

Principal Component Analysis

Announcements

- Exams still being graded
- HW 6 due next Thursday
- Final project
 - Begin forming groups. Group size is 2-4.
 - Start thinking about the project topic.
 - See project guidelines – linked from Canvas home page
 - As part of HW7, you will be asked to submit a project proposal.

Dimension(ality) Reduction

- The transformation of data to a lower-dimensional space such that meaningful information is retained
- Feature extraction vs. feature selection

new features are
functions of old
features

select a few of the
existing features

- Motivations:

Reduce computational cost

Extract interpretable features

Improve performance of a supervised learning algorithm.

Eliminate redundant/noisy features to avoid overfitting

Feature selection: Find
important original features

Visualization

PCA

- Feature extraction
- Unsupervised
- Linear
- Objective: squared reconstruction error

PCA

- Given $\mathbf{x}_1, \dots, \mathbf{x}_n$, the idea behind PCA is to approximate

$$\mathbf{x}_i \approx \boldsymbol{\mu} + \mathbf{A}\boldsymbol{\theta}_i$$

$\begin{matrix} & \nearrow & \uparrow & \nwarrow \\ d \times 1 & & d \times k & k \times 1 \end{matrix}$

where

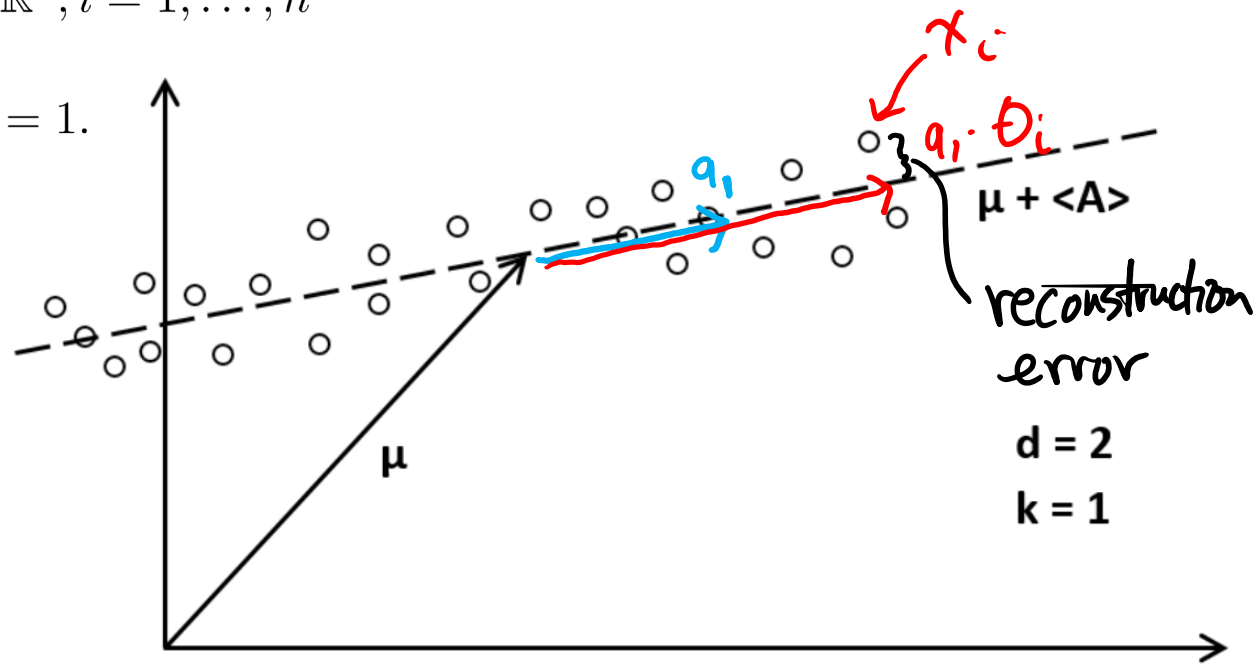
$$\boldsymbol{\mu} \in \mathbb{R}^d$$

$$\mathbf{A} \in \mathcal{A}_k := \{\mathbf{A} \in \mathbb{R}^{d \times k} \mid \mathbf{A}^T \mathbf{A} = \mathbf{I}_{k \times k}\}$$

$$\boldsymbol{\theta}_i \in \mathbb{R}^k, i = 1, \dots, n$$

- Example:** $d = 2, k = 1$.

$$\mathbf{A} = \begin{bmatrix} a_1 \end{bmatrix} \in \mathbb{R}^2$$

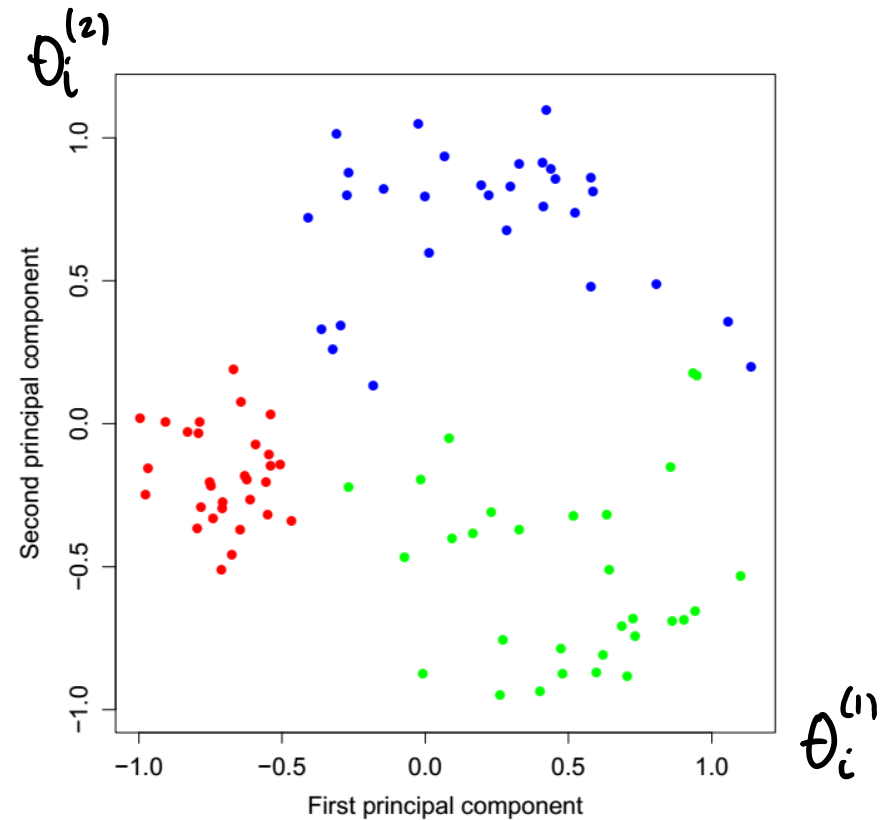
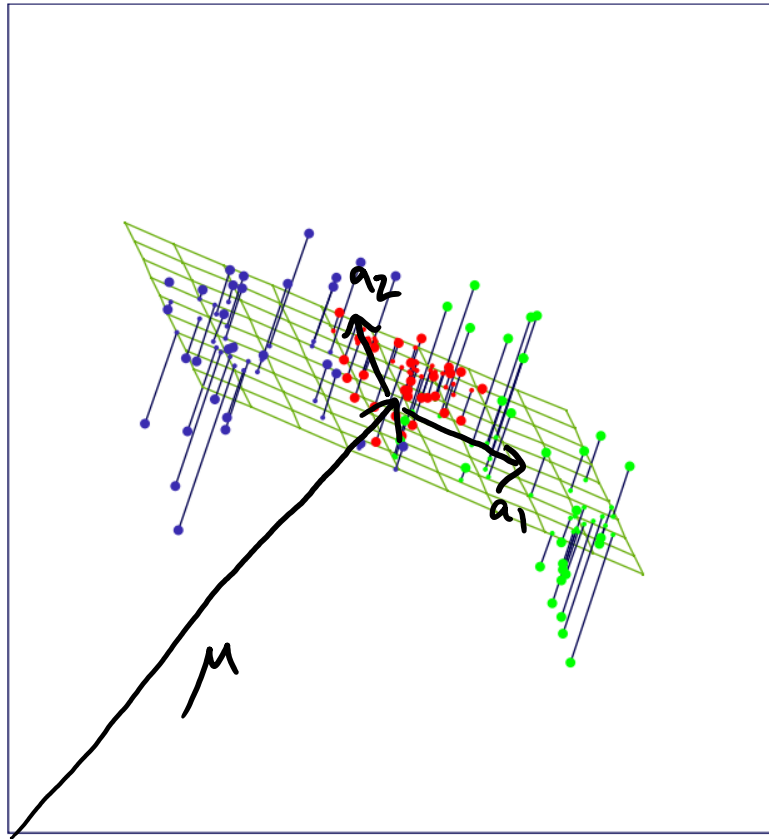


PCA

$$A = \begin{bmatrix} a_1 & a_2 \end{bmatrix}, \theta_i = \begin{bmatrix} \theta_i^{(1)} \\ \theta_i^{(2)} \end{bmatrix}$$

(3x2)

- Example: $d = 3, k = 2$



$$x_i \approx \mu + A \theta_i = \mu + \theta_i^{(1)} \cdot a_1 + \theta_i^{(2)} a_2$$

PCA

- Mathematically, we define $\mu, A, \theta_1, \dots, \theta_n$ to be the solution of

$$\begin{array}{l} \min \\ \mu \in \mathbb{R}^d, \theta_1, \dots, \theta_n \in \mathbb{R}^k \\ A \in \mathbb{R}^{d \times k}: A^T A = I \end{array} \quad \sum_{i=1}^n \|x_i - (\mu + A\theta_i)\|^2$$

- PCA gives the least squares rank- k linear approximation to the data set.
- The solution is given in terms of the spectral (or eigenvalue) decomposition of the sample covariance matrix

$$S = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

$$\bar{x} = \frac{1}{n} \sum x_i$$

Poll

What do we know about the sample covariance matrix

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T?$$

Select the best answer.

- (A) diagonal
- (B) positive semi-definite
- (C) positive definite
- (D) B and C

$$\begin{aligned} \mathbf{z}^T \mathbf{S} \mathbf{z} &= \frac{1}{n} \sum \underbrace{\mathbf{z}^T (\mathbf{x}_i - \bar{\mathbf{x}})}_{\text{scalar}} \underbrace{(\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{z}}_{\text{scalar}} \\ &= \frac{1}{n} \sum_{i=1}^n [\mathbf{z}^T (\mathbf{x}_i - \bar{\mathbf{x}})]^2 \\ &\geq 0 \end{aligned}$$

PCA Solution

- Denote the eigenvalue decomposition of S

$$S = U \Lambda U^T \quad \text{where} \quad \Lambda = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_d \end{bmatrix}$$
$$U = \begin{bmatrix} u_1 & \dots & u_d \end{bmatrix}, \quad U^T U = U U^T = I, \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0.$$

- A solution to PCA is

$$\mu = \bar{x}$$
$$A = \begin{bmatrix} u_1 & \dots & u_k \end{bmatrix}, \quad \theta_i = A^T (x_i - \bar{x})$$

Terminology and Concepts

- Principal components

$$\begin{aligned}\theta_i^{(j)} &= j^{\text{th}} \text{ principal component of } x_i \\ &= u_j^T (x_i - \bar{x})\end{aligned}$$

- Principal eigenvectors/directions

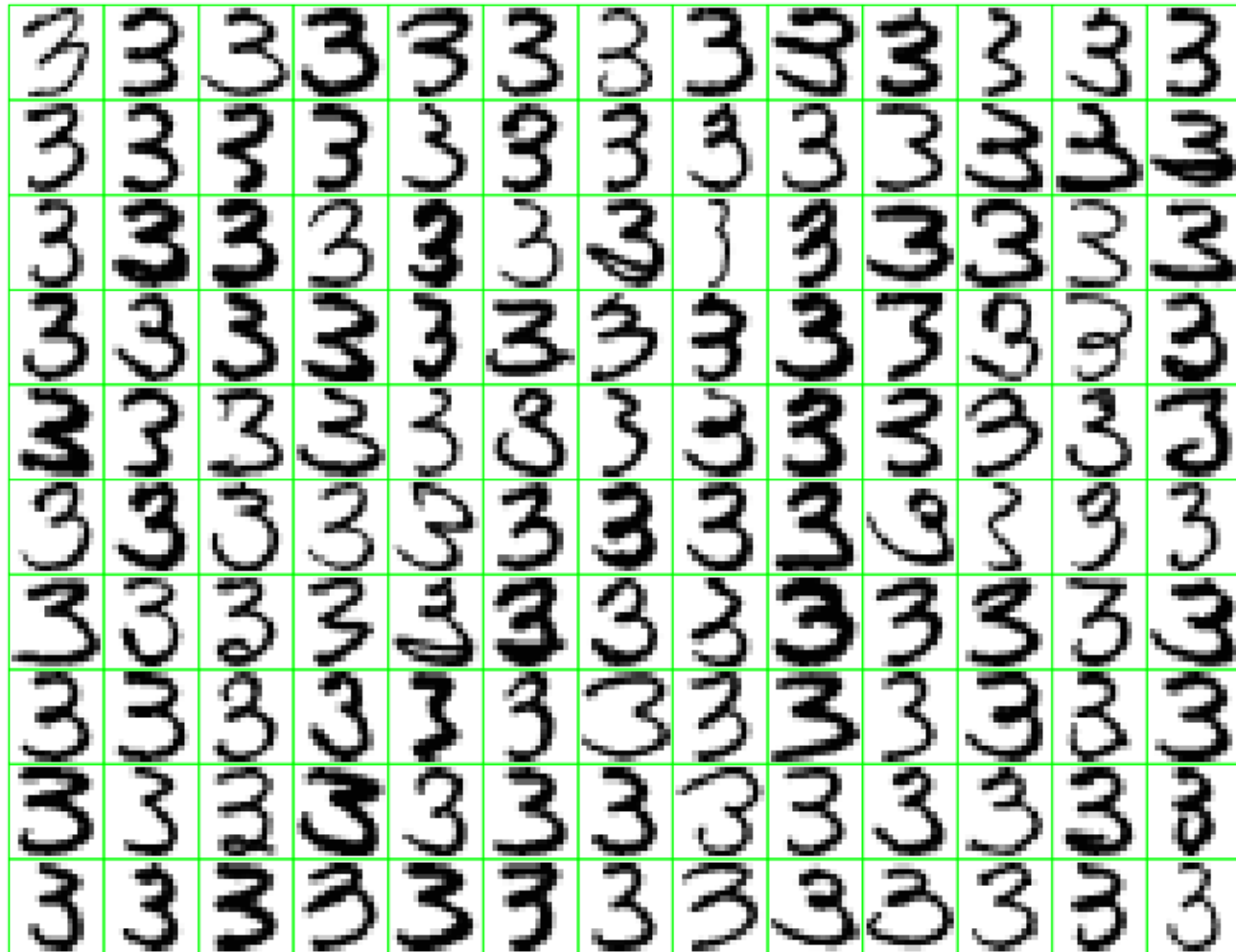
$$u_j = j^{\text{th}} \text{ principal eigenvector}$$

- Reconstruction of x_i

$$\hat{x}_i = \mu + A\theta_i = \bar{x} + AA^T(x_i - \bar{x})$$

Digits Example

- Training data



Digits Example

- Reconstruction

$k = 2$



$k = 10$



$k = 100$



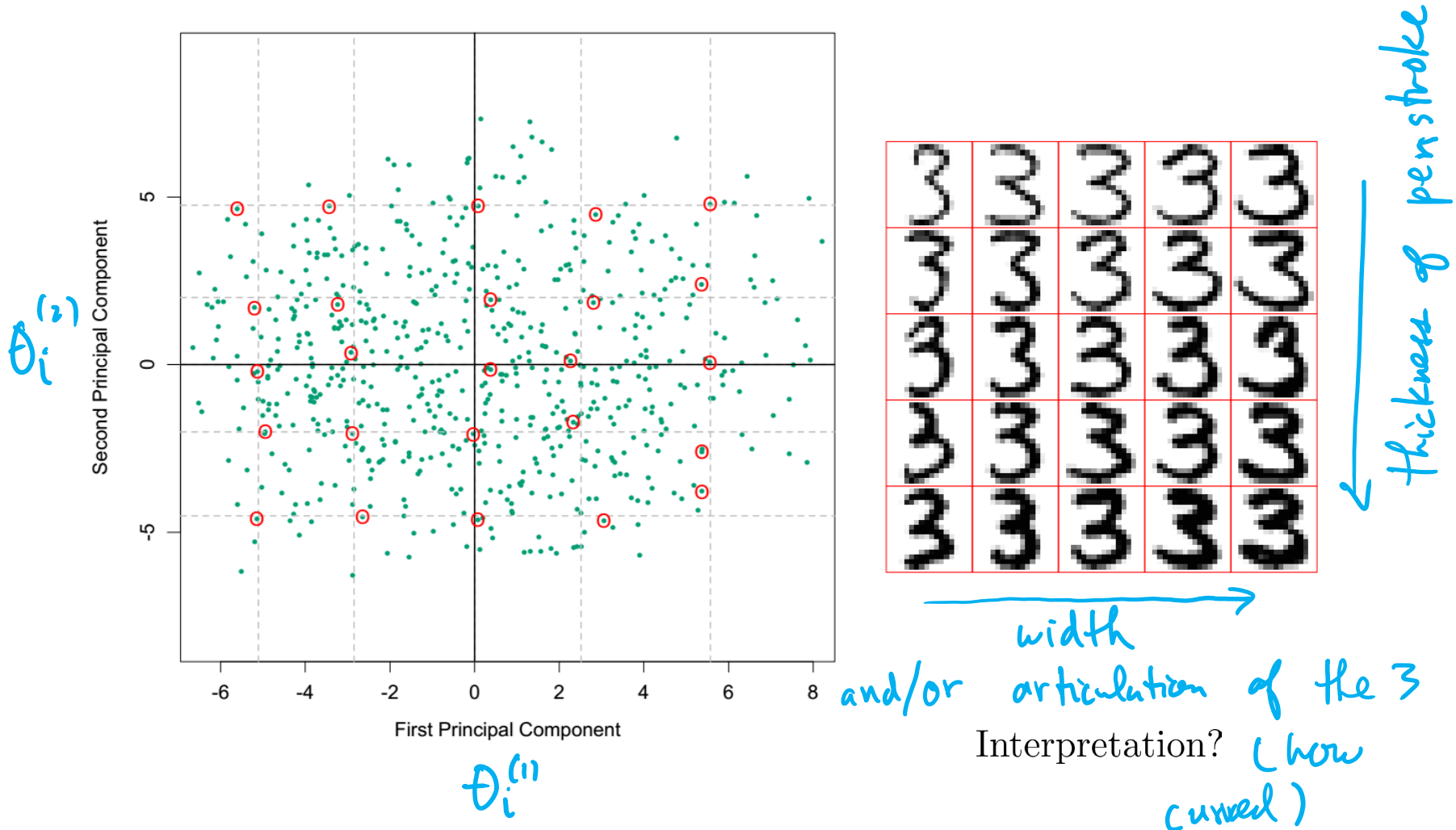
$k = 506$



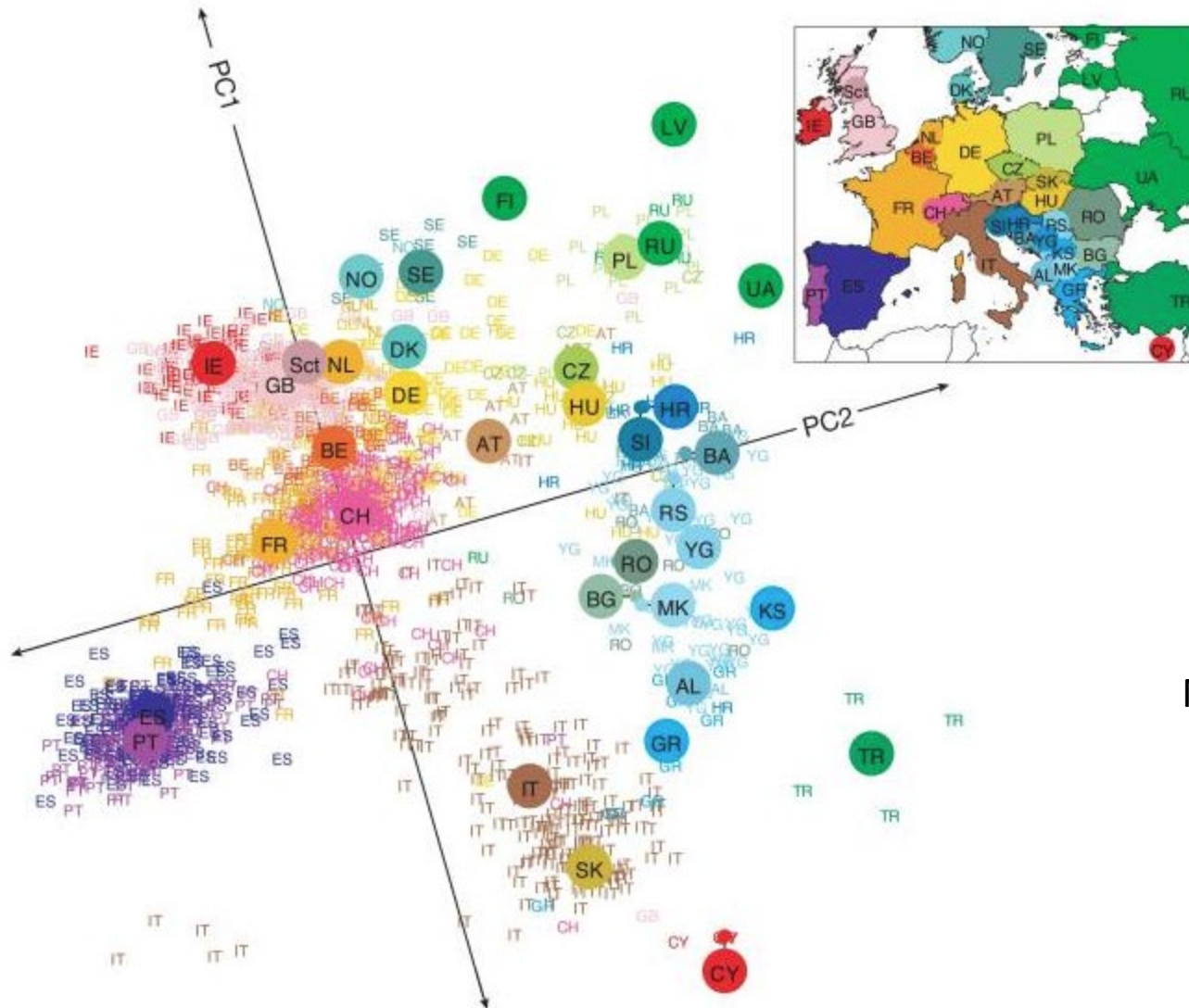
$$\hat{x}_i = \mu + \theta_i^{(1)} u_1 + \theta_i^{(2)} u_2 = \boxed{3} + \theta_i^{(1)} \boxed{3} + \theta_i^{(2)} \boxed{3}$$

Digits Example

First two PCs



Genotype Data



Novembre et al. (2008)

Poll

True or False: If we first center the data (that is, we subtract the mean \bar{x} from each x_i) before applying PCA, the principal components and principal eigenvectors do not change.

(A) True

(B) False

(C) Not enough information

$$\tilde{x}_i = x_i - \bar{x}$$

$$\tilde{S} = \frac{1}{n} \sum_{i=1}^n (\tilde{x}_i - 0)(\tilde{x}_i - 0)^T$$

$$= \frac{1}{n} \sum (x_i - \bar{x})(x_i - \bar{x})^T$$

$$= S$$

$$\tilde{\theta}_i = \tilde{A}^T (\tilde{x}_i - 0) = A^T (x_i - \bar{x}) = \theta_i$$

Connection to Projections

- Suppose $\bar{\mathbf{x}} = \mathbf{0}$ (a common preprocessing step).
- Then the rank- k approximation to \mathbf{x}_i is

$$\hat{\mathbf{x}}_i = \mu + A\boldsymbol{\theta}_i = \bar{\mathbf{x}} + \underbrace{AA^T}_{\text{"A(A^T A)^{-1}A^T"}}(\mathbf{x}_i - \bar{\mathbf{x}})$$

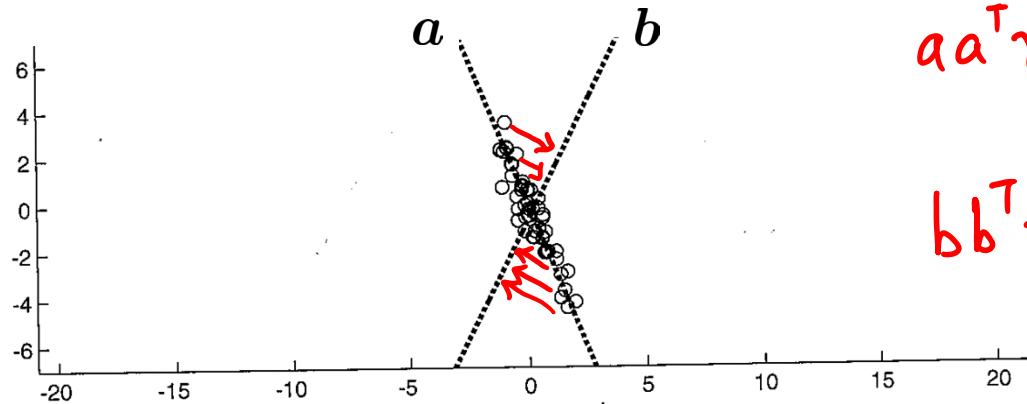
$$= A\boldsymbol{\theta}_i = \underbrace{AA^T}_{\text{projection onto colspan(A) =: <A>}}\mathbf{x}_i$$

$$= \sum_{j=1}^k \theta_i^{(j)} \mathbf{u}_j$$

" $A(A^T A)^{-1}A^T$
 \uparrow
 projection onto
 $\text{colspan}(A) =: \langle A \rangle$

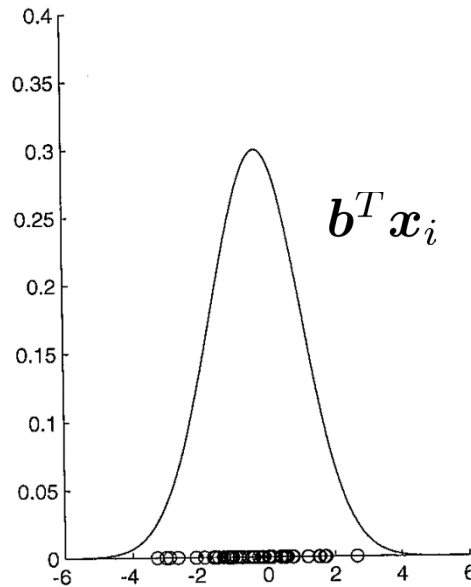
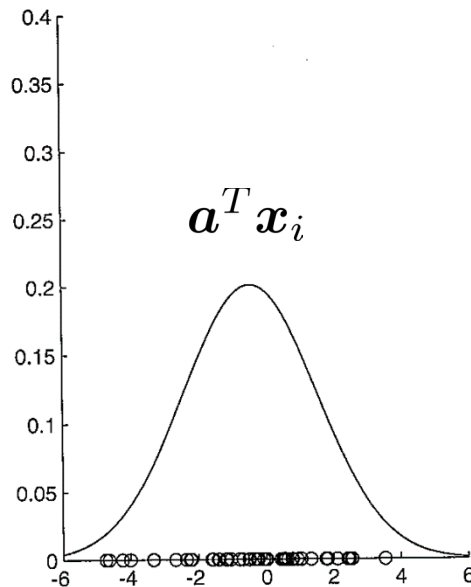
- Intuition:
 - $\hat{\mathbf{x}}_i$ is the projection of \mathbf{x}_i onto $\langle \mathbf{A} \rangle$.
 - Columns of \mathbf{A} define a k dimensional coordinate system for $\langle \mathbf{A} \rangle$.
 - $\boldsymbol{\theta}_i = \mathbf{A}^T \mathbf{x}_i$ are the coordinates of $\hat{\mathbf{x}}_i$ in the subspace

Maximum Variance Projections



$$aa^T x_i = a \underbrace{(a^T x_i)}_{\theta_i}$$

$$bb^T x_i = b \underbrace{(b^T x_i)}_{\theta_i}$$



Maximum Variance Perspective

- Suppose $\bar{\mathbf{x}} = \mathbf{0}$.
- What is the unit vector $\mathbf{a}_1 \in \mathbb{R}^d$ ($\|\mathbf{a}_1\| = 1$) for which the sample variance of

$$\theta^{(1)} = \mathbf{a}_1^T X$$

is maximized?

- The sample mean of $\mathbf{a}_1^T X$ is

$$\frac{1}{n} \sum a_1^T x_i = a_1^T \bar{x} = 0$$

- The sample variance of $\theta^{(1)}$ is therefore

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \left(a_1^T x_i - 0 \right)^2 &= \frac{1}{n} \sum (a_1^T x_i) (x_i^T a_1) \\ &= a_1^T \left(\frac{1}{n} \sum x_i x_i^T \right) a_1 = a_1^T S a_1 \end{aligned}$$

Maximum Variance Perspective

- We can express the sample variance of $\theta^{(1)}$ as

$$\text{var}(\theta^{(1)}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{a}_1^T \mathbf{x}_i)^2 = \mathbf{a}_1^T \mathbf{S} \mathbf{a}_1$$

- The solution of

$$\max_{\mathbf{a}_1: \|\mathbf{a}_1\|=1} \mathbf{a}_1^T \mathbf{S} \mathbf{a}_1$$

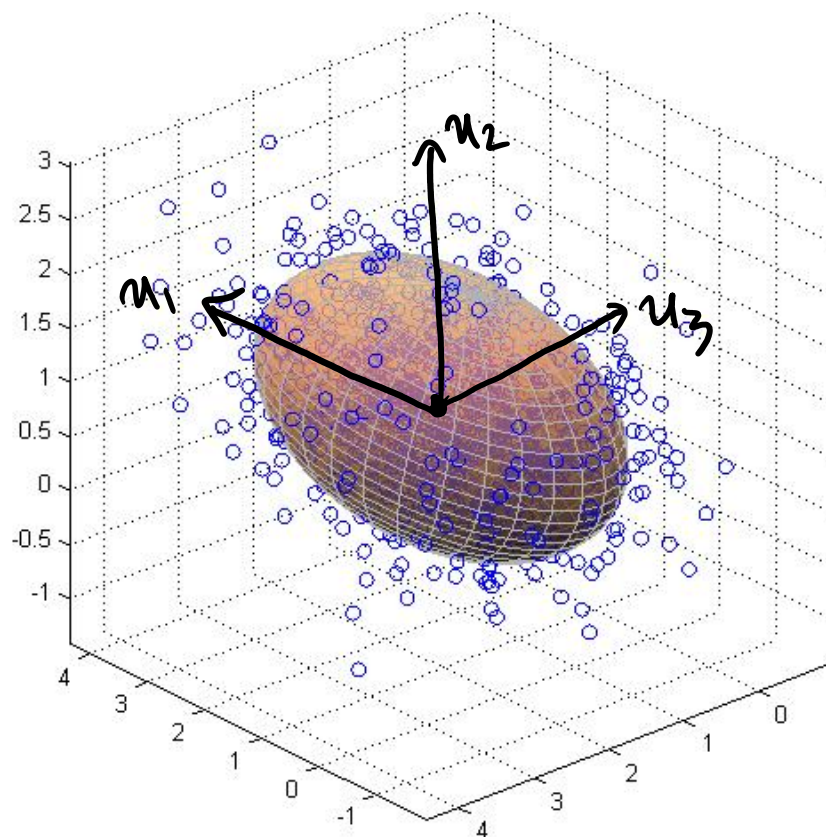
is

\mathbf{u}_1 = the "largest eigenvector" of \mathbf{S}

Maximum Variance Perspective

- More generally, we have the following result:
- **Theorem:** Let $\theta^{(k)} = \mathbf{a}_k^T \mathbf{X}$ and $\text{var}(\theta^{(k)}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{a}_k^T \mathbf{x}_i)^2$. A vector \mathbf{a}_k that maximizes $\text{var}(\theta^{(k)})$ subject to
 - $\|\mathbf{a}_k\| = 1$
 - $\mathbf{a}_k \perp \mathbf{u}_1, \dots, \mathbf{u}_{k-1}$is $\mathbf{a}_k = \mathbf{u}_k$.
- What is the variance of $\theta^{(k)}$?

$$\begin{aligned} \mathbf{a}_k^T \mathbf{S} \mathbf{a}_k &= \mathbf{u}_k^T (\mathbf{U} \mathbf{\Lambda} \mathbf{U}^T) \mathbf{u}_k \\ &= \mathbf{e}_k^T \mathbf{\Lambda} \mathbf{e}_k = \lambda_k \end{aligned}$$



Selecting k

- It can be shown that the optimal objective function value is

$$\min_{\mu, A, \theta_i} \sum \|x_i - \mu - A\theta_i\|^2 = n(\lambda_{k+1} + \dots + \lambda_d)$$

- When $k = 0$, this specializes to

$$\min_{\mu} \sum \|x_i - \mu\|^2 = n(\lambda_1 + \dots + \lambda_d)$$

which we call the *total variation* of the data.

- One heuristic for choosing k is to select the smallest k such that

$$\frac{\lambda_1 + \dots + \lambda_k}{\lambda_1 + \dots + \lambda_d} \geq .95$$

\nearrow % total variation explained by 1st k PCs.

Connection to SVD

- Assume $\bar{\mathbf{x}} = \mathbf{0}$

- Data matrix ($d \times n$)

$$\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_n].$$

- Let the singular value decomposition of \mathbf{X} be

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

where \mathbf{U} ($d \times d$) and \mathbf{V} ($n \times n$) are orthogonal matrices, and $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_{\min\{d,n\}})$ is $d \times n$.

- Then

principal eigenvectors = left singular vectors

$$\lambda_j = \begin{cases} \frac{1}{n} \sigma_j^2, & \text{if } j \leq \min(d, n) \\ 0, & \text{ow.} \end{cases}$$

Extensions

- Kernel PCA
- Streaming PCA
- Supervised PCA
- Robust PCA
- Sparse PCA
- Probabilistic PCA
- Nonnegative matrix factorization
- ...