

# Learning Tools for Big Data Analytics

---

*Georgios B. Giannakis*

**Acknowledgments:** Profs. G. Mateos and K. Slavakis

NSF 1343860, 1442686, and MURI-FA9550-10-1-0567

# Growing data torrent

**\$600** to buy a disk drive that can store all of the world's music

**5 billion** mobile phones in use in 2010

**30 billion** pieces of content shared on Facebook every month

**40%** projected growth in global data generated per year vs. **5%** growth in global IT spending



**Source:** McKinsey Global Institute, "Big Data: The next frontier for innovation, competition, and productivity," May 2011.



# Big Data: Capturing its value

**\$300 billion**

potential annual value to US health care—more than double the total annual health care spending in Spain

**€250 billion**

potential annual value to Europe's public sector administration—more than GDP of Greece

**\$600 billion**

potential annual consumer surplus from using personal location data globally

**60%** potential increase in retailers' operating margins possible with big data

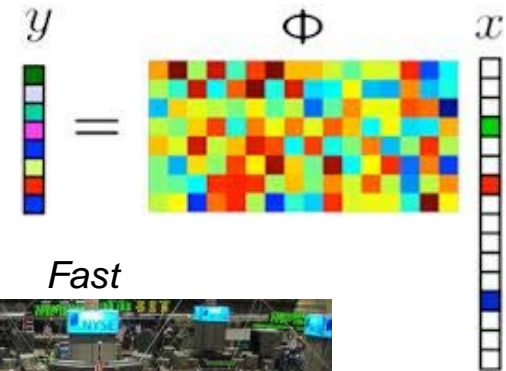


**Source:** McKinsey Global Institute, “Big Data: The next frontier for innovation, competition, and productivity,” May 2011.

# Challenges

## ❑ Sheer volume of data

- Decentralized and parallel processing
- Security and privacy measures



## ❑ Modern massive datasets involve many attributes

- Parsimonious models to ease interpretability and enhance learning performance

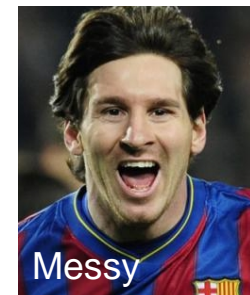
## ❑ Real-time streaming data

- Online processing
- Quick-rough answer vs. slow-accurate answer?



## ❑ Outliers and misses

- Robust imputation approaches



# Opportunities

Big tensor data models and factorizations

High-dimensional statistical SP

Network data visualization

## **Theoretical and Statistical Foundations of Big Data Analytics**

Resource tradeoffs

Pursuit of low-dimensional structure

Analysis of multi-relational data

Common principles across networks

Scalable online, decentralized optimization

Information processing over graphs

Randomized algorithms

## **Algorithms and Implementation Platforms to Learn from Massive Datasets**

Graph SP

Convergence and performance guarantees

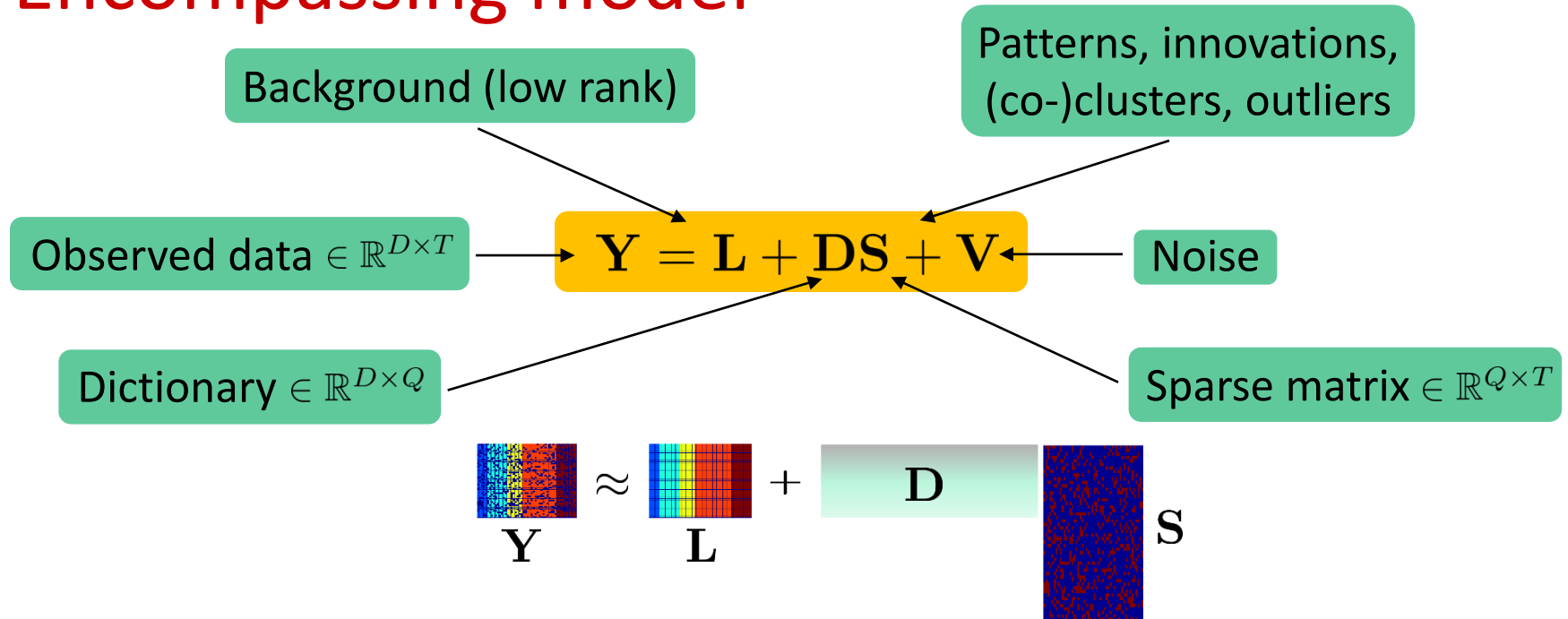
Novel architectures for large-scale data analytics

Robustness to outliers and missing data

# Roadmap

- ❑ Context and motivation
- ❑ Critical Big Data tasks
  - Encompassing and parsimonious data modeling
  - Dimensionality reduction
  - Data cleansing, anomaly detection, and inference
- ❑ Randomized learning via data sketching
- ❑ Conclusions and future research directions

# Encompassing model



- ❑ Subset  $\Omega \subset \{1, \dots, D\} \times \{1, \dots, T\}$  of observations and projection operator

$$[\mathcal{P}_\Omega(\mathbf{Y})]_{ij} = \begin{cases} [\mathbf{Y}]_{ij}, & \text{if } (i, j) \in \Omega \\ 0, & \text{o.w.} \end{cases}$$

allow for misses

- ❑ Large-scale data  $D \gg$  and/or  $T \gg$
- ❑ Any of  $\{\mathbf{L}, \mathbf{D}, \mathbf{S}\}$  unknown



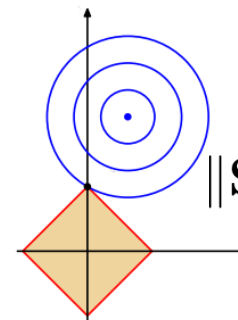
# Subsumed paradigms

## □ Structure leveraging criterion

$$\min_{\{ \}} \frac{1}{2} \| \mathbf{Y} \|_F^2$$



Nuclear norm:  $\| \mathbf{L} \|_* := \sum_{j=1}^{\text{rank}(\mathbf{L})} \sigma_j(\mathbf{L})$   
 $\{\sigma_j(\mathbf{L})\}_{j=1}^{\text{rank}(\mathbf{L})}$ : singular val. of  $\mathbf{L}$



$\ell_1$ -norm

$$\| \mathbf{S} \|_1 := \sum_{q,t} |s_{q,t}|$$

(With or without misses)

- $\mathbf{L} = \mathbf{0}, \mathbf{D}$  known  $\Rightarrow$  Compressive sampling (CS) [Candes-Tao '05]
- $\mathbf{L} = \mathbf{0} \Rightarrow$  Dictionary learning (DL) [Olshausen-Field '97]
- $\mathbf{L} = \mathbf{0}, [\mathbf{D}]_{ij} \geq 0, [\mathbf{S}]_{ij} \geq 0 \Rightarrow$  Non-negative matrix factorization (NMF) [Lee-Seung '99]
- $\mathbf{D} = \mathbf{I}_D \Rightarrow$  Principal component pursuit (PCP) [Candes et al '11]
- $\mathbf{S} = \mathbf{0}, \text{rank}(\mathbf{L}) \leq \rho \Rightarrow$  Principal component analysis (PCA) [Pearson 1901]



# PCA formulations

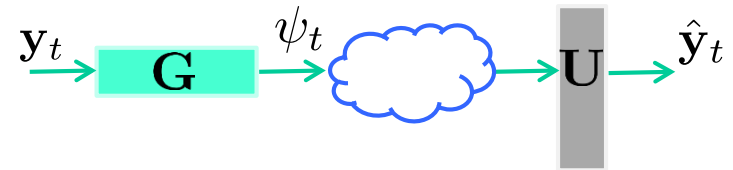
□ Training data  $\{\mathbf{y}_t \in \mathbb{R}^D\}_{t=1}^T$      $\hat{\mathbf{C}}_{yy} := (1/T) \sum_{t=1}^T \mathbf{y}_t \mathbf{y}_t^\top$

□ Minimum reconstruction error

➤ Compression  $\mathbf{G} \in \mathbb{R}^{d \times D}$

➤ Reconstruction  $\mathbf{U} \in \mathbb{R}^{D \times d}$

$$d \ll D$$



$$\min_{\mathbf{U}, \mathbf{G}} \sum_{t=1}^T \|\mathbf{y}_t - \mathbf{U} \mathbf{G} \mathbf{y}_t\|_2^2, \quad \text{s.to. } \mathbf{U}^\top \mathbf{U} = \mathbf{I}_d$$

□ Component analysis model  $\mathbf{y}_t = \mathbf{U} \boldsymbol{\psi}_t + \boldsymbol{\varepsilon}_t$

$$\min_{\mathbf{U}, \boldsymbol{\psi}_t} \sum_{t=1}^T \|\mathbf{y}_t - \mathbf{U} \boldsymbol{\psi}_t\|_2^2, \quad \text{s.to. } \mathbf{U}^\top \mathbf{U} = \mathbf{I}_d$$



**Solution:**  $\hat{\mathbf{U}}_d = d\text{-evecs}(\hat{\mathbf{C}}_{yy}), \quad \hat{\mathbf{G}} = \hat{\mathbf{U}}_d^\top, \quad \hat{\boldsymbol{\psi}}_t = \hat{\mathbf{U}}_d^\top \mathbf{y}_t$

# Dual and kernel PCA

□ SVD:  $\underbrace{\mathbf{Y}}_{D \times T} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$

$T \gg D$   $\mathbf{Y} \mathbf{Y}^\top = \mathbf{U} \mathbf{\Sigma}^2 \mathbf{U}^\top \in \mathbb{R}^{D \times D} \quad \mathcal{O}(TD^2)$

$D \gg T$   $\mathbf{Y}^\top \mathbf{Y} = \mathbf{V} \mathbf{\Sigma}^2 \mathbf{V}^\top \in \mathbb{R}^{T \times T} \quad \mathcal{O}(DT^2)$

Gram matrix

$\hat{\mathbf{U}}_d = \mathbf{Y} \hat{\mathbf{V}}_d \hat{\mathbf{\Sigma}}_d^{-1}$

$\mathbf{y}_t \rightarrow \hat{\mathbf{U}}_d^\top \mathbf{y}_t = \hat{\mathbf{\Sigma}}_d^{-1} \hat{\mathbf{V}}_d^\top \mathbf{Y}^\top \mathbf{y}_t \rightarrow \hat{\psi}_t$

Inner products

$\hat{\mathbf{U}}_d \hat{\psi}_t = \mathbf{Y} \hat{\mathbf{V}}_d \hat{\mathbf{\Sigma}}_d^{-1} \hat{\psi}_t \rightarrow \hat{\mathbf{y}}_t$

**Q.** What if approximating low-dim space not a hyperplane?

**A1.** Stretch it to become linear: Kernel PCA; e.g., [Scholkopf-Smola'01]

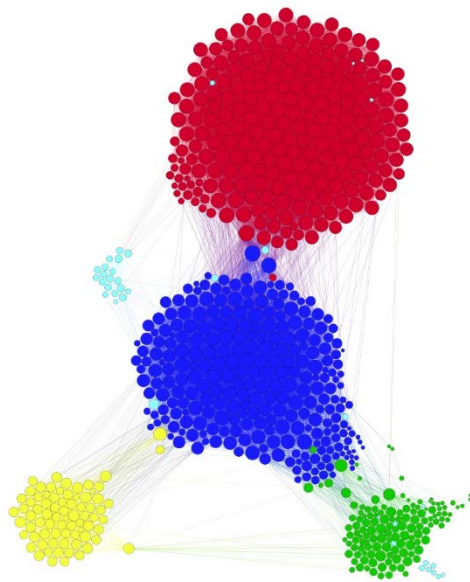
➤ maps  $\mathbf{y}_t$  to  $\varphi(\mathbf{y}_t)$ , and leverages dual PCA in high-dim spaces

**A2.** General (non)linear models; e.g., union of hyperplanes, or, locally linear

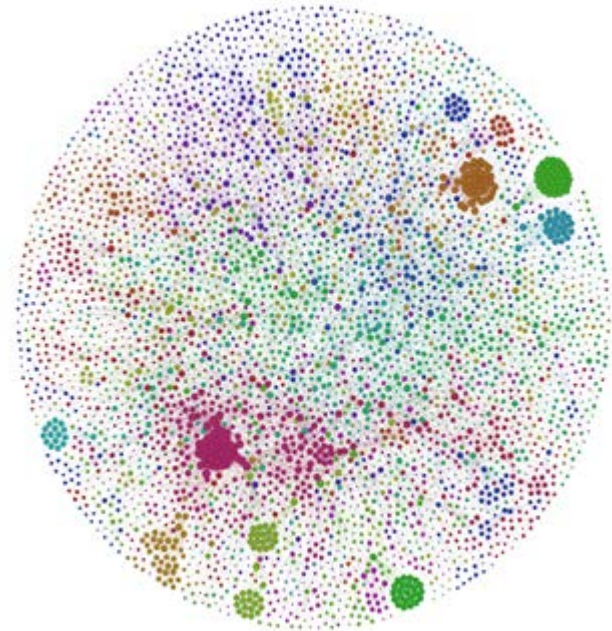
➤ tangential hyperplanes

# Identification of network communities

- ❑ Kernel PCA instrumental for partitioning of **large** graphs (**spectral clustering**)
  - Relies on graph Laplacian to capture nodal correlations



Facebook egonet  
744 nodes, 30,023 edges



arXiv collaboration network (General Relativity)  
4,158 nodes, 13,422 edges

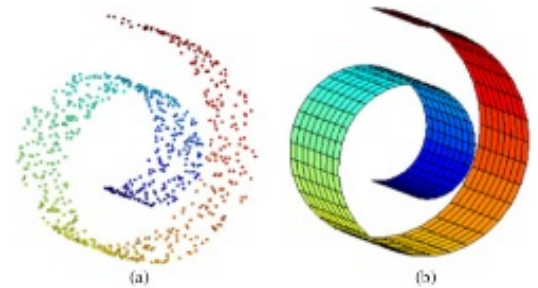
- ❑ For  $D \gg$  random sketching and validation reduces complexity to  $\mathcal{O}(d)$

# Local linear embedding

□ For each  $\mathbf{y}_t$  find neighborhood  $\{\mathbf{y}_{t'}\}_{t' \in \mathcal{N}_t}$ , e.g., k-nearest neighbors

□ Weight matrix captures local affine relations

$$\min_{\mathbf{W} := [\mathbf{w}_1, \dots, \mathbf{w}_T] \in \mathbb{R}^{T \times T} \atop \mathbf{W}^\top \mathbf{1}_T = \mathbf{0}_T} \sum_{t=1}^T \left\| \mathbf{y}_t - \sum_{t' \in \mathcal{N}_t} w_{t't} \mathbf{y}_{t'} \right\|^2$$



Sparse  $\mathcal{N}_t$  and  $\mathbf{W}$  [Elhamifar-Vidal'11]

□ Identify **low-dimensional** vectors **preserving** local geometry [Saul-Roweis'03]

$$\min_{\substack{\Psi := [\psi_1, \dots, \psi_T] \in \mathbb{R}^{d \times T} \\ \Psi \Psi^\top = \mathbf{I}_d \\ \Psi \mathbf{1}_T = \mathbf{0}_d}} \left\{ \sum_{t=1}^T \left\| \psi_t - \sum_{t'=1}^T w_{t't} \psi_{t'} \right\|^2 = \text{trace} \left[ \Psi (\mathbf{I}_T - \mathbf{W}) (\mathbf{I}_T - \mathbf{W})^\top \Psi^\top \right] \right\}$$



**Solution:** The rows of  $\Psi$  are the  $d$  minor, excluding  $\mathbf{1}_T$ ,  
evecs  $\left[ (\mathbf{I}_T - \mathbf{W}) (\mathbf{I}_T - \mathbf{W})^\top \right]$



# Dictionary learning

□ Solve for **dictionary**  $\mathbf{D}$  and **sparse**  $\mathbf{S}$ : 
$$\min_{\substack{\mathbf{D} \in \mathfrak{D} \\ \mathbf{S} \in \mathbb{R}^{Q \times T}}} \underbrace{\frac{1}{2} \|\mathcal{P}_\Omega(\mathbf{Y} - \mathbf{D}\mathbf{S})\|_F^2 + \lambda_1 \sum_{t=1}^T \|\mathbf{s}_t\|_1}_{=:\mathcal{L}(\mathbf{D}, \mathbf{S})}$$

$$\mathfrak{D} := \{\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_Q] : \|\mathbf{d}_q\| \leq 1, \forall q\} \quad Q \geq D$$

$$\begin{cases} \mathbf{S}_{k+1} \in \arg \min_{\mathbf{S} \in \mathbb{R}^{Q \times T}} \overbrace{\frac{1}{2} \|\mathbf{Y} - \mathbf{D}_k \mathbf{S}\|_F^2 + \lambda_1 \|\mathbf{S}\|_1}^{\mathcal{L}(\mathbf{D}_k, \mathbf{S})} \\ \mathbf{D}_{k+1} \in \arg \min_{\mathbf{D} \in \mathfrak{D}} \|\mathbf{Y} - \mathbf{D} \mathbf{S}_{k+1}\|_F^2 = \arg \min_{\mathbf{D} \in \mathfrak{D}} \mathcal{L}(\mathbf{D}, \mathbf{S}_{k+1}) \end{cases}$$

(Lasso task; sparse coding)

(Constrained LS task)

□ Alternating minimization; both  $\mathcal{L}(\mathbf{D}_k, \cdot)$  and  $\mathcal{L}(\cdot, \mathbf{S}_{k+1})$  are convex

□ Under conditions,  $(\mathbf{D}_k, \mathbf{S}_k)_{k=0}^\infty$  converges to a stationary point of  $\mathcal{L}$  [Tseng'01]

# Joint DL-LLE paradigm

$$\min_{\substack{\mathbf{D}_y, \mathbf{S}, \mathbf{S}^\top \mathbf{1}_Q = \mathbf{1}_T \\ \mathbf{W}, \mathbf{W}^\top \mathbf{1}_Q = \mathbf{1}_Q \\ \text{diag}(\mathbf{W}) = \mathbf{0}}} \underbrace{\frac{1}{2} \|\mathcal{P}_\Omega(\mathbf{Y} - \mathbf{D}_y \mathbf{S})\|_F^2}_{\text{DL fit}} + \underbrace{\frac{\lambda_y}{2} \|\mathbf{D}_y - \mathbf{D}_y \mathbf{W}\|_F^2}_{\text{LLE fit}} + \underbrace{\sum_{t=1}^T \lambda_{st} \|\mathbf{s}_t\|_1 + \sum_{q=1}^Q \lambda_{wq} \|\mathbf{w}_q\|_1}_{\text{Sparsity regularization}}$$

- Dictionary morphs data to a smooth basis; reduces noise and complexity

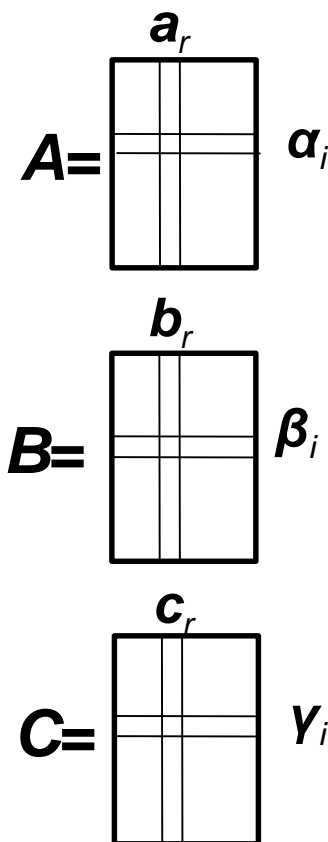
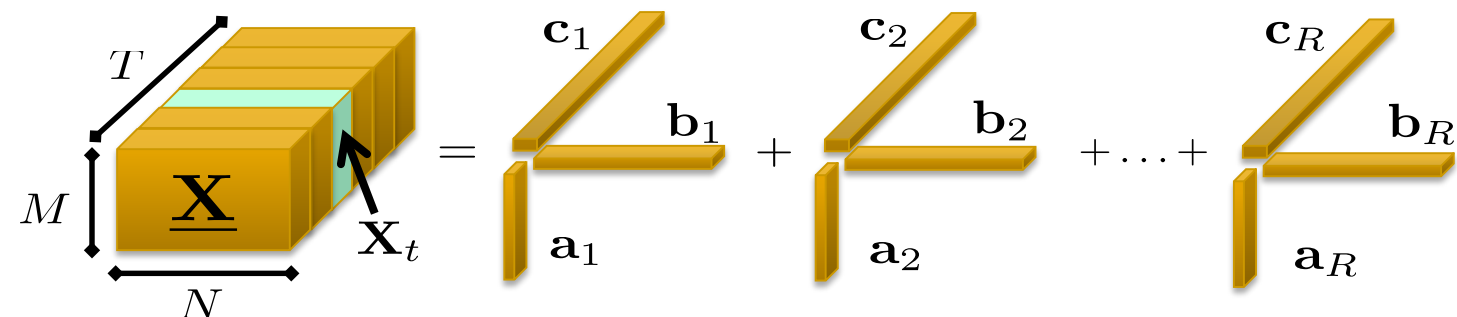


- IDpaiting for data-affine geometry-preserving for 50% missing entries; PSNR 32dB

# From low-rank matrices to tensors

- Data cube  $\underline{\mathbf{X}} \in \mathbb{R}^{M \times N \times T}$ , e.g., sub-sampled MRI frames

$$\mathbf{Y}_t^\Omega \approx \mathcal{F}_{\Omega_t}(\mathbf{X}_t)$$



- **PARAFAC** decomposition per slab  $t$  [Harshman '70]

$$\mathbf{X}_t = \sum_{r=1}^R \gamma_{t,r} \mathbf{a}_r \mathbf{b}_r^\top = \mathbf{A} \text{diag}(\gamma_t) \mathbf{B}^\top$$

- Tensor subspace comprises  $R$  rank-one matrices  $\{\mathbf{a}_r \mathbf{b}_r^\top\}_{r=1}^R$

**Goal:** Given streaming  $\mathbf{Y}_t^\Omega \approx \mathcal{F}_{\Omega_t}(\mathbf{A} \text{diag}(\gamma_t) \mathbf{B}^\top)$ , learn the subspace matrices  $(\mathbf{A}, \mathbf{B})$  recursively, and impute possible misses of  $\mathbf{Y}_t$

# Online tensor subspace learning

- Image domain low tensor rank  $\mathbf{Y}_t^\Omega \approx \mathcal{F}_{\Omega_t}(\mathbf{A} \text{diag}(\gamma_t) \mathbf{B}^\top)$

$$(\hat{\mathbf{A}}_t, \hat{\mathbf{B}}_t) = \arg \min_{\mathbf{A}, \mathbf{B}} \frac{1}{t} \sum_{\tau=1}^t \min_{\gamma_\tau} \left\{ \|\mathbf{Y}_\tau^\Omega - \mathcal{F}_{\Omega_\tau}(\mathbf{A} \text{diag}(\gamma_\tau) \mathbf{B}^\top)\|_F^2 + \frac{\lambda}{2} \|\gamma_\tau\|^2 \right\} \\ + \frac{\lambda}{2t} (\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2)$$

- Tikhonov regularization promotes low rank

**Proposition [Bazerque-GG '13]:** With  $[\boldsymbol{\sigma}]_r = \|\mathbf{a}_r\| \|\mathbf{b}_r\| \|\mathbf{c}_r\|$

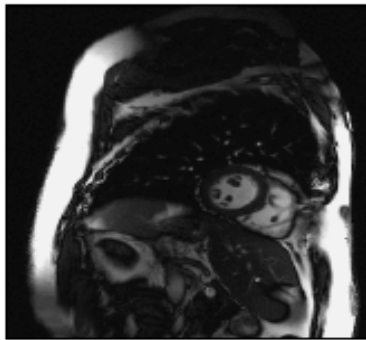
$$\|\boldsymbol{\sigma}\|_{2/3}^{2/3} := \arg \min_{\{\mathbf{A}, \mathbf{B}, \mathbf{C}\}} (\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2 + \|\mathbf{C}\|_F^2)$$

- Stochastic alternating minimization; parallelizable across bases
- Real-time reconstruction (FFT per iteration)  $\hat{\mathbf{X}}_t = \hat{\mathbf{A}}_t \text{diag}(\hat{\gamma}_t) \hat{\mathbf{B}}_t^\top$

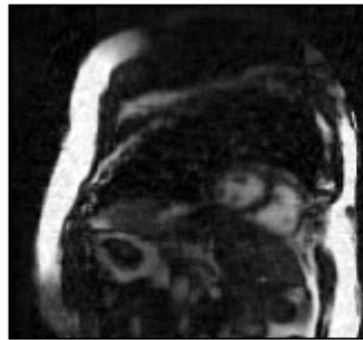


# Dynamic cardiac MRI test

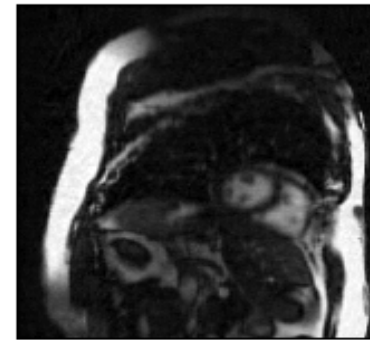
- *in vivo* dataset: 256 k-space 200x256 frames



Ground-truth frame



$R=100$ , 90% misses



$R=150$ , 75% misses



Sampling trajectory

- Potential for accelerating MRI at high spatio-temporal resolution
- Low-rank  $\mathcal{F}_{\Omega_t}(\mathbf{X}_t)$  plus  $\mathcal{F}_{\Omega_t}(\mathbf{D}\mathbf{S}_t)$  can also capture motion effects

# Roadmap

- ❑ Context and motivation
- ❑ Critical Big Data tasks
- ❑ Randomized learning via data sketching
  - Johnson-Lindenstrauss lemma
  - Randomized linear regression
  - Randomized clustering
- ❑ Conclusions and future research directions

# Randomized linear algebra

- ❑ **Basic tools:** Random sampling and random projections
- ❑ **Attractive features**
  - Reduced dimensionality to lower complexity with Big Data
  - Rigorous error analysis at reduced dimension

**Ordinary least-squares (LS)**    Given  $\mathbf{y} \in \mathbb{R}^D$ ,  $\mathbf{X} \in \mathbb{R}^{D \times p}$

$$\boldsymbol{\theta}_{\text{LS}} := \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^p} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2$$

$$\text{If } \text{rank}(\mathbf{X}) = p \implies \boldsymbol{\theta}_{\text{LS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

- ❑ SVD incurs complexity  $\mathcal{O}(Dp^2)$ . **Q:** What if  $D \gg p$ ?

# Randomized LS for linear regression

- LS estimate using (pre-conditioned) random projection matrix  $\mathbf{R}_{d \times D}$

$$\check{\boldsymbol{\theta}}_{\text{LS}} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^p} \left\| \overbrace{\underbrace{\boldsymbol{\Gamma}_d \mathbf{S}_d}_{\mathbf{R}_2} \underbrace{\mathbf{H}_D \boldsymbol{\Delta}_D}_{\mathbf{R}_1}}^{\mathbf{R}} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \right\|_2^2$$

- Random diagonal w/  $[\boldsymbol{\Delta}_D]_{ii} \in \{1, -1\} \sim \text{Ber}(1/2)$  and Hadamard matrix

$$\mathbf{H}_D = \frac{1}{\sqrt{D}} \begin{bmatrix} \mathbf{H}_{D/2} & \mathbf{H}_{D/2} \\ \mathbf{H}_{D/2} & -\mathbf{H}_{D/2} \end{bmatrix}, \quad \mathbf{H}_2 := \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

- subsets of data obtained by uniform sampling/scaling via  $\mathbf{S}_d, \boldsymbol{\Gamma}_d$  yield LS estimates of “comparable quality”

- Select reduced dimension  $d = \mathcal{O}(p \log p \cdot \log D + \epsilon^{-1} D \log p)$

- Complexity reduced from  $\mathcal{O}(Dp^2)$  to  $o(Dp^2)$



# Johnson-Lindenstrauss lemma

- The “workhorse” for proofs involving random projections

**JL lemma:** If  $0 < \epsilon < 1$ , integer  $T$ , and reduced dimension satisfies

$$d \geq 4(\epsilon^2/2 - \epsilon^3/3)^{-1} \ln T$$

then for any  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_T] \in \mathbb{R}^{D \times T}$  there exists a mapping  $f : \mathbb{R}^D \rightarrow \mathbb{R}^d$  s.t.

$$(1 - \epsilon) \|\mathbf{y}_{t_1} - \mathbf{y}_{t_2}\|^2 \leq \|f(\mathbf{y}_{t_1}) - f(\mathbf{y}_{t_2})\|^2 \leq (1 + \epsilon) \|\mathbf{y}_{t_1} - \mathbf{y}_{t_2}\|^2 \quad (\star)$$

**Almost preserves pairwise distances!**

- If  $f(\mathbf{y}) := d^{-1/2} \mathbf{R} \mathbf{y}$  with i.i.d.  $\mathcal{N}(0, 1)$  entries of  $\mathbf{R}$  and reduced dimension  $d \geq 4(\epsilon^2/2 - \epsilon^3/3)^{-1} \ln T + \mathcal{O}(\log \log T)$ , then  $(\star)$  holds w.h.p. [Indyk-Motwani'98]
- If  $f(\mathbf{y}) := d^{-1/2} \mathbf{R} \mathbf{y}$  with i.i.d. uniform over  $\{+1, -1\}$  entries of  $\mathbf{R}$  and reduced dimension as in JL lemma, then  $(\star)$  holds w.h.p. [Achlioptas'01]

# Performance of randomized LS

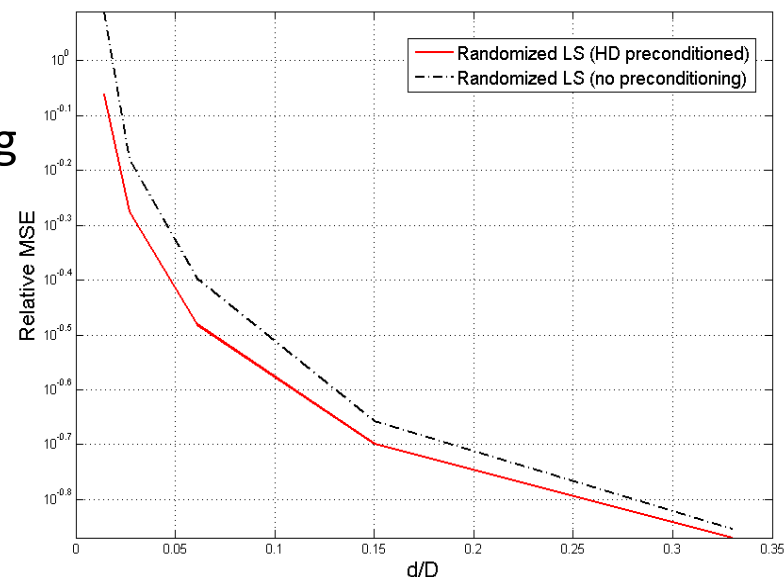
**Theorem** For any  $\epsilon > 0$ , if  $d = \mathcal{O}(p \log p / \epsilon^2)$ , then w.h.p.

$$\|\mathbf{y} - \mathbf{X}\check{\boldsymbol{\theta}}_{\text{LS}}\|_2 \leq (1 + \epsilon) \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}_{\text{LS}}\|_2$$
$$\|\boldsymbol{\theta}_{\text{LS}} - \check{\boldsymbol{\theta}}_{\text{LS}}\|_2 \leq \sqrt{\epsilon} \kappa(\mathbf{X}) \sqrt{\gamma^{-2} - 1} \|\boldsymbol{\theta}_{\text{LS}}\|_2$$

$\kappa(\mathbf{X})$  condition number of  $\mathbf{X}$ ; and  $\gamma = \|\hat{\mathbf{y}}\|_2 / \|\mathbf{y}\|_2$

## □ Uniform sampling vs Hadamard preconditioning

- $D = 10,000$  and  $p = 50$
- Performance depends on  $\mathbf{X}$

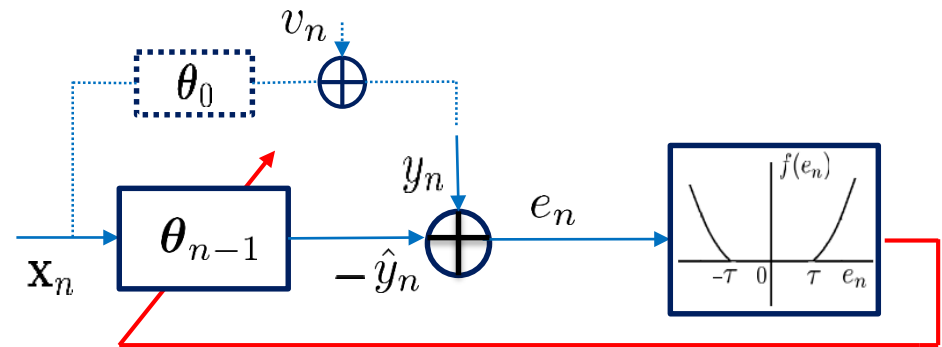


# Online censoring for large-scale regression

❑ **Key idea:** Sequentially test and update RLS estimates only for informative data

❑ Adaptive censoring (AC) rule

$$c_n = \begin{cases} 1, & \frac{|y_n - \mathbf{x}_n^T \boldsymbol{\theta}_{n-1}|}{\sigma} \leq \tau \\ 0, & \text{otherwise.} \end{cases}$$



❑ Criterion reveals “causal” support vectors (SVs)

$$f_n(\boldsymbol{\theta}) = f(e_n) := \begin{cases} \frac{e_n^2}{2} - \frac{\tau^2 \sigma^2}{2} & |e_n| > \tau \sigma \\ 0 & |e_n| \leq \tau \sigma \end{cases}$$

❑ Threshold controls avg. data reduction:  $\tau \approx Q^{-1}(\frac{1}{2}(1 - \frac{d}{D}))$ ,  $D \gg p$

# Censoring algorithms and performance

- ❑ AC least mean-squares (LMS)

$$\hat{\boldsymbol{\theta}}_n = \hat{\boldsymbol{\theta}}_{n-1} + \mu(1 - c_n) \mathbf{x}_n (y_n - \mathbf{x}_n^T \hat{\boldsymbol{\theta}}_{n-1})$$

- ❑ AC recursive least-squares (RLS) at complexity  $\mathcal{O}(dp^2)$

$$\begin{aligned} \hat{\boldsymbol{\theta}}_n &= \hat{\boldsymbol{\theta}}_{n-1} + (1 - c_n) \frac{1}{n} \hat{\mathbf{C}}_n \mathbf{x}_n (y_n - \mathbf{x}_n^T \hat{\boldsymbol{\theta}}_{n-1}) \\ \hat{\mathbf{C}}_n &= \frac{n}{n-1} \left[ \hat{\mathbf{C}}_{n-1} - (1 - c_n) \hat{\mathbf{C}}_{n-1} \mathbf{x}_n \mathbf{x}_n^T \hat{\mathbf{C}}_{n-1} \left( n - 1 + \mathbf{x}_n^T \hat{\mathbf{C}}_{n-1} \mathbf{x}_n \right)^{-1} \right] \end{aligned}$$

**Proposition: AC-RLS**  $\frac{1}{n} \text{tr}(\mathbf{R}_x^{-1}) \sigma^2 \leq \mathbb{E} \left[ \|\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0\|_2^2 \right] \leq \frac{1}{n} \frac{\text{tr}(\mathbf{R}_x^{-1}) \sigma^2}{2Q(\tau)} \quad \forall n \geq k$

**AC-LMS**  $\mathbb{E} \left[ \|\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0\|_2^2 \right] \leq \frac{\exp(4L^2/\alpha^2)}{n^2} \left( \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\|_2^2 + \frac{\Delta}{L^2} \right) + 8 \frac{\Delta}{\alpha^2} \frac{\log n}{n}$

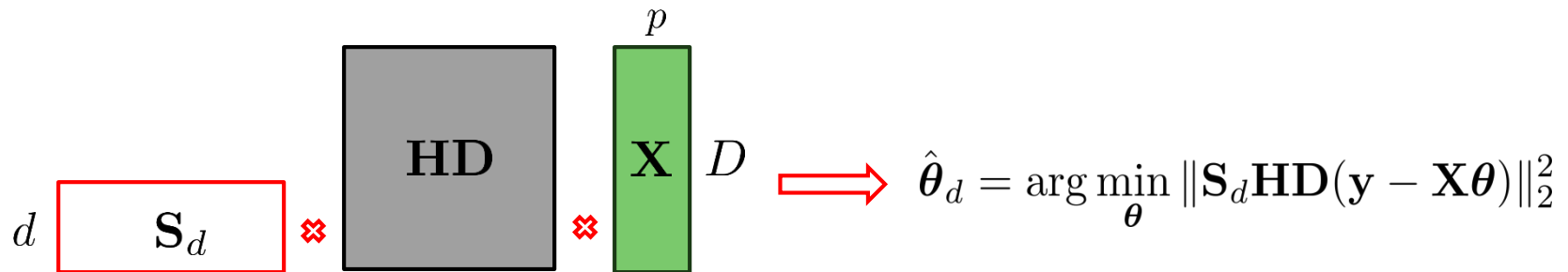
- ❑ AC Kalman Filtering and Smoothing for “tracking with a budget”



# Censoring vis-a-vis random projections

## ❑ Random projections for linear regression [Mahoney '11]

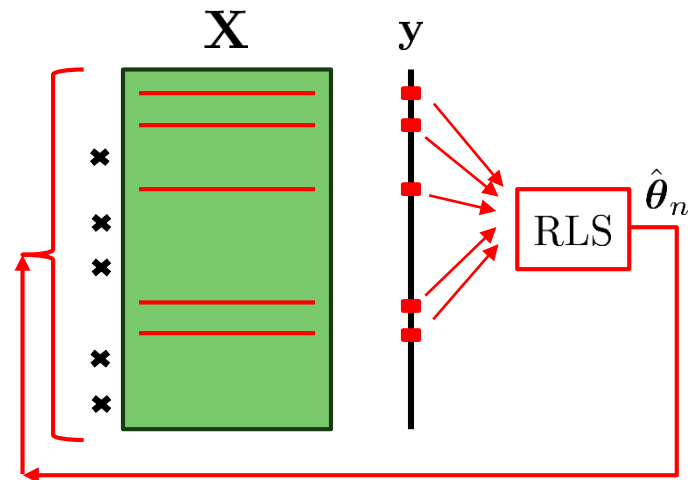
- **Data-agnostic** reduction decoupled from LS solution



The diagram illustrates the process of random projections for linear regression. It shows a sequence of operations: a matrix  $S_d$  of size  $d \times D$  is multiplied (indicated by a red double cross symbol) by a matrix  $HD$  of size  $D \times D$ , which is then multiplied by a matrix  $X$  of size  $D \times p$ . The result is a matrix of size  $d \times p$ . This is followed by a red arrow pointing to the equation  $\hat{\theta}_d = \arg \min_{\theta} \|S_d HD(y - X\theta)\|_2^2$ .

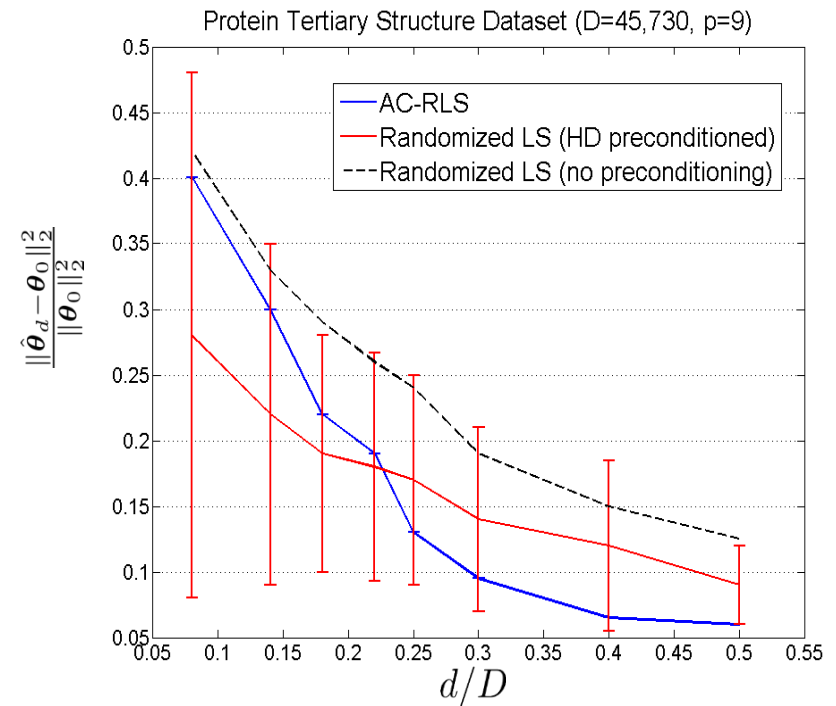
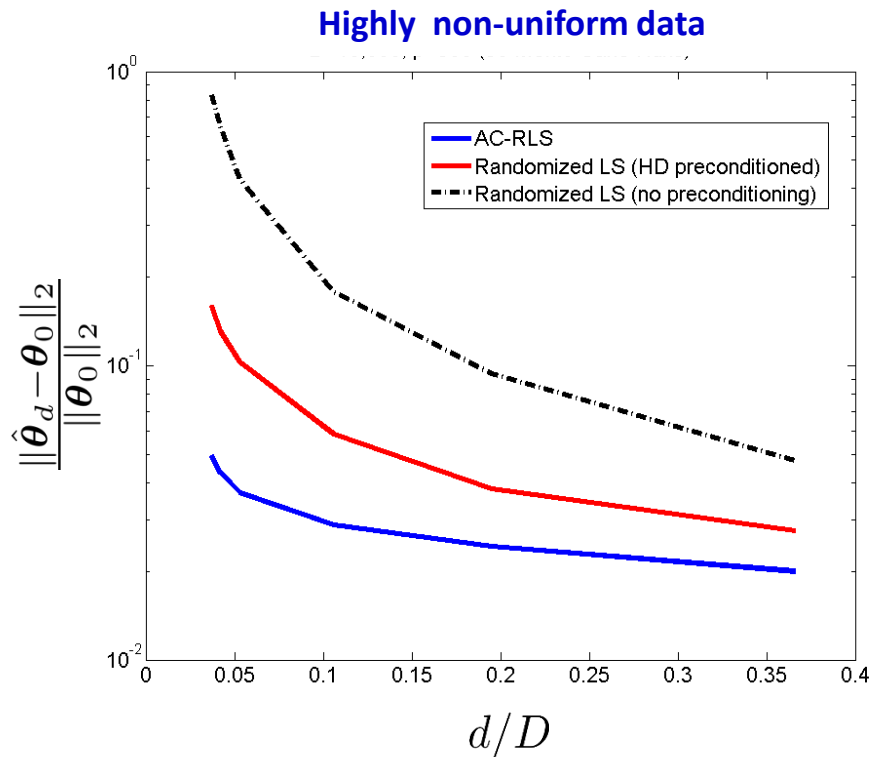
## ❑ Adaptive censoring (AC-RLS)

- **Data-driven** measurement selection
- Suitable also for streaming data
- Minimal memory requirements



# Performance comparison

❑ **Synthetic:**  $D=10,000$ ,  $p=300$  (50 MC runs); **Real data:**  $\theta_0, \sigma$  estimated from full set



- ❑ AC-RLS outperforms alternatives at comparable complexity
- ❑ Robustness to uniform data (all rows of  $\mathbf{X}$  equally “important”)

# Big data clustering

- ❑ Given  $\{\mathbf{y}_t\}_{t=1}^T$  with  $\dim(\mathbf{y}_t) = D \gg$  assign them to clusters
- ❑ **Key idea:** Reduce dimensionality via random projections
- ❑ **Desiderata:** Preserve the pairwise data distances in lower dimensions

$$\mathbf{Y} := [\mathbf{y}_1, \dots, \mathbf{y}_T]$$

## Feature extraction

- ❑ Construct  $d \ll D$  “combined” features (e.g., via  $\mathbf{R}\mathbf{Y}$ )
- ❑ Apply  $K$ -means to  $d$ -space

## Feature selection

- ❑ Select  $d \ll D$  of input features (rows of  $\mathbf{Y}$ )
- ❑ Apply  $K$ -means to  $d$ -space

# Random sketching and validation (SkeVa)

□ Randomly select  $d \ll D$  “informative” dimensions

□ **Algorithm** For  $r = 1, \dots, R_{\max}$

❖ **Sketch**  $d \ll D$  dimensions:  $\mathbf{X} \rightarrow \check{\mathbf{X}}^{(r)} \in \mathbb{R}^{d \times N}$

❖ Run k-means on  $\check{\mathbf{X}}^{(r)} \rightarrow \{\check{\mathcal{C}}_k^{(r)}\}_{k=1}^K, \{\check{\mathbf{c}}_k^{(r)}\}_{k=1}^K$

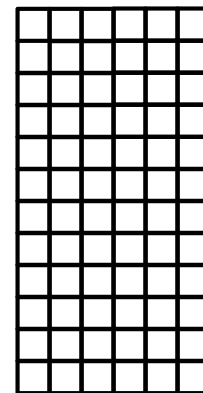
❖ Re-sketch  $d' \leq D - d$  dimensions  $\rightarrow \check{\mathbf{X}}^{(r')} \in \mathbb{R}^{d' \times N}$

❖ Augment centroids  $\bar{\mathbf{c}}_k^{(r)} := [\check{\mathbf{c}}_k^{(r)\top}, \check{\mathbf{c}}_k^{(r')\top}]^\top \quad \forall k, \check{\mathbf{c}}_k^{(r')} = \frac{1}{|\check{\mathcal{C}}_k^{(r)}|} \sum_{\check{\mathbf{x}}_n^{(r)} \in \check{\mathcal{C}}_k^{(r)}} \check{\mathbf{x}}_n^{(r')}$

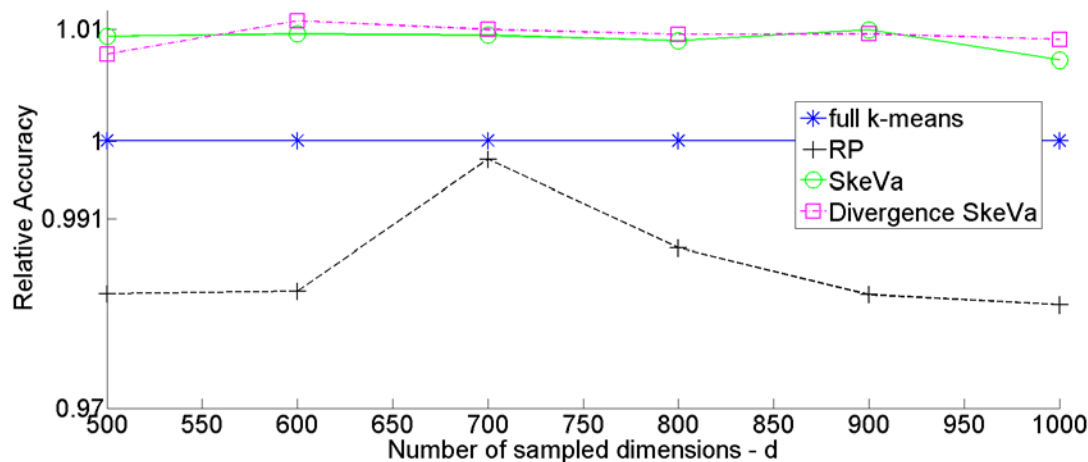
❖ **Validate** using consensus set  $\mathcal{S}^{(r)} = \{\mathbf{x}_n | \check{\mathbf{x}}_n^r \in \check{\mathcal{C}}_{k_1}^{(r)}, \bar{\mathbf{x}}_n^r \in \bar{\mathcal{C}}_{k_2}^{(r)}, \text{ and } k_1 = k_2\}$

➤  $r^* = \operatorname{argmax}_r f(\mathcal{S}^{(r)})$

□ Similar approaches possible for  $N \gg$  □ Sequential and kernel variants available

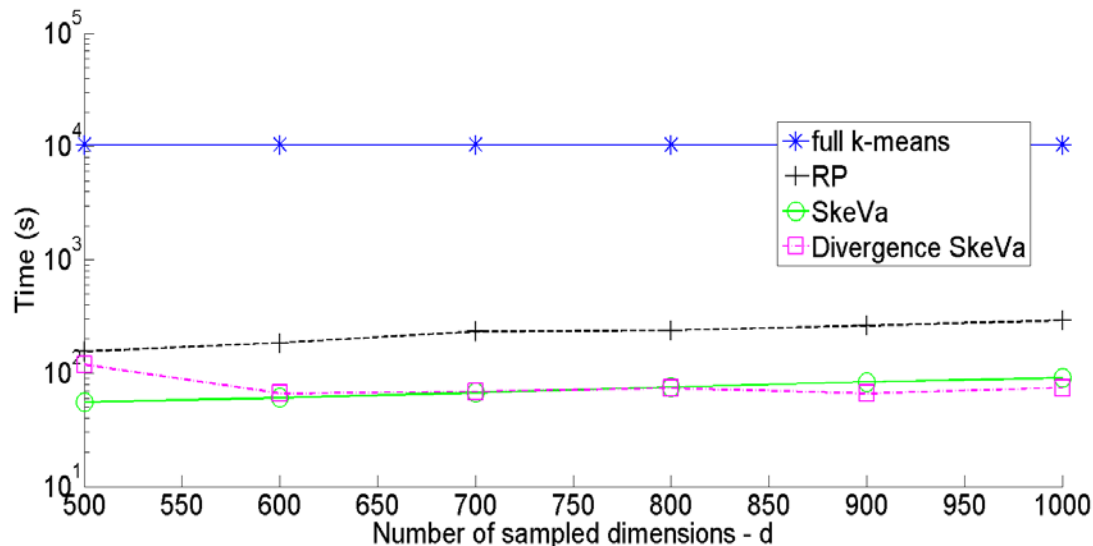


# RP versus SkeVA comparisons



**KDDb** dataset (subset)

**$D = 2,990,384$ ,  $T = 10,000$ ,  $K = 2$**



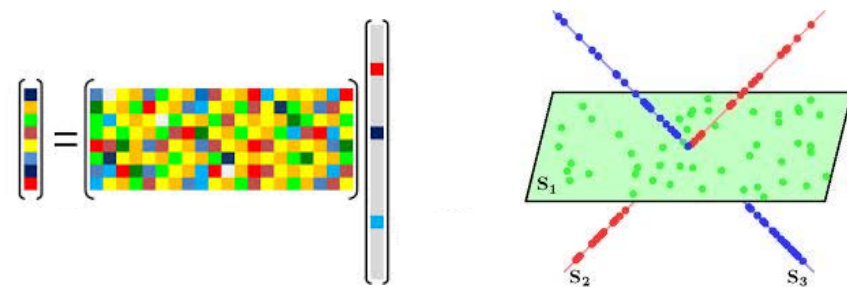
RP: [Boutsidis et al '13]

SkeVa: Sketch and validate

# Closing comments

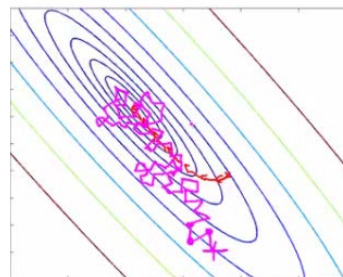
## □ Big Data modeling and tasks

- Dimensionality reduction
- Succinct representations
- Vectors, matrices, and tensors



## □ Learning algorithms

- Data sketching via random projections
- Streaming, parallel, decentralized



## □ Implementation platforms

- Scalable computing platforms
- Analytics in the cloud

