_k \mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}^k)\end{aligned}\tag{13}$$

- $\mathbf{X}^k = \mathcal{D}_\tau(\mathbf{Y}^{k-1})$ is similar to a shrinkage step
- τ is the shrinkage threshold
- $\mathbf{Y}^k = \mathbf{Y}^{k-1} + \delta_k \mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}^k)$ is the update

Choice of δ_k

“Heuristic Choice” by Cai-Candes-Shen:

$$\delta = 1.2 \frac{n_1 n_2}{m} \quad (14)$$

- $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$, $m = |\Omega|$.
- Incoherence assumption: If $\text{rank}(\mathbf{A})$ is not too large, then with high probability

$$(1 - \epsilon)p \|\mathbf{A}\|_F^2 \leq \|\mathcal{P}_\Omega(\mathbf{A})\|_F^2 \leq (1 + \epsilon)p \|\mathbf{A}\|_F^2, \quad p = m/(n_1 n_2)$$

- As part of the convergence proof, the algorithm will converge when

$$0 < \delta < 2 \frac{\|\mathbf{X}^* - \mathbf{X}^k\|_F^2}{\|\mathcal{P}_\Omega(\mathbf{X}^* - \mathbf{X}^k)\|_F^2}.$$

- Put $\mathbf{A} = \mathbf{X}^* - \mathbf{X}^k$, then if $\epsilon = 1/4$, $\delta \leq 1.6p^{-1}$. They take $1.2p^{-1}$.

Stopping Criteria

Another “heuristic choice”:

$$\frac{\|\mathcal{P}_\Omega(\mathbf{X}^k - \mathbf{M})\|_F}{\|\mathcal{P}_\Omega(\mathbf{M})\|_F} \leq \epsilon \quad (15)$$

- At optimal point, we must have $\mathcal{P}_\Omega(\mathbf{X} - \mathbf{M}) = 0$.
- Under suitable assumptions: $\|\mathcal{P}_\Omega(\mathbf{M})\|_F^2 \approx p \|\mathcal{P}_\Omega(\mathbf{M})\|_F^2$
- So roughly speaking: $\|\mathcal{P}_\Omega(\mathbf{X}^k - \mathbf{M})\|_F^2 \approx p \|\mathcal{P}_\Omega(\mathbf{X}^k - \mathbf{M})\|_F^2$
- Hence,

$$\frac{\|\mathcal{P}_\Omega(\mathbf{X}^k - \mathbf{M})\|_F}{\|\mathcal{P}_\Omega(\mathbf{M})\|_F} \approx \frac{\|\mathbf{X}^k - \mathbf{M}\|_F}{\|\mathbf{M}\|_F}.$$

Algorithm

Algorithm 1: Singular Value Thresholding (SVT) Algorithm

Input: sampled set Ω and sampled entries $\mathcal{P}_\Omega(\mathbf{M})$, step size δ , tolerance ϵ , parameter τ , increment ℓ , and maximum iteration count k_{\max}

Output: \mathbf{X}^{opt}

Description: Recover a low-rank matrix \mathbf{M} from a subset of sampled entries

```
1  Set  $\mathbf{Y}^0 = k_0 \delta \mathcal{P}_\Omega(\mathbf{M})$  ( $k_0$  is defined in (5.3))
2  Set  $r_0 = 0$ 
3  for  $k = 1$  to  $k_{\max}$ 
4      Set  $s_k = r_{k-1} + 1$ 
5      repeat
6          Compute  $[\mathbf{U}^{k-1}, \mathbf{\Sigma}^{k-1}, \mathbf{V}^{k-1}]_{s_k}$ 
7          Set  $s_k = s_k + \ell$ 
8      until  $\sigma_{s_k-\ell}^{k-1} \leq \tau$ 
9      Set  $r_k = \max\{j : \sigma_j^{k-1} > \tau\}$ 
10     Set  $\mathbf{X}^k = \sum_{j=1}^{r_k} (\sigma_j^{k-1} - \tau) \mathbf{u}_j^{k-1} \mathbf{v}_j^{k-1}$ 
11
12     if  $\|\mathcal{P}_\Omega(\mathbf{X}^k - \mathbf{M})\|_F / \|\mathcal{P}_\Omega \mathbf{M}\|_F \leq \epsilon$  then break
13
14     Set  $\mathbf{Y}_{ij}^k = \begin{cases} 0 & \text{if } (i, j) \notin \Omega, \\ \mathbf{Y}_{ij}^{k-1} + \delta(M_{ij} - \mathbf{X}_{ij}^k) & \text{if } (i, j) \in \Omega \end{cases}$ 
15 end for  $k$ 
16 Set  $\mathbf{X}^{\text{opt}} = \mathbf{X}^k$ 
```


Results

Unknown M				Computational results		
size ($n \times n$)	rank (r)	m/d_r	m/n^2	time(s)	# iters	relative error
$1,000 \times 1,000$	10	6	0.12	23	117	1.64×10^{-4}
	50	4	0.39	196	114	1.59×10^{-4}
	100	3	0.57	501	129	1.68×10^{-4}
$5,000 \times 5,000$	10	6	0.024	147	123	1.73×10^{-4}
	50	5	0.10	950	108	1.61×10^{-4}
	100	4	0.158	3,339	123	1.72×10^{-4}
$10,000 \times 10,000$	10	6	0.012	281	123	1.73×10^{-4}
	50	5	0.050	2,096	110	1.65×10^{-4}
	100	4	0.080	7,059	127	1.79×10^{-4}
$20,000 \times 20,000$	10	6	0.006	588	124	1.73×10^{-4}
	50	5	0.025	4,581	111	1.66×10^{-4}
$30,000 \times 30,000$	10	6	0.004	1,030	125	1.73×10^{-4}

Table 1: Experimental results for matrix completion. The rank r is the rank of the unknown matrix M , m/d_r is the ratio between the number of sampled entries and the number of degrees of freedom in an $n \times n$ matrix of rank r (oversampling ratio), and m/n^2 is the fraction of observed entries. All the computational results on the right are averaged over five runs.

Extension

Robust PCA:

$$(\hat{\mathbf{A}}, \hat{\mathbf{E}}) = \arg \min_{\mathbf{A}, \mathbf{E}} \|\mathbf{A}\|_* + \lambda \|\mathbf{E}\|_1, \quad \mathbf{A} + \mathbf{E} = \mathbf{D} \quad (16)$$

- \mathbf{D} is data matrix
- \mathbf{A} is low rank: background video frames
- \mathbf{E} is sparse: foreground objects

To solve the problem, translate into

$$(\hat{\mathbf{A}}, \hat{\mathbf{E}}) = \arg \min_{\mathbf{A}, \mathbf{E}} \mu \|\mathbf{A}\|_* + \lambda \|\mathbf{E}\|_1 + \frac{1}{2} \|\mathbf{D} - \mathbf{A} - \mathbf{E}\|_F^2 \quad (17)$$

Algorithm

Algorithm 1: Robust PCA via Proximal Gradient with Continuation

- 1: **Input:** Observation matrix $D \in \mathbb{R}^{m \times n}$, weight λ .
 - 2: $A_0, A_{-1} \leftarrow 0, E_0, E_{-1} \leftarrow 0, t_0, t_{-1} \leftarrow 1, \mu_0 \leftarrow .99\|D\|_{2,2}, \bar{\mu} \leftarrow 10^{-5}\mu_0$.
 - 3: **while** not converged
 - 4: $\tilde{A}_k \leftarrow A_k + \frac{t_{k-1}-1}{t_k} (A_k - A_{k-1}), \tilde{E}_k \leftarrow E_k + \frac{t_{k-1}-1}{t_k} (E_k - E_{k-1})$.
 - 5: $Y_k^A \leftarrow \tilde{A}_k - \frac{1}{2} \left(\tilde{A}_k + \tilde{E}_k - D \right)$.
 - 6: $(U, S, V) \leftarrow \text{svd}(Y_k^A), A_{k+1} \leftarrow U \left[S - \frac{\mu}{2} \mathbf{I} \right]_+ V^*$.
 - 7: $Y_k^E \leftarrow \tilde{E}_k - \frac{1}{2} \left(\tilde{A}_k + \tilde{E}_k - D \right)$.
 - 8: $E_{k+1} \leftarrow \text{sign}[Y_k^E] \circ \left[|Y_k^E| - \frac{\lambda\mu}{2} \mathbf{1} \right]_+$.
 - 9: $t_{k+1} \leftarrow \frac{1+\sqrt{1+4t_k^2}}{2}, \mu \leftarrow \max(.9\mu, \bar{\mu})$.
 - 10: **end while**
 - 11: **Output:** A, E .
-

Example

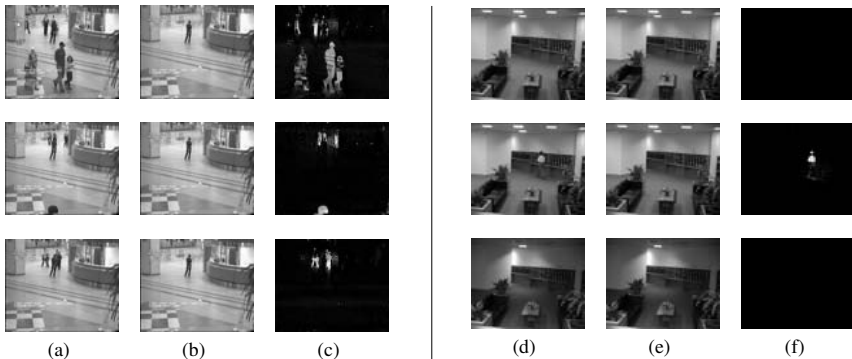


Figure 2: **Background modeling.** (a) Video sequence of a scene in an airport. The size of each frame is 72×88 pixels, and a total of 200 frames were used. (b) Static background recovered by our algorithm. (c) Sparse error recovered by our algorithm represents activity in the frame. (d) Video sequence of a lobby scene with changing illumination. The size of each frame is 64×80 pixels, and a total of 550 frames were used. (e) Static background recovered by our algorithm. (f) Sparse error. The background is correctly recovered even when the illumination in the room changes drastically in the frame on the last row.

Reference

- Jian-Feng Cai, Emmanuel J. Candes, Zuowei Shen, A Singular Value Thresholding Algorithm for Matrix Completion, arXiv:0810.3286
- Emmanuel J. Candes, Benjamin Recht, Exact Matrix Completion via Convex Optimization, arXiv:0805.4471
- John Wright, Yigang Peng, Yi Ma, Arvind Ganesh, and Shankar Rao, Robust Principal Component Analysis: Exact Recovery of Corrupted Low-Rank Matrices by Convex Optimization, arXiv:0905.0233