



# Improved Distributed Principal Component Analysis

Maria-Florina Balcan<sup>†</sup>, Vandana Kanchanapally<sup>‡</sup>, Yingyu Liang<sup>◇</sup>, and David Woodruff<sup>\*</sup>

<sup>†</sup>ninamf@cs.cmu.edu, <sup>‡</sup>vvandana@gatech.edu, <sup>◇</sup>yingyul@cs.princeton.edu, <sup>\*</sup>dpwoodru@us.ibm.com

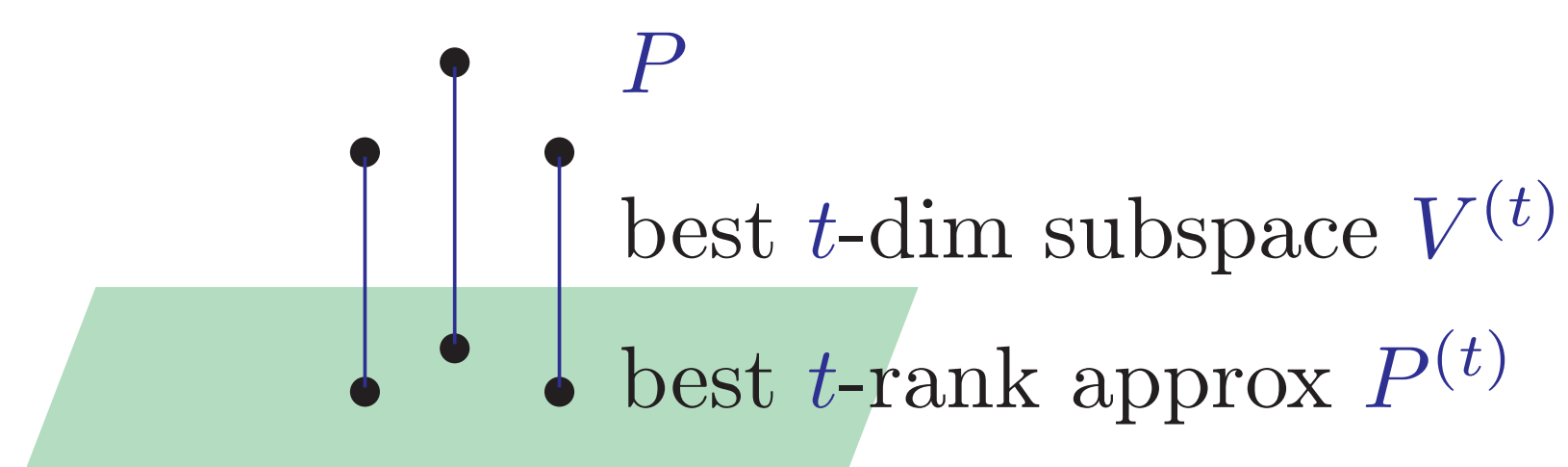


## Motivation

### Principal Component Analysis

**Given:** data matrix  $P$  with zero mean

Traditional Goal: to summarize/visualize the data in low dim, find  $t$ -dim subspace  $V$  to  $\min \|P - PVV^T\|_F^2$



More Important: **PCA pre-processes for downstream applications**

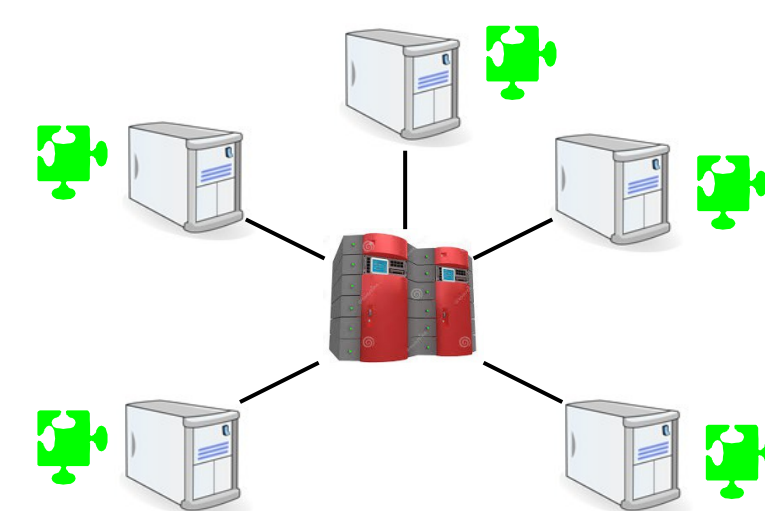
e.g. k-means, regression, kernel methods, ...

**Our Goal:** find  $t$ -dim subspace  $V$  s.t. for downstream applications, any  $\alpha$ -approx for  $PVV^T$  is  $(1 + \epsilon)\alpha$ -approx for  $P$

### Distributed Data Model

PCA is most useful for large-scale data collected distributedly

- $s$  servers linked to a coordinator
- Local datasets  $P_1, P_2, \dots, P_s$ : points in dim  $d$
- Global dataset  $P = \bigcup_i P_i$



### Challenges in Distributed PCA

1. Arbitrary local datasets: no distributional or low rank assumption
2. Relative error guarantees for downstream applications
3. Low communication/computation cost

## Our Results

### disPCA algorithm

- Relative  $(1 + \epsilon)$  error for downstream applications
- Holds for **arbitrary** data and a wide family of problems: includes low rank approx,  $k$ -means, LDA, NNMF, ...
- Low communication cost:  $O(skd/\epsilon^2)$  words

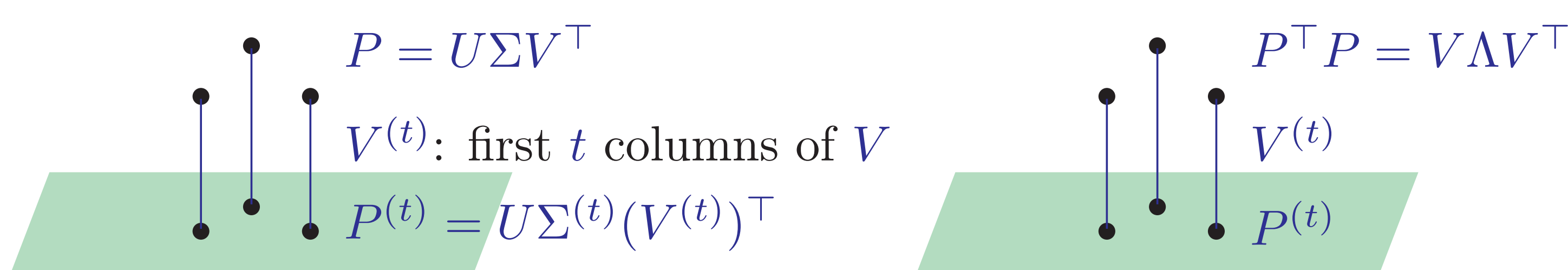
### Fast disPCA algorithm

- Easy plug-in of randomized techniques
- Merely comprise the quality and communication
- Speedup of **orders of magnitude**

## Review: Classic Non-Distributed PCA Algorithm

### 1. SVD on Data

### 2. Factorize Covariance $S = P^T P$



For a matrix  $A$ , let  $A^{(t)}$  denote the first  $t$  columns of  $A$

## Distributed PCA Algorithm

### Algorithm disPCA

1. Each server: SVD  $P_i = U_i \Sigma_i V_i^T$ ;  
send  $\Sigma_i^{(t)}$  and  $V_i^{(t)}$  to coordinator
2. Coordinator: estimate  $S = \sum_i V_i^{(t)} \Sigma_i^{(t)} \Sigma_i^{(t)} (V_i^{(t)})^T$ ;  
factorize  $S = V\Lambda V^T$  and get  $V^{(t)}$

### Equivalent Form of Algorithm disPCA

Let  $Y_i = \Sigma_i^{(t)} (V_i^{(t)})^T$ , then  $S = \sum_i Y_i^T Y_i$ .

Factorizing  $S$  is equivalent to SVD on  $Y = [Y_1; Y_2; \dots; Y_s]$ .

$$P = \begin{bmatrix} P_1 \\ \vdots \\ P_s \end{bmatrix} \xrightarrow{\text{Local PCA}} \begin{bmatrix} \Sigma_1^{(t)} (V_1^{(t)})^T \\ \vdots \\ \Sigma_s^{(t)} (V_s^{(t)})^T \end{bmatrix} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_s \end{bmatrix} = Y \xrightarrow{\text{Global PCA}} V^{(t)}$$

## Guarantee of disPCA

### Main Theorem (for $k$ -means):

Suppose  $t = O(k/\epsilon^2)$ . Then any  $\alpha$ -approx for  $\tilde{P} = PV^{(t)}(V^{(t)})^T$  is  $(1 + \epsilon)\alpha$ -approx for  $P$ .

- When  $t = O(rk/\epsilon^2)$ , the same holds for  $r$ -Subspace  $k$ -Clustering  
 $\min_{\mathcal{L}} \sum_{p \in P} \min_{L \in \mathcal{L}} d^2(p, L)$   
where  $\mathcal{L} = \{L_j\}_{j=1}^k$  is a set of  $k$  linear/affine subspaces of dim  $r$ .

## Analysis

### I. Close Projection Property of SVD:

Let  $A = U\Sigma W^T$  be its SVD, and  $\hat{A} = AW^{(t)}(W^{(t)})^T$ . For any  $k$ -dim subspace  $X$ ,  $AXX^T$  and  $\hat{A}XX^T$  are close.

### II. Close Projection Property of disPCA:

For any  $k$ -dim subspace  $X$ ,  $PXX^T$  and  $\tilde{P}XX^T$  are close.

$$P \xrightarrow{\text{project to } V^{(t)}} \tilde{P}$$

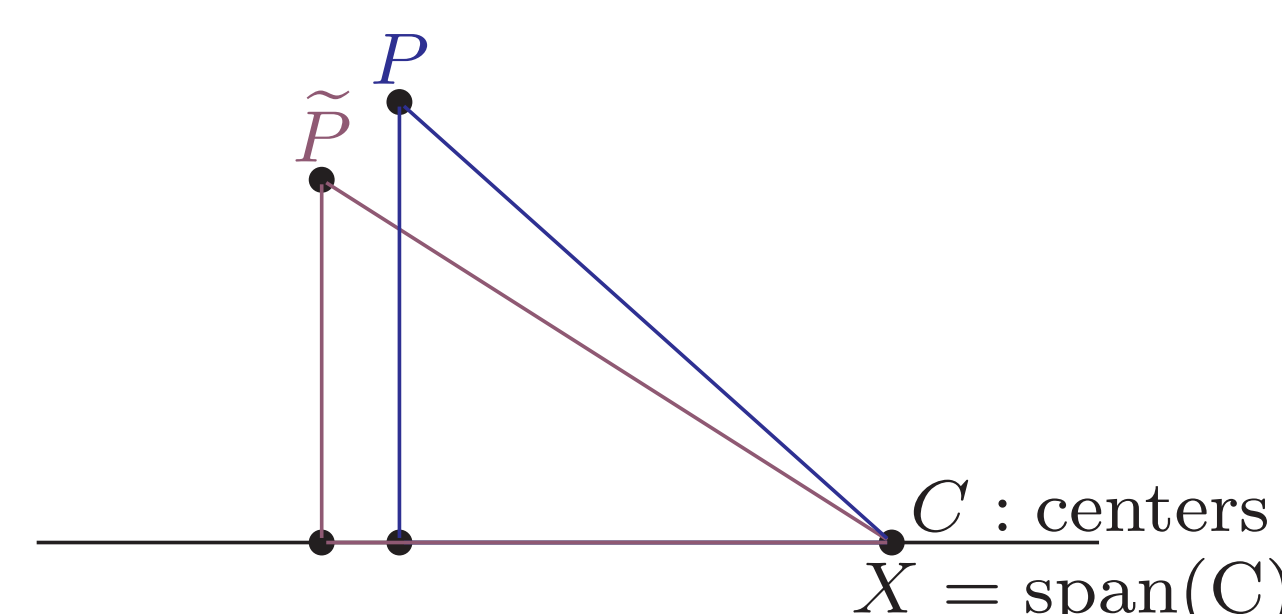
$$P_i \xrightarrow{\text{project to } V_i^{(t)}} \tilde{P}_i \xrightarrow{\text{project to } V^{(t)}} \tilde{P}_i$$

$$\|(P - \tilde{P})XX^T\|_F^2 \leq \sum_i \|(P_i - \tilde{P}_i)XX^T\|_F^2 + 2\|(\tilde{P} - \tilde{P})XX^T\|_F^2 + 2\|(\tilde{P} - \tilde{P})XX^T\|_F^2$$

each term can be bounded by **I** or its generalization

### III. Cost Decomposition:

Decompose the costs of  $\tilde{P}$  and  $P$  and show that they are close



$$\left. \begin{aligned} \text{cost}(P) &= d^2(P, PXX^T) + d^2(PXX^T, C) \\ \text{cost}(\tilde{P}) &= d^2(\tilde{P}, \tilde{P}XX^T) + d^2(\tilde{P}XX^T, C) \end{aligned} \right\} \begin{array}{l} \text{difference in each part bounded by} \\ \text{II and weak triangle inequality} \end{array}$$

## Fast Distributed PCA

### Techniques: subspace embedding; randomized SVD

### Subspace Embedding

Property of embedding matrix  $H$ :

- With prob  $\geq 99/100$ , for any  $y$

$$\|HPy\| = (1 \pm \epsilon)\|Py\|$$

such that **disPCA**( $HP$ ) almost equivalent to **disPCA**( $P$ )

- $HP$  much smaller than  $P$ , computed in time  $O(\text{nnz}(P))$ ; particularly suitable for sparse data
- $HP$  can be computed locally

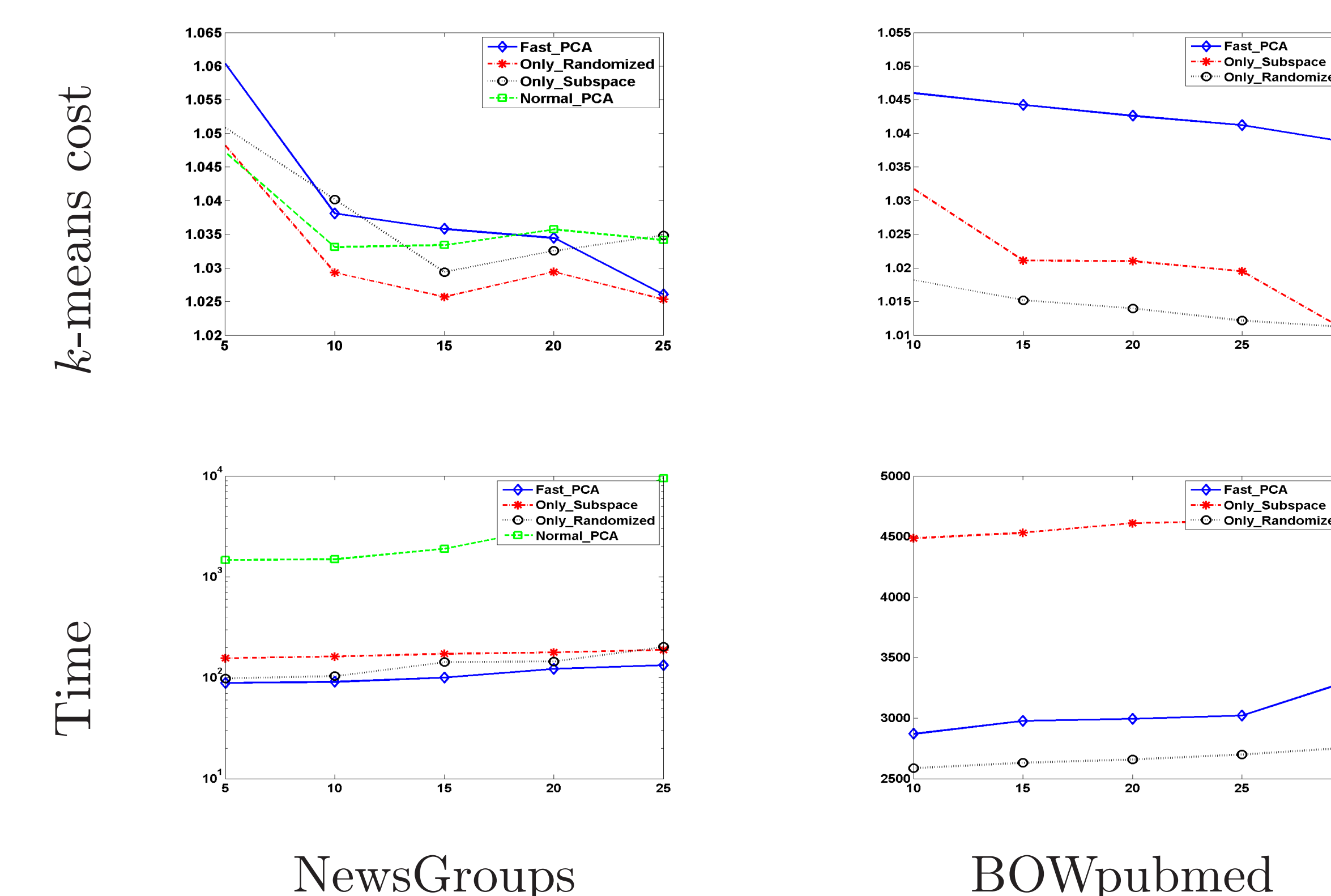
Challenge: boost the constant success prob to high prob

- Solved by our new cross validation style technique
- Dimension and sparsity of  $HP$  independent of the success prob

## Experimental Results

### Performance: cost increase < 5%; $\times 10$ to $\times 100$ speedup

### $k$ -Means Clustering: $k$ -means cost/time vs dimension



### Principal Component Regression: PCR cost/time vs dimension

