

# Sparse Principal Components and Subspaces

Concepts, Theory, and Computation

**Vincent Q. Vu**

Department of Statistics  
The Ohio State University

October 2, 2013

*This talk is based on joint work with...*



**Jing Lei**  
Carnegie Mellon U.



**Juhee Cho**  
U. Wisconsin-Madison



**Karl Rohe**  
U. Wisconsin-Madison

# Outline

- Background on PCA and high-dimensions
- Sparsity of the leading eigenvector
- Consistent estimation and minimax theory
- Sparse principal subspaces
- Computationally tractable estimation

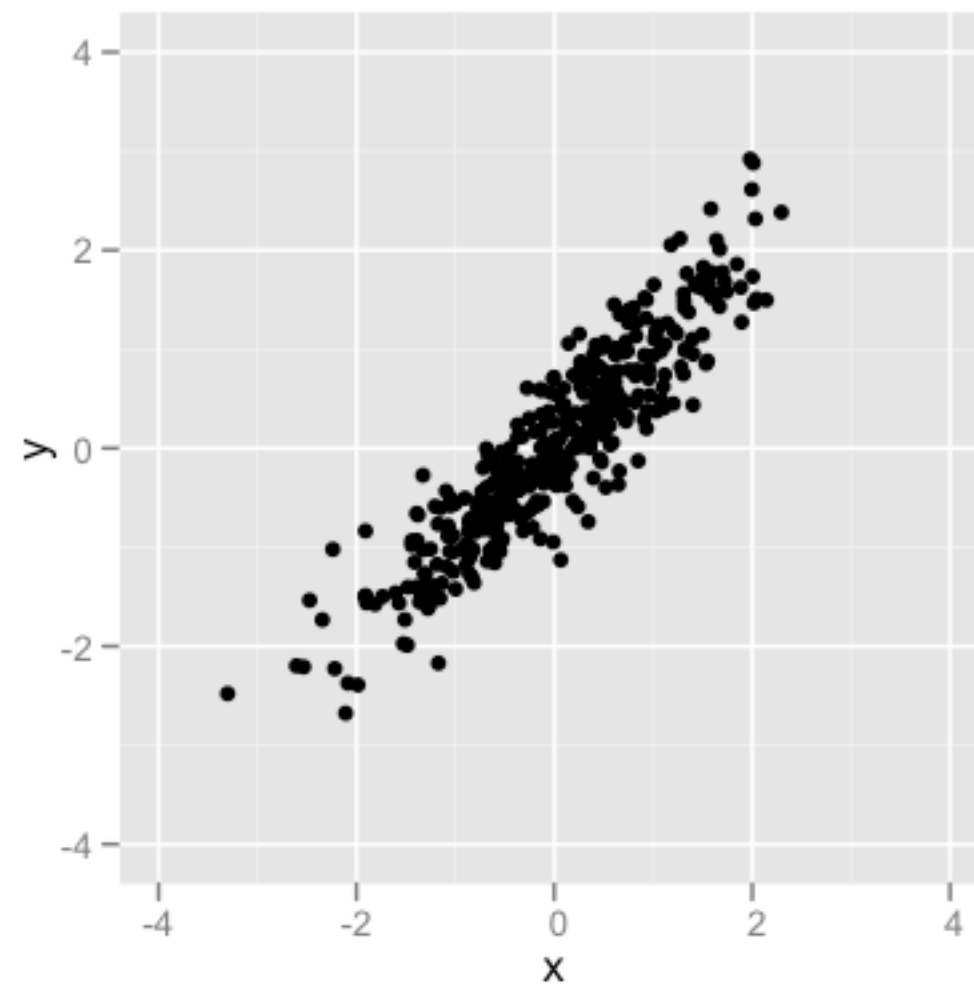
# High-Dimensional PCA

*“In many physical, statistical, and biological investigations it is desirable to represent a system of points in ... higher dimensioned space by the ‘best fitting’ straight line or plane.”*

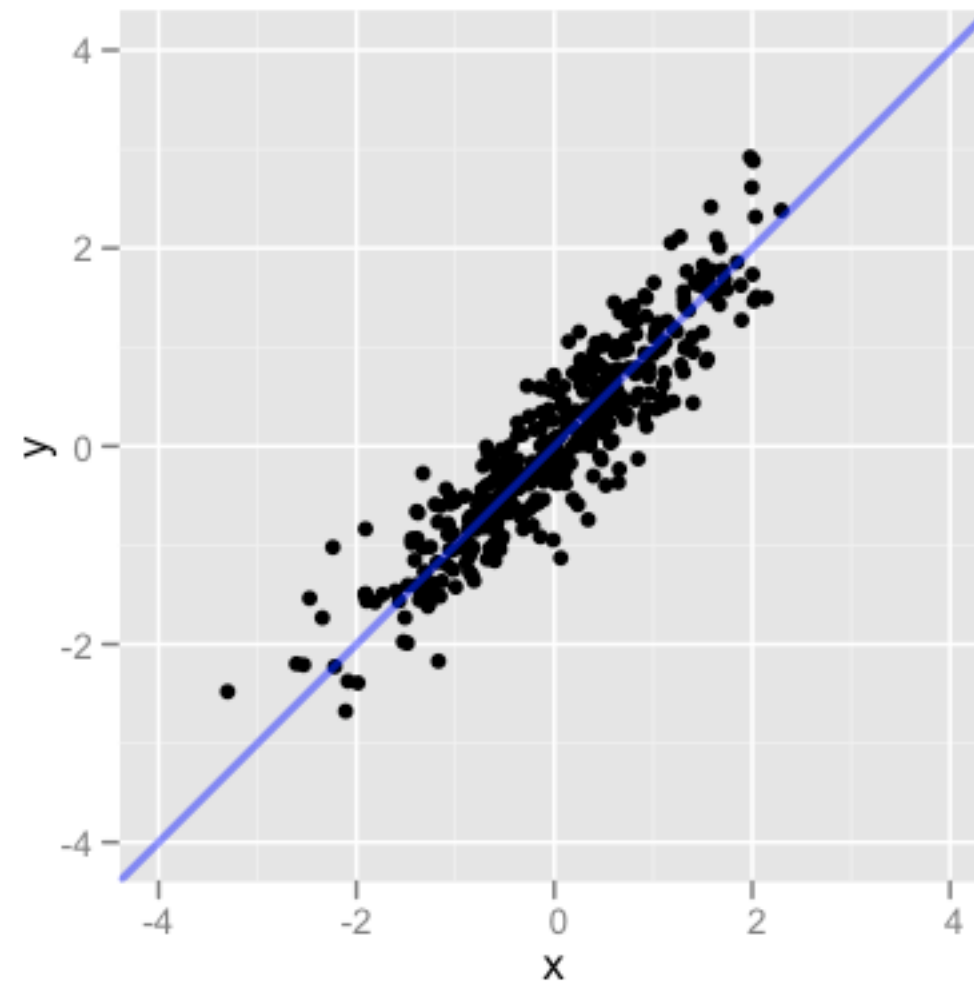
**– Karl Pearson (1901)**

On lines and planes of closest fit to systems of points in space

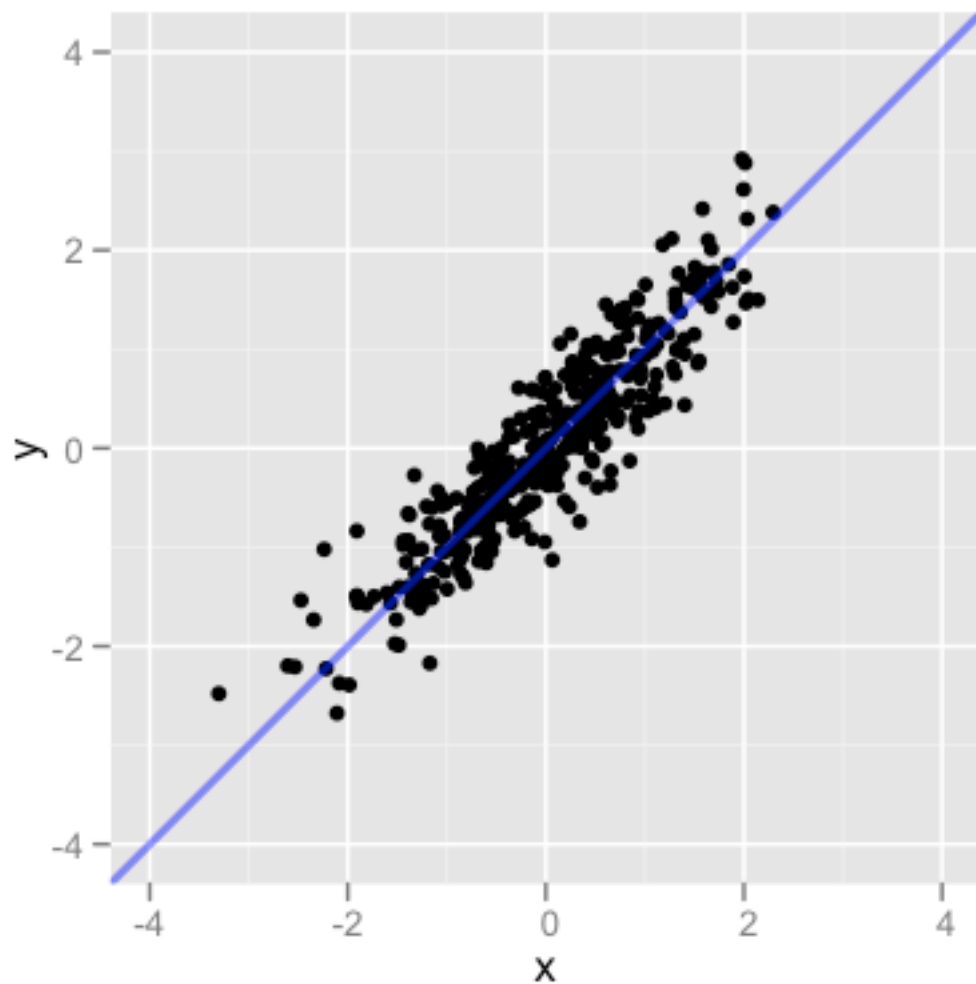
# Principal Components Analysis



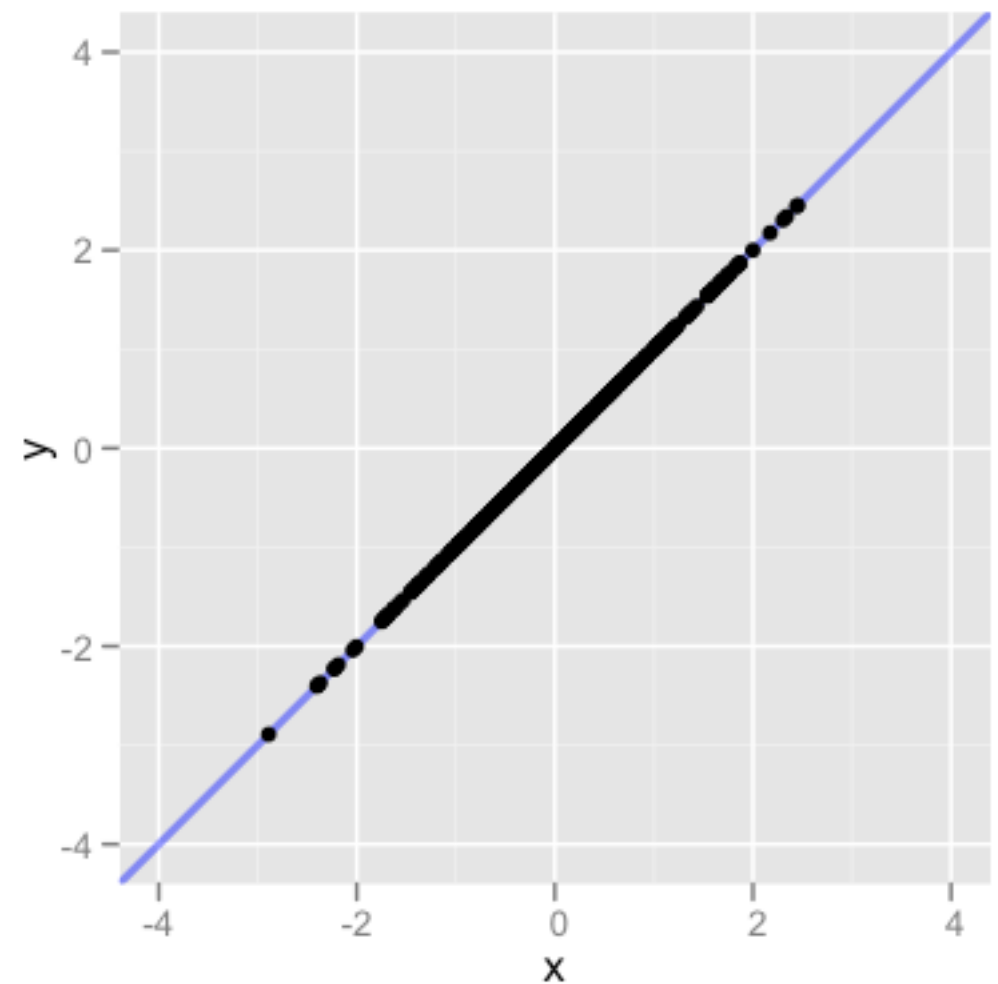
# Principal Components Analysis



# Principal Components Analysis

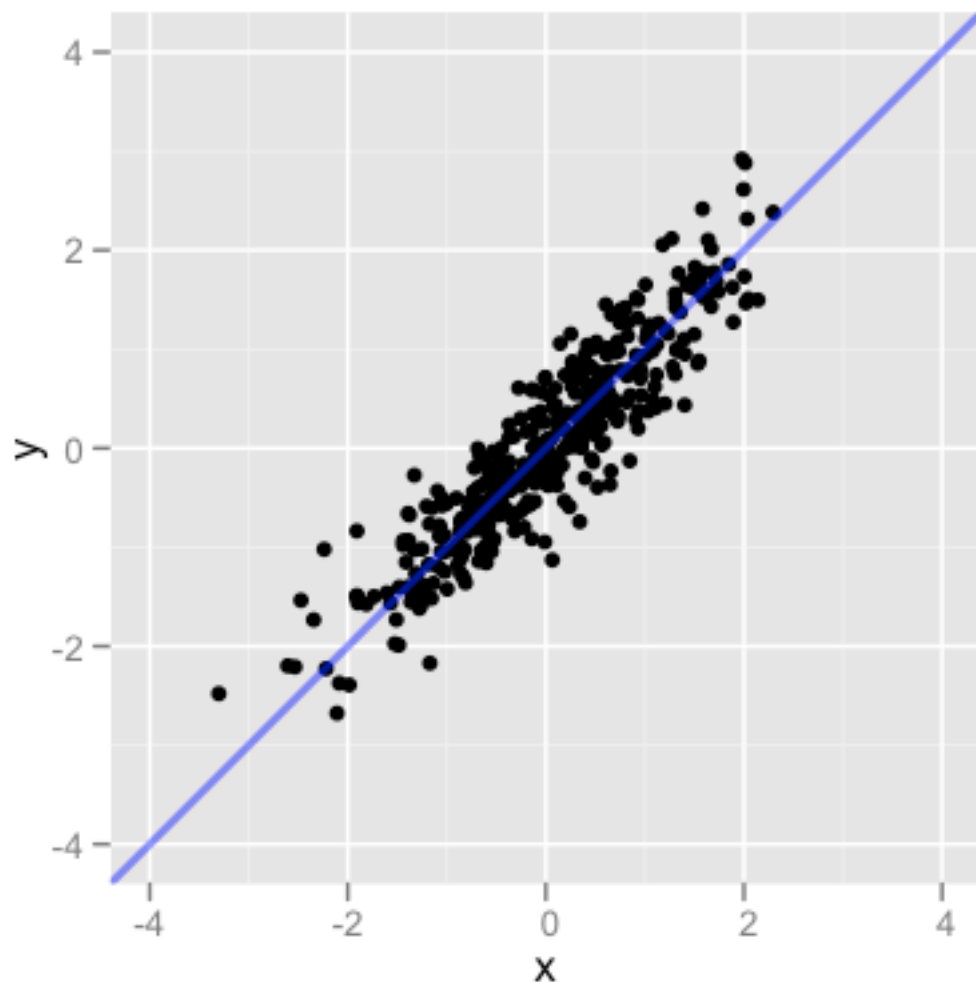


$\approx$



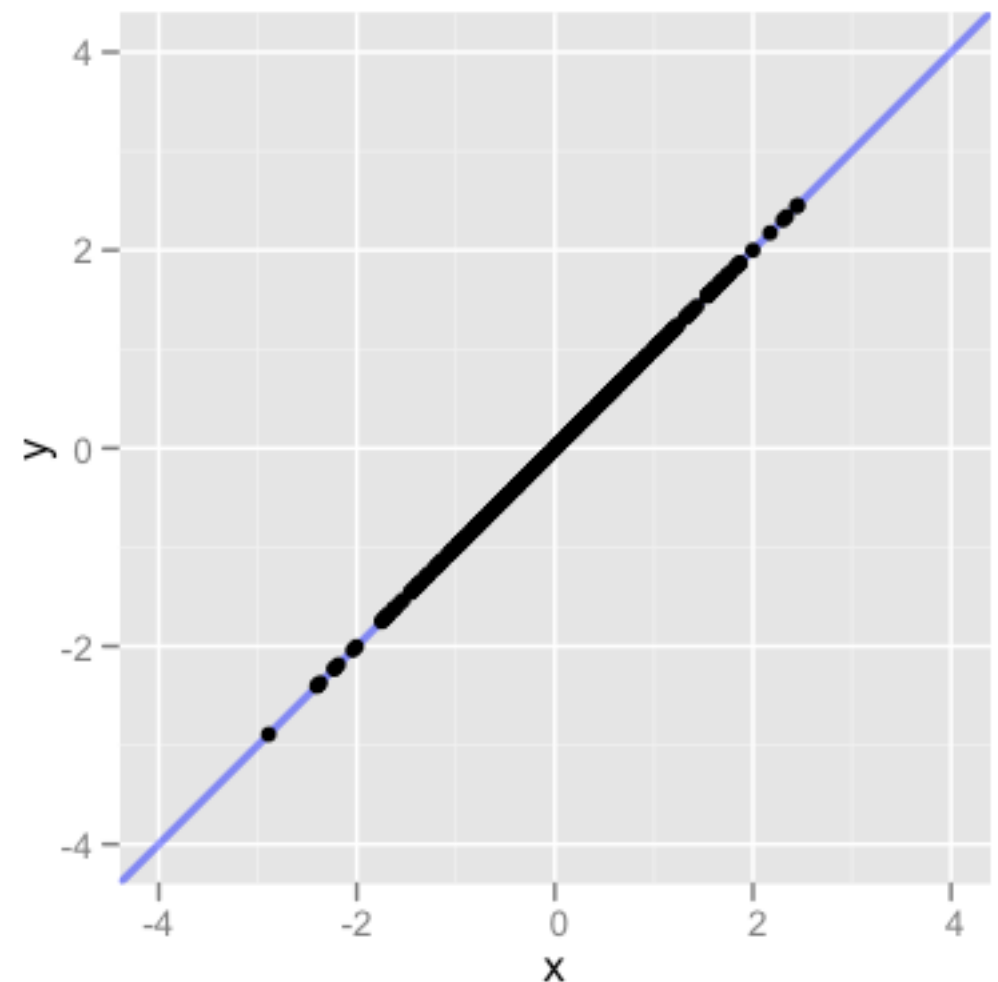


# Principal Components Analysis



original data

$\approx$



lower-dimensional  
projection

# Principal Components Analysis

- Suppose  $\{X_1, X_2, \dots, X_n\}$  is a dataset of i.i.d. observations on **p** variables
- **p** is *large*, so **PCA** could be used for dimension reduction

# Principal Components Analysis

“Optimal” dimension reduction is determined by **eigenvectors** of the **population covariance matrix**:

$$\Sigma \equiv \mathbb{E}(X X^T)$$

(assume  $\mathbb{E}X = 0$  to simplify presentation)

# Principal Components Analysis

## Eigendecomposition

$$\Sigma = V\Lambda V^T = \lambda_1 v_1 v_1^T + \cdots + \lambda_p v_p v_p^T$$

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p), \lambda_1 \geq \cdots \geq \lambda_p \geq 0 \quad (\text{eigenvalues})$$

$$V = (v_1, \dots, v_p), V^T V = I_p \quad (\text{eigenvectors})$$

# Principal Components Analysis

## Eigendecomposition

$$\Sigma = V\Lambda V^T = \lambda_1 v_1 v_1^T + \cdots + \lambda_p v_p v_p^T$$

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p), \lambda_1 \geq \cdots \geq \lambda_p \geq 0 \quad (\text{eigenvalues})$$

$$V = (v_1, \dots, v_p), V^T V = I_p \quad (\text{eigenvectors})$$

## Optimal projector

$$\Pi_1 = v_1 v_1^T \quad (\text{rank-1 projector})$$

$$\Pi_d = V_d V_d^T, V_d = (v_1, \dots, v_d) \quad (\text{rank-d projector})$$

# Classical PCA

Estimate eigenvectors by eigendecomposition  
of **sample covariance matrix**:

$$\hat{\Sigma} \equiv \frac{1}{n} \sum_{i=1}^n X_i X_i^T$$

# Classical PCA

Estimate eigenvectors by eigendecomposition of **sample covariance matrix**:

$$\hat{\Sigma} \equiv \frac{1}{n} \sum_{i=1}^n X_i X_i^T$$

Standard PCA estimator:

$$\hat{V}_d = (\hat{v}_1, \dots, \hat{v}_d), \quad \hat{\Pi}_d = \hat{V}_d \hat{V}_d^T$$

# Classical PCA

Estimate eigenvectors by eigendecomposition of **sample covariance matrix**:

$$\hat{\Sigma} \equiv \frac{1}{n} \sum_{i=1}^n X_i X_i^T$$

Standard PCA estimator:

$$\hat{V}_d = (\hat{v}_1, \dots, \hat{v}_d), \quad \hat{\Pi}_d = \hat{V}_d \hat{V}_d^T$$

Standard theory for **p** fixed and **n**  $\rightarrow \infty$ :

$$\hat{\Pi}_d \rightarrow \Pi_d \text{ a.s. if } \lambda_d - \lambda_{d+1} > 0$$



# High-Dimensional PCA: Challenges

# High-Dimensional PCA: Challenges

- In modern applications, e.g. neuroimaging, genetics:  $\mathbf{p} \approx \mathbf{n}$  and often  $\mathbf{p} \gg \mathbf{n}$

# High-Dimensional PCA: Challenges

- In modern applications, e.g. neuroimaging, genetics:  $\mathbf{p} \approx \mathbf{n}$  and often  $\mathbf{p} \gg \mathbf{n}$
- **Accuracy:** Standard PCA estimator can be inconsistent (Johnstone & Lu 2009):

$$\hat{v}_1^T v_1 \approx 0 \quad (\text{when } p/n \rightarrow c > 0, \lambda_1 - \lambda_2 \rightarrow c' > 0)$$

# High-Dimensional PCA: Challenges

- In modern applications, e.g. neuroimaging, genetics:  $\mathbf{p} \approx \mathbf{n}$  and often  $\mathbf{p} \gg \mathbf{n}$
- **Accuracy:** Standard PCA estimator can be inconsistent (Johnstone & Lu 2009):  
$$\hat{v}_1^T v_1 \approx 0 \quad (\text{when } p/n \rightarrow c > 0, \lambda_1 - \lambda_2 \rightarrow c' > 0)$$
- **Interpretability:** PCA difficult to interpret when estimated projector depends on many variables

# PCA in High-Dimensions

- Unconstrained estimation generally inconsistent
- Need additional **structural constraints** to have consistency
- What structural constraints make sense?

# Sparsity

# Sparsity

- Few variables have large effects – most others negligible

# Sparsity

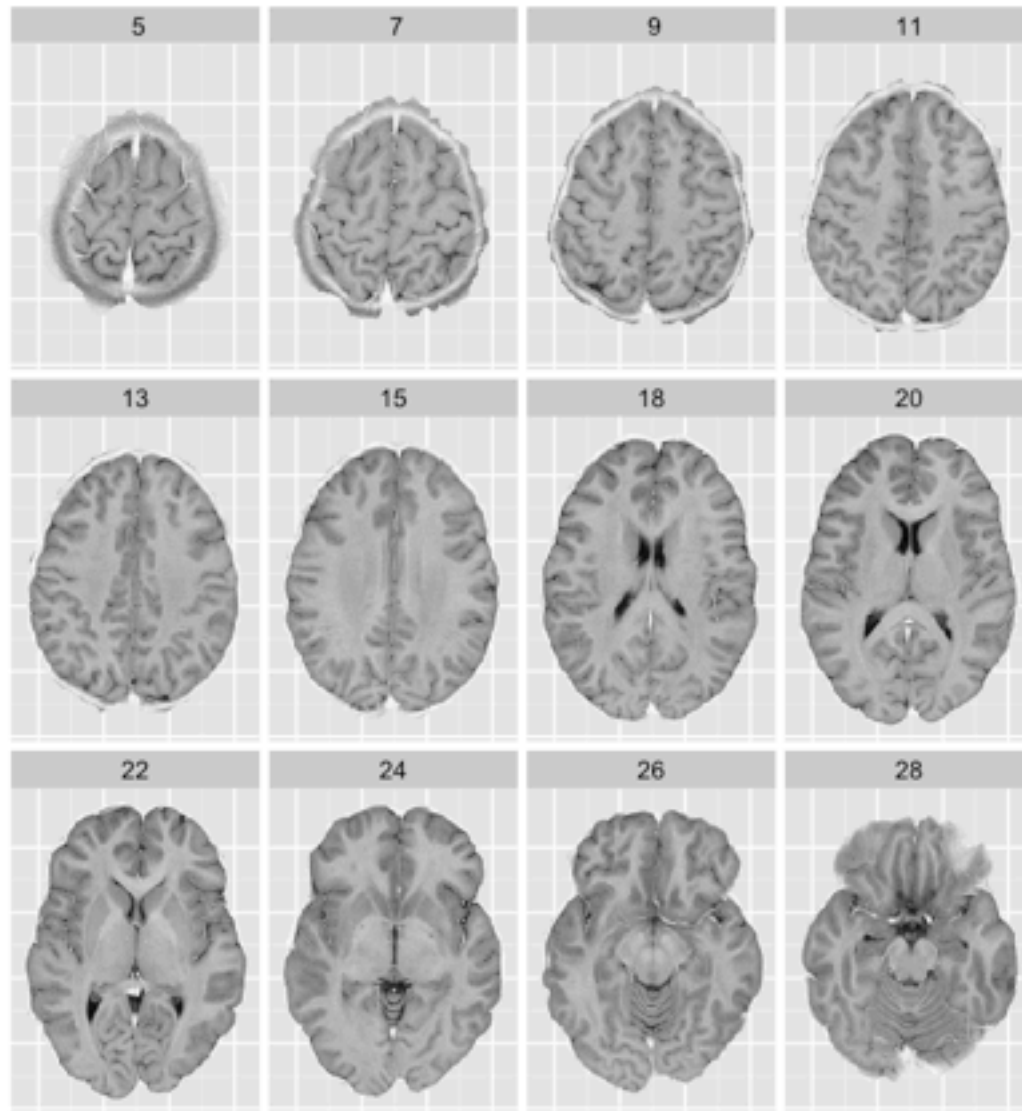
- Few variables have large effects – most others negligible
- Sometimes appropriate after *sparsifying* transformation
  - smoothness, localization, or periodicity correspond to sparsity in known bases: e.g. wavelets or Fourier



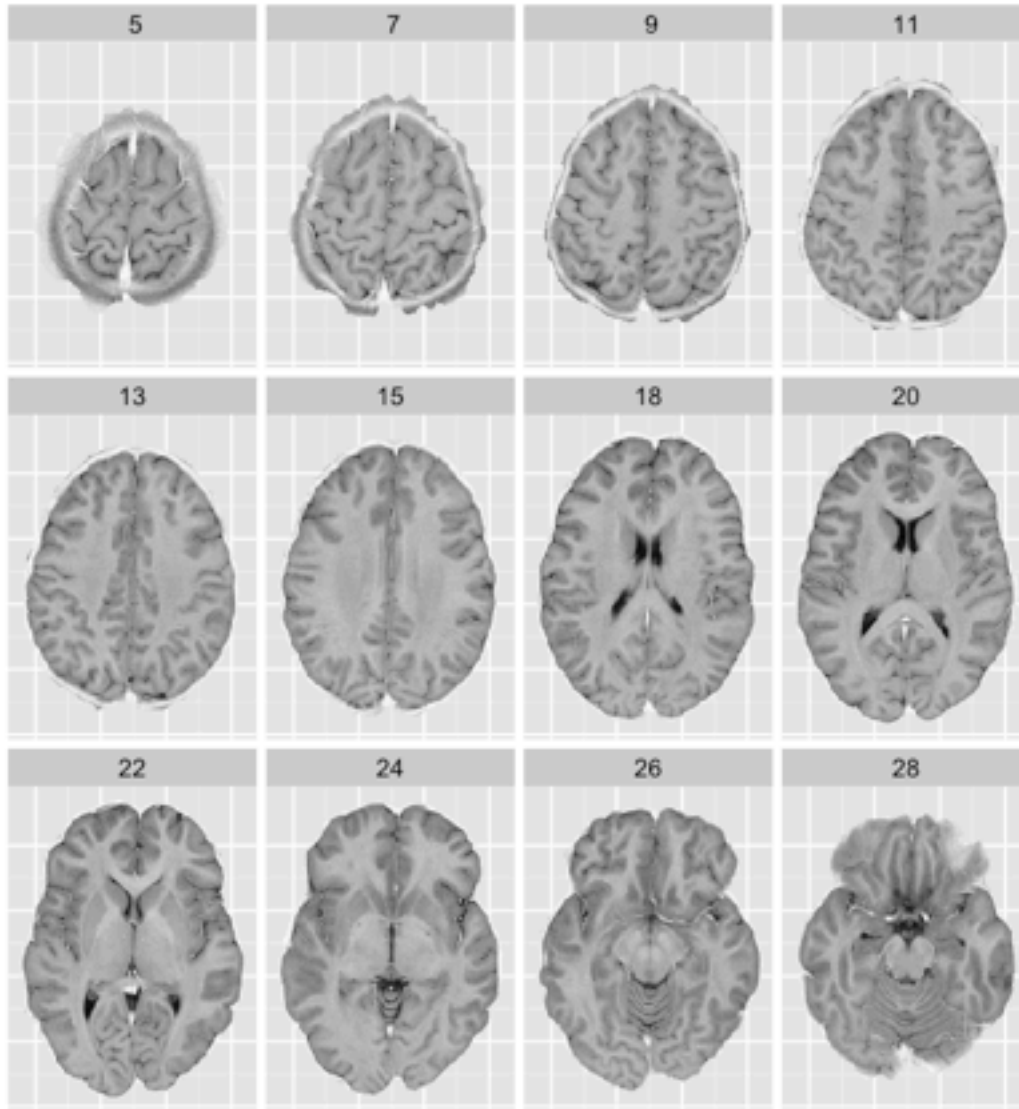
# Sparsity

- Few variables have large effects – most others negligible
- Sometimes appropriate after *sparsifying* transformation
  - smoothness, localization, or periodicity correspond to sparsity in known bases: e.g. wavelets or Fourier
- Can make estimation feasible **and** enhance interpretability

# Example: fMRI

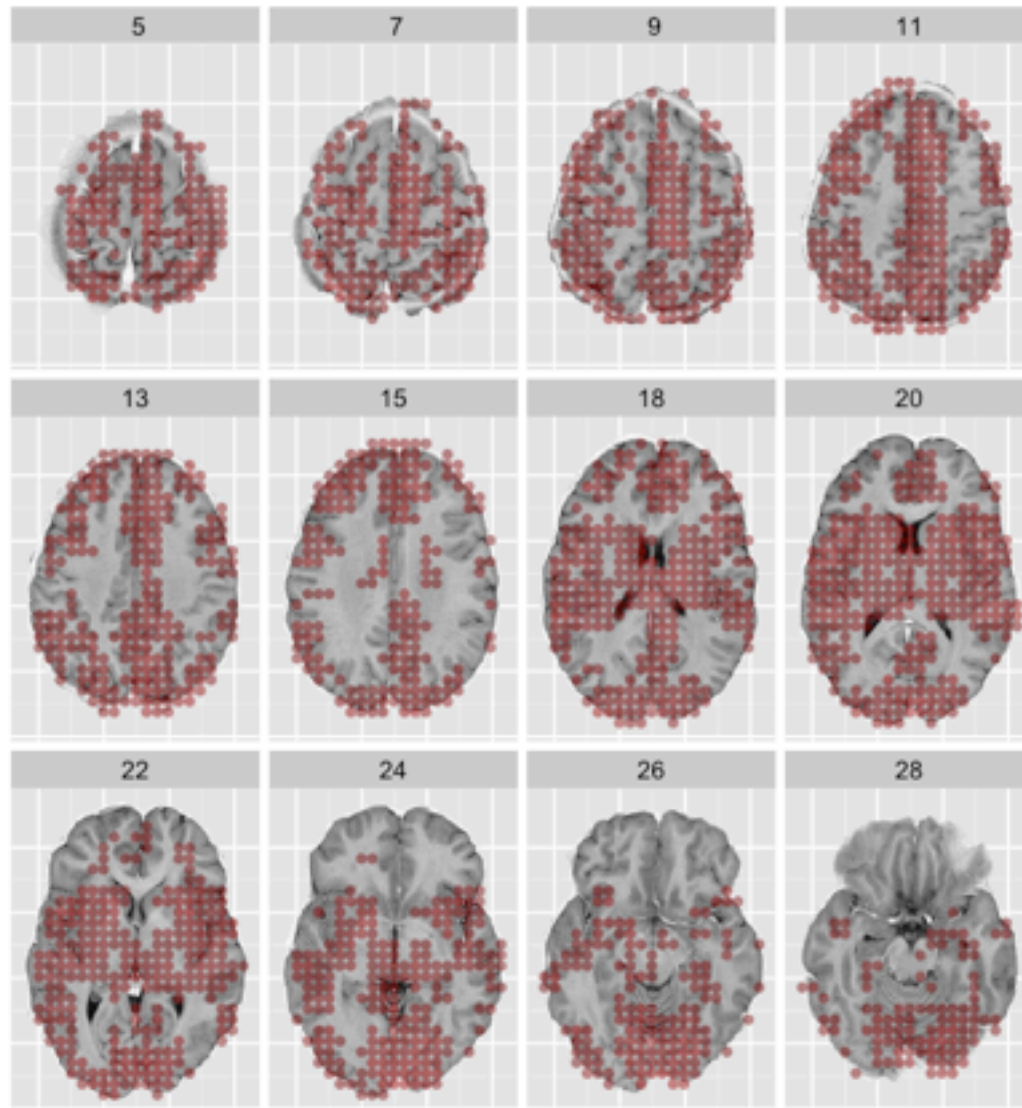


# Example: fMRI



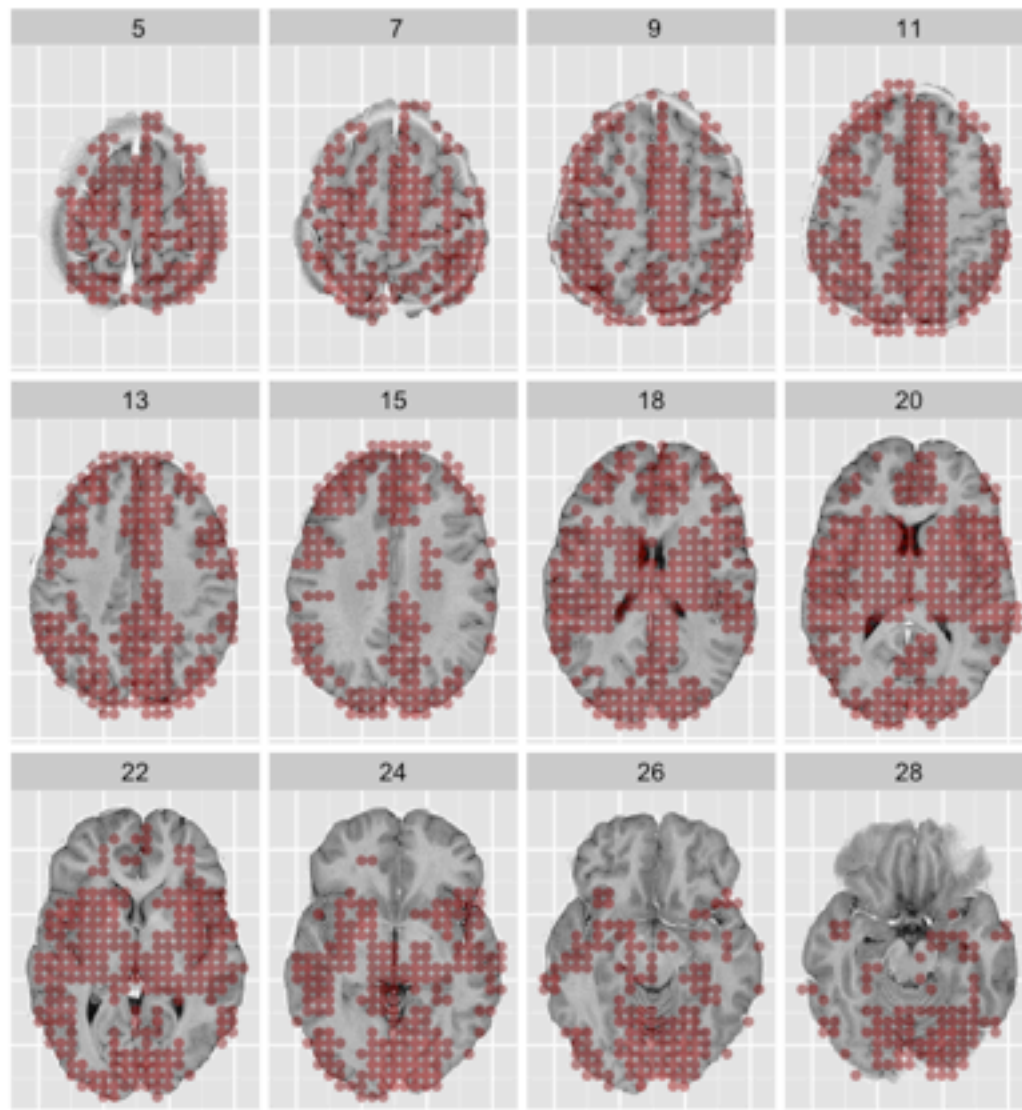
- $\mathbf{p} \approx 10,000 \sim 40,000$

# Example: fMRI



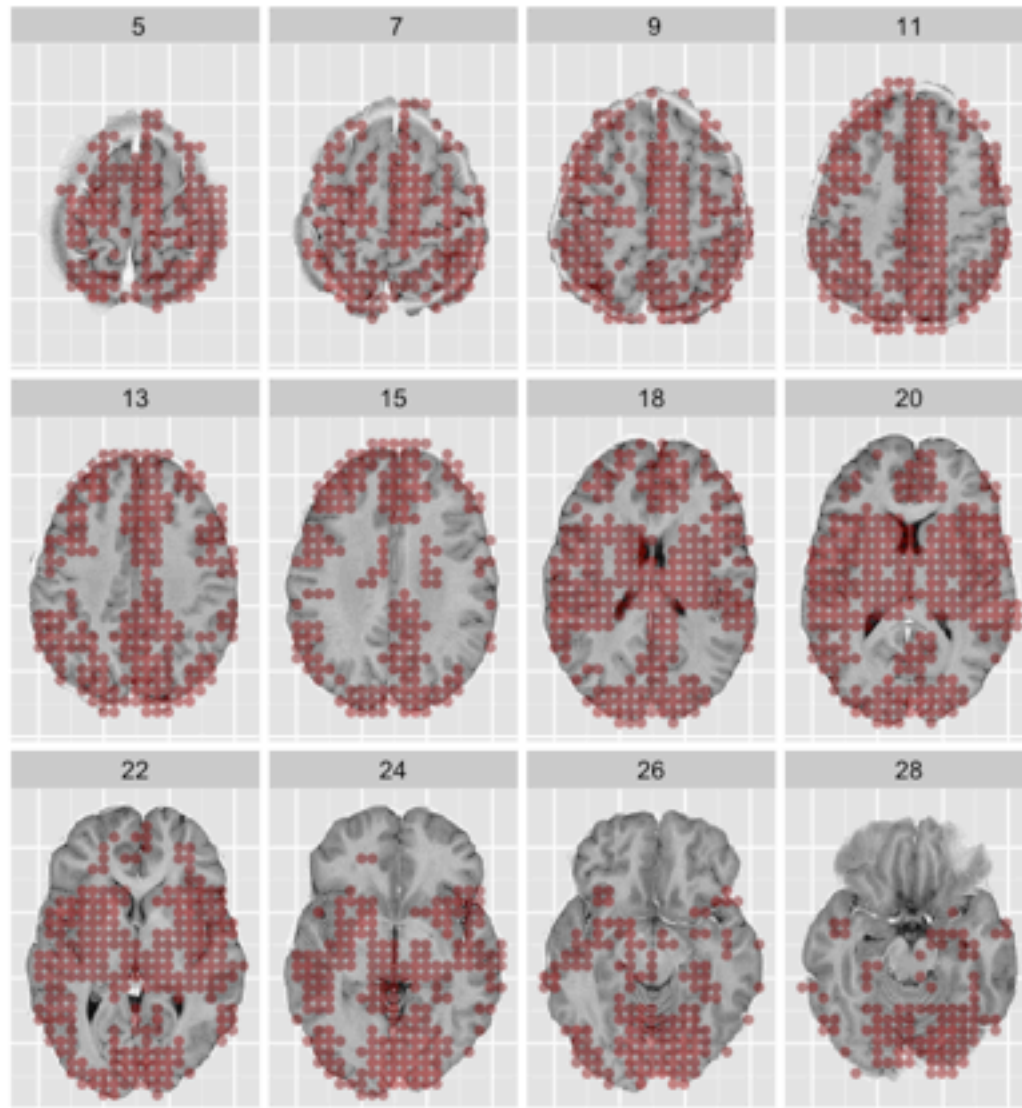
- $p \approx 10,000 \sim 40,000$
- “Interesting” activity spatially localized

# Example: fMRI



- $p \approx 10,000 \sim 40,000$
- “Interesting” activity spatially localized
- Locations not known in advance

# Example: fMRI



- $p \approx 10,000 \sim 40,000$
- “Interesting” activity spatially localized
- Locations not known in advance
- Sparsity = spatial localization

# Sparse PCA

- Many methods proposed over last 10 years:

*Joliffe, et al. (2003); Zou, et al. (2006); d'Aspremont, et al. (2007); Shen and Huang (2008); Johnstone and Lu (2009); Witten, et al. (2009); Journée et al. (2010); and many more*

# Sparse PCA

- Many methods proposed over last 10 years:

*Joliffe, et al. (2003); Zou, et al. (2006); d'Aspremont, et al. (2007); Shen and Huang (2008); Johnstone and Lu (2009); Witten, et al. (2009); Journée et al. (2010); and many more*

- Mostly algorithmic proposals



# Sparse PCA

- Many methods proposed over last 10 years:

*Joliffe, et al. (2003); Zou, et al. (2006); d'Aspremont, et al. (2007); Shen and Huang (2008); Johnstone and Lu (2009); Witten, et al. (2009); Journée et al. (2010); and many more*

- Mostly algorithmic proposals
- Few theoretical guarantees on statistical error – strong assumptions (spiked covariance model)

# Sparse PCA framework

( $d=1$  case)

# Sparse PCA framework

( $d=1$  case)

- **Goal:** Estimate principal eigenvector  $v_1$

# Sparse PCA framework

( $d=1$  case)

- **Goal:** Estimate principal eigenvector  $v_1$
- Identifiability condition:  $\lambda_1 > \lambda_2$

# Sparse PCA framework

( $d=1$  case)

- **Goal:** Estimate principal eigenvector  $v_1$
- Identifiability condition:  $\lambda_1 > \lambda_2$
- Sparsity assumption:  $\|v_1\|_0 \leq R_0 \ll p$

# Sparse PCA framework

( $d=1$  case)

- **Goal:** Estimate principal eigenvector  $v_1$
- Identifiability condition:  $\lambda_1 > \lambda_2$
- Sparsity assumption:  $\|v_1\|_0 \leq R_0 \ll p$
- Covariance assumption –
  - Spiked model:  $\Sigma = (\lambda_1 - \lambda_2)v_1v_1^T + \lambda_2I_p$
  - General model:  $\Sigma = \lambda_1v_1v_1^T + \dots + \lambda_pv_pv_p^T$

# Sparse PCA framework

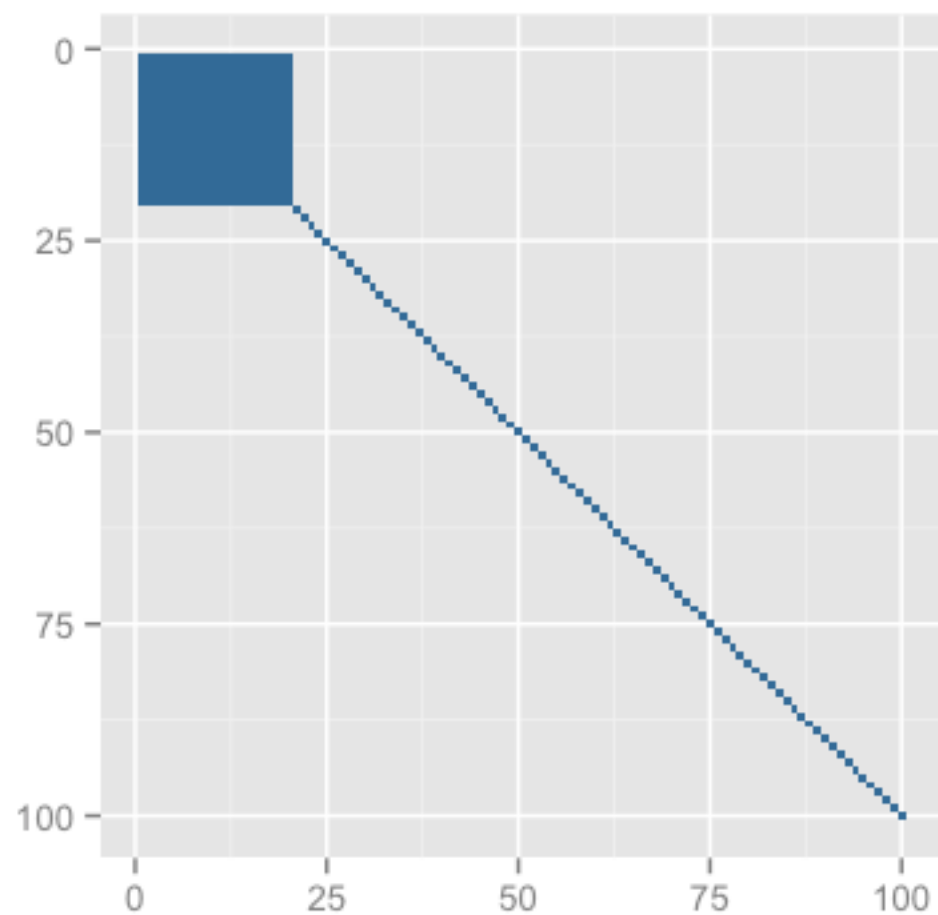
( $d=1$  case)

- **Goal:** Estimate principal eigenvector  $v_1$
- Identifiability condition:  $\lambda_1 > \lambda_2$
- Sparsity assumption:  $\|v_1\|_0 \leq R_0 \ll p$
- Covariance assumption –
  - Spiked model:  $\Sigma = (\lambda_1 - \lambda_2)v_1v_1^T + \lambda_2I_p$
  - General model:  $\Sigma = \lambda_1v_1v_1^T + \dots + \lambda_pv_pv_p^T$
- i.i.d. sub-Gaussian data

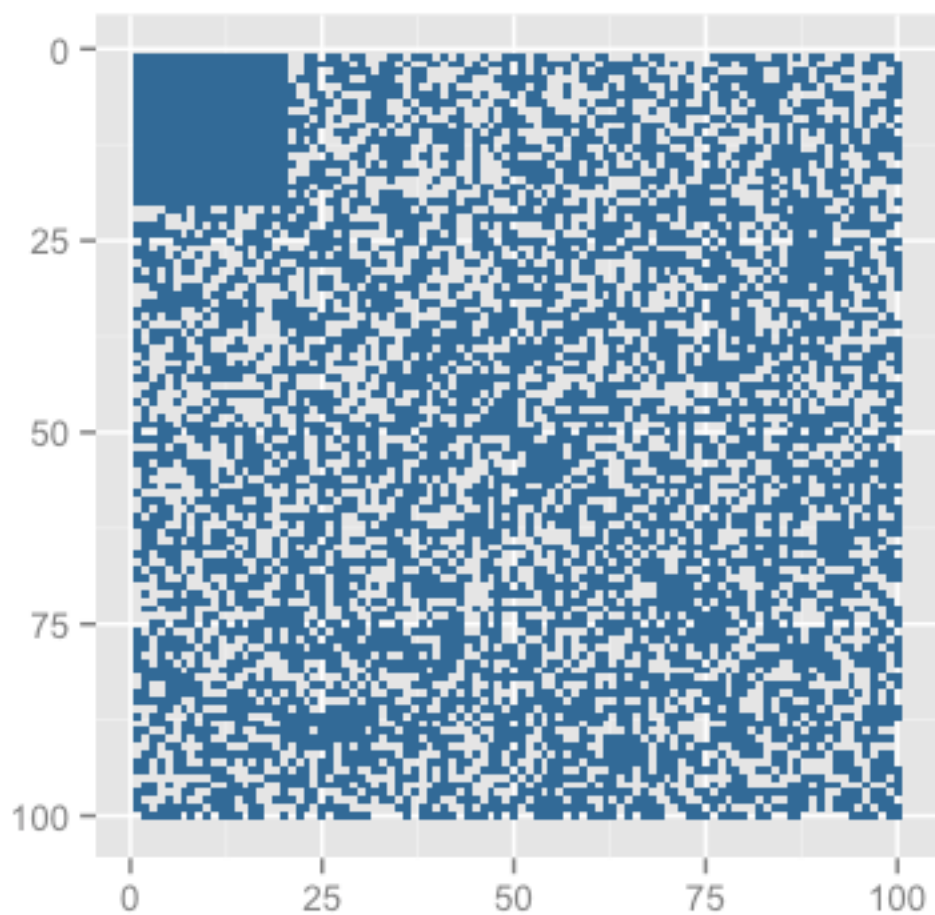
# ***Spiked Model vs General Model***

Locations of large nonzero entries:  $|\Sigma(i, j)| \geq 0.01$

$$\|v_1\|_0 = 20$$



Spiked model



General model



# How does sparsity help PCA?

# How does sparsity help PCA?

- **Questions**

- How well can we estimate  $v_1$  if sparsity assumed?
- How do we estimate  $v_1$  if sparsity assumed?

# How does sparsity help PCA?

- **Questions**

- How well can we estimate  $v_1$  if sparsity assumed?
- How do we estimate  $v_1$  if sparsity assumed?

- **Intuition** – Estimation is easy if

- $R_0$  small and  $\lambda_1 - \lambda_2$  large

# How does sparsity help PCA?

- **Questions**

- How well can we estimate  $v_1$  if sparsity assumed?
- How do we estimate  $v_1$  if sparsity assumed?

- **Intuition** – Estimation is easy if

- $R_0$  small and  $\lambda_1 - \lambda_2$  large

- Under **spiked model** (Johnstone & Lu 2003/9)  
give a consistent estimator of  $v_1$  when  $p/n \rightarrow c$

# Minimax theory

( $d = 1$  case)

# Minimax Framework

Find  $f(n, p, R_0, \lambda_1, \lambda_2)$  such that

$$f(n, p, R_0, \lambda_1, \lambda_2) \lesssim \sup_{\Sigma} \mathbb{E} \|\hat{v}_1 - v_1\|_2^2$$

for all estimators  $\hat{v}_1$  and a  
particular estimator  $\hat{v}_1$  such that

$$\sup_{\Sigma} \mathbb{E} \|\hat{v}_1 - v_1\|_2^2 \lesssim f(n, p, R_0, \lambda_1, \lambda_2)$$

for all covariance matrices in the sparse PCA  
model.

# Existing results for $d=1$

Under the spiked model Birnbaum et al. (2013) and Ma (2013) show that (**roughly**)

$$f(n, p, R_0, \lambda_1, \lambda_2) \sim \frac{R_0}{n} \cdot \frac{\lambda_1 \lambda_2}{(\lambda_1 - \lambda_2)^2} \cdot \log p$$

where the estimator is a thresholded power method (**up to  $\log n$  factor**)

# Existing results for $d=1$

Under the spiked model Birnbaum et al. (2013) and Ma (2013) show that (**roughly**)

$$f(n, p, R_0, \lambda_1, \lambda_2) \sim \frac{R_0}{n} \cdot \frac{\lambda_1 \lambda_2}{(\lambda_1 - \lambda_2)^2} \cdot \log p$$

where the estimator is a thresholded power method (**up to  $\log n$  factor**)

**Can we close the  $\log$  term gap?**  
**What about the general model?**



# Minimax Optimal Rate

**Theorem** (*V and Lei, 2013*)

Under the **general model**, the minimax error rate of estimating  $\mathbf{v}_1$  is

$$f(n, p, R_0, \lambda_1, \lambda_2) \asymp \frac{R_0}{n} \cdot \frac{\lambda_1 \lambda_2}{(\lambda_1 - \lambda_2)^2} \cdot \log p$$

where the **exact rate** is achieved by

$$\hat{v}_1 = \arg \max_{\|v\|_2=1, \|v\|_0 \leq R_0} v^T \hat{\Sigma} v$$

# Consequences

- **Good news**

# Consequences

- **Good news**

- Exact minimax rate in  $(n, p, R_0, \lambda_1, \lambda_2)$  for general model

# Consequences

- **Good news**

- Exact minimax rate in  $(n, p, R_0, \lambda_1, \lambda_2)$  for general model
- Extensions to  $L_q$  sparsity provide first consistency result for  $L_q$  constrained/penalized PCA (Joliffe et al 2003, Shen and Huang 2008, Witten et al 2009)

# Consequences

- **Good news**

- Exact minimax rate in  $(n, p, R_0, \lambda_1, \lambda_2)$  for general model
- Extensions to  $L_q$  sparsity provide first consistency result for  $L_q$  constrained/penalized PCA (Joliffe et al 2003, Shen and Huang 2008, Witten et al 2009)

- **Bad news**

# Consequences

- **Good news**

- Exact minimax rate in  $(n, p, R_0, \lambda_1, \lambda_2)$  for general model
- Extensions to  $L_q$  sparsity provide first consistency result for  $L_q$  constrained/penalized PCA (Joliffe et al 2003, Shen and Huang 2008, Witten et al 2009)

- **Bad news**

- Estimator is computationally intractable (**NP**-hard in **p**)

# Consequences

- **Good news**

- Exact minimax rate in  $(n, p, R_0, \lambda_1, \lambda_2)$  for general model
- Extensions to  $L_q$  sparsity provide first consistency result for  $L_q$  constrained/penalized PCA (Joliffe et al 2003, Shen and Huang 2008, Witten et al 2009)

- **Bad news**

- Estimator is computationally intractable (**NP**-hard in **p**)
- Estimation not possible if  $\lambda_1 \approx \lambda_2$

# Multiple Eigenvectors?

- Most sparse PCA methods only estimate single eigenvectors, and are extended to multiple eigenvectors by iterative deflation



# Multiple Eigenvectors?

- Most sparse PCA methods only estimate single eigenvectors, and are extended to multiple eigenvectors by iterative deflation
- Iterative deflation methods are heuristic and can be suboptimal (Mackey 2009)

# Multiple Eigenvectors?

- Most sparse PCA methods only estimate single eigenvectors, and are extended to multiple eigenvectors by iterative deflation
- Iterative deflation methods are heuristic and can be suboptimal (Mackey 2009)
- If  $\lambda_1 \approx \lambda_2$ , then it makes less sense to think about *distinct* eigenvectors

# Sparse Principal Subspaces

( $d \geq 1$  case)

# Sparse Principal Subspaces

- **Identifiability** – If  $\lambda_1 = \lambda_2 = \dots = \lambda_d$  then impossible to distinguish  $V_d$  and  $V_d Q$  from the data for any orthogonal  $Q$ .

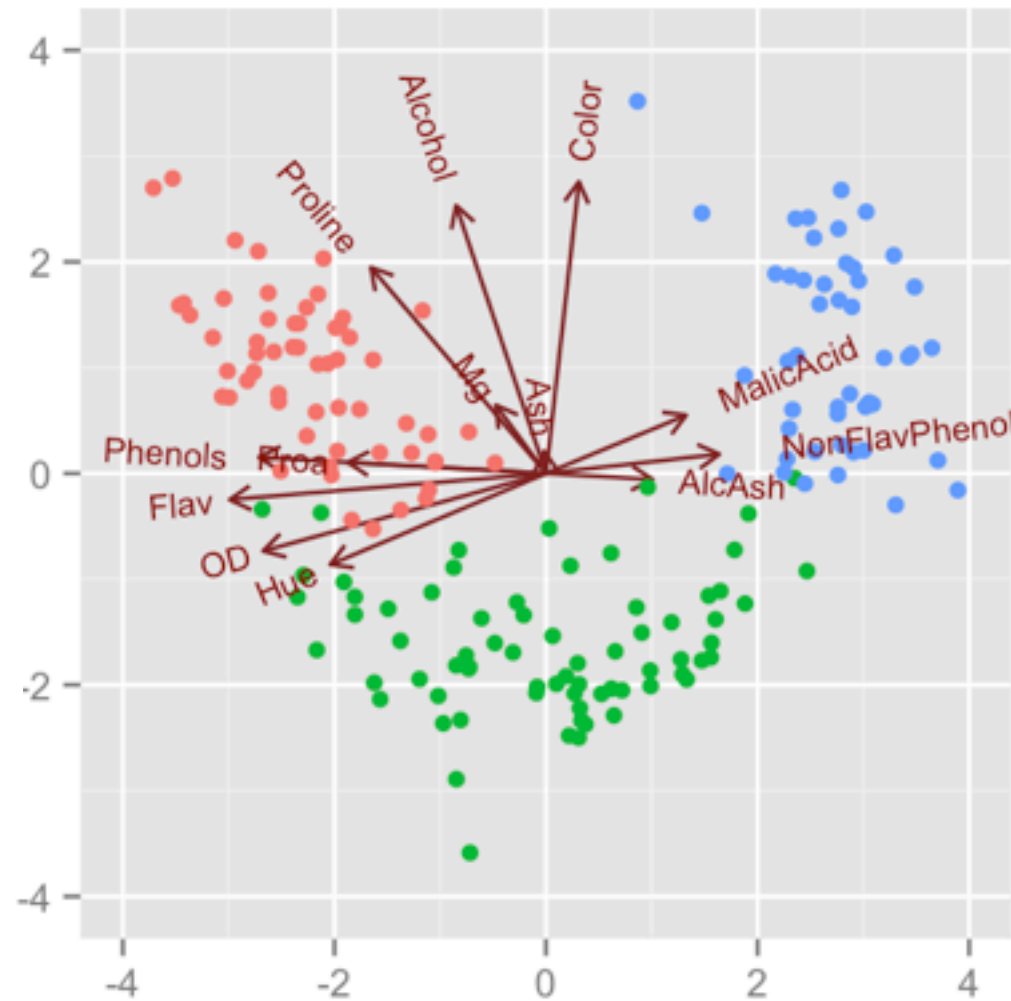
# Sparse Principal Subspaces

- **Identifiability** – If  $\lambda_1 = \lambda_2 = \dots = \lambda_d$  then impossible to distinguish  $V_d$  and  $V_d Q$  from the data for any orthogonal  $Q$ .
- **Sparsity** – How to extend sparsity to subspaces? Good notion of sparsity should be **rotation invariant**

# Sparse Principal Subspaces

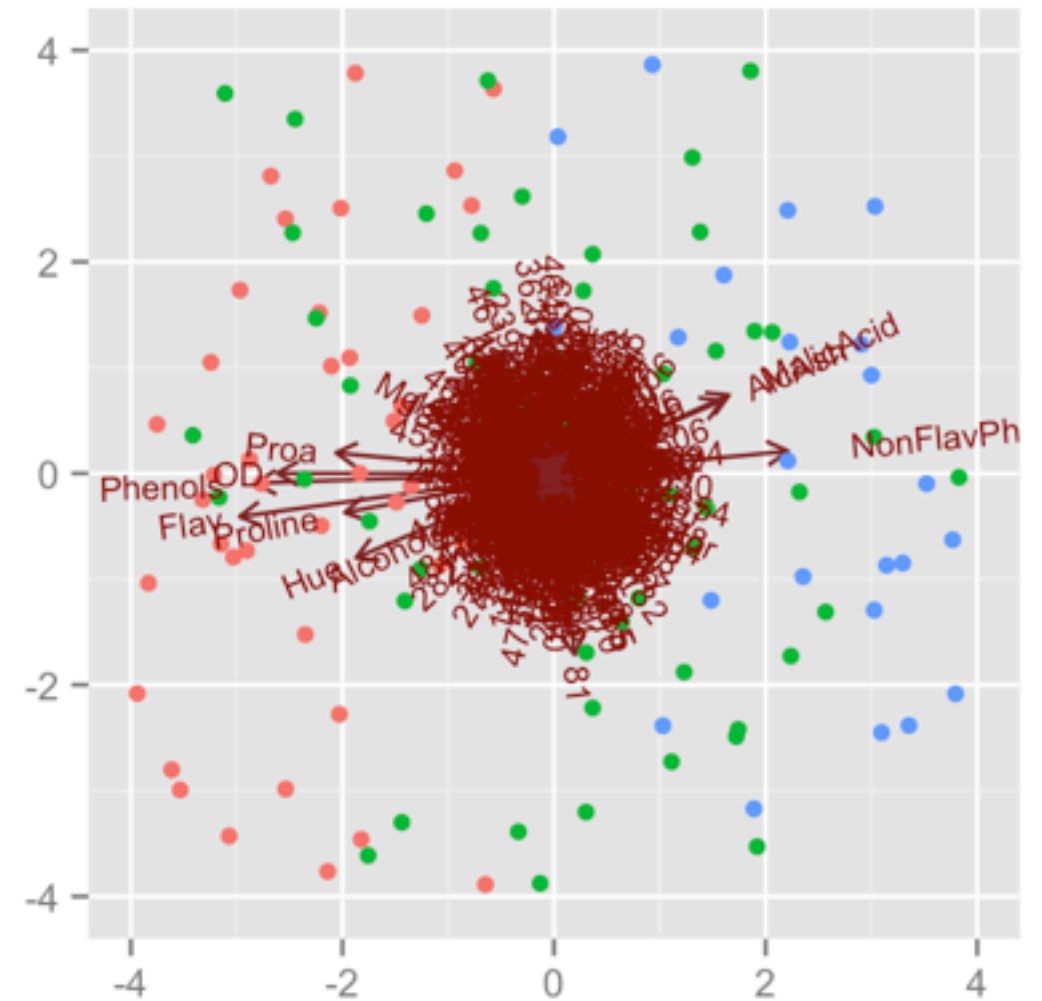
- **Identifiability** – If  $\lambda_1 = \lambda_2 = \dots = \lambda_d$  then impossible to distinguish  $V_d$  and  $V_d Q$  from the data for any orthogonal  $Q$ .
- **Sparsity** – How to extend sparsity to subspaces? Good notion of sparsity should be **rotation invariant**
- **Intuitively** – A subspace is sparse if its projector depends on a small number of variables

## Sparse



Projection depends on  
13 variables

## Not sparse



Projection depends on  
500 variables

# Row sparsity

- Matrix (2,0)-norm – for any  $p \times d$  matrix

$$\|V\|_{2,0} = \# \text{ of nonzero rows in } V$$

- Row sparsity:

$$\|V_d\|_{2,0} \leq R_0 \ll p, \quad V_d = (v_1, \dots, v_d)$$

- Row sparsity is rotation invariant – for any orthogonal  $Q$ :

$$\|V_d\|_{2,0} = \|V_d Q\|_{2,0}$$



# Subspace distance

- Measure distance between two subspaces with *canonical angles*

$$\|\sin \Theta(\hat{V}_d, V_d)\|_F^2 = \frac{1}{2} \|\hat{V}_d \hat{V}_d^T - V_d V_d^T\|_F^2$$

- Sum of squares of sines of canonical angles
- If  $d=1$ , equivalent to squared Euclidean distance

# Minimax Optimal Rate ( $d \geq 1$ )

**Theorem** (V and Lei, 2013)

Under the **general model**, the minimax error rate of estimating  $\mathbf{V}_d$  is

$$\min_{\hat{V}_d} \max_{\Sigma} \mathbb{E} \|\sin \Theta(\hat{V}_d, V_d)\|_F^2 \asymp R_0 \cdot \frac{\lambda_1 \lambda_{d+1}}{n(\lambda_d - \lambda_{d+1})^2} \cdot (d + \log p)$$

where the **exact rate** is achieved by

$$\hat{V}_d = \arg \max_{V^T V = I_d, \|V\|_{2,0} \leq R_0} \text{trace}(V^T \hat{\Sigma} V)$$

*Rate independently obtained by Cai et al. (2013) for Gaussian spiked model*

$$\begin{array}{l} \text{optimal} \\ \text{minimax} \\ \text{rate} \end{array} \asymp R_0 \cdot \frac{\lambda_1 \lambda_{d+1}}{n(\lambda_d - \lambda_{d+1})^2} \cdot (d + \log p)$$

optimal  
minimax  
rate

$$\asymp R_0 \cdot \frac{\lambda_1 \lambda_{d+1}}{n(\lambda_d - \lambda_{d+1})^2} \cdot (d + \log p)$$



# of active variables

optimal  
minimax  
rate

$$\asymp R_0 \cdot \frac{\lambda_1 \lambda_{d+1}}{n(\lambda_d - \lambda_{d+1})^2} \cdot (d + \log p)$$



# of active variables



effective noise variance

optimal  
minimax  
rate

$$\asymp R_0 \cdot \frac{\lambda_1 \lambda_{d+1}}{n(\lambda_d - \lambda_{d+1})^2} \cdot (d + \log p)$$



# of active variables



effective noise variance



estimation error

optimal  
minimax  
rate

$$\asymp R_0 \cdot \frac{\lambda_1 \lambda_{d+1}}{n(\lambda_d - \lambda_{d+1})^2} \cdot (d + \log p)$$



# of active variables



effective noise variance



estimation error



selection error

# Consequences

- **Good news**



# Consequences

- **Good news**
  - Exact minimax rate for **general model**

# Consequences

- **Good news**
  - Exact minimax rate for **general model**
  - **Sparsity enables** estimation in high-dimensions

# Consequences

- **Good news**
  - Exact minimax rate for **general model**
  - **Sparsity enables** estimation in high-dimensions
  - Extensions to  $L_q$  (weak) row sparsity

# Consequences

- **Good news**

- Exact minimax rate for **general model**
- **Sparsity enables** estimation in high-dimensions
- Extensions to  $L_q$  (weak) row sparsity

- **Bad news**

# Consequences

- **Good news**

- Exact minimax rate for **general model**
- **Sparsity enables** estimation in high-dimensions
- Extensions to  $L_q$  (weak) row sparsity

- **Bad news**

- Estimator is computationally intractable (**NP**-hard)

# Computation

# Sparse PCA computation

- Almost all formulations of sparse PCA involve:

# Sparse PCA computation

- Almost all formulations of sparse PCA involve:
  - Non-convex optimization problems with no statistical guarantees for local optima



# Sparse PCA computation

- Almost all formulations of sparse PCA involve:
  - Non-convex optimization problems with no statistical guarantees for local optima
  - Strong assumptions and sensitivity to initial value

# Sparse PCA computation

- Almost all formulations of sparse PCA involve:
  - Non-convex optimization problems with no statistical guarantees for local optima
  - Strong assumptions and sensitivity to initial value
- Is there a **polynomial time** method with **strong statistical guarantees** for the general model?

# Minimax Optimal Estimator – but NP-Hard

$$\arg \max_V \quad \text{trace}(V^T \hat{\Sigma} V) - \lambda \|V\|_{2,0}$$

subject to  $V^T V = I_d$

# Minimax Optimal Estimator – but NP-Hard

$$\arg \max_V \quad \text{trace}(V^T \hat{\Sigma} V) - \lambda \|V\|_{2,0}$$

subject to  $V^T V = I_d$

or, equivalently,

$$\arg \max_Z \quad \text{trace}(\hat{\Sigma} Z) - \lambda \|Z\|_{2,0}$$

subject to  $Z \in \{Z : Z \text{ is a rank-}d \text{ projector}\}$

# Difficulties

$$\arg \max_Z \quad \text{trace}(\hat{\Sigma}Z) - \lambda \|Z\|_{2,0}$$

subject to  $Z \in \{Z : Z \text{ is a rank-}d \text{ projector}\}$

# Difficulties

$$\arg \max_Z \quad \text{trace}(\hat{\Sigma}Z) - \lambda \|Z\|_{2,0}$$

subject to  $Z \in \{Z : Z \text{ is a rank-}d \text{ projector}\}$

- **Penalty is non-convex**

# Difficulties

$$\arg \max_Z \quad \text{trace}(\hat{\Sigma}Z) - \lambda \|Z\|_{2,0}$$

subject to  $Z \in \{Z : Z \text{ is a rank-}d \text{ projector}\}$

- Penalty is non-convex
- Constraint set is non-convex

# Difficulties

$$\arg \max_Z \quad \text{trace}(\hat{\Sigma}Z) - \lambda \|Z\|_{2,0}$$

subject to  $Z \in \{Z : Z \text{ is a rank-}d \text{ projector}\}$

- Penalty is non-convex
- Constraint set is non-convex
- Solution? **Use convex hulls!**



# Convex Relaxation

$$\arg \max_Z \quad \text{trace}(\hat{\Sigma}Z) - \lambda \|Z\|_{2,0}$$

subject to  $Z \in \{Z : Z \text{ is a rank-}d \text{ projector}\}$

# Convex Relaxation

$$\arg \max_Z \quad \text{trace}(\hat{\Sigma}Z) - \lambda \|Z\|_1$$

subject to  $Z \in \{Z : Z \text{ is a rank-}d \text{ projector}\}$

- Convex penalty function – entrywise L1 norm

# Convex Relaxation

$$\arg \max_Z \quad \text{trace}(\hat{\Sigma}Z) - \lambda \|Z\|_1$$

subject to  $Z \in \{Z : 0 \preceq Z \preceq I \text{ and } \text{trace}(Z) = d\}$

- Convex penalty function – entrywise L1 norm
- Convex constraint set – **The Fantope**

# Convex Relaxation

$$\arg \max_Z \quad \text{trace}(\hat{\Sigma}Z) - \lambda \|Z\|_1$$

subject to  $Z \in \{Z : 0 \preceq Z \preceq I \text{ and } \text{trace}(Z) = d\}$

- Convex penalty function – entrywise L1 norm
- Convex constraint set – **The Fantope**
- Amazing fact – Fantope is convex hull of rank- $d$  projectors (see Overton & Womersley, 1992)

# Fantope Projection and Selection

$$\begin{aligned} & \arg \max_Z \quad \text{trace}(\hat{\Sigma}Z) - \lambda \|Z\|_1 \\ & \text{subject to } Z \in \{Z : 0 \preceq Z \preceq I \text{ and } \text{trace}(Z) = d\} \end{aligned}$$

- Equivalent to a semidefinite program (SDP)
- $d=I$  case proposed by d'Aspremont et al. 2007
- Avoids orthogonality/deflation issues by ***directly*** estimating projector

# Fantope Projection and Selection

- Computable in polynomial time

# Fantope Projection and Selection

- Computable in polynomial time
- Alternating Direction Method of Multipliers (ADMM) algorithm has two main steps:

# Fantope Projection and Selection

- Computable in polynomial time
- Alternating Direction Method of Multipliers (ADMM) algorithm has two main steps:
  - Fantope Projection – have exact analytic solution



# Fantope Projection and Selection

- Computable in polynomial time
- Alternating Direction Method of Multipliers (ADMM) algorithm has two main steps:
  - Fantope Projection – have exact analytic solution
  - Element-wise soft-thresholding (***selection***)

# Guarantee for FPS

## **Theorem** (VLCR 2013)

Under the **general model**, assume the principal subspace is  $R_0$  row-sparse. If regularization parameter chosen appropriately\*, then FPS estimate  $\hat{Z}$  satisfies

$$\|\hat{Z} - V_d V_d^T\|_F^2 \lesssim R_0^2 \cdot \frac{\lambda_1 \lambda_{d+1}}{n(\lambda_d - \lambda_{d+1})^2} \cdot \log p$$

with high probability, regardless of its rank.

$$* \quad \lambda \asymp \sqrt{(\lambda_1 \lambda_{d+1} \log p)/n}$$

# FPS is *near*-Optimal

Recall minimax optimal rate ( $d=1$ ) vs FPS rate:

$R_0 \cdot \frac{\lambda_1 \lambda_{d+1}}{n(\lambda_d - \lambda_{d+1})^2} \cdot \log p$	$R_0^2 \cdot \frac{\lambda_1 \lambda_{d+1}}{n(\lambda_d - \lambda_{d+1})^2} \cdot \log p$
minimax optimal	FPS

FPS rate is off by factor of  $R_0$

# FPS is *near*-Optimal

Recall minimax optimal rate ( $d=1$ ) vs FPS rate:

$$R_0 \cdot \frac{\lambda_1 \lambda_{d+1}}{n(\lambda_d - \lambda_{d+1})^2} \cdot \log p \quad R_0^2 \cdot \frac{\lambda_1 \lambda_{d+1}}{n(\lambda_d - \lambda_{d+1})^2} \cdot \log p$$

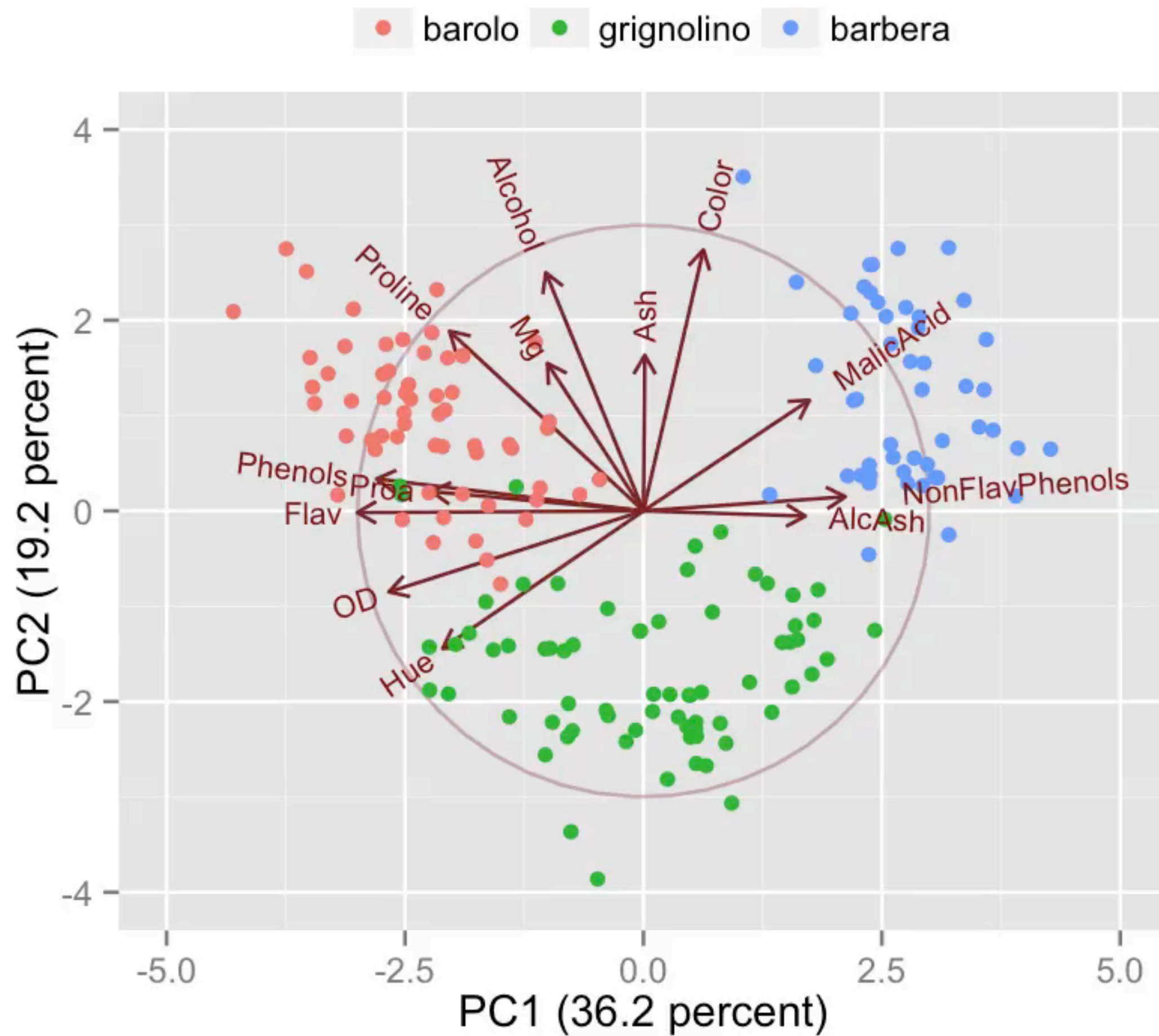
minimax optimal FPS

FPS rate is off by factor of  $R_0$

When  $d=1$ ,  **$R_0$**  factor maybe unavoidable for **any polynomial time** algorithm in a hypothesis testing framework (Berthet & Rigollet 2013)

# Small illustration

- Data on  **$n=178$**  wines grown over a decade in the same region of Italy
- 3 different cultivars: Barolo, Grignolino, Barbera
- Measurements on  **$p=13$**  constituents
- Will show  **$d=2$**  subspace estimated by FPS over a range of regularization parameter values

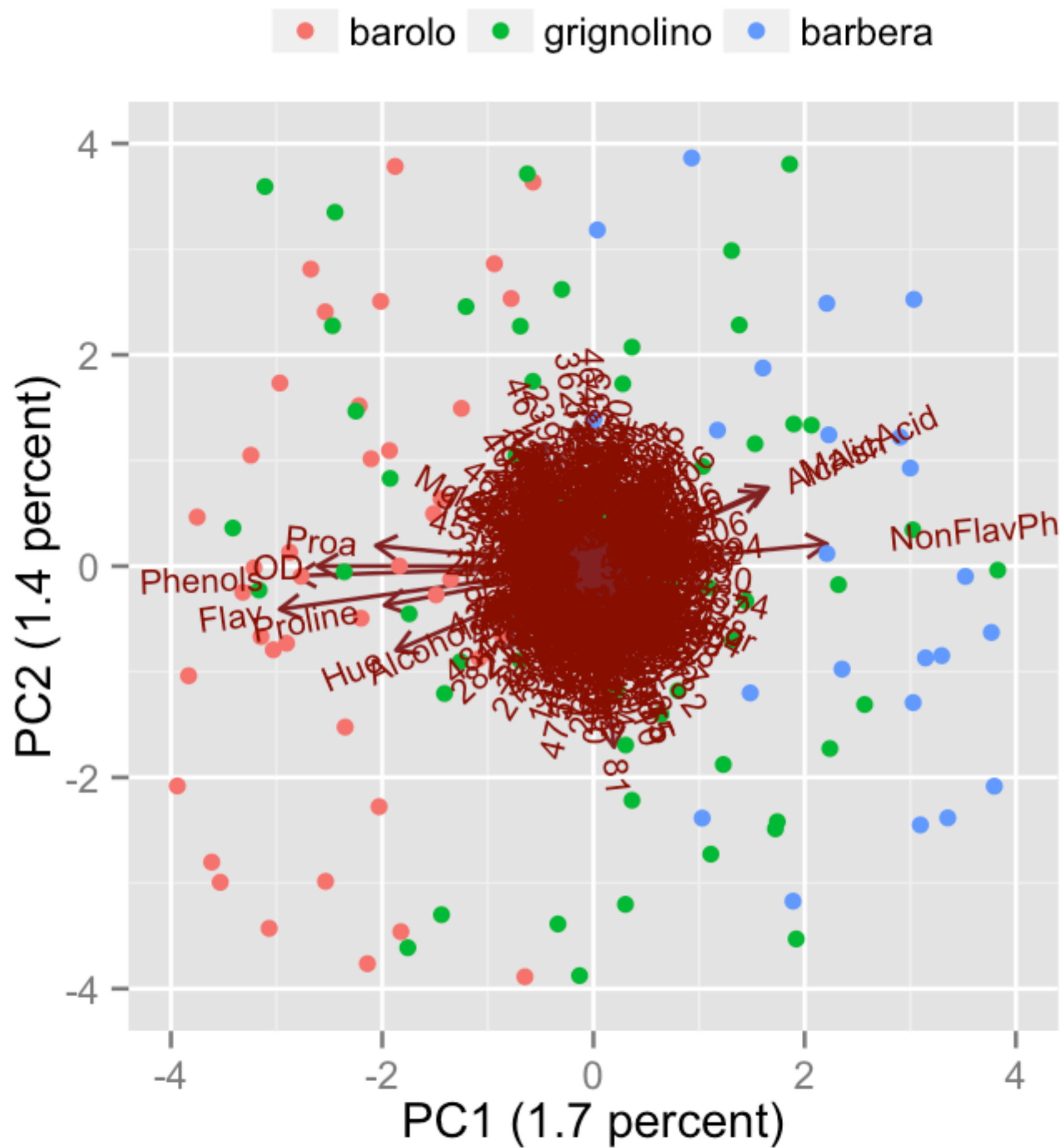


Movie not shown

# Synthetic illustration

- Dataset ***synthetically enlarged*** by adding 487 noise variables by randomly, independently copying and permuting the real variables – result **n=178, p=500**
- Does FPS recover the 13 real variables?
- Does FPS projection reveal the 3 clusters?





Movie not shown

# Summary

- Sparsity helps both estimation **accuracy** and **interpretation** of PCA in high dimensions
- Row sparsity is a rotation invariant notion of **subspace sparsity**
- Minimax rates reveal **gap** between computationally tractable and optimal (NP-hard) procedures
- Convex relaxation (FPS) is **near-optimal**

# Ongoing work

- Fast algorithm and computational insights for FPS to enable processing of larger scale data
- Is FPS rate optimal among polynomial time methods?

**Thank you!**

# References

- Vu & Lei (2012) “Minimax rates of estimation for sparse PCA in high dimensions.” AISTATS
- Vu & Lei (2013) “Minimax sparse principal subspace estimation in high dimensions.” Annals of Statistics, *to appear*
- Vu, Cho, Lei & Rohe (2013) “Fantope projection and selection.” *manuscript in preparation; preliminary report to appear in NIPS*