

# Efficiently Implementing Sparsity in Learning

M. Magdon-Ismail

Rensselaer Polytechnic Institute

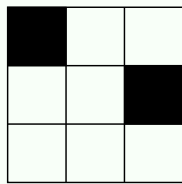
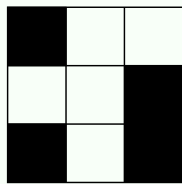
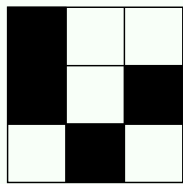


(Joint Work)

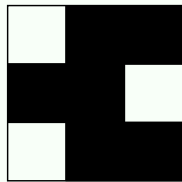
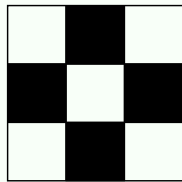
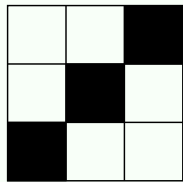


December 9, 2013.

# Out-of-Sample is What Counts

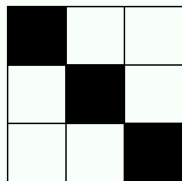


NO



YES

- A **pattern** exists
- We **don't know it**
- We **have data** to learn it
- Tested on **new cases**



?

# Data

## Data Matrix

$d$  dimensions ☹️

$n$  data points 😊

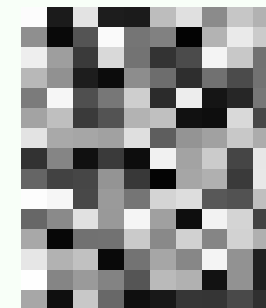
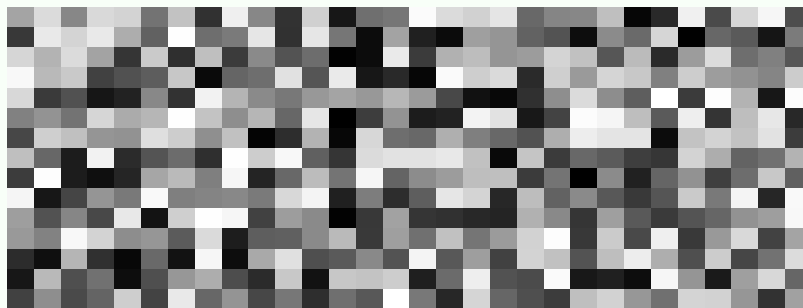
name	age	debt	income	...	hair	weight	sex
John	21yrs	-\$10K	\$65K	...	black	175lbs	M
Joe	74yrs	-\$100K	\$25K	...	blonde	275lbs	M
Jane	27yrs	-\$20K	\$85K	...	blonde	135lbs	F
⋮							
Jen	37yrs	-\$400K	\$105K	...	brun	155lbs	F

## Response Matrix

credit?	limit	risk
✓	2K	high
✗	0	—
✓	10K	low
⋮		
✓	15K	high

$$X \in \mathbb{R}^{n \times d}$$

$$Y \in \mathbb{R}^{n \times \omega}$$



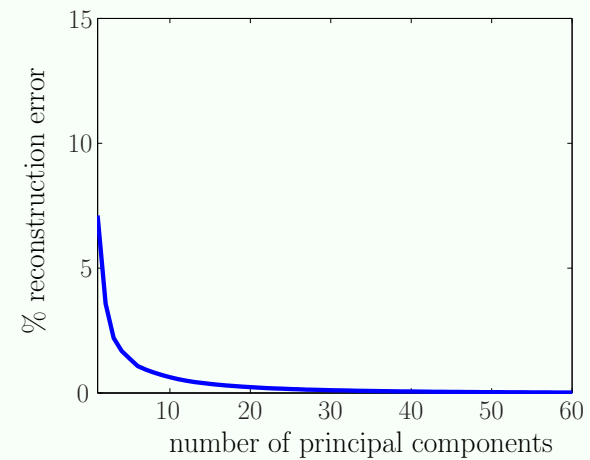
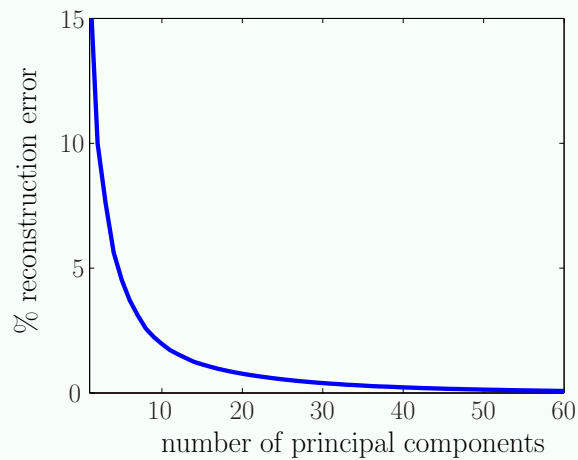
# More Beautiful Data



$$X \in \mathbb{R}^{231 \times 174}$$



$$Y \in \mathbb{R}^{231 \times 166}$$



---

# Throwing Out Unnecessary Features is Good

Sparsity: represent your solution **using only a few features**.

‘Sparse’ solutions generalize to out-of-sample better – less *overfitting*.

Sparse solutions are easier to interpret – few important features.

Computations are more efficient.

**Problem:** How to find the few relevant features *quickly*.

# PCA, $K$ -means, Linear Regression



$$k = 20$$

$$r = 2k$$

## PCA



Exact



Approx, fast (relative error)



Sparse, approx, fast (relative error)

## $K$ -Means



Exact



Sparse, approx, fast (relative error)

## Regression



Exact

$$\left[ \begin{array}{c} \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \\ \phantom{\text{ }} \end{array} \right] =$$



top- $k$  PCA regression

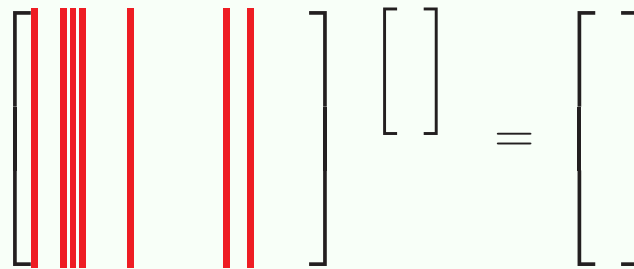


Fast-sparse regression (additive error)

# Sparsity

Represent your solution using **only a few** ...

**Example:** linear regression


$$\left[ \begin{array}{c} | \\ | \\ | \\ | \\ | \\ | \end{array} \right] \left[ \begin{array}{c} \\ \\ \\ \\ \\ \end{array} \right] = \left[ \begin{array}{c} \\ \\ \\ \\ \\ \end{array} \right]$$

$$X\mathbf{w} = \mathbf{y}$$

$\mathbf{y}$  is an optimal linear combination of **only a few** columns in  $X$ .

(sparse regression; regularization ( $\|\mathbf{w}\|_0 \leq k$ ); feature subset selection; ...)

# Singular Value Decomposition (SVD)

$$X = \begin{bmatrix} U_k & U_{d-k} \end{bmatrix} \begin{bmatrix} \Sigma_k & 0 \\ 0 & \Sigma_{d-k} \end{bmatrix} \begin{bmatrix} V_k^T \\ V_{d-k}^T \end{bmatrix} \quad O(nd^2)$$
$$\begin{array}{ccc} U & \Sigma & V^T \\ (n \times d) & (d \times d) & (d \times d) \end{array}$$

$$\begin{aligned} X_k &= U_k \Sigma_k V_k^T \\ &= X V_k V_k^T \end{aligned}$$

$X_k$  is the best rank- $k$  approximation to  $X$ .

Reconstruction of  $X$  using **only a few deg. of freedom**.



X



X<sub>20</sub>



X<sub>40</sub>



X<sub>60</sub>

$V_k$  is an orthonormal basis for the best  $k$ -dimensional subspace of the row space of  $X$ .



# Fast Approximate SVD

- 1:  $Z = XR$
  - 2:  $Q = \text{QR.FACTORIZE}(Z)$
  - 3:  $\hat{V}_k \leftarrow \text{SVD}_k(Q^T X)$
- $R \sim \mathcal{N}(d \times r), Z \in \mathbb{R}^{n \times r}$

**Theorem.** Let  $r = \lceil k(1 + \frac{1}{\epsilon}) \rceil$  and  $E = X - X\hat{V}_k\hat{V}_k^T$ . Then,

$$\mathbb{E} [\| E \|] \leq (1 + \epsilon) \| X - X_k \|$$

running time is  $O(ndk) = o(\text{SVD})$

[BDM, FOCS 2011]

# $V_k$ and Sparsity

Important “dimensions” of  $V_k^T$  are important for  $X$

$$\left[ \begin{array}{|c|} \hline \times s_1 \\ \hline \times s_2 \\ \hline \end{array} \quad \begin{array}{|c|} \hline \times s_3 \\ \hline \times s_4 \\ \hline \end{array} \quad \begin{array}{|c|} \hline \times s_5 \\ \hline \end{array} \right] \longrightarrow \left[ \hat{V}_k^T \in \mathbb{R}^{k \times r} \right]$$

$V_k^T$

The sampled  $r$  columns are “good” if

$$I = V_k^T V_k \approx \hat{V}_k^T \hat{V}_k.$$

Sampling schemes: [Largest norm \(Jolliffe, 1972\)](#);  
[Randomized norm sampling \(Rudelson, 1999; RudelsonVershynin, 2007\)](#);  
[Greedy \(Batson et al, 2009; BDM, 2011\)](#).

# Sparse PCA – Algorithm

- 1: Choose a few columns  $C$  of  $X$ ;  $C \in \mathbb{R}^{n \times r}$ .
- 2: Find the best rank- $k$  approximation of  $X$  in the span of  $C$ ,  $X_{C,k}$ .
- 3: Compute the  $SVD_k$  of

$$X_{C,k} = U_{C,k} \Sigma_{C,k} V_{C,k}^T.$$

4:

$$Z = X V_{C,k}.$$

Each feature in  $Z$  is a mixture of **only the few** original  $r$  feature dimensions in  $C$ .

$$\|X - X V_{C,k} V_{C,k}^T\| \leq \|X - X_{C,k} V_{C,k} V_{C,k}^T\| = \|X - X_{C,k}\| \leq \left(1 + O\left(\frac{2k}{r}\right)\right) \|X - X_k\|.$$

[BDM, FOCS 2011]

# Sparse PCA

$k = 20$



$k = 40$



$k = 60$



Dense PCA



Sparse PCA,  $r = 2k$

**Theorem.** One can construct, in  $o(\text{SVD})$ ,  $k$  features that are  $r$ -sparse,  $r = O(k)$ , that are as good as exact dense top- $k$  PCA-features.

# Clustering: $K$ -Means

Full, slow

Fast, sparse



3 clusters

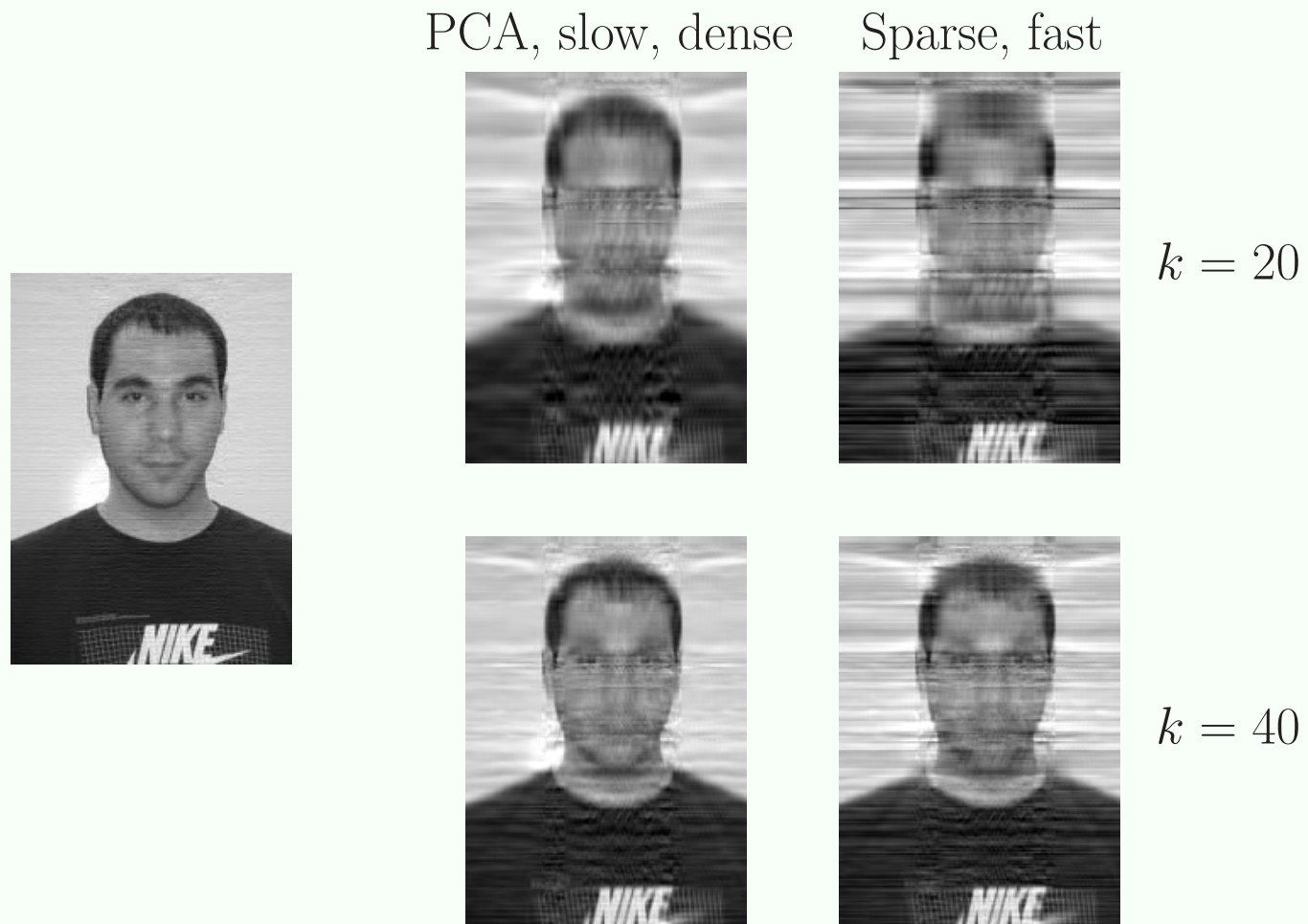


4 Clusters

**Theorem.** There is a subset of features of size  $O(\#clusters)$  which produces nearly the optimal partition (within a constant factor). One can quickly produce features with a log-approximation factor.

[BDM,2013]

# Fast Regression using Few Important Features



**Theorem.** Can find  $O(k)$  pure features which performs as well top- $k$  PCA-regression (additive error controlled by  $\|X - X_k\|_F / \sigma_k$ ).

[BDM,2013]

---

# The Proofs

All the algorithms use the sparsifier of  $V_k^T$  in [BDM,FOCS2011].

1. Choose columns of  $V_k^T$  to preserve its singular values.
2. Ensure that the selected columns preserve the structural properties of the objective with respect to the columns of  $X$  that are sampled.
3. Use dual set sparsification algorithms to accomplish (2).

# THANKS!

- **Data compression (PCA):**  
quick and reveals few important features
- **Unsupervised clustering:**  
quick and reveals few important features
- **Supervised Regression:**  
quick and reveals few important features



**Few features:** easy to interpret; better generalizers; faster computations.

