

Principal Component Analysis for Distributed Data

David Woodruff
IBM Almaden

Based on works with Ken Clarkson, Ravi Kannan, and Santosh Vempala

Outline

1. What is low rank approximation?
2. How do we solve it offline?
3. How do we solve it in a distributed setting?

Low rank approximation

- A is an $n \times d$ matrix
 - Think of n points in \mathbb{R}^d
- E.g., A is a customer-product matrix
 - $A_{i,j}$ = how many times customer i purchased item j
- A is typically well-approximated by low rank matrix
 - E.g., high rank because of noise
- **Goal:** find a low rank matrix approximating A
 - Easy to store, data more interpretable

What is a good low rank approximation?

Singular Value

Any matrix A

$$A_k = \operatorname{argmin}_{\text{rank } k \text{ matrices } B} \|A - B\|_F$$

- U

The rows of V_k are
the top k **principal
components**

- R

$$\begin{pmatrix} \mathbf{A} \end{pmatrix} = \begin{pmatrix} \mathbf{U}_k \end{pmatrix} \begin{pmatrix} \Sigma_k \end{pmatrix} \begin{pmatrix} \mathbf{V}_k \end{pmatrix} + \begin{pmatrix} \mathbf{E} \end{pmatrix}$$

Low rank approximation

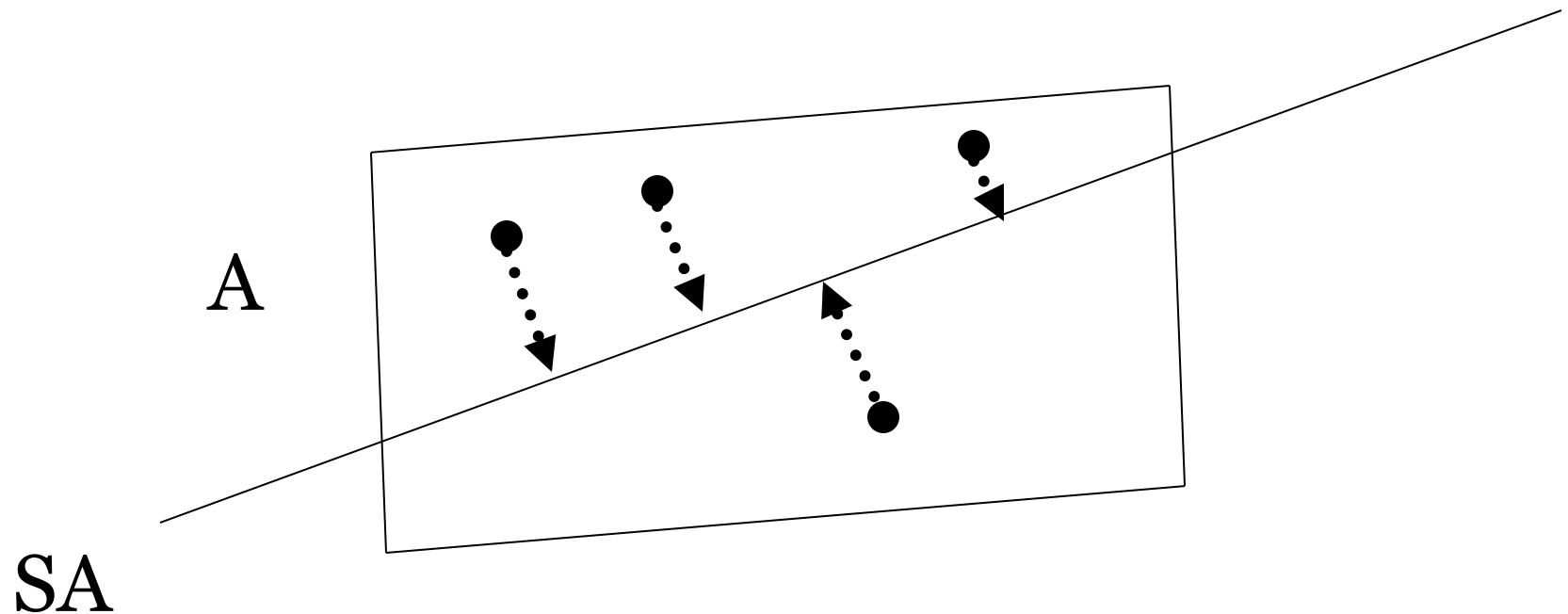
- **Goal:** output a rank k matrix A' , so that

$$\|A - A'\|_F \leq (1 + \varepsilon) \|A - A_k\|_F$$

- Can do this in $\text{nnz}(A) + (n+d) \cdot \text{poly}(k/\varepsilon)$ time [S,CW]
 - $\text{nnz}(A)$ is number of non-zero entries of A

Solution to low-rank approximation [S]

- Given $n \times d$ input matrix A
- Compute S^*A using a sketching matrix S with $k/\epsilon \ll n$ rows. S^*A takes random linear combinations of rows of A



- Project rows of A onto SA , then find best rank- k approximation to points inside of SA .

What is the matrix S ?

- S can be a $k/\epsilon \times n$ matrix of i.i.d. normal random variables
- [S] S can be a $k/\epsilon \times n$ Fast Johnson Lindenstrauss Matrix
 - Uses Fast Fourier Transform
- [CW] S can be a $\text{poly}(k/\epsilon) \times n$ CountSketch matrix

$$\begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 & 0 & -1 & 0 \\ 0 & -1 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$S \odot A$ can be computed in $\text{nnz}(A)$ time!

Caveat: projecting the points onto SA is slow

- Current algorithm:
 1. Compute S^*A
 2. Project each of the rows onto S^*A
 3. Find best rank- k approximation of projected points inside of rowspace of S^*A
- Bottleneck is step 2
- [CW] Approximate the projection
 - Fast algorithm for approximate regression
$$\min_{\text{rank-}k \ X} \|X(SA) - A\|_F^2$$
 - $\text{nnz}(A) + (n+d) \cdot \text{poly}(k/\epsilon)$ time

Distributed low rank approximation

- *We have fast algorithms, but can they be made to work in a distributed setting?*
- Matrix A distributed among s servers
- For $t = 1, \dots, s$, we get a customer-product matrix from the t -th shop stored in server t . Server t 's matrix = A^t
- Customer-product matrix $A = A^1 + A^2 + \dots + A^s$
- More general than row-partition model in which each customer shops in only one shop

Communication cost of low rank approximation

- **Input:** $n \times d$ matrix A stored on s servers
 - Server t has $n \times d$ matrix A^t
 - $A = A^1 + A^2 + \dots + A^s$
- **Output:** Server t has $n \times d$ matrix C^t satisfying
 - $C = C^1 + C^2 + \dots + C^s$ has rank at most k
 - $|A-C|_F \cdot (1+\varepsilon) |A-A_k|_F$
 - Application: distributed clustering
- **Resources:** Each server is polynomial time, linear space, communication is $O(1)$ rounds. Bound the total number of words communicated
- [KVW]: $O(skd/\varepsilon)$ communication, independent of n

Protocol

- Designate one machine the Central Processor (CP)
- Let L_t be the local data matrix on server t
- **Problems:**
 - Can't output $A^t U U^T$ since rank too large
 - Could communicate $A^t U$ to CP, then CP computes SVD of $\sum_t A^t U U^T = A U U^T$
 - But communicating $A^t U$ depends on n
- CP sends U to each server
- Server t computes $A^t U$

Approximate SVD lemma

- Problem reduces to
 - Server t has $n \times r$ matrix B^t
 - $B = \sum_t B^t$
 - CP outputs top k principal components

Communication independent of n !

■ Approximate SVD

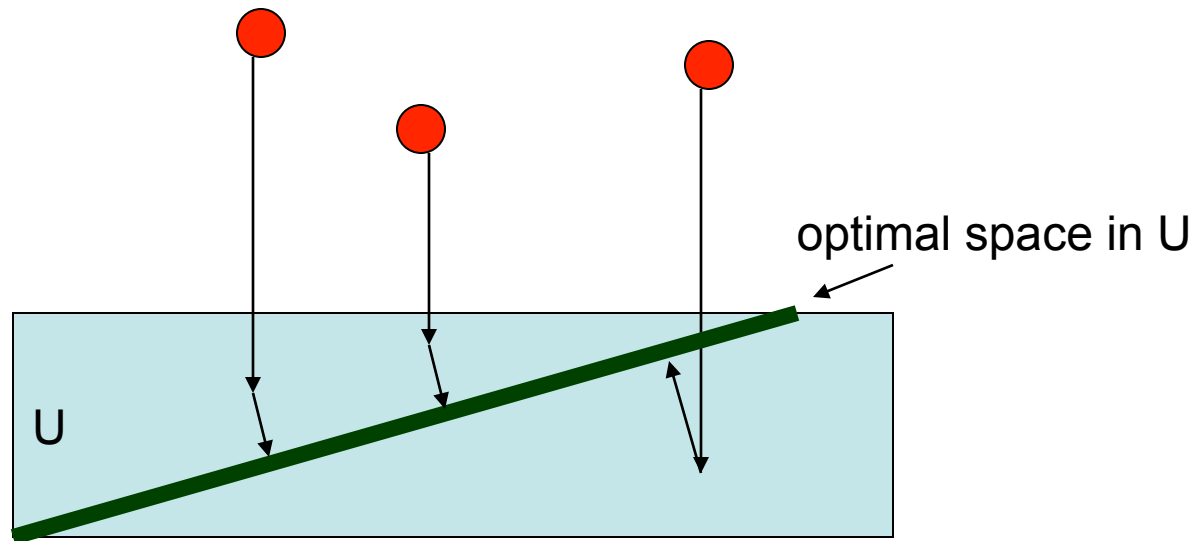
- If $W^T \in \mathbb{R}^{k \times r}$ is the matrix of top k principal components of PB , where P is a random $r/\epsilon^2 \times n$ matrix,

$$\|B - BW W^T\|_F \leq (1 + \epsilon) \|B - B_k\|_F$$

- CP sends P to every server
- Server t sends PB^t to CP who computes $PB = \sum_t PB^t$
- CP computes W , sends everyone W

The protocol

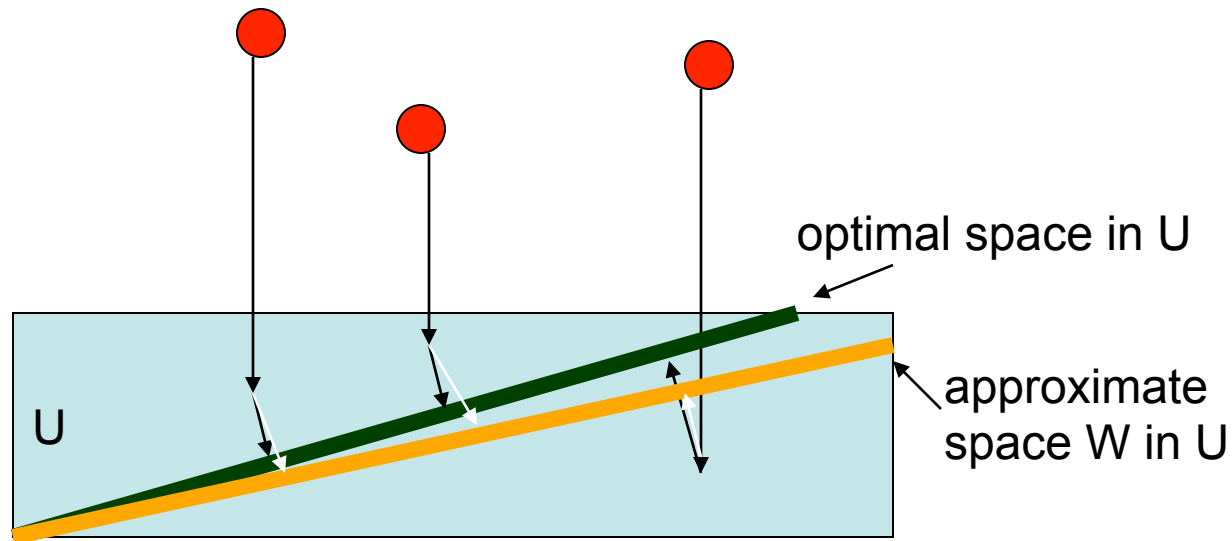
- Phase 1:
- Learn an orthonormal basis U for row space of SA



$$\text{cost} \cdot (1+\epsilon) |A - A_k|_F$$

The protocol

- Phase 2:
- Find an approximately optimal space W inside of U



$$\text{cost} \cdot (1+\epsilon)^2 |A - A_k|_F$$

Conclusion

- $O(\text{sdk}/\varepsilon)$ communication protocol for low rank approximation
- A bit sloppy with words vs. bits but can be dealt with
- Almost matching $\Omega(\text{sdk})$ bit lower bound
 - Can be strengthened to $\Omega(\text{sdk}/\varepsilon)$ in one-way model
 - Can we remove the one-way restriction?
- Communication cost of other optimization problems?
 - Linear programming
 - Frequency moments
 - Matching
 - etc.