## Lecture 8 Principal Component Analysis

#### Luigi Freda

#### ALCOR Lab DIAG University of Rome "La Sapienza"

December 13, 2016

#### Eigen Decomposition

- Eigenvalues and Eigenvectors
- Some Properties

#### Singular Value Decomposition

- Singular Values
- Connection with Eigenvalues
- Singular Value Decomposition



#### Eigen Decomposition

- Eigenvalues and Eigenvectors
- Some Properties

#### 2 Singular Value Decomposition

- Singular Values
- Connection with Eigenvalues
- Singular Value Decomposition
- Principal Component Analysis
   Discovering Latent Factors

# Eigenvalues

given a matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$ 

• a nonzero vector  $\mathbf{v} \in \mathbb{C}^n$  is said to be its (right) eigenvector if

 $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$ 

for some scalar  $\lambda \in \mathbb{C}$ ;  $\lambda$  is called an **eigenvalue** of **A** 

- the set of all eigenvalues of A is called its spectrum and it's denoted by  $\sigma(A)$
- the Matlab command

[V,D] = eig(A)

produces a diagonal matrix  ${\boldsymbol{\mathsf{D}}}$  of eigenvalues and a matrix  ${\boldsymbol{\mathsf{V}}}$  such that

$$AV = VD$$
 with  $D = diag(\lambda_i)$ 

if A is diagonalizable then V is full-rank (i.e. rank(V) = n) and the columns of V correspond to n linearly independent eigenvectors, i.e. V = {v<sub>1</sub>, ..., v<sub>n</sub>}; in this case one has

$$\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{V}^{-1} = \sum_{i=1}^{n} \lambda_i \mathbf{v}_i \mathbf{v}_i^T$$

if A = TAT<sup>-1</sup>, where T is a nonsingular matrix, then A and A are called similar matrices; at the previous point, A and D are similar matrices

#### 1 Eigen Decomposition

- Eigenvalues and Eigenvectors
- Some Properties

#### Singular Value Decomposition

- Singular Values
- Connection with Eigenvalues
- Singular Value Decomposition

# Principal Component Analysis Discovering Latent Factors

#### Eigenvalues Real Matrix

given a real matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ 

• all its eigenvalues  $\sigma(\mathbf{A})$  are the roots of the characteristic polynomial equation

$$\det(\lambda \mathbf{I} - \mathbf{A}) = 0$$

given that A is a real matrix, if λ<sup>\*</sup> ∈ C is an eigenvalue then its conjugate λ̄<sup>\*</sup> ∈ C is also an eigenvalue, i.e., σ(A) = σ̄(A)

$$\det(\lambda^*\mathbf{I} - \mathbf{A}) = 0 \implies \overline{\det(\lambda^*\mathbf{I} - \mathbf{A})} = \det(\overline{\lambda^*\mathbf{I}} - \overline{\mathbf{A}}) = \det(\overline{\lambda^*\mathbf{I}} - \mathbf{A}) = 0$$

•  $\sigma(\mathbf{A}) = \sigma(\mathbf{A}^T)$  since

$$det(\lambda \mathbf{I} - \mathbf{A}) = det((\lambda \mathbf{I} - \mathbf{A})^{T}) = det(\lambda \mathbf{I} - \mathbf{A}^{T})$$

• if  $\tilde{\mathbf{A}} = \mathbf{T}\mathbf{A}\mathbf{T}^{-1}$ , where **T** is a nonsingular matrix, then  $\sigma(\tilde{\mathbf{A}}) = \sigma(\mathbf{A})$  since  $\det(\lambda \mathbf{I} - \tilde{\mathbf{A}}) = \det(\lambda \mathbf{T}\mathbf{T}^{-1} - \mathbf{T}\mathbf{A}\mathbf{T}^{-1}) = \det(\mathbf{T}(\lambda \mathbf{I} - \mathbf{A})\mathbf{T}^{-1}) =$   $= \det(\mathbf{T})\det(\lambda \mathbf{I} - \mathbf{A})\det(\mathbf{T}^{-1}) = \det(\lambda \mathbf{I} - \mathbf{A}) \qquad (\det \mathbf{I} = \det(\mathbf{T}^{-1})\det \mathbf{T} = 1)$ 

(a)

o i

given a real matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ 

• eigenvectors of A associated to different eigenvalues are linearly independent

$$\mathbf{A}\mathbf{v}_{i} = \lambda_{i}\mathbf{v}_{i}, \quad \lambda_{1} \neq \lambda_{2}, \quad \mathbf{v}_{1} \neq \mathbf{v}_{2}$$

$$\alpha_{1}\mathbf{v}_{1} + \alpha_{2}\mathbf{v}_{2} = 0 \implies \mathbf{A}(\alpha_{1}\mathbf{v}_{1} + \alpha_{2}\mathbf{v}_{2}) = 0 \implies \alpha_{1}\lambda_{1}\mathbf{v}_{1} + \alpha_{2}\lambda_{2}\mathbf{v}_{2} = 0$$
since  $\lambda_{1}(\alpha_{1}\mathbf{v}_{1} + \alpha_{2}\mathbf{v}_{2}) = 0 \implies \alpha_{2}(\lambda_{2} - \lambda_{1})\mathbf{v}_{2} = 0 \implies \alpha_{2} = 0 \implies \alpha_{1} = 0$ 
if  $|\sigma(\mathbf{A})| = n$ , i.e.  $\sigma(\mathbf{A}) = \{\lambda_{1}, \lambda_{2}, ..., \lambda_{n}\}$  with  $\lambda_{i} \neq \lambda_{j}$  for  $i \neq j$ , then  $\mathbf{A}$  is diagonalizable

Image: A matrix of the second seco

э

given a real matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ 

• in general  $\sigma(\mathbf{A}) = \{\lambda_1, ..., \lambda_m\}$  with  $m \leq n$  and  $\lambda_i \neq \lambda_j$  for  $i \neq j$ , and one has

$$\det(\lambda \mathbf{I} - \mathbf{A}) = (\lambda - \lambda_1)^{\mu_a(\lambda_1)} (\lambda - \lambda_2)^{\mu_a(\lambda_2)} \dots (\lambda - \lambda_m)^{\mu_a(\lambda_m)} \qquad \sum_{i=1}^m \mu_a(\lambda_i) = n$$

where  $\mu_a(\lambda_i)$  is the algebraic multiplicity of  $\lambda_i$ 

• in the general case one can always find the Jordan canonical form

$$\mathbf{A} = \mathbf{V} \mathbf{J} \mathbf{V}^{-1} \qquad \mathbf{J} = \begin{bmatrix} \mathbf{J}_1 & & & \\ & \mathbf{J}_2 & & \\ & & \ddots & \\ & & \ddots & \mathbf{J}_p \end{bmatrix} \qquad \mathbf{J}_i = \begin{bmatrix} \lambda_i & 1 & & \\ & \lambda_i & \ddots & \\ & & \ddots & 1 \\ & & & \ddots & 1 \\ & & & & \lambda_i \end{bmatrix}$$

Image: A matrix of the second seco

m

given a real matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ 

• let 
$$E(\lambda_i) \triangleq \ker(\mathbf{A} - \lambda_i \mathbf{I}) = \{\mathbf{v} \in \mathbb{R}^n : (\mathbf{A} - \lambda_i \mathbf{I})\mathbf{v} = 0\}$$
 be the **eigenspace** of  $\lambda_i$ 

• let  $\mu_g(\lambda_i) \triangleq \dim(E(\lambda_i))$  be the geometric multiplicity of  $\lambda_i$   $(\mu_g(\lambda_i) \le \mu_a(\lambda_i))$ 

• if 
$$E(\lambda_i) = \text{span}\{\mathbf{v}_{i1}, ..., \mathbf{v}_{i\mu_i}\}$$
, with  $\mu_i \triangleq \mu_g(\lambda_i)$ , then

$$\mathbf{AV}_i = \mathbf{V}_i \mathbf{D}_i$$
 with  $\mathbf{V}_i = [\mathbf{v}_{i1}, ..., \mathbf{v}_{i\mu_i}] \in \mathbb{R}^{n \times \mu_i}$  and  $\mathbf{D}_i = \lambda_i \mathbf{I}_{\mu_i}$ 

• if 
$$\mu_g(\lambda_i) = \mu_a(\lambda_i)$$
 for each  $i \in \{1, ..., m\}$  then

$$\mathbb{R}^n = E(\lambda_1) \oplus ... \oplus E(\lambda_m)$$
  $\sum_{i=1}^n \mu_g(\lambda_i) = n$ 

#### and A is diagonalizable

Image: A matrix of the second seco

3

given a real symmetric matrix  $\mathbf{S} \in \mathbb{R}^{n \times n}$  (i.e.  $\mathbf{S} = \mathbf{S}^{T}$ )

- all eigenvalues of **S** are real, i.e.  $\sigma(\mathbf{S}) \subset \mathbb{R}$
- eigenvectors  $v_1, v_2$  of **S** associated to different eigenvalues are linearly independent and orthogonal

$$\begin{aligned} \mathbf{S}\mathbf{v}_{i} &= \lambda_{i}\mathbf{v}_{i}, \ \lambda_{1} \neq \lambda_{2}, \ \mathbf{v}_{1} \neq \mathbf{v}_{2}, \ \mathbf{v}_{2}^{T}(\mathbf{S}\mathbf{v}_{1} - \lambda_{1}\mathbf{v}_{1}) = 0 \\ \implies (\mathbf{S}\mathbf{v}_{2})^{T}\mathbf{v}_{1} - \lambda_{1}\mathbf{v}_{2}^{T}\mathbf{v}_{1} = 0 \implies \mathbf{v}_{2}^{T}\mathbf{v}_{1}(\lambda_{2} - \lambda_{1}) = 0 \implies \mathbf{v}_{2}^{T}\mathbf{v}_{1} = 0 \end{aligned}$$

• S is diagonalizable and one can find an orthonormal basis V such that  $V^T = V^{-1}$ 

$$\mathbf{S} = \mathbf{V}\mathbf{D}\mathbf{V}^{\mathsf{T}} = \sum_{i=1}^{n} \lambda_{i} \mathbf{v}_{i} \mathbf{v}_{i}^{\mathsf{T}}$$

• if  $\mathbf{S} > 0$  ( $\mathbf{S} \ge 0$ ) then  $\lambda_i > 0$  ( $\lambda_i \ge 0$ ) for  $i \in \{1, ..., m\}$ 

N.B.: the above results can be applied for instance to covariance matrices  $\boldsymbol{\Sigma}$  which are symmetric and positive definite

#### 1 Eigen Decomposition

- Eigenvalues and Eigenvectors
- Some Properties

#### Singular Value Decomposition

- Singular Values
- Connection with Eigenvalues
- Singular Value Decomposition

# Principal Component Analysis Discovering Latent Factors

idea: the singular values generalize the notion of eigenvalues to any kind of matrix

given a matrix  $\mathbf{X} \in \mathbb{R}^{N \times D}$ 

• this can be decomposed as follows

$$\mathbf{X}_{N \times D} = \mathbf{U}_{N \times N} \mathbf{S}_{N \times D} \mathbf{V}^{T}_{D \times D} = \sum_{i=1}^{r} \sigma_{i} \mathbf{u}_{i} \mathbf{v}_{i}^{T}$$

- $\mathbf{U} \in \mathbb{R}^{N \times N}$  is orthonormal,  $\mathbf{U}^T \mathbf{U} = \mathbf{U} \mathbf{U}^T = \mathbf{I}_N$
- $\mathbf{V} \in \mathbb{R}^{D \times D}$  is orthonormal,  $\mathbf{V}^T \mathbf{V} = \mathbf{V} \mathbf{V}^T = \mathbf{I}_D$
- $\sigma_1 \ge \sigma_2 \ge ... \ge \sigma_r \ge 0$  are the singular values,  $r = \min(N, D)$  and  $\mathbf{S} \in \mathbb{R}^{N \times D}$

$$\mathbf{S} = \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_D \\ \hline & \mathbf{0}_{N-D} \end{bmatrix} \text{ if } N > D \text{ or } \mathbf{S} = \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_N \end{bmatrix} \mathbf{0}_{D-N} \end{bmatrix} \text{ if } N \le D$$

#### 1 Eigen Decomposition

- Eigenvalues and Eigenvectors
- Some Properties

#### 2 Singular Value Decomposition

Singular Values

#### • Connection with Eigenvalues

• Singular Value Decomposition

# Principal Component Analysis Discovering Latent Factors

Connection with Eigenvalues

given a matrix  $\mathbf{X} \in \mathbb{R}^{N \times D}$ 

• one has 
$$\mathbf{X}_{N \times D} = \bigcup_{N \times N} \sum_{N \times D} \mathbf{V}_{D \times D}^{T} = \sum_{i=1}^{r} \sigma_{i} \mathbf{u}_{i} \mathbf{v}_{i}^{T}$$
 where  $r = \min(N, D)$ 

•  $\mathbf{U} \in \mathbb{R}^{N \times N}$  is orthonormal,  $\mathbf{U}^T \mathbf{U} = \mathbf{U} \mathbf{U}^T = \mathbf{I}_N$ 

•  $\mathbf{V} \in \mathbb{R}^{D \times D}$  is orthonormal,  $\mathbf{V}^T \mathbf{V} = \mathbf{V} \mathbf{V}^T = \mathbf{I}_D$ 

•  $\mathbf{S} \in \mathbb{R}^{N \times D}$  contains the singular values  $\sigma_i$  on the main diagonal and 0s elsewhere

#### we have

$$\mathbf{X}^{\mathsf{T}}\mathbf{X} = (\mathbf{V}\mathbf{S}^{\mathsf{T}}\mathbf{U}^{\mathsf{T}})(\mathbf{U}\mathbf{S}\mathbf{V}^{\mathsf{T}}) = \mathbf{V}\mathbf{S}^{\mathsf{T}}\mathbf{S}\mathbf{V}^{\mathsf{T}} = \mathbf{V}\mathbf{D}\mathbf{V}^{\mathsf{T}}$$

where  $\mathbf{D} \triangleq \mathbf{S}^T \mathbf{S}$  is a diagonal matrix containing the squared singular values  $\sigma_i^2$  (and possibly some additional zeros) on the main diagonal and 0s elsewhere

- since  $(\mathbf{X}^{\mathsf{T}}\mathbf{X})\mathbf{V} = \mathbf{V}\mathbf{D}$ , then **V** contains the eigenvectors of  $\mathbf{X}^{\mathsf{T}}\mathbf{X}$
- the eigenvalues of  $\mathbf{X}^T \mathbf{X}$  are the squared singular values  $\sigma_i^2$ , contained in  $\mathbf{D}$
- the columns of V are called the right singular vectors of X

Connection with Eigenvalues

given a matrix  $\mathbf{X} \in \mathbb{R}^{N \times D}$ 

• one has 
$$\mathbf{X}_{N \times D} = \bigcup_{N \times N} \sum_{N \times D} \mathbf{V}_{D \times D}^{T} = \sum_{i=1}^{r} \sigma_{i} \mathbf{u}_{i} \mathbf{v}_{i}^{T}$$
 where  $r = \min(N, D)$ 

•  $\mathbf{U} \in \mathbb{R}^{N \times N}$  is orthonormal,  $\mathbf{U}^T \mathbf{U} = \mathbf{U} \mathbf{U}^T = \mathbf{I}_N$ 

•  $\mathbf{V} \in \mathbb{R}^{D \times D}$  is orthonormal,  $\mathbf{V}^T \mathbf{V} = \mathbf{V} \mathbf{V}^T = \mathbf{I}_D$ 

•  $\mathbf{S} \in \mathbb{R}^{N \times D}$  contains the singular values  $\sigma_i$  on the main diagonal and 0s elsewhere

#### we have

$$\mathbf{X}\mathbf{X}^{\mathsf{T}} = (\mathbf{U}\mathbf{S}\mathbf{V}^{\mathsf{T}})(\mathbf{V}\mathbf{S}^{\mathsf{T}}\mathbf{U}^{\mathsf{T}}) = \mathbf{U}\mathbf{S}\mathbf{S}^{\mathsf{T}}\mathbf{U}^{\mathsf{T}} = \mathbf{U}\hat{\mathbf{D}}\mathbf{U}^{\mathsf{T}}$$

where  $\hat{\mathbf{D}} \triangleq \mathbf{SS}^{T}$  is a diagonal matrix containing the squared singular values  $\sigma_{i}^{2}$  (and possibly some additional zeros) on the main diagonal and 0s elsewhere

- since  $(XX^{T})U = U\hat{D}$ , then U contains the eigenvectors of  $XX^{T}$
- the eigenvalues of  $XX^T$  are the squared singular values  $\sigma_i^2$ , contained in  $\hat{D}$
- the columns of **U** are called the left singular vectors of **X**

#### 1 Eigen Decomposition

- Eigenvalues and Eigenvectors
- Some Properties

#### Singular Value Decomposition

- Singular Values
- Connection with Eigenvalues
- Singular Value Decomposition

# Principal Component Analysis Discovering Latent Factors

given a matrix  $\mathbf{X} \in \mathbb{R}^{N \times D}$  with N > D

• this can be decomposed as follows

$$\mathbf{X}_{N \times D} = \mathbf{U}_{N \times N} \mathbf{S}_{N \times D} \mathbf{V}^{T}_{D \times D} = \sum_{i=1}^{D} \sigma_{i} \mathbf{u}_{i} \mathbf{v}_{i}^{T}$$
$$\mathbf{S} = \begin{bmatrix} \sigma_{1} & & \\ & \ddots & \\ & & \sigma_{D} \\ \hline & & \mathbf{0}_{N-D} \end{bmatrix}$$

• the last N - D columns of **U** are irrelevant, since they will be multiplied by zero



given a matrix  $\mathbf{X} \in \mathbb{R}^{N \times D}$  with N > D

- this can be decomposed as  $\mathbf{X}_{N \times D} = \bigcup_{N \times N} \mathbf{S}_{N \times D} \mathbf{V}_{D \times D}^{T} = \sum_{i=1}^{D} \sigma_{i} \mathbf{u}_{i} \mathbf{v}_{i}^{T}$
- since the last N D columns of **U** are irrelevant, we can directly neglect them and consider the economy sized SVD

$$\mathbf{X}_{N\times D} = \underbrace{\mathbf{\hat{U}}}_{N\times D} \underbrace{\mathbf{\hat{S}}}_{D\times D} \underbrace{\mathbf{V}^{\mathsf{T}}}_{D\times D} = \sum_{i=1}^{D} \sigma_{i} \mathbf{u}_{i} \mathbf{v}_{i}^{\mathsf{T}}$$

where  $\hat{\mathbf{U}}^T \hat{\mathbf{U}} = \mathbf{I}_D$ ,  $\mathbf{V}^T \mathbf{V} = \mathbf{V} \mathbf{V}^T = \mathbf{I}_D$  and  $\hat{\mathbf{S}} = \text{diag}(\sigma_1, ..., \sigma_D)$ 



given a matrix  $\mathbf{X} \in \mathbb{R}^{N \times D}$  with N > D

- this can be decomposed as  $\mathbf{X}_{N \times D} = \bigcup_{N \times N} \mathbf{S}_{N \times D} \mathbf{V}_{D \times D}^{T} = \sum_{i=1}^{D} \sigma_{i} \mathbf{u}_{i} \mathbf{v}_{i}^{T}$
- a rank *L* approximation of **X** considers only the first *L* < *D* columns of **U** and the first *L* singular values
- we can then consider the truncated SVD

$$\underbrace{\mathbf{X}_{L}}_{N \times D} = \underbrace{\mathbf{U}_{L}}_{N \times L} \underbrace{\mathbf{S}_{L}}_{L \times L} \underbrace{\mathbf{V}_{L}}_{D \times L} = \sum_{i=1}^{L} \sigma_{i} \mathbf{u}_{i} \mathbf{v}_{i}^{T}$$

where  $\mathbf{U}_{L}^{T}\mathbf{U}_{L} = \mathbf{I}_{L}$ ,  $\mathbf{V}_{L}^{T}\mathbf{V}_{L} = \mathbf{I}_{D}$  and  $\mathbf{S}_{L} = \text{diag}(\sigma_{1}, ..., \sigma_{L})$ 



given a matrix  $\mathbf{X} \in \mathbb{R}^{N \times D}$  with N > D

• it is possible to show that  $\mathbf{X}_{L} = \mathbf{U}_{L}\mathbf{S}_{L}\mathbf{V}_{L}^{T}$  is the **best rank** *L* approximation matrix which solve

$$\mathbf{X}_L = \underset{\hat{\mathbf{X}}}{\operatorname{argmin}} \|\mathbf{X} - \hat{\mathbf{X}}\|_F$$
 with the constraint  $\operatorname{rank}(\hat{\mathbf{X}}) = L$ 

where  $\|\mathbf{A}\|_{F}$  denotes the Frobenious norm

$$\|\mathbf{A}\|_{F} = \left(\sum_{i=1}^{m}\sum_{j=1}^{n}a_{ij}^{2}\right)^{1/2} = \left(\mathsf{trace}(\mathbf{A}^{\mathsf{T}}\mathbf{A})\right)^{1/2}$$

• by replacing  $\hat{\mathbf{X}} = \mathbf{X}_L$  one obtains

$$\|\mathbf{X} - \mathbf{X}_L\|_F = \|\mathbf{U}(\mathbf{S} - \begin{bmatrix} \mathbf{S}_L & \\ & \mathbf{0} \end{bmatrix})\mathbf{V}^T\|_F = \|\mathbf{S} - \begin{bmatrix} \mathbf{S}_L & \\ & \mathbf{0} \end{bmatrix}\|_F = \sum_{i=L+1}^r \sigma_i$$

#### 1 Eigen Decomposition

- Eigenvalues and Eigenvectors
- Some Properties

#### 2 Singular Value Decomposition

- Singular Values
- Connection with Eigenvalues
- Singular Value Decomposition

#### Optimize Principal Component Analysis

Discovering Latent Factors

# Discovering Latent Factors

An example



- dimensionality reduction: it is often useful to reduce the dimensionality by projecting the data to a lower dimensional subspace which captures the "essence" of the data
- for example in the above plot:
  - the 2d approximation is quite good, most points lie close to this subspace
  - projecting points onto the red line 1d approx is a rather poor approximation

# Discovering Latent Factors

Motivations



- although data may appear high dimensional, there may only be a small number of degrees of variability, corresponding to latent factors
- low dimensional representations, when used as input to statistical models, often result in **better predictive accuracy** (focusing on the "essence")
- low dimensional representations are useful for **enabling fast** nearest neighbor **searches**
- two dimensional projections are very useful for visualizing high dimensional data

• the most common approach to dimensionality reduction is called **Principal Components Analysis** or **PCA** 

• given  $\mathbf{x}_i \in \mathbb{R}^D$ , we compute an approximation  $\hat{\mathbf{x}}_i = \mathbf{W}\mathbf{z}_i$  with

- $\mathbf{z}_i \in \mathbb{R}^L$  where L < D (dimensionality reduction)
- $\mathbf{W} \in \mathbb{R}^{D \times L}$  and  $\mathbf{W}^T \mathbf{W} = \mathbf{I}$  (orthogonal  $\mathbf{W}$ )

so as to minimize the reconstruction error

$$J(\mathbf{W}, \mathbf{Z}) = \frac{1}{N} \sum_{i=1}^{N} \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 = \frac{1}{N} \sum_{i=1}^{N} \|\mathbf{x}_i - \mathbf{W}\mathbf{z}_i\|^2$$

• since  $\mathbf{W}^T \mathbf{W} = \mathbf{I}$  one has  $\|\mathbf{x}\|_2 = (\mathbf{x}^T \mathbf{x})^{1/2} = (\mathbf{z}^T \mathbf{W}^T \mathbf{W} \mathbf{z})^{1/2} = (\mathbf{z}^T \mathbf{z})^{1/2} = \|\mathbf{z}\|_2$ 

- in 3D or 2D W can represent part of a rotation matrix
- this can be thought of as an unsupervised version of (multi-output) linear regression, where we observe the high-dimensional response y = x, but not the low-dimensional "cause" z

• the objective function (reconstruction error)

$$J(\mathbf{W}, \mathbf{Z}) = \frac{1}{N} \sum_{i=1}^{N} \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 = \frac{1}{N} \sum_{i=1}^{N} \|\mathbf{x}_i - \mathbf{W}\mathbf{z}_i\|^2$$

is equivalent to

$$J(\mathbf{W}, \mathbf{Z}) = \|\mathbf{X}^T - \mathbf{W}\mathbf{Z}^T\|_F^2$$

where  $\|\mathbf{A}\|_{F}$  denotes the **Frobenious norm** 

$$\|\mathbf{A}\|_{F} = \left(\sum_{i=1}^{m} \sum_{j=1}^{n} a_{ij}^{2}\right)^{1/2} = \left(\mathsf{trace}(\mathbf{A}^{\mathsf{T}}\mathbf{A})\right)^{1/2}$$

we have to minimize the objective function

$$J(\mathbf{W}, \mathbf{Z}) = rac{1}{N} \sum_{i=1}^{N} \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 = rac{1}{N} \sum_{i=1}^{N} \|\mathbf{x}_i - \mathbf{W}\mathbf{z}_i\|^2$$

or equivalently

$$J(\mathbf{W}, \mathbf{Z}) = \|\mathbf{X}^{T} - \mathbf{W}\mathbf{Z}^{T}\|_{F}^{2}$$

subject to  $\mathbf{W} \in \mathbb{R}^{D \times L}$  and  $\mathbf{W}^T \mathbf{W} = \mathbf{I}$ 

• it is possible to prove that the **optimal solution** is obtained by setting  $\hat{W} = V_L$  $\hat{Z} = XW$   $(z_i = \hat{W}^T x_i = V_L^T x_i)$ 

where  $\mathbf{X}_{L} = \mathbf{U}_{L} \mathbf{S}_{L} \mathbf{V}_{L}^{T}$  is the rank L truncated SVD of **X** 

- we have  $\mathbf{X}^T \mathbf{X} = \mathbf{V} \mathbf{S}^2 \mathbf{V}^T$  hence  $(\mathbf{X}^T \mathbf{X}) \mathbf{V}_L = \mathbf{V}_L \mathbf{S}_L^2$  since  $\mathbf{V} = [\mathbf{V}_L | *]$
- $V_L$  contains the *L* eigenvectors (principal directions) with the largest eigenvalues of the empirical covariance matrix  $\Sigma = \frac{1}{N} X^T X = \frac{1}{N} \sum_{i=1}^{N} x_i x_i^T$
- $\mathbf{z}_i = \mathbf{V}_L^T \mathbf{x}_i$  is the orthogonal projection of  $\mathbf{x}_i$  on the eigenvectors in  $\mathbf{V}_L$

• by replacing  $\mathbf{Z} = \mathbf{X}\mathbf{W}$  in the objective function and recalling that  $\mathbf{W}^{\mathsf{T}}\mathbf{W} = \mathbf{I}$ 

$$J(\mathbf{Z}) = \|\mathbf{X}^{T} - \mathbf{W}\mathbf{Z}^{T}\|_{F}^{2} = \operatorname{trace}\left((\mathbf{X}^{T} - \mathbf{W}\mathbf{Z}^{T})^{T}(\mathbf{X}^{T} - \mathbf{W}\mathbf{Z}^{T})\right) =$$

$$= \operatorname{trace}(\mathbf{X}\mathbf{X}^{T} - \mathbf{X}\mathbf{W}\mathbf{Z}^{T} - \mathbf{Z}\mathbf{W}^{T}\mathbf{X}^{T} + \mathbf{Z}\mathbf{Z}^{T}) = \operatorname{trace}(\mathbf{X}^{T}\mathbf{X}) - \operatorname{trace}(\mathbf{W}^{T}\mathbf{X}^{T}\mathbf{X}\mathbf{W}) =$$
$$= \operatorname{trace}(\mathbf{X}^{T}\mathbf{X}) - \operatorname{trace}(\mathbf{Z}^{T}\mathbf{Z}) = \operatorname{trace}(N\Sigma_{X}) - \operatorname{trace}(N\Sigma_{Z})$$

where we used the following facts: trace(**AB**) = trace(**BA**), trace(**A**<sup>T</sup>) = trace(**A**)

 minimizing the reconstruction error is equivalent to maximizing the variance of the projected data z<sub>i</sub> = Wx<sub>i</sub>

- the principal directions are the ones along which the data shows maximal variance
- this means that PCA can be "misled" by directions in which the variance is high merely because of the measurement scale
- in the figure below, the vertical axis (weight) uses a large range than the horizontal axis (height), resulting in a line that looks somewhat "unnatural"
- it is therefore standard practice to standardize the data first

$$x_{ij} 
ightarrow rac{(x_{ij}-\mu_j)}{\sigma_j}$$

where  $\mu_j = \frac{1}{N} \sum_{i=1}^N x_{ij}$  and  $\sigma_j^2 = \frac{1}{N-1} \sum_{i=1}^N (x_{ij} - \mu_j)^2$ 





- left: the mean and the first three PC basis vectors (eigendigits) based on 25 images of the digit 3 (from the MNIST dataset)
- right: reconstruction of an image based on 2, 10, 100 and all the basis vectors

- ∢ 🗇 እ



- low rank approximations to an image
- top left: the original image is of size 200  $\times$  320, so has rank 200
- subsequent images have ranks 2, 5, and 20.

• Kevin Murphy's book

< □ > < ---->

2