

Low-rank Matrix Completion from Noisy Entries

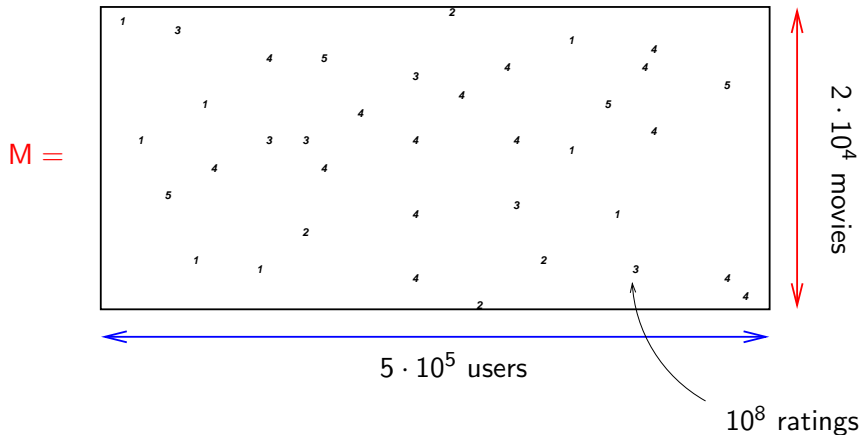
Sewoong Oh

Joint work with Raghunandan Keshavan and Andrea Montanari
Stanford University

Forty-Seventh Allerton Conference
October 1, 2009

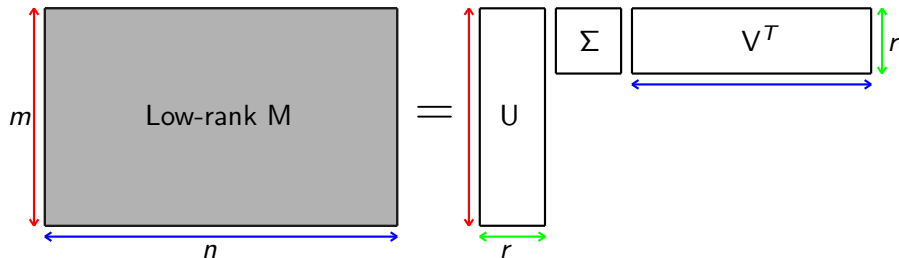
Motivating Example

- Netflix Challenge



The Model

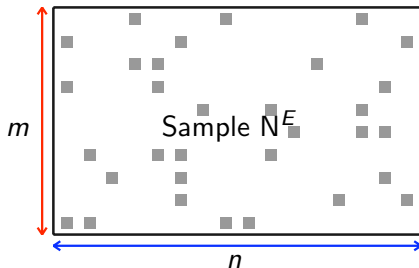
Matrix Completion Problem



1. Low-rank matrix M
2. $N = M + Z$
3. Uniformly random sample E

$$N_{ij}^E = \begin{cases} M_{ij} + Z_{ij} & \text{if } (i,j) \in E, \\ 0 & \text{otherwise.} \end{cases}$$

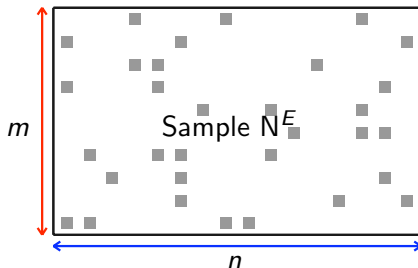
Matrix Completion Problem



1. Low-rank matrix M
2. $N = M + Z$
3. Uniformly random sample E

$$N_{ij}^E = \begin{cases} M_{ij} + Z_{ij} & \text{if } (i,j) \in E, \\ 0 & \text{otherwise.} \end{cases}$$

Matrix Completion Problem



Goal : Estimation $\hat{M}(E, N^E)$ that minimizes

$$\text{RMSE} \equiv \left(\frac{1}{mn} \sum_{i,j} (M_{ij} - \hat{M}_{ij})^2 \right)^{1/2}.$$

Matrix Completion Problem

- Characterized by 4 parameters:
 - ▶ Data size $m \times n$, assume fixed $\alpha \equiv m/n$
 - ▶ Rank r
 - ▶ Sample size $|E|$
 - ▶ Noise Z^E
{ Running example : $Z_{ij} \sim \text{i.i.d. } N(0, \sigma_z^2)$ }

Q. Is there an efficient algorithm with performance guarantee?

$$\text{RMSE}_{\text{ALG}} \leq F(n, r, |E|, Z^E)$$

Matrix Completion Problem

- Characterized by 4 parameters:
 - ▶ Data size $m \times n$, assume fixed $\alpha \equiv m/n$
 - ▶ Rank r
 - ▶ Sample size $|E|$
 - ▶ Noise Z^E
{ Running example : $Z_{ij} \sim \text{i.i.d. } N(0, \sigma_z^2)$ }

Q. Is there an efficient algorithm with performance guarantee?

$$\text{RMSE}_{\text{ALG}} \leq F(n, r, |E|, Z^E)$$

Matrix Completion Problem

- Characterized by 4 parameters:
 - ▶ Data size $m \times n$, assume fixed $\alpha \equiv m/n$
 - ▶ Rank r
 - ▶ Sample size $|E|$
 - ▶ Noise Z^E
{ Running example : $Z_{ij} \sim \text{i.i.d. } N(0, \sigma_z^2)$ }

Q. Is there an efficient algorithm with performance guarantee?

$$\text{RMSE}_{\text{OptSpace}} \leq C \frac{n\sqrt{r}}{|E|} \|Z^E\|_2$$

$\left(C \sigma_z \sqrt{rn \log n / |E|}, \text{ for Gaussian} \right)$

Matrix Completion Problem

- Characterized by 4 parameters:
 - ▶ Data size $m \times n$, assume fixed $\alpha \equiv m/n$
 - ▶ Rank r
 - ▶ Sample size $|E|$
 - ▶ Noise Z^E
- { Running example : $Z_{ij} \sim \text{i.i.d. } N(0, \sigma_z^2)$ }

Q. Is there an efficient algorithm with performance guarantee?

$$\begin{aligned} \text{RMSE}_{\text{OptSpace}} &\leq C \frac{n\sqrt{r}}{|E|} \|Z^E\|_2 \\ &\quad \left(C \sigma_z \sqrt{rn \log n / |E|}, \text{ for Gaussian} \right) \\ \text{RMSE}_{\text{Oracle}} &\simeq \frac{1}{\sqrt{|E|}} \|Z^E\|_F \\ &\quad \left(\sigma_z \sqrt{2rn / |E|}, \text{ for Gaussian} \right) \end{aligned}$$

Matrix Completion Problem

- Characterized by 4 parameters:
 - ▶ Data size $m \times n$, assume fixed $\alpha \equiv m/n$
 - ▶ Rank r
 - ▶ Sample size $|E|$
 - ▶ Noise Z^E
{ Running example : $Z_{ij} \sim \text{i.i.d. } N(0, \sigma_z^2)$ }

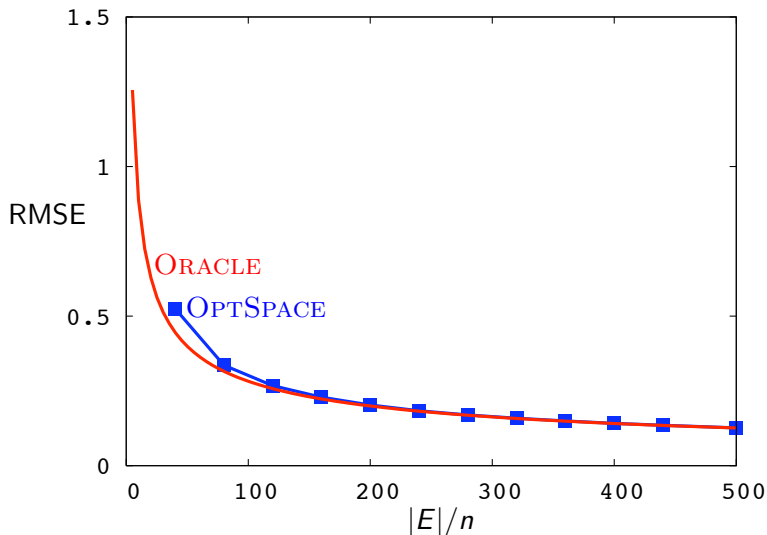
Q. Is there an efficient algorithm with performance guarantee?

$$\begin{aligned} \text{RMSE}_{\text{OptSpace}} &\leq C \frac{n\sqrt{r}}{|E|} \|Z^E\|_2 \\ &\quad \left(C \sigma_z \sqrt{rn \log n / |E|}, \text{ for Gaussian} \right) \\ \text{RMSE}_{\text{Oracle}} &\simeq \frac{1}{\sqrt{|E|}} \|Z^E\|_F \\ &\quad \left(\sigma_z \sqrt{2rn / |E|}, \text{ for Gaussian} \right) \end{aligned}$$

OPTSPACE is Near-optimal

Numerical Simulation Results

- Fixed $n = 500, r = 4, \sigma_z = 1$



Main Contribution

OPTSPACE

1. Complexity: Low complexity
2. Theory: Near-optimal performance guarantee
3. Practice: Numerical simulations

The Algorithm and Main Theorems

Naïve Approach

$$N^E = \sum_{k=1}^n x_k \sigma_k y_k^T$$

Rank- r projection :

$$\mathcal{P}_r(N^E) \equiv \frac{mn}{|E|} \sum_{k=1}^r x_k \sigma_k y_k^T$$

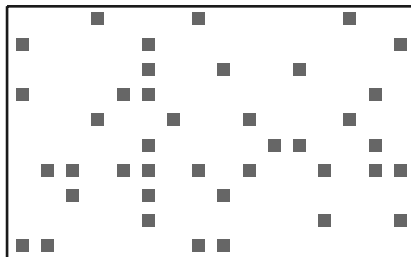
Naïve Approach Fails

- Define : $\deg(\text{row}_i) \equiv \#$ of samples in row i .
- For $|E| = O(n)$, there exists a row with degree $\Omega(\log n / (\log \log n))$.
- *spurious* singular values of $\Omega(\sqrt{\log n / (\log \log n)})$.

Trimming

- Solution : Trimming

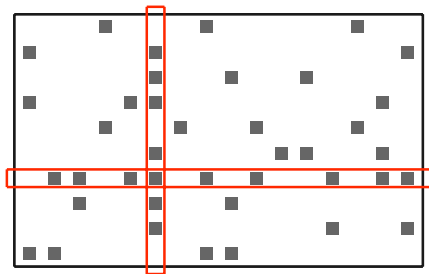
$$\tilde{N}_{ij}^E = \begin{cases} 0 & \text{if } \deg(\text{row}_i) > 2\mathbb{E}[\deg(\text{row}_i)] , \\ 0 & \text{if } \deg(\text{col}_j) > 2\mathbb{E}[\deg(\text{col}_j)] , \\ N_{ij}^E & \text{otherwise.} \end{cases}$$



Trimming

- Solution : Trimming

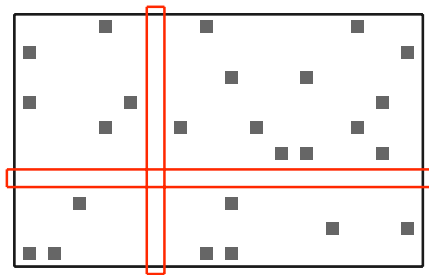
$$\tilde{N}_{ij}^E = \begin{cases} 0 & \text{if } \deg(\text{row}_i) > 2\mathbb{E}[\deg(\text{row}_i)] , \\ 0 & \text{if } \deg(\text{col}_j) > 2\mathbb{E}[\deg(\text{col}_j)] , \\ N_{ij}^E & \text{otherwise.} \end{cases}$$



Trimming

- Solution : Trimming

$$\tilde{N}_{ij}^E = \begin{cases} 0 & \text{if } \deg(\text{row}_i) > 2\mathbb{E}[\deg(\text{row}_i)] , \\ 0 & \text{if } \deg(\text{col}_j) > 2\mathbb{E}[\deg(\text{col}_j)] , \\ N_{ij}^E & \text{otherwise.} \end{cases}$$



The Algorithm

OPTSPACE

Input : sample positions E , sample values N^E , rank r

Output : estimation \hat{M}

- 1: Trim N^E , and let \tilde{N}^E be the output;
 - 2: Compute rank- r projection $\mathcal{P}_r(\tilde{N}^E) = X_0 S_0 Y_0^T$;
 - 3:
-

Main Result

Theorem (Keshavan, Montanari, Oh, 2009 Thm. 1.1)

Let M be an $n \times n$ matrix of rank- r bounded by M_{\max} . Then, w.h.p., rank- r projection achieves

$$\text{RMSE} \leq CM_{\max} \sqrt{\frac{nr}{|E|}} + C' \frac{n\sqrt{r}}{|E|} \|Z^E\|_2.$$

$$\left(\text{Example: } CM_{\max} \sqrt{\frac{nr}{|E|}} + C' \sigma_z \sqrt{\frac{rn \log n}{|E|}} \right)$$

The Algorithm

OPTSPACE

Input : sample positions E , sample values N^E , rank r

Output : estimation \hat{M}

- 1: Trim N^E , and let \tilde{N}^E be the output;
 - 2: Compute rank- r projection $\mathcal{P}_r(\tilde{N}^E) = X_0 S_0 Y_0^T$;
 - 3: Minimize RMSE by gradient descent starting at (X_0, S_0, Y_0) .
-

Main Result

Theorem (Keshavan, Montanari, Oh, 2009 Thm. 1.1)

Let M be an $n \times n$ matrix of rank- r bounded by M_{\max} . Then, w.h.p., rank- r projection achieves

$$\text{RMSE} \leq C M_{\max} \sqrt{nr/|E|} + C' \|Z^E\|_2 n\sqrt{r}/|E|.$$

Theorem (Keshavan, Montanari, Oh, 2009 Thm. 1.2)

Let M be an $n \times n$ rank- r *incoherent* matrix with $\sigma_1(M)/\sigma_r(M) = O(1)$. If $|E| \geq C n r \max\{r, \log n\}$, then, w.h.p., OPTSPACE achieves

$$\text{RMSE} \leq C'' \frac{n\sqrt{r}}{|E|} \|Z^E\|_2,$$

provided that the RHS is smaller than $\sigma_r(M)$.

$$\left(\text{Example: } C'' \sigma_z \sqrt{r n \log n / |E|} \right)$$

Comparison: Theory

Theorem (Candés, Plan, 2009)

Assume *strongly incoherent* matrix M . If $|E| \geq C r n (\log n)^6$ then SEMIDEFINITE PROGRAMMING achieves, w.h.p.,

$$\text{RMSE} \leq C' \sqrt{\frac{n}{|E|}} \|Z^E\|_F + C'' \frac{1}{n} \|Z^E\|_F.$$

$$\left(\text{Example: } C' \sigma_z \sqrt{n} + C'' \sigma_z \frac{\sqrt{|E|}}{n} \right)$$

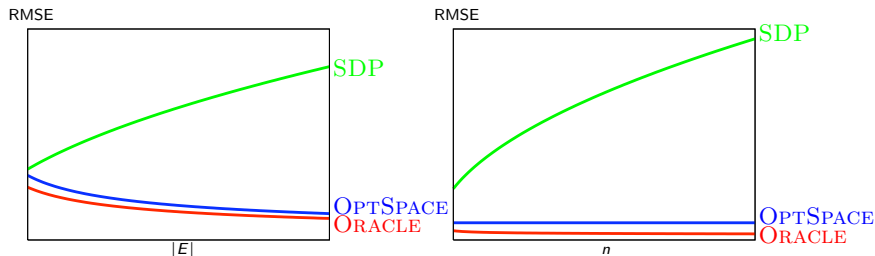
Comparison: Theory

When Z is i.i.d $N(0, \sigma_z^2)$,

$$\text{ORACLE: RMSE} \simeq C\sigma_z \sqrt{\frac{rn}{|E|}}$$

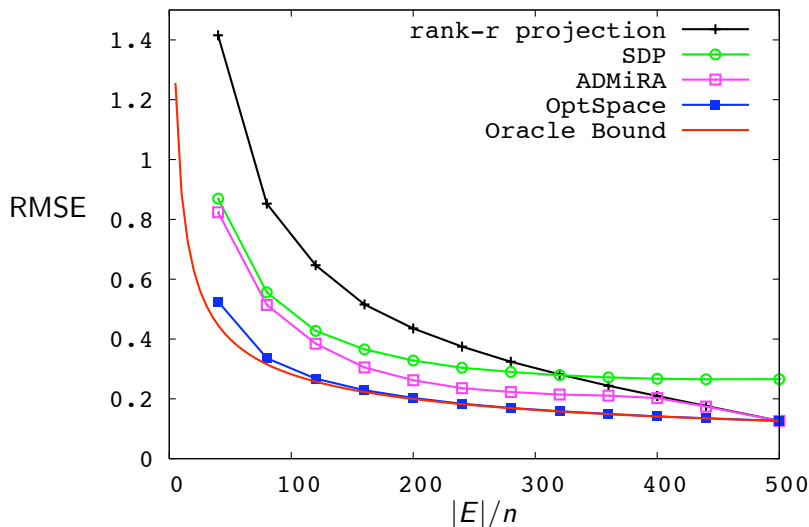
$$\text{OPTSPACE: RMSE} \leq C'\sigma_z \sqrt{\frac{rn \log n}{|E|}}$$

$$\text{SDP: RMSE} \leq C''\sigma_z \left\{ \sqrt{n} + \frac{\sqrt{|E|}}{n} \right\}$$



Comparison

- Fixed $n = 500, r = 4, \sigma_z = 1$, example from [Candés, Plan, 2009]



Conclusion

OPTSPACE

1. Complexity: Low complexity
2. Theory: Near-optimal performance guarantee
3. Practice: Numerical simulations

Thank you!

Conclusion

OPTSPACE

1. Complexity: Low complexity
2. Theory: Near-optimal performance guarantee
3. Practice: Numerical simulations

Thank you!