

# Lecture 13: Principal Component Analysis

Statistical Learning and Data Mining

Xingye Qiao

Department of Mathematical Sciences  
Binghamton University

E-mail: [qiao@math.binghamton.edu](mailto:qiao@math.binghamton.edu)

Read: ELSII Ch. 14.5, ISLR 10.2 & 10.4, and SLS 8.2

# Outline

- 1 Interpretations & Uses
- 2 Models & Optimization Problems
- 3 Solution via the SVD
- 4 Amount of Variance Explained
- 5 Real Example
- 6 Extensions

# From Supervised Learning to Unsupervised Learning

- Supervised Learning: regression and classification. Response  $Y$  is used in training and the goal is to predict  $Y$
- Unsupervised Learning: No information about the response is used. The goal is to understand the  $X$  data.

# The next section would be .....

## 1 Interpretations & Uses

- Data Visualization
- Pattern Recognition
- Dimension Reduction

## 2 Models & Optimization Problems

## 3 Solution via the SVD

## 4 Amount of Variance Explained

## 5 Real Example

## 6 Extensions

The next section would be .....

- 1 Interpretations & Uses
  - Data Visualization
  - Pattern Recognition
  - Dimension Reduction
- 2 Models & Optimization Problems
- 3 Solution via the SVD
- 4 Amount of Variance Explained
- 5 Real Example
- 6 Extensions

# Data Visualization

- Suppose there are variables  $X_1, \dots, X_p$ . To visualize the data, one may draw pairwise scatterplots. But there are  $p(p-1)/2$  such plots.
- Data lie in  $p$ -dimensional space, but not all the dimensions are interesting.
- Solution: find a low-dimensional representation of the data that captures as much of the **information** as possible.
- PCA seeks a small number of dimensions that are **as interesting as possible**, where interesting-ness is measured by the amount that the observations **vary along each dimension**.

# Data Exploration - Multivariate data

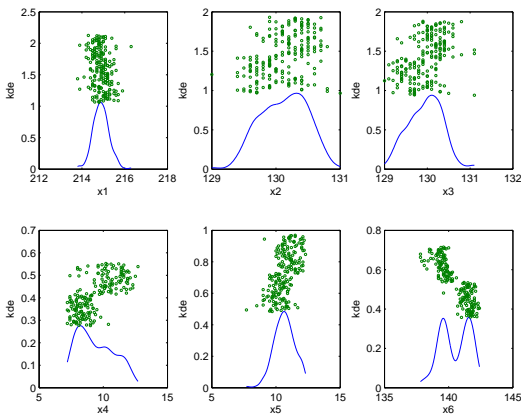
## Dimension $p = 6$ example – Swiss bank notes

- $n = 200$  Swiss bank notes (See Fig. 1.1, Härdle and Simar)
- Each note (obs.) has  $p = 6$  measurements (variables).
- Additional information: first half are genuine; the other half are counterfeit.
- Visualization of 6-dim'l data?
- Can use 6 KDEs overlaid with jitterplot for each measurements (variables)
- jitterplot: heights of dots (y value) are random for visualization. The x value represents the realized value of the data point.

## Swiss bank notes - Marginal KDEs

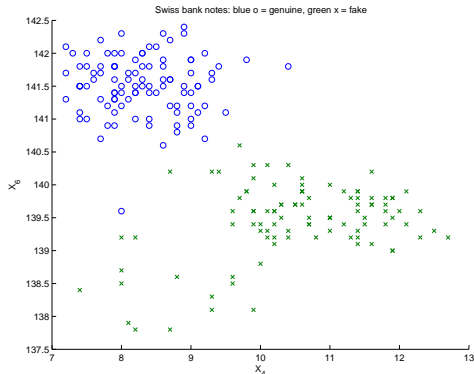
Marginal KDEs overlaid with jitterplot for each of 6 variables.

- Informative, realistic when  $p$  is small
- No information about association between variables.

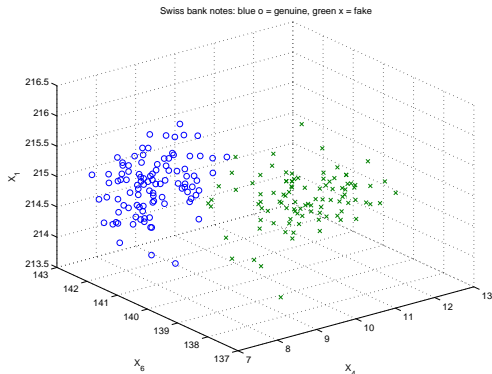




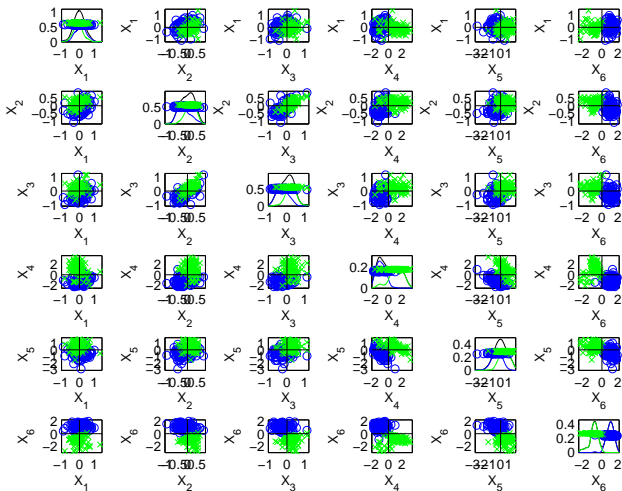
- Variable pair best visualized by scatterplot, e.g.  $X_4$  vs  $X_6$ .
- Understood as point clouds, which empirically representing the distribution
- $\binom{6}{2} = 15$  many pairs to choose from. 15 plots?



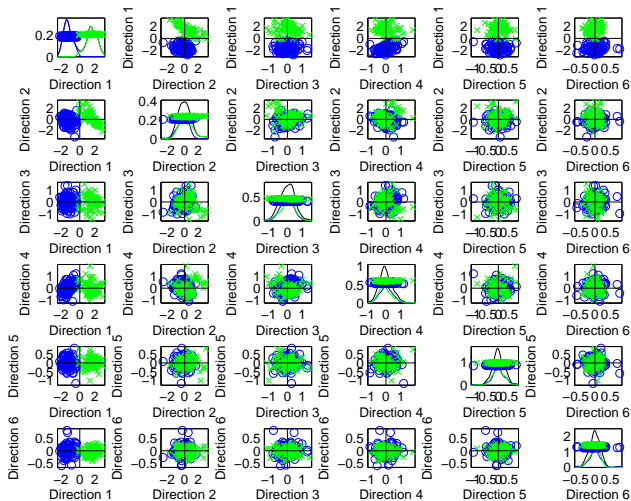
- Scatters of three variables can also be informative
- But only if software allows to rotate the axes.
- Otherwise, the 3D scatterplot is just a 2D scatterplot of two linear combinations of the three variables.
- Angle matters.



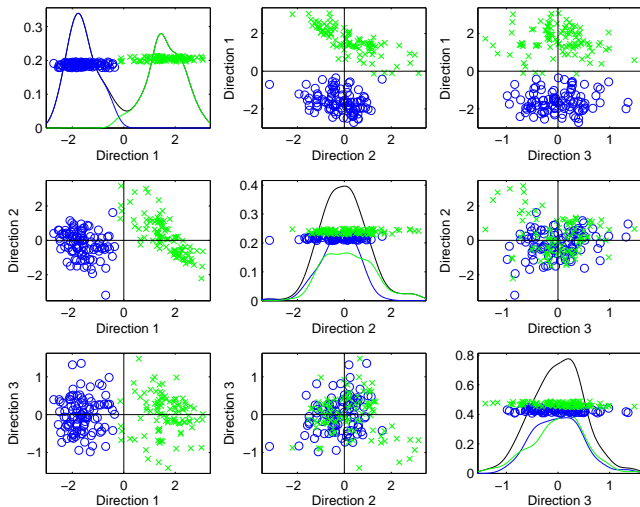
- A traditional, yet powerful, tool is to construct a matrix of scatterplots. - Too busy with  $p = 6$ .



■ Better to visualize with principal component scores.



- With principal component scores, we can focus on **fewer** combinations (it's called Dimension Reduction)



## The next section would be .....

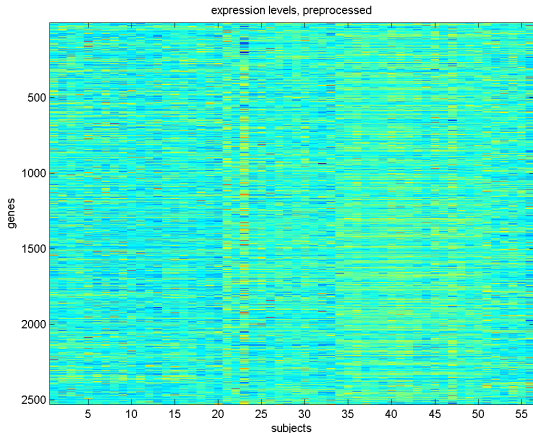
- 1 Interpretations & Uses
  - Data Visualization
  - Pattern Recognition
  - Dimension Reduction
- 2 Models & Optimization Problems
- 3 Solution via the SVD
- 4 Amount of Variance Explained
- 5 Real Example
- 6 Extensions

# Pattern Recognition

PCA can sometimes help discover previously unknown patterns, and help learn from labelled data.

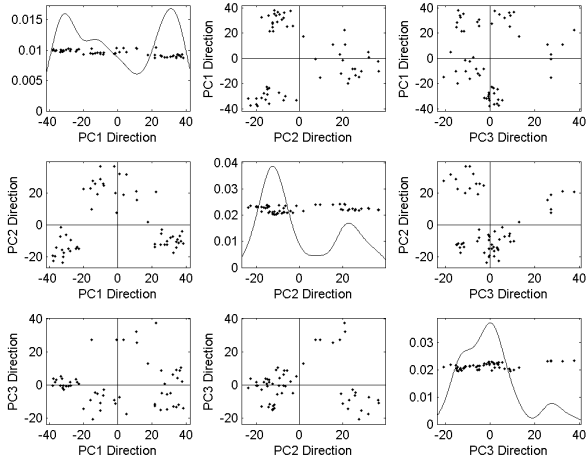
## Example: mRNA expression profiling

- Bhattacharjee et al (2001) PNAS
- Preprocessed gene expressions with  $d = 2530$  genes and  $n = 56$  subjects with lung cancer.
- Subgroup for different types of lung cancers?



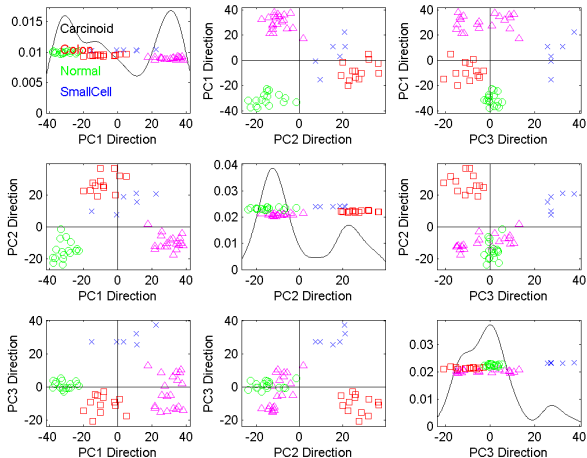


# mRNA expression profiling



# mRNA expression profiling

## ■ Color by true subgroups



Successfully capture the major pattern in the data: **the black, red, green, blue and cyan observations that are near each other in the high-dimensional space remain nearby in these two-dimensional representations.**

The next section would be .....

- 1 Interpretations & Uses
  - Data Visualization
  - Pattern Recognition
  - Dimension Reduction
- 2 Models & Optimization Problems
- 3 Solution via the SVD
- 4 Amount of Variance Explained
- 5 Real Example
- 6 Extensions

## Dimension Reduction

- Lastly, PCA is a way of dimension reduction.
- For example, as in principal component regression, we simply use principal components as predictors in a regression model in place of the original larger set of variables.

# The next section would be .....

- 1 Interpretations & Uses
- 2 Models & Optimization Problems
  - Matrix Factorization
  - Covariance
- 3 Solution via the SVD
- 4 Amount of Variance Explained
- 5 Real Example
- 6 Extensions

# The next section would be .....

- 1 Interpretations & Uses
- 2 Models & Optimization Problems
  - Matrix Factorization
  - Covariance
- 3 Solution via the SVD
- 4 Amount of Variance Explained
- 5 Real Example
- 6 Extensions

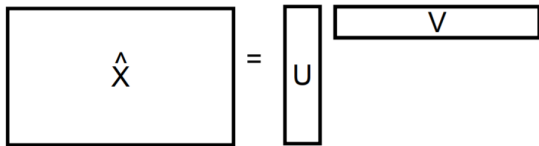
## Matrix factorization

- Given matrix  $X$ , we seek to find matrices  $U$  and  $V$  such that  $X \approx UV := \hat{X}$ .

$$\hat{X} = UV$$

- It is ideal that  $U$  has very few columns.
- Why are low-rank approximations important?
  - Intuitively, if matrix is low rank, then the observations can be explained by linear combinations of **few** underlying factors
  - Want to know which factors control the observations





A diagram illustrating the PCA equation  $\hat{X} = UV$ . On the left is a large rectangle containing the symbol  $\hat{X}$ . To its right is an equals sign. Further right is a tall, narrow vertical rectangle containing the letter  $U$ . To the right of  $U$  is a horizontal rectangle containing the letter  $V$ .

- Imagine that  $X$  is the **centered** data matrix where the  $i$ th column  $X_{(i)}$  is the  $i$ th observation.
- $X \approx UV$  means that we seek to find  $U$  and  $V$  so that  $X_{(i)} \approx \sum_{j=1}^q v_{ji} U_j$  where  $U_j$  is the  $j$ th column of  $U$ .

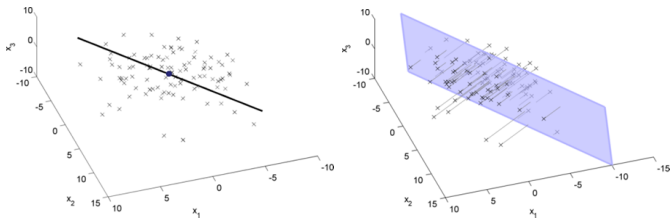
To resolve identifiable issue, we may require  $U^T U = \mathbb{I}$ .

PCA can be viewed as the following matrix factorization / matrix approximation problem.

$$(U, V) = \underset{U, V, U^T U = \mathbb{I}}{\operatorname{argmin}} \|X - UV\|_F^2 = \underset{U, V, U^T U = \mathbb{I}}{\operatorname{argmin}} \sum_{i=1}^n \|X_{(i)} - \sum_{j=1}^q v_{ji} U_j\|_2^2$$

where  $\|A\|_F^2 = \sum_{i=1}^n \sum_{j=1}^n A_{ij}^2$  is the Frobenius norm of matrix  $A$

# Geometric understanding of PCA



A 3D point cloud.

Mean is removed

left: best 1-d approximation ( $q = 1$ )

right: best 2-d approximation ( $q = 2$ )

Next, another formulation of PCA using eigen-decomposition of covariance matrix.

# The next section would be .....

- 1 Interpretations & Uses
- 2 Models & Optimization Problems
  - Matrix Factorization
  - Covariance
- 3 Solution via the SVD
- 4 Amount of Variance Explained
- 5 Real Example
- 6 Extensions

## Linear dimension reduction

For a random vector  $\mathbf{X} \in \mathbb{R}^p$ , consider reducing the dimension from  $p$  to  $d$ , i.e.,  $p$  variables  $(X_1, \dots, X_p)^T$  to a set of *most interesting*  $d$  variables. Here,  $1 \leq d \leq p$ .

- Best subset?
- Linear dimension reduction: Construct  $d$  variables  $Z_1, \dots, Z_d$  as linear combinations of  $X_1, \dots, X_p$ , i.e.

$$Z_i = a_{i1}X_1 + \dots + a_{ip}X_p = \mathbf{a}_i' \mathbf{X} \quad (i = 1, \dots, d),$$

with  $\mathbf{a}_i \in \mathbb{R}^p$ .

Linear dimension reduction seeks a sequence of such  $Z_i$ , or equivalently a sequence of  $\mathbf{a}_i$ , where the random variables  $Z_i$ 's are most **important** among all choices.

# Principal Component Analysis

Require  $\|\mathbf{a}_1\| = 1$  and  $\langle \mathbf{a}_i, \mathbf{a}_j \rangle = 0$ . Thus the problem is to find an interesting set of (orthogonal) **direction vectors**  $\{\mathbf{a}_i : i = 1, \dots, p\}$ , where the projection scores of  $\mathbf{X}$  onto  $\mathbf{a}_i$  are useful.

PCA aims for a set of direction vectors which lead to **maximal variances of the projected random variables**.

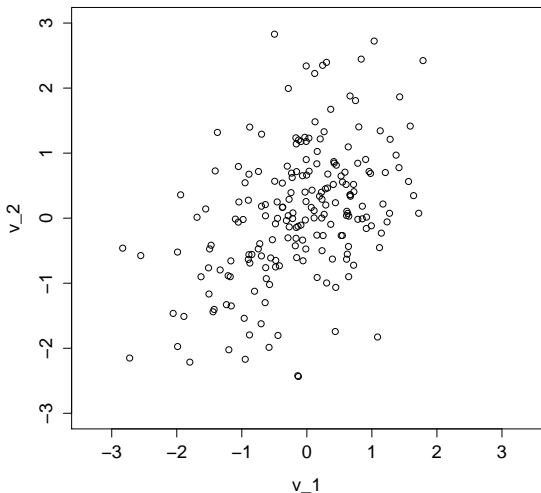
Take  $d = 1$ . PCA for the distribution of  $\mathbf{X}$  finds  $\mathbf{a}_1$  such that

$$\mathbf{a}_1 = \operatorname{argmax}_{\mathbf{a} \in \mathbb{R}^p, \|\mathbf{a}\|=1} \operatorname{Var}(Z_1(\mathbf{a})) \left( = \operatorname{argmax}_{\mathbf{a} \in \mathbb{R}^p, \|\mathbf{a}\|=1} \mathbf{a}' \operatorname{Var}(\mathbf{X}) \mathbf{a} \right),$$

where  $Z_1(\mathbf{a}) = a_1 X_1 + \dots + a_p X_p = \mathbf{a}' \mathbf{X}$ .

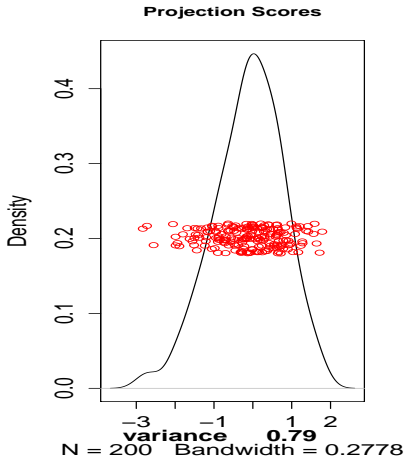
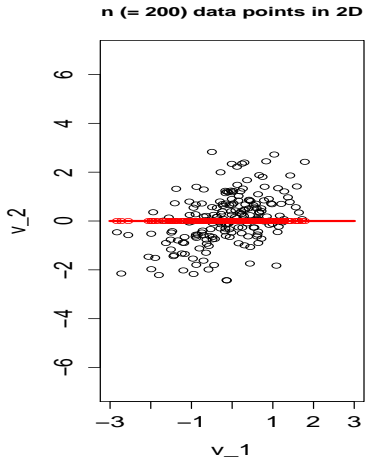
# Geometric understanding of PCA for point cloud

**n (= 200) data points in 2D**



PCA is best understood with a point cloud. Take a look at this 2D example.

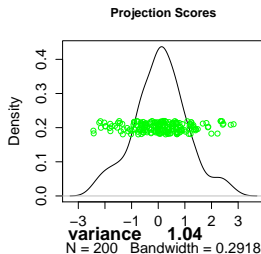
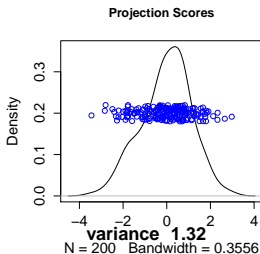
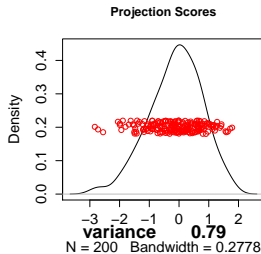
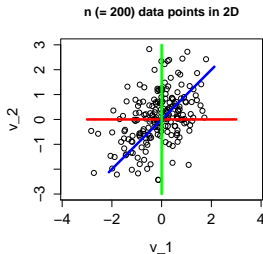
Take  $a = (1, 0)'$ .





## Which one is better?

Take  $\mathbf{a} = (1, 0)'$ ,  $(0, 1)'$ ,  $(1/\sqrt{2}, 1/\sqrt{2})'$ .



## Formulation of population PCA-1

Suppose a random vector  $\mathbf{X}$  with mean  $\mu$ , covariance  $\Sigma$  (not necessarily normal).

The first principal component (PC) direction vector is the unit vector  $\mathbf{u}_1 \in \mathbb{R}^p$  that maximizes the variance of  $\mathbf{u}_1' \mathbf{X}$  among all unit vectors, i.e.,

$$\mathbf{u}_1 = \operatorname{argmax}_{\mathbf{u} \in \mathbb{R}^p, \|\mathbf{u}\|=1} \operatorname{Var}(\mathbf{u}' \mathbf{X}).$$

- $\mathbf{u}_1 = (u_{11}, \dots, u_{1p})'$  is the first PC **direction vector**, sometimes called *loading vector*.
- $u_{11}, \dots, u_{1p}$  are loadings of the 1st PC.
- $Z_1 = u_{11}X_1 + \dots + u_{1p}X_p = \mathbf{u}_1' \mathbf{X}$  is the first PC score or the first principal component (it's a random variable).
- $\lambda_1 = \operatorname{Var}(\mathbf{u}_1' \mathbf{X}) = \operatorname{Var}(Z_1)$  is the variance explained by the first PC.

## Formulation of population PCA-2

The second PC direction is the unit vector  $\mathbf{u}_2 \in \mathbb{R}^p$  that

- can maximize the variance of  $\mathbf{u}_2' \mathbf{X}$ ;
- among directions orthogonal to the first PC direction  $\mathbf{u}_1$ .

That is,

$$\mathbf{u}_2 = \underset{\mathbf{u} \in \mathbb{R}^p, \|\mathbf{u}\|=1, \mathbf{u}'\mathbf{u}_1=0}{\operatorname{argmax}} \operatorname{Var}(\mathbf{u}'\mathbf{X}).$$

- $\mathbf{u}_2 = (u_{21}, \dots, u_{2p})'$  is the second PC direction vector, and is the vector of the 2nd set of loadings.
- $Z_2 = \mathbf{u}_2' \mathbf{X}$  is the second principal component.
- $\lambda_2 = \operatorname{Var}(Z_2)$  is the variance explained by the second PC, and  $\lambda_1 \geq \lambda_2$ .
- $\operatorname{Corr}(Z_1, Z_2) = 0$ .

## Formulation of population PCA-(3,4,...p)

Given the first  $k - 1$  PC directions  $\mathbf{u}_1, \dots, \mathbf{u}_{k-1}$ , the  $k$ th PC direction is the unit vector  $\mathbf{u}_k \in \mathbb{R}^p$  that

- maximizes the variance of  $\mathbf{u}'_k \mathbf{X}$ ;
- among those orthogonal to the 1st to the  $(k - 1)$ th PC directions  $\mathbf{u}_j$  ( $j = 1, \dots, k - 1$ )

That is,

$$\mathbf{u}_k = \underset{\substack{\mathbf{u} \in \mathbb{R}^p, \|\mathbf{u}\|=1 \\ \mathbf{u}'\mathbf{u}_j=0, j=1, \dots, k-1}}{\operatorname{argmax}} \operatorname{Var}(\mathbf{u}'\mathbf{X}).$$

- $\mathbf{u}_k = (u_{k1}, \dots, u_{kp})'$  is the  $k$ th PC direction vector, and is the vector of the  $k$ th loadings.
- $Z_k = \mathbf{u}'_k \mathbf{X}$  is the  $k$ th principal component.
- $\lambda_k = \operatorname{Var}(Z_k)$  is the variance explained by the  $k$  PC score, and  $\lambda_1 \geq \dots \geq \lambda_{k-1} \geq \lambda_k$ .
- $\operatorname{Corr}(Z_i, Z_j) = 0$  for all  $i \neq j \leq k$ .

## Relation to eigen-decomposition of $\Sigma$

Recall the eigen-decomposition of the symmetric positive definite  $\Sigma = \mathbf{U}\mathbf{\Lambda}\mathbf{U}'$  with

- $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_p]$  orthogonal matrix
- $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_p)$  with  $\lambda_1 \geq \dots \geq \lambda_p$ ,
- $\Sigma \mathbf{u}_i = \lambda_i \mathbf{u}_i$ .

In next two slides we show that:

- 1 The  $k$ th eigenvector  $\mathbf{u}_k$  is the  $k$ th PC direction vector.
- 2 The  $k$ th eigenvalue  $\lambda_k$  is the variance explained by the  $k$ th principal component.
- 3 PC directions are both orthogonal  $\mathbf{u}_i' \mathbf{u}_j = 0$  ( $i \neq j$ ) and  $\Sigma$ -orthogonal

$$\mathbf{u}_i' \Sigma \mathbf{u}_j = 0 \iff \text{Cov}(Z_i, Z_j) = 0 \quad (i \neq j).$$

## Gradient

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ . Define

$$\nabla f(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_d} \end{pmatrix}_{d \times 1}$$

Facts:



$$\nabla_{\mathbf{x}}(\mathbf{c}'\mathbf{x}) = \frac{\partial(\mathbf{c}'\mathbf{x})}{\partial \mathbf{x}} = \mathbf{c}$$



$$\nabla_{\mathbf{x}}(\mathbf{x}'\mathbf{A}\mathbf{x}) = \frac{\partial(\mathbf{x}'\mathbf{A}\mathbf{x})}{\partial \mathbf{x}} = 2\mathbf{A}$$

for symmetric  $\mathbf{A}$

## Relation to eigen-decomposition of $\Sigma$

The first PC direction maximizes  $\text{Var}(\mathbf{u}'\mathbf{X})$  with the constraint  $\mathbf{u}'\mathbf{u} = 1$ . Using Lagrange multiplier  $\lambda$ , it is the same as finding a stationary point of

$$\begin{aligned}\Phi(\mathbf{u}, \lambda) &= \text{Var}(\mathbf{u}'\mathbf{X}) - \lambda(\mathbf{u}'\mathbf{u} - 1) \\ &= \mathbf{u}'\Sigma\mathbf{u} - \lambda(\mathbf{u}'\mathbf{u} - 1).\end{aligned}$$

The stationary point solves the following:

$$\frac{1}{2} \frac{\partial}{\partial \mathbf{u}} \Phi(\mathbf{u}, \lambda) = \Sigma\mathbf{u} - \lambda\mathbf{u} = \mathbf{0},$$

which leads to

$$\lambda = \mathbf{u}'\Sigma\mathbf{u} = \text{Var}(\mathbf{u}'\mathbf{X}), \quad \Sigma\mathbf{u} = \lambda\mathbf{u}. \quad (1)$$

Recall: any eigenvector-eigenvalue pair  $(\mathbf{u}_i, \lambda_i)$ ,  $(i = 1, \dots, p)$  satisfies the second eq. in (1). It is clear that the first PC direction is the first eigenvector  $\mathbf{u}_1$ , as it gives the largest variance  $\lambda_1 = \mathbf{u}_1'\Sigma\mathbf{u}_1 = \text{Var}(\mathbf{u}_1'\mathbf{X}) \geq \lambda_j$  ( $j > 1$ ).

## Relation to eigen-decomposition of $\Sigma$

For the  $k$ th PC direction, we form a Lagrangian function

$$\Phi(\mathbf{u}, \lambda, \gamma_1^k) = \mathbf{u}'\Sigma\mathbf{u} - \lambda(\mathbf{u}'\mathbf{u} - 1) - \sum_{j=1}^{k-1} 2\gamma_j\mathbf{u}'_j\mathbf{u},$$

given the first  $k - 1$  PC directions. The derivative of  $\Phi$ , equated to zero, is then

$$\begin{aligned}\frac{1}{2} \frac{\partial}{\partial \mathbf{u}} \Phi(\mathbf{u}, \lambda, \gamma_1^k) &= \Sigma\mathbf{u} - \lambda\mathbf{u} - \sum_{j=1}^{k-1} \gamma_j\mathbf{u}_j = \mathbf{0}, \\ \frac{\partial}{\partial \gamma_j} \Phi(\mathbf{u}, \lambda, \gamma_1^k) &= \mathbf{u}'_j\mathbf{u} = 0.\end{aligned}\tag{2}$$

We have  $\gamma_j = \mathbf{u}'_j\Sigma\mathbf{u} = 0$  (since  $\Sigma\mathbf{u}_j = \lambda_j\mathbf{u}_j$ ), thus

$$\lambda = \mathbf{u}'\Sigma\mathbf{u} = \text{Var}(\mathbf{u}'\mathbf{X}), \quad \Sigma\mathbf{u} = \lambda\mathbf{u}.\tag{3}$$

The  $k$ th to the last eigen-pairs  $(\mathbf{u}_i, \lambda_i)$ ,  $(i = k, \dots, p)$  all satisfy both (3) and (2). Thus, the  $k$ th PC direction is  $\mathbf{u}_k$ , as it gives the largest variance  $\lambda_k = \mathbf{u}'_k\Sigma\mathbf{u}_k$  among the remaining eigen-pairs.



## Computation of PCA

PCA is either computed using **eigenvalue decomposition** of  $\mathbf{S} = \frac{1}{n-1} \tilde{\mathbf{X}} \tilde{\mathbf{X}}'$  or using the **singular value decomposition** of  $\tilde{\mathbf{X}}$ .

### Eigen-decomposition of $\mathbf{S}$

For  $\mathbf{S} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}'$ ,

- 1 PC directions  $\mathbf{u}_k$  (eigenvectors)
- 2 Variance of PC (scores)  $\lambda_k$  (eigenvalues)
- 3 Matrix of centered principal component scores

$$\mathbf{U}' \tilde{\mathbf{X}} = \mathbf{Z} = \begin{bmatrix} \mathbf{z}_{(1)} \\ \vdots \\ \mathbf{z}_{(p)} \end{bmatrix}.$$

# The next section would be .....

- 1 Interpretations & Uses
- 2 Models & Optimization Problems
- 3 Solution via the SVD
  - Properties of the SVD
  - SVD & PCs
  - PC loadings/directions and PC scores
- 4 Amount of Variance Explained
- 5 Real Example
- 6 Extensions

## The next section would be .....

- 1 Interpretations & Uses
- 2 Models & Optimization Problems
- 3 Solution via the SVD
  - Properties of the SVD
    - SVD & PCs
    - PC loadings/directions and PC scores
- 4 Amount of Variance Explained
- 5 Real Example
- 6 Extensions

# Computation of PCA

## Singular value decomposition (SVD) of $\tilde{\mathbf{X}}$

The singular value decomposition (SVD) of  $p \times n$  matrix  $\tilde{\mathbf{X}}$  has the form

$$\tilde{\mathbf{X}} = \mathbf{U}\mathbf{D}\mathbf{V}'.$$

- The left singular vectors  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_p]_{p \times p}$  and the right singular vectors  $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_p]_{n \times p}$  are orthogonal ( $\mathbf{U}'\mathbf{U} = \mathbb{I}_p$ ,  $\mathbf{V}'\mathbf{V} = \mathbb{I}_p$ ).
- The columns of  $\mathbf{U}$  span the column space of  $\tilde{\mathbf{X}}$ ; the columns of  $\mathbf{V}$  (which are  $n$ -vectors) span the row space.
- $\mathbf{D} = \text{diag}(d_1, \dots, d_p)$ ,  $d_1 \geq d_2 \geq \dots \geq d_p \geq 0$  are the singular values of  $\tilde{\mathbf{X}}$ .

## The next section would be .....

- 1 Interpretations & Uses
- 2 Models & Optimization Problems
- 3 Solution via the SVD
  - Properties of the SVD
  - SVD & PCs
  - PC loadings/directions and PC scores
- 4 Amount of Variance Explained
- 5 Real Example
- 6 Extensions

## SVD and Eigen-decomposition Connection

If SVD of  $\tilde{\mathbf{X}} = \mathbf{U}\mathbf{D}\mathbf{V}'$ , then

$$\mathbf{S} = \frac{1}{n-1} \tilde{\mathbf{X}} \tilde{\mathbf{X}}' = \frac{1}{n-1} \mathbf{U} \mathbf{D} \mathbf{V}' \mathbf{V} \mathbf{D} \mathbf{U}' = \mathbf{U} \text{diag}\left(\frac{1}{n-1} d_j^2\right) \mathbf{U}'$$

- 1 PC directions  $\mathbf{u}_k$  (left singular vectors)
- 2 Variance of PC (scores) =  $\frac{1}{n-1} d_j^2$  (scaled singular values<sup>2</sup>)
- 3 Matrix of principal component scores (scaled right singular vectors)

$$\begin{bmatrix} \mathbf{z}_{(1)} \\ \vdots \\ \mathbf{z}_{(p)} \end{bmatrix} = \mathbf{Z} = \mathbf{U}' \tilde{\mathbf{X}} = \mathbf{D} \mathbf{V}' = \begin{bmatrix} d_1 \mathbf{v}'_1 \\ \vdots \\ d_p \mathbf{v}'_p \end{bmatrix}$$

NOTE: we are working with the centered  $\tilde{\mathbf{X}}$  here, not  $\mathbf{X}$ !!

## PCA in R

The standard data format is the  $n \times p$  data frame or matrix  $x$ .  
To perform PCA by eigen decomposition:

```
spr <- princomp(x)
U <- spr$loadings
L <- (spr$sdev)^2
Z <- spr$scores
```

To perform PCA by singular value decomposition

```
gpr <- prcomp(x)
U <- gpr$rotation
L <- (gpr$sdev)^2
Z <- gpr$x
```

## Scaling? Correlation PCA

- PCA is not scale invariant.
- SOMETIMES, good idea to do normalization.
- Correlation matrix of a random vector  $\mathbf{X}$  is given by

$$\mathbf{R} = \mathbf{D}_{\Sigma}^{-\frac{1}{2}} \Sigma \mathbf{D}_{\Sigma}^{-\frac{1}{2}},$$

where  $\mathbf{D}_{\Sigma}$  is the  $p \times p$  diagonal matrix consisting of diagonal elements of  $\Sigma$ .

- Correlation PCA: PC directions obtained by eigen-decomposition of  $\mathbf{R} = \mathbf{U}_R \Lambda_R \mathbf{U}_R'$ .
- Preferred if measurements are not commensurate (e.g.  $X_1$  = household income,  $X_2$  = years in school).



## The next section would be .....

- 1 Interpretations & Uses
- 2 Models & Optimization Problems
- 3 Solution via the SVD
  - Properties of the SVD
  - SVD & PCs
  - PC loadings/directions and PC scores
- 4 Amount of Variance Explained
- 5 Real Example
- 6 Extensions

## PC loadings and PC scores

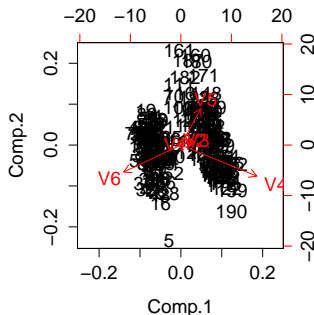
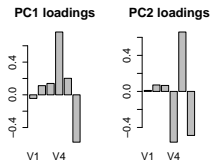
- For some students, it is often confusing between PC loadings and PC scores.
- PC direction  $\mathbf{U}_j$ :  $j$ th column of  $\mathbf{U}$ , a  $p$ -dimensional vector.
- PC loadings: elements of  $\mathbf{U}_j$ , measuring contributions from different dimensions (variables) to the  $j$ th principal component
- PC scores: inner products of  $\mathbf{x}_i^T \mathbf{U}_j$ ,  $i = 1, \dots, n$ , coordinates of obs.  $i$  in the new coordinate system spanned by  $\mathbf{U}_j$ 's

## Which variables are most responsible for the principal components?

- Check loadings of principal component directions.
- **Biplot** - scatterplot of **PC1** and **PC2** scores, overlaid with  $p$  vectors each representing the loadings of the first two PC directions.

In the Swiss Bank Note Data, the loadings are

	Comp.1	Comp.2	Comp.3	Comp.4	Comp.5	Comp.6
V1			-0.326	0.562	0.753	
V2	0.112		-0.259	0.455	-0.347	-0.767
V3	0.139		-0.345	0.415	-0.535	0.632
V4	0.768	-0.563	-0.218	-0.186		
V5	0.202	0.659	-0.557	-0.451	0.102	
V6	-0.579	-0.489	-0.592	-0.258		



Recall that the scatter plot of PC1+PC2 is a visualization after rotation and projection.

Hence **red** vectors can be viewed as the rotated and projected coordinate direction vectors. For example, V1 is the rotated and projected  $(1, 0, 0, 0, \dots)'$

## The next section would be .....

- 1 Interpretations & Uses
- 2 Models & Optimization Problems
- 3 Solution via the SVD
- 4 Amount of Variance Explained**
- 5 Real Example
- 6 Extensions

## How many components to keep? (1)

“How much of the variation within the data have PCs explained?”

- 1** Total variation in  $\mathbf{X}$  is the sum of all marginal (sample) variances

$$\begin{aligned}\sum_{k=1}^p \text{Var}(\{x_{ki} : i = 1, \dots, n\}) &= \text{Trace}(\mathbf{S}) = \text{Trace}(\hat{\Lambda}) \\ &= \sum_{k=1}^p \hat{\lambda}_k = \sum_{k=1}^p \text{Var}(\{z_{(k)i} : i = 1, \dots, n\}).\end{aligned}$$

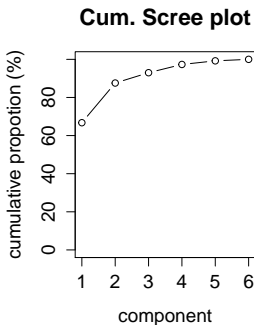
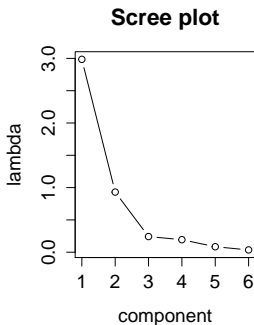
- 2** (Sample) variance of the  $k$ th PC:

$$\text{Var}(\{z_{(k)i} : i = 1, \dots, n\}) = \hat{\lambda}_k$$

- 3** Total variance in the 1st to the  $k$ th PCs:  $\hat{\lambda}_1 + \dots + \hat{\lambda}_k$ .

## How many components to keep? (1)

- 1 In scree plot ( $k, \hat{\lambda}_k$ ), we look for an elbow.
- 2 In cumulative scree plot (proportion of variance explained,  $(k, \frac{\sum_{j=1}^k \hat{\lambda}_j}{\sum_{j=1}^p \hat{\lambda}_j})$ ), use 90% as a cutoff.



## How many components to keep? (2)

- 1 Kaiser's rule of thumb: Retain PCs 1– $k$  satisfying

$$\lambda_k > \bar{\lambda} = \frac{1}{p} \sum_{j=1}^p \lambda_j.$$

Tends to choose fewer components.

- 2 Likelihood ratio testing on null hypothesis

$$H_0(k) : \lambda_{k+1} = \cdots = \lambda_p,$$

The first  $k$  components will be retained if  $H_0(k)$  is not rejected at a specified level.



## The next section would be .....

- 1 Interpretations & Uses
- 2 Models & Optimization Problems
- 3 Solution via the SVD
- 4 Amount of Variance Explained
- 5 Real Example**
- 6 Extensions

# PCA for Olivetti Faces data

## Olivetti Faces data

- Obtained from <http://www.cs.nyu.edu/~roweis/data.html>.
- Grayscale faces 8 bit [0-255], a few (10) images of several (40) different people.
- 400 total images, 64x64 size.
- From the Olivetti database at ATT.



## Images as data

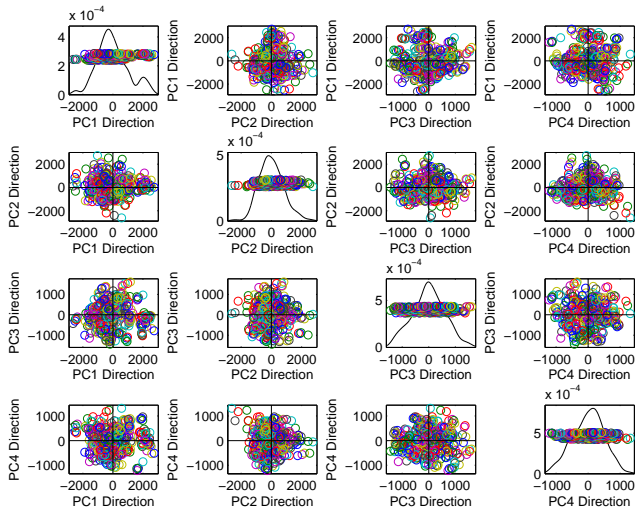
An image is a matrix-valued datum. In Olivetti Faces data, the matrix is of size  $64 \times 64$ , with each pixel having values between  $[0-255]$ . The matrix, corresponding one observation, is vectorized (vec'd) by stacking each column into one long vector of size  $d = 4096 = 64 \times 64$ .

So,  $\mathbf{x}_1$  is a  $d \times 1$  vector corresponding to

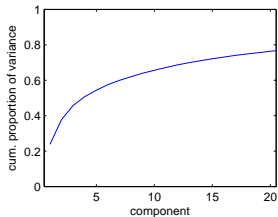
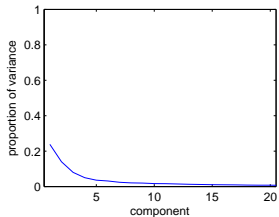
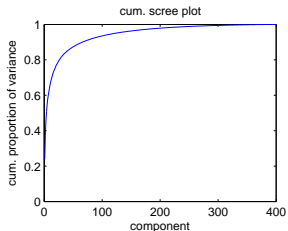
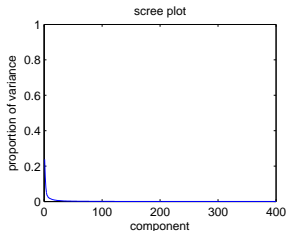


PCA is applied to the data matrix  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ . Tall and skinny data.

# Olivetti Faces data–Major components

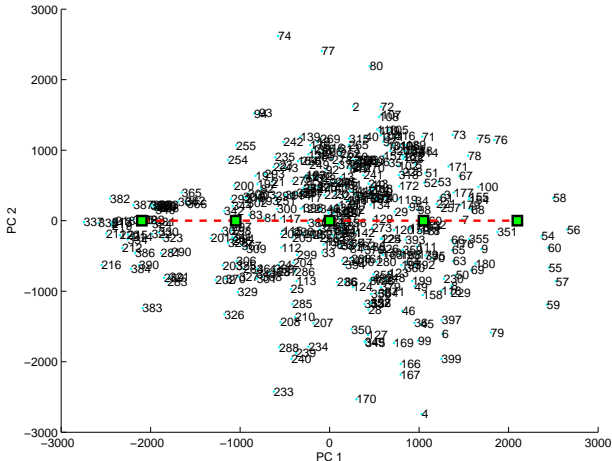


# Olivetti Faces data–Scree plots



# Olivetti Faces data–Interpretation

Examine the mode of variation by walking along the PC direction through the mean. Shown here are  $\pm 1, 2$  standard deviations of  $Z_{(1)}$  apart from the mean in the direction of PC1.



# Olivetti Faces data–Interpretation (Eigenfaces)

PC walk.  $\pm 2$  std along PC dir. (Top—PC 1, Mid—PC 2, Bottom—PC 3)



PC1  $\sim$  darker to lighter face

PC2  $\sim$  feminine to masculine face

PC3  $\sim$  oval to rectangle face

## How to walk along the PC direction?

- Vectorized data in  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ , with mean  $\bar{\mathbf{x}}$ .
- Compute  $(\mathbf{u}_j, \lambda_j)$ : the  $j$ th PC direction and PC variance.
- Walk along PC- $j$  direction and examine at position  $s = \pm 2, \pm 1, 0$  by
  - 1 Reconstruction at  $s$ :  $\mathbf{w}_s = \bar{\mathbf{x}} \pm s\sqrt{\lambda_j}\mathbf{u}_j$
  - 2 Convert to image by reshaping the  $4096 \times 1$  vector  $\mathbf{w}_s$  into  $64 \times 64$  matrix  $\mathbf{W}_s$ .

Next, reconstruct the original face using PCs.



## Approximation to the original data matrix

Recall the matrix factorization viewpoint of PCA:  $X \approx UV$ .

$$\mathbf{x}_i = \bar{\mathbf{x}} + \sum_{j=1}^p z_{(j)i} \mathbf{u}_j, \quad (i = 1, \dots, n)$$

Approximation of the original observation  $\mathbf{x}_i$  by the first  $m < p$  principal components:

$$\hat{\mathbf{x}}_i = \bar{\mathbf{x}} + \sum_{j=1}^m z_{(j)i} \mathbf{u}_j,$$

- The larger  $m$ , the better approximation by  $\hat{\mathbf{x}}_i$ .
- The smaller  $m$ , the more succinct dimension reduction of  $\mathbf{X}$ .

See some mathematical explanations in the next page.

## Olivetti Faces data–Reconstruction of original data

Recall

1  $\tilde{\mathbf{X}} = \mathbf{U}\mathbf{D}\mathbf{V}'$

2  $\mathbf{Z} = \mathbf{U}'\tilde{\mathbf{X}} = \mathbf{D}\mathbf{V}'$

3  $\tilde{\mathbf{X}} = \mathbf{U}\mathbf{Z}$

4  $\tilde{\mathbf{x}}_i = \mathbf{U}\mathbf{z}_i = \sum_{j=1}^p \mathbf{u}_j z_{(j)i}$

In a coordinate system with  $\{\mathbf{u}_i, i = 1, \dots, n\}$  as the  $p$  basis vectors,  $z_{(j)i}$  is the  $j$ th coordinate for the  $i$ th observation

$$\tilde{\mathbf{x}}_i = \mathbf{x}_i - \bar{\mathbf{x}}.$$

Hence

$$\mathbf{x}_i = \bar{\mathbf{x}} + \sum_{j=1}^p \mathbf{u}_j z_{(j)i}$$

## Reconstruction of original face

Observation index  $i = 5$ .

5th face. from top left to bottom right: (mean, 1, 5, 10) & (20, 50, 100, 400) PCs



## Reconstruction of original face

Observation index  $i = 19$ .

19th face. from top left to bottom right: (mean, 1, 5, 10) & (20, 50, 100, 400) PCs



## Reconstruction of original face

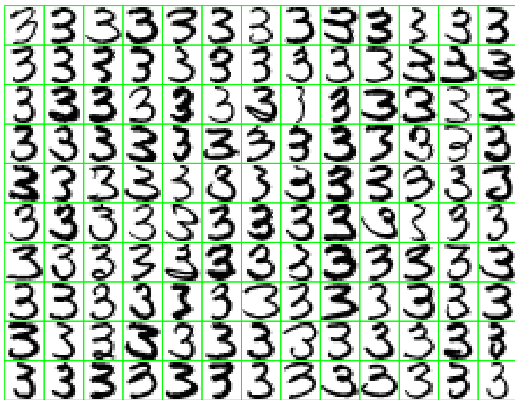
Observation index  $i = 100$ .

100th face. from top left to bottom right: (mean, 1, 5, 10) & (20, 50, 100, 400) PCs

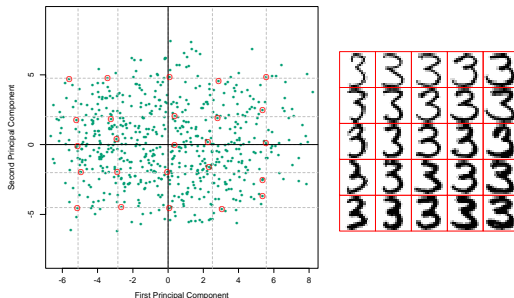


- Human eyes require  $> 50$  principal components to see resemblance between  $\hat{\mathbf{x}}_i$  and  $\mathbf{x}_i$ .
- Corresponds to about 90 percent of variance explained in PCs.
- 50 is still much smaller than 4096!
- Subjective and heuristic decision on “how many components to use”
- Reconstruction by PCA most useful and meaningful when
  - each datum is visually represented (rather than being just numbers).
  - for example: images, functions, shapes.

## Handwritten Digits



**FIGURE 14.22.** A sample of 130 handwritten 3's shows a variety of writing styles.



**FIGURE 14.23.** (Left panel:) the first two principal components of the handwritten threes. The circled points are the closest projected images to the vertices of a grid, defined by the marginal quantiles of the principal components. (Right panel:) The images corresponding to the circled points. These show the nature of the first two principal components.

$$\begin{aligned}
 \hat{f}(\lambda) &= \bar{x} + \lambda_1 v_1 + \lambda_2 v_2 \\
 &= \text{3} + \lambda_1 \cdot \text{3} + \lambda_2 \cdot \text{3}.
 \end{aligned}$$



## PCA as a mean of dimension reduction

- We can use only the first  $d$  PCs to approximately represent the data. Instead of  $\mathbf{X}_{p \times n}$ , we store the data as  $\mathbf{Z}_{d \times n}$ .
- However,
  - 1 Unsupervised learning (no information on  $Y$ ).
  - 2 Hard to interpret. Each PC (new variable) is a linear combination of  $p$  variables.
  - 3 Eigen-decomposition/SVD are problematic when  $p \gg n$ .

## The next section would be .....

- 1 Interpretations & Uses
- 2 Models & Optimization Problems
- 3 Solution via the SVD
- 4 Amount of Variance Explained
- 5 Real Example
- 6 Extensions**

## Matrix Completion

The matrix-completion problem has attracted a lot of attention, largely as a result of the celebrated Netflix Prize competition.

		Item			
		W	X	Y	Z
User	A		4.5	2.0	
	B	4.0		3.5	
	C		5.0		2.0
	D		3.5	4.0	1.0

Rating Matrix

=

A	1.2	0.8
B	1.4	0.9
C	1.5	1.0
D	1.2	0.8

User Matrix

X

	W	X	Y	Z
A	1.5	1.2	1.0	0.8
B	1.7	0.6	1.1	0.4

Item Matrix

Foundation of collaborative filtering and recommendation system.

Candès and Tao (2009), Mazumder et al. (2010):

$$\min_M \frac{1}{2} \|(X - M)_\Omega\|_F^2 + \lambda \|M\|_*$$

where  $\Omega$  is the set of available entries and  $\|M\|_*$  is the nuclear norm which is the sum of the singular values of  $M$

Rennie and Srebro (2005):

$$\min_{A,B} \frac{1}{2} \|(X - AB^T)_\Omega\|_F^2 + \lambda (\|A\|_F^2 + \|B\|_F^2)$$

where  $A$  and  $B$  have  $r$  columns.

## Sparse PCA

Goal: PC directions should be sparse (many zero loadings)

Why? Better interpretation

Shen and Huang (2006): suppose rank = 1

$$\operatorname{argmin}_{\mathbf{u}, \mathbf{v}} \|\mathbf{X} - \mathbf{u} \mathbf{v}^T\|_F^2 + \lambda \|\mathbf{u}\|_1$$

subject to  $\|\mathbf{v}\|_2 = 1$

Zou, Hastie and Tibshirani (2006):

$$\operatorname{argmin}_{\mathbf{u}, \mathbf{v}} \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\theta} \mathbf{u}^T \mathbf{x}_i\|_2^2 + \lambda \|\mathbf{u}\|_2 + \lambda_1 \|\mathbf{u}\|_1$$

subject to  $\|\boldsymbol{\theta}\|_2 = 1$

## Functional PCA

- Functional extension of PCA.
- FPCA. Suppose we observe functions  $X_1(\cdot), X_2(\cdot), \dots, X_n(\cdot)$ . We want to find an orthonormal basis  $\phi_1(\cdot), \dots, \phi_K(\cdot)$  such that

$$\sum_{i=1}^n \|X_i - \sum_{k=1}^K \langle X_i, \phi_k \rangle \phi_k\|^2$$

is minimized.

- Once such a basis is found, we can replace each curve  $X_i$  by  $\sum_{k=1}^K \langle X_i, \phi_k \rangle \phi_k$  as a good approximation.
- This means instead of working with infinitely dimensional curves  $X_i$ , we can work with  $K$ -dimensional vectors  $(\langle X_i, \phi_1 \rangle, \dots, \langle X_i, \phi_K \rangle)^\top$ .

See Ramsay, J. and Silverman, B. (1997). Functional Data Analysis, Springer, New York.

## PCA for Functional Data In Practice

- Each  $X(\cdot)$  is observed at  $p$  times and stored as a  $p$ -dimensional vector
- $n$  curves are organized as a  $n \times p$  data matrix.
- Apply the regular PCA
- The  $p$  dimensional PC direction vector is converted to the eigenfunction  $\phi_j(\cdot)$
- The PC scores are  $\langle X_i, \phi_k \rangle$

## Kernel PCA

- Goal: re-express PCA using inner products.
- Recall  $\tilde{\mathbf{X}} = \mathbf{U}\mathbf{D}\mathbf{V}'$ .
- Here  $\mathbf{U}$  are the loadings and  $\mathbf{Z} = \mathbf{D}\mathbf{V}'$  is the PC scores.
- Let  $\mathbf{K} = \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \in \mathbb{R}^{n \times n}$ . Then  $\mathbf{K} = \mathbf{V}\mathbf{D}^2\mathbf{V}'$
- Conclusion: PCA = eigen decomposition of  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} = (\mathbb{I} - \mathbf{M})^T \mathbf{X}^T \mathbf{X} (\mathbb{I} - \mathbf{M})$
- Kernel PCA: eigen decomposition of  $(\mathbb{I} - \mathbf{M})^T \mathbf{K} (\mathbb{I} - \mathbf{M})$ , where  $\mathbf{K}$  is the kernel matrix.

Note that there is no loading matrix (what is the PC direction vector in this case anyway?)



# Supervised Dimension Reduction

- Partial Least Squares:
  - Best dimension reduction of cross-covariance between  $\mathbf{X}$  and  $\mathbf{Y}$  such that factors are orthogonal to  $\mathbf{X}$ .
- Canonical Correlations Analysis:
  - Best dimension reduction of cross-covariance between  $\mathbf{X}$  and  $\mathbf{Y}$  such that bi-projection is orthogonal to  $\mathbf{X}$  or  $\mathbf{Y}$ .
- Linear Discriminant Analysis (classification):
  - Best dimension reduction of between class covariance matrix relative to within-class covariance.