Testing Matrix Rank, Optimally

Hongyang Zhang, CMU

Nina Balcan (CMU)



Yi Li (NTU)



David P. Woodruff (CMU)



Era of Big Data



2

Era of Big Data



*Zettabyte = 35, 000,000,000,000,000,000,000 bytes

3

Netflix Challenge (low-rank data)









Action

	AND DOWN.		18. C 16	
omedy	Historical		Cartoon	
	8	8	1	1
	8	8	1	1
	4	4	2	2
?	4	4	2	2
	1	1	4	4
	1	1	4	4
	0	0	8	8
	0	0	8	8

The user-movie matrix is of low rank *r*.

Need to know *r* in many applications, e.g., in matrix completion.



Question:

• What is the true rank given a few observations? (Testing of Rank)

Outline

Motivation and examples

Our settings

- Our algorithms
- Our hardness results
- Conclusions

Property Testing of Rank

Goal: Design non-adaptive sampling scheme to distinguish H_0 from H_1 with sample complexity independent of size of matrix *n*, or even optimally



Previous Work



n

Algorithm [KS'03]:

- Randomly permute the row and column indices of the matrix
- Sample top $O(\frac{d}{\epsilon}) \times O(\frac{d}{\epsilon})$ submatrix X
- If $rank(X) \le d$, output " H_0 "; otherwise, output " H_1 "

Sample complexity: $O(\frac{d^2}{\epsilon^2})$

Problem of the algorithm ($d=1, A \in H_1$)

Outline

- Motivation and examples
- Our settings
- Our algorithms
- Our hardness results
- Conclusions



Sample complexity: $\tilde{O}(\frac{d^2}{\epsilon})$

n

Why does the Algorithm Work?

Proof Idea: Starting from an empty-by-empty matrix, augment the matrix until finding a $(d + 1) \times (d + 1)$ full-rank matrix



•
$$r:=\min_{A_{[d/\varepsilon]\times[d/\varepsilon]\setminus Q}} rank(A_{[d/\varepsilon]\times[d/\varepsilon]})$$

• If $r \le d$, output " H_0 "; otherwise, output " H_1 "

When
$$A \in H_0$$
, $r \leq d$

Only need to show when $A \in H_1$, r > d

Oracle Lemma:

For any *t*, we can augment any $t \times t$ full-rank matrix to a $(t + 1) \times (t + 1)$ full-rank matrix by an augmentation entry in the sampled region, as long as $A \in H_1$ and $t \leq d$.



d

ε

Challenge: Though the oracle lemma guarantees we find a $(t + 1) \times (t + 1)$ fullrank matrix from hindsight, the algorithm may not know it is full-rank due to unobservations. (Let's say d=1)

$$\begin{bmatrix} A_{r,c'} & A_{r,c} \\ A_{r',c'} & ? \end{bmatrix} = \begin{bmatrix} 0 & \neq 0 \\ A_{r',c'} & ? \end{bmatrix}$$

Case I: $A_{r,c'} = 0$

 d/ε

Oracle Lemma:

For any *t*, we can augment any $t \times t$ full-rank matrix to a $(t + 1) \times (t + 1)$ full-rank matrix by an augmentation entry in the sampled region, as long as $A \in H_1$ and $t \leq d$.



 $\frac{d}{\varepsilon}$

Challenge: Though the oracle lemma guarantees we find a $(t + 1) \times (t + 1)$ fullrank matrix from hindsight, the algorithm may not know it is full-rank due to unobservations.

$$\begin{bmatrix} A_{r,c'} & A_{r,c} \\ A_{r',c'} & ? \end{bmatrix} = \begin{bmatrix} 0 & \neq 0 \\ \neq 0 & ? \end{bmatrix}$$

Case I: $A_{r,c'} = 0$

Oracle Lemma:

For any *t*, we can augment any $t \times t$ full-rank matrix to a $(t + 1) \times (t + 1)$ full-rank matrix by an augmentation entry in the sampled region, as long as $A \in H_1$ and $t \leq d$.



 $\frac{d}{\varepsilon}$

Challenge: Though the oracle lemma guarantees we find a $(t + 1) \times (t + 1)$ fullrank matrix from hindsight, the algorithm may not know it is full-rank due to unobservations.

$$\begin{bmatrix} A_{r,c'} & A_{r,c} \\ A_{r',c'} & ? \end{bmatrix} = \begin{bmatrix} \neq 0 & A_{r,c} \\ A_{r',c'} & ? \end{bmatrix}$$

Case II: $A_{r,c'} \neq 0$

 d/ε

Oracle Lemma:

For any *t*, we can augment any $t \times t$ full-rank matrix to a $(t + 1) \times (t + 1)$ full-rank matrix by an augmentation entry in the sampled region, as long as $A \in H_1$ and $t \leq d$.





 $\frac{d}{\varepsilon}$

Challenge: Though the oracle lemma guarantees we find a $(t + 1) \times (t + 1)$ fullrank matrix from hindsight, the algorithm may not know it is full-rank due to unobservations.

$$\begin{bmatrix} A_{r^0,c'} & A_{r^0,c^0} \\ A_{r,c'} & ? \end{bmatrix} = \begin{bmatrix} 0 & A_{r^0,c^0} \\ \neq 0 & ? \end{bmatrix}$$

Case I: $A_{r^0,c'} = 0$

 d/ε

Oracle Lemma:

For any *t*, we can augment any $t \times t$ full-rank matrix to a $(t + 1) \times (t + 1)$ full-rank matrix by an augmentation entry in the sampled region, as long as $A \in H_1$ and $t \leq d$.

.

 $\frac{d}{\varepsilon}$



Challenge: Though the oracle lemma guarantees we find a $(t + 1) \times (t + 1)$ fullrank matrix from hindsight, the algorithm may not know it is full-rank due to unobservations.

$$\begin{bmatrix} A_{r^0,c'} & A_{r^0,c^0} \\ A_{r,c'} & ? \end{bmatrix} = \begin{bmatrix} 0 & \neq 0 \\ \neq 0 & ? \end{bmatrix}$$

Case I: $A_{r^0,c'} = 0$

 d/ε

Oracle Lemma:

For any *t*, we can augment any $t \times t$ full-rank matrix to a $(t + 1) \times (t + 1)$ full-rank matrix by an augmentation entry in the sampled region, as long as $A \in H_1$ and $t \leq d$.



 $\frac{d}{\varepsilon}$

Challenge: Though the oracle lemma guarantees we find a $(t + 1) \times (t + 1)$ fullrank matrix from hindsight, the algorithm may not know it is full-rank due to unobservations.

$$\begin{bmatrix} A_{r^{0},c'} & A_{r^{0},c^{0}} \\ A_{r,c'} & ? \end{bmatrix} = \begin{bmatrix} \neq 0 & A_{r^{0},c^{0}} \\ A_{r,c'} & ? \end{bmatrix}$$

Case II: $A_{r^0,c'} \neq 0$

 d/ε

Oracle Lemma:

For any *t*, we can augment any $t \times t$ full-rank matrix to a $(t + 1) \times (t + 1)$ full-rank matrix by an augmentation entry in the sampled region, as long as $A \in H_1$ and $t \leq d$.



d

Е

Challenge: Though the oracle lemma guarantees we find a $(t + 1) \times (t + 1)$ fullrank matrix from hindsight, the algorithm may not know it is full-rank due to unobservations.

$$\begin{bmatrix} A_{r^0,c'} & A_{r^0,c} \\ A_{r,c'} & A_{r,c} \end{bmatrix}$$

fully observed

 d/ε

A Polynomial-Time Algorithm

✓ × ✓



Our algorithm:

- Randomly permute the row and column indices of the matrix
- Sample the block region Q
- $r := \min_{A_{[d/\varepsilon] \times [d/\varepsilon] \setminus Q}} rank(A_{[d/\varepsilon] \times [d/\varepsilon]})$
- If $r \le d$, output " H_0 "; otherwise, output " H_1 "

Other Extensions

Theorem (Stable Rank Upper Bounds):

There is an algorithm which tests the stable rank by a constant success probability with $\tilde{O}(d^3/\epsilon^4)$ samples.



$$\operatorname{srank}(A) = \frac{\|A\|_F^2}{\sigma_1^2(A)}$$

stable version of matrix rank

Outline

- Motivation and examples
- Our settings
- Our algorithms
- Our hardness results
- Conclusions

Hardness ---- Testing Matrix Rank

Our positive result: $\tilde{O}(\frac{d^2}{s})$ samples

Theorem (Lower Bounds):

Any non-adaptive algorithm with constant success probability requires at least $\tilde{\Omega}(d^2/\varepsilon)$ samples over reals and finite fields.

 $U, V \sim \mathcal{G}(n, d) \qquad G \sim \mathcal{G}(n, n)$

 d_{TV} distance is a constant with small samples

 UV^T







Outline

- Motivation and examples
- Our settings
- Our algorithms
- Our hardness results
- Conclusions

Conclusions

- Property testing of rank
 - Polynomial-time algorithm
 - Rebasing argument
 - Sample complexity $\tilde{O}(d^2/\varepsilon)$
 - Other extensions
- Hardness results
 - Sample complexity $\widetilde{\Omega}(d^2/\varepsilon)$



Thank You

Structure

