

# Data Science

## Principal Component Analysis

Linear latent decomposition

Stéphane Marchand-Maillet

Department of Computer Science



UNIVERSITÉ  
DE GENÈVE

FACULTÉ DES SCIENCES



Master en Sciences Informatiques - Autumn semester

# Table of contents

Motivation

Formalization

- Variance maximization

- Error minimization

PCA for dimension reduction

- Latent space

- Approximation

Practical PCA

Alternatives

- Linear AutoEncoder

- Further formulations

# What is this lecture about?

- ★ Base representations may not be optimal (to be defined)
- ★ **Latent models** promise to exhibit the underlying (latent) factors that drive the process in question
- ★ This initial (but fundamental) definition of latent factors uses statistical correlation
- ⇒ It exhibits **linear** latent factors
- ⇒ Enables “simplification” of the data by sound decimation
- ★ Also: we will study several interpretations

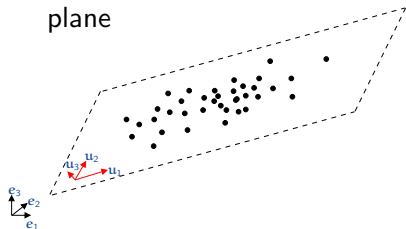
**Reading:** [1] (chap 12) and [3] (chap 3 and 10.3)

# Intuition

Given  $\mathcal{X} \subset \Omega$ , we wish to decompose  $\Omega$  into subspaces such that the projection of  $\mathcal{X}$  onto these subspaces retains the most “information”.

Q: What information should we consider?

- ★ Say  $\mathcal{X}$  is almost “contained” into a 2D plane in a 3D space
- ★ A relevant choice for our subspace is to choose a basis  $\{\mathbf{u}_1, \mathbf{u}_2\}$  for the plane



Q: What characterizes  $\{\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3\}$ ?

- ⇒ the fact that the data varies most along these directions
- ⇒ the fact that the data varies least orthogonally to these directions ( $\mathbf{u}_3$ )

# Formalization

Given  $\mathcal{X}$ ,  $\{\mathbf{u}_i\}_{i \in \llbracket D \rrbracket}$  is a new orthonormal basis of  $\mathbb{R}^D$ . The **Principal Component**  $\mathbf{u}_1$  is chosen such that the variance of the data projected over  $\mathbf{u}_1$  is maximum.  $\mathbf{u}_2$  is chosen using  $\text{Proj}_{\mathbf{u}_1^\perp}(\mathcal{X})$ .

## Model

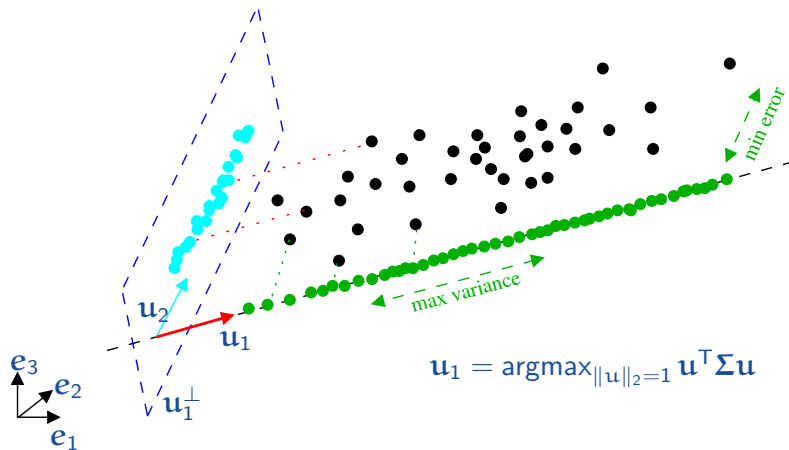
Given  $\mathcal{X}$ , the sample mean  $(\bar{\mathbf{x}})$  and the variance of the data projected over  $\mathbf{u}_1$  are

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \quad \text{and} \quad v_{\mathbf{u}_1} = \frac{1}{N} \sum_{i=1}^N (\mathbf{u}_1^T \mathbf{x}_i - \mathbf{u}_1^T \bar{\mathbf{x}})^2 = \mathbf{u}_1^T \boldsymbol{\Sigma} \mathbf{u}_1$$

where  $\boldsymbol{\Sigma} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$  is the data covariance matrix. Therefore

$$\mathbf{u}_1 = \underset{\mathbf{u}^T \mathbf{u} = 1}{\operatorname{argmax}} \mathbf{u}^T \boldsymbol{\Sigma} \mathbf{u}$$

# Intuition



# Formalization

$$\mathbf{u}_1 = \underset{\mathbf{u}^T \mathbf{u} = 1}{\operatorname{argmax}} \mathbf{u}^T \boldsymbol{\Sigma} \mathbf{u} \quad \Rightarrow \quad J(\mathbf{u}) = \mathbf{u}^T \boldsymbol{\Sigma} \mathbf{u} + \lambda(1 - \mathbf{u}^T \mathbf{u})$$

So that

$$\frac{\partial J(\mathbf{u})}{\partial \mathbf{u}} \Big|_{\mathbf{u}=\mathbf{u}_1} = 0 \quad \text{and} \quad \frac{\partial J(\mathbf{u})}{\partial \lambda} \Big|_{\lambda=\lambda_1} = 0$$

Hence

$$\boldsymbol{\Sigma} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1 \quad \Rightarrow \quad v_{\mathbf{u}_1} = \mathbf{u}_1^T \boldsymbol{\Sigma} \mathbf{u}_1 = \lambda_1$$

$\Rightarrow (\mathbf{u}_1, \lambda_1)$  is an eigenpair of the covariance matrix  $\boldsymbol{\Sigma}$

$\Rightarrow$  continuing with the decimation process, we obtain the set of **Principal Components** as the eigenpairs  $\{(\mathbf{u}_i, \lambda_i)\}_{i \in \llbracket D \rrbracket}$  of the covariance matrix  $\boldsymbol{\Sigma}$  of data  $\mathcal{X}$

$$\boldsymbol{\Lambda} = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_D \end{pmatrix} = \begin{pmatrix} | & & | \\ \mathbf{u}_1 & \cdots & \mathbf{u}_D \\ | & & | \end{pmatrix}^T \boldsymbol{\Sigma} \begin{pmatrix} | & & | \\ \mathbf{u}_1 & \cdots & \mathbf{u}_D \\ | & & | \end{pmatrix} = \mathbf{U}^T \boldsymbol{\Sigma} \mathbf{U}$$

# Formalization

- ★ The variance  $v_{\mathbf{u}_1}$  is expressed as a sum of squares

$$v_{\mathbf{u}_1} = \frac{1}{N} \sum_{i=1}^N (\mathbf{u}_1^T (\mathbf{x}_i - \bar{\mathbf{x}}))^2$$

⇒ To maximize  $v_{\mathbf{u}_1}$ , terms  $\mathbf{u}_1^T (\mathbf{x}_i - \bar{\mathbf{x}})$  should be collectively maximized

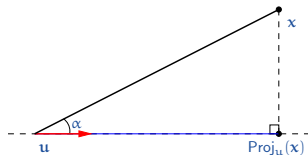
⇒ Since  $(\mathbf{x}_i - \bar{\mathbf{x}})$  is fixed, Pythagoras tells us it is equivalent to minimize the distance to the axis of projection (**approximation error**)

⇒ A Principal Component is a **quadratic regression** over the data

$$\mathbf{u}_1 = \underset{\mathbf{u}^T \mathbf{u} = 1}{\operatorname{argmin}} \sum_{i=1}^N \|(\mathbf{x}_i - \bar{\mathbf{x}}) - [\mathbf{u}^T (\mathbf{x}_i - \bar{\mathbf{x}})] \mathbf{u}\|_2^2 \Rightarrow \Sigma \mathbf{u}_1 = \lambda_1 \mathbf{u}_1 \quad \textcircled{\varphi}$$

Maximize Projection Variance

⇔ Minimize Approximation Error





# Physical interpretation

- ★ Consider a physical system  $\mathcal{X}$  with masses  $m_i = 1$  at positions  $\mathbf{x}_i$
- ★ The **inertia** of the system w.r.t  $\mathbf{a} \in \Omega$  is  $I_{\mathbf{a}}(\mathcal{X}) = \sum_{i=1}^N d^2(\mathbf{a}, \mathbf{x}_i)$
- ★ Huygens theorem tells us that if  $\mathbf{g} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$  then

$$I_{\mathbf{a}}(\mathcal{X}) = d^2(\mathbf{a}, \mathbf{g}) + I_{\mathbf{g}}(\mathcal{X}) \quad (\text{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}X)^2)$$

- ★ and if  $\mathcal{F}$  is a subspace of  $\Omega$  going thru  $\mathbf{g}$  then

$$I_{\mathcal{F}}(\mathcal{X}) = \sum_{i=1}^N d^2(\mathcal{F}, \mathbf{x}_i) \quad \text{where} \quad d(\mathcal{F}, \mathbf{x}_i) = \|\mathbf{x}_i - \text{Proj}_{\mathcal{F}}(\mathbf{x}_i)\|$$

$\Rightarrow$  A Principal Component is a **subspace of least inertia w.r.t  $\mathcal{X}$**

# Structure of the latent space

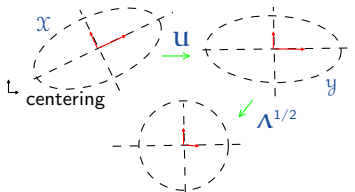
The **latent space** with basis  $\{\mathbf{u}_1, \dots, \mathbf{u}_D\}$  has the following properties:

- ★ By construction  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D$  and  $\lambda \stackrel{\text{def}}{=} \sum_{d=1}^D \lambda_d = \text{Tr}(\Sigma)$
- ★  $v_{\mathbf{u}_d} = \lambda_d \Rightarrow \lambda$  represents the total variance  $\Rightarrow$  latent features are decorrelated  $\mathbf{u}_d^T \mathbf{u}_{d'} = 0$  ( $\Lambda$  is the diagonal **latent covariance matrix**)
- ★ The basis  $\{\mathbf{u}_1, \dots, \mathbf{u}_D\}$  induces **latent coordinates**  $\mathbf{y}_i$  for the data:

$$\mathbf{y}_i(d) = \langle \mathbf{u}_d, \mathbf{x}_i - \bar{\mathbf{x}} \rangle = \mathbf{u}_d^T (\mathbf{x}_i - \bar{\mathbf{x}}) \quad \text{so that} \quad \mathbf{y}_i = \mathbf{U}^T (\mathbf{x}_i - \bar{\mathbf{x}})$$

Ⓢ The transform is linear and composed of:

- Centering on  $\bar{\mathbf{x}}$
- Rotation using  $\mathbf{U}$
- Scaling using  $\Lambda^{1/2}$  (whitening)

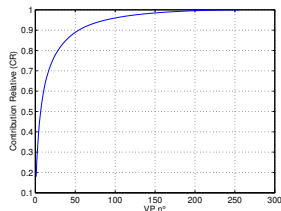
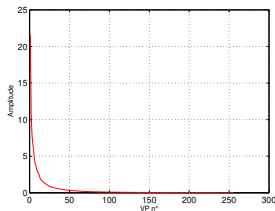


# Structure of the model

- ★ The latent space preserves the variance (of the centered data)
- ⇒ the underlying data model is  $\mathbf{X}_i \sim \mathbf{f}_X = \mathcal{N}(\bar{\mathbf{x}}, \Sigma)$
- ⇒ PCA will **not** be relevant for non-Gaussian data (e.g clustered)
- ⚠ **Important:** So far PCA is an **exact** transform

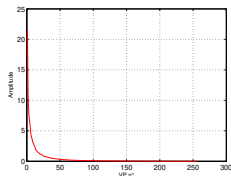
$$\mathbf{y}_i = \mathbf{U}^T(\mathbf{x}_i - \bar{\mathbf{x}}) \quad \text{so that} \quad \mathbf{U}\mathbf{y}_i = \mathbf{U}\mathbf{U}^T(\mathbf{x}_i - \bar{\mathbf{x}}) = \mathbf{x}_i - \bar{\mathbf{x}}$$

- ⇒ at this stage, one purpose is to study the spectrum  $\{\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D\}$  of the data

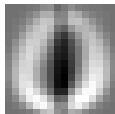


# Decomposition via PCA

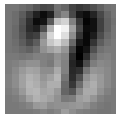
Ⓢ MNIST partial dataset:  $N = 7291$  images  $16 \times 16$  (8bits)  $\Rightarrow D = 256$



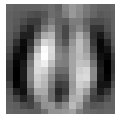
CP n°1

 $cr(\Delta_1)=18\%$ 

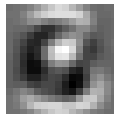
CP n°2

 $cr(\Delta_2)=27\%$ 

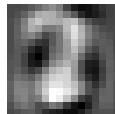
CP n°3

 $cr(\Delta_3)=33\%$ 

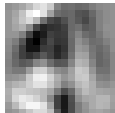
CP n°4

 $cr(\Delta_4)=39\%$ 

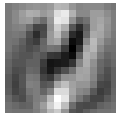
CP n°5

 $cr(\Delta_5)=44\%$ 

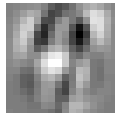
CP n°6

 $cr(\Delta_6)=48\%$ 

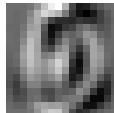
CP n°7

 $cr(\Delta_7)=51\%$ 

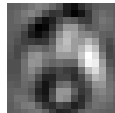
CP n°8

 $cr(\Delta_8)=54\%$ 

CP n°9

 $cr(\Delta_9)=57\%$ 

CP n°10

 $cr(\Delta_{10})=59\%$

# Approximation via PCA

**Low-rank approximation** from the Eckart and Young theorem:

If  $\Sigma = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$  and for  $K < D$  define  $\Sigma_K \stackrel{\text{def}}{=} \sum_{d=1}^K \lambda_d \mathbf{u}_d \mathbf{u}_d^T$  then

$$\operatorname{argmin}_{\operatorname{rank}(\mathbf{S})=K} \|\Sigma - \mathbf{S}\|_F^2 = \Sigma_K \quad \text{and} \quad \|\Sigma - \Sigma_K\|_F^2 = \sum_{d=K+1}^D \lambda_d$$

$\Rightarrow \Sigma_K$  is the closest  $K$ -rank matrix to  $\Sigma$

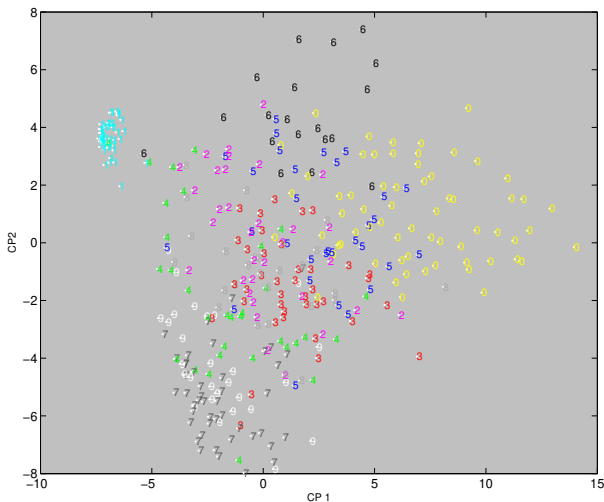
Truncation (of  $\mathbf{U}$  and  $\mathbf{\Lambda}$ )

$$\Sigma_K = \begin{pmatrix} | & & | \\ \mathbf{u}_1 & \cdots & \mathbf{u}_K \\ | & & | \end{pmatrix} \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_K \end{pmatrix} \begin{pmatrix} | & & | \\ \mathbf{u}_1 & \cdots & \mathbf{u}_K \\ | & & | \end{pmatrix}^T = \mathbf{u}_K \mathbf{\Lambda}_K \mathbf{u}_K^T$$

$$\triangle \Rightarrow \tilde{\mathbf{y}}_i = \mathbf{u}_K^T \mathbf{x}_i \in \mathbb{R}^K \text{ and } \tilde{\mathbf{x}}_i = \mathbf{u}_K \mathbf{u}_K^T \mathbf{x}_i + \bar{\mathbf{x}}$$

# Visualization via PCA

Ⓢ MNIST partial dataset:  $K = 2 \Rightarrow \tilde{\mathbf{y}}_i \in \mathbb{R}^2$



# Geometry of PCA

The quality of reconstruction can be measured by

- ★ the **relative contribution** of each dimension to the variance

$$c_d = \frac{\lambda_d}{\sum_k \lambda_k}$$

⇒ depends on the distribution of the spectrum

- ★ the **projection ratio** of each data  $\mathbf{x}_i$  over a latent factor

$$\rho_d(\mathbf{x}_i) = \frac{\langle \mathbf{u}_d, \mathbf{x}_i \rangle^2}{\|\mathbf{x}_i\|^2} = \cos^2(\angle(\mathbf{u}_d, \mathbf{x}_i))$$

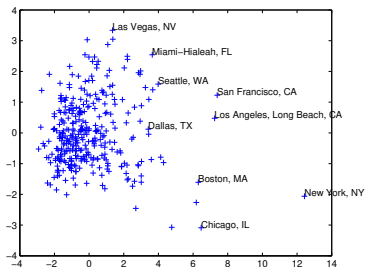
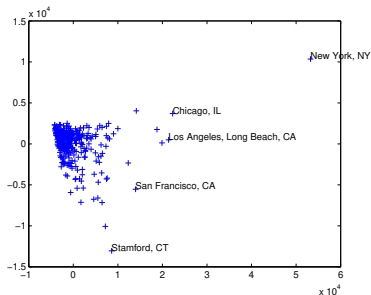
⇒ the closer  $\rho_d(\mathbf{x}_i)$  is to 1, the more  $\mathbf{x}_i$  lies on  $\mathbf{u}_d$

⇒ The above can be grouped (summed) to evaluate wrt a subspace  $\{\mathbf{u}_1, \mathbf{u}_2, \dots\}$

# Geometry of PCA

Every original data comes with its unit (scale), that we can estimate via  $\sigma_d^2$  the sample variance along original dimension  $d$ . PCA is more effective if all scales are similar.

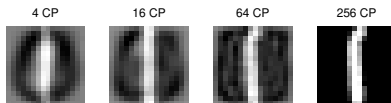
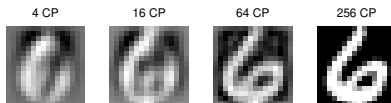
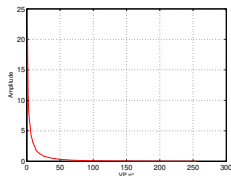
$\Rightarrow$  we create the scaling matrix  $\mathbf{S} = \text{diag}[\sigma_1^2, \dots, \sigma_d^2]$  and we define the metric  $d_S^2(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^\top \mathbf{S}^{-1} (\mathbf{x} - \mathbf{y}) \Rightarrow$  in that metric space, the covariance matrix  $\Sigma_S$  is also rescaled (into the **correlation matrix**) and used as a base for PCA.





# Approximation via PCA

Ⓢ MNIST partial dataset:  $N = 7291$  images  $16 \times 16$  (8bits)  $\Rightarrow D = 256$

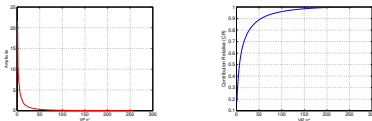


# Practical PCA

Given  $\mathcal{X} \in \Omega$

- ★ Compute the sample mean  $\bar{\mathbf{x}}$
- ★ Center the data  $\mathbf{x}_i \leftarrow (\mathbf{x}_i - \bar{\mathbf{x}})$  and form centered data matrix  $\mathbf{X}$
- ★  $\Sigma = \frac{1}{N} \mathbf{X} \mathbf{X}^T$  and  $\Sigma = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$  ←  $\Delta$  Exact transform so far
- ★ Select the number of components  $K$
- ★ Define  $\mathbf{U}_K$ ,  $\mathbf{\Lambda}_K$  and compute  $\{\tilde{\mathbf{y}}_i\}_{i \in \llbracket N \rrbracket}$  and/or  $\{\tilde{\mathbf{x}}_i\}_{i \in \llbracket N \rrbracket}$

## Choice of $K$



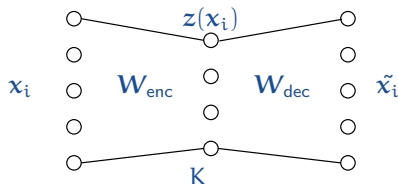
1.  $K = d$ , the target dimension  $\Rightarrow \tilde{\mathbf{y}}_i \in \mathbb{R}^d$
2. Require  $\text{Var}(\mathcal{Y}) = \tau \cdot \text{Var}(\mathcal{X}) \Rightarrow K$  such that  $\frac{\sum_{d=1}^K \lambda_d}{\sum_{d=1}^D \lambda_d} \geq \tau$
3. Train  $K$  such that  $\mathcal{L}(\mathcal{X}) \leq \varepsilon$  (e.g  $\mathcal{L}(\mathcal{X}) = \sum_i \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|^2$ )

# PCA and linear AutoEncoders

A linear AE is a simple structure

- ★  $\mathbf{W}_{\text{enc}}$  and  $\mathbf{W}_{\text{dec}}$  : encoder and decoder weights

$$\mathbf{z}(\mathbf{x}_i) = \mathbf{W}_{\text{enc}} \mathbf{x}_i \quad \tilde{\mathbf{x}}_i = \mathbf{W}_{\text{dec}} \mathbf{z}(\mathbf{x}_i)$$



We optimize the weight matrices

$$\theta^* = \underset{\mathbf{W}_1, \mathbf{W}_2}{\operatorname{argmin}} \sum_{i=1}^N \frac{1}{2} \sum_{d=1}^D (\mathbf{x}_i[d] - \tilde{\mathbf{x}}_i[d])^2 = \underset{\operatorname{rank}(\mathbf{W})=K}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{W}\mathbf{Z}\|_F^2$$

Using again the Eckart and Young theorem with  $\mathbf{X} = \mathbf{U}\Psi\mathbf{V}^T$  then the solution is  $\mathbf{W}_{\text{dec}}\mathbf{Z} = \mathbf{U}_K\Psi_K\mathbf{V}_K^T$ .

Setting  $\mathbf{W}_{\text{dec}} = \mathbf{U}_K\Psi_K$ , then clearly  $\mathbf{W}_{\text{enc}} = \Psi_K^{-1}\mathbf{U}_K^T$

$$\mathbf{W}_{\text{enc}}\mathbf{X} = \mathbf{Z} = \mathbf{V}^T = \mathbf{V}^T(\mathbf{X}^T\mathbf{X})^{-1}(\mathbf{X}^T\mathbf{X}) = \mathbf{V}^T(\mathbf{V}\Psi\Psi\mathbf{V}^T)^{-1}(\mathbf{V}\Psi\mathbf{U}^T)\mathbf{X} = \Psi_K^{-1}\mathbf{U}_K^T\mathbf{X}$$



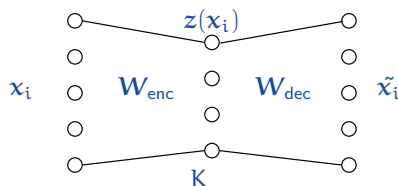
Note:  $\mathbf{W}_{\text{dec}} = \mathbf{W}_{\text{dec}} = \text{Id}_D$  cannot be a solution if  $K < D$  and  $\mathbf{W}_{\text{enc}}$  is not the inverse of  $\mathbf{W}_{\text{dec}}$

# PCA and linear AutoEncoders

A linear AE is a simple structure

- ★  $\mathbf{W}_{\text{enc}}$  and  $\mathbf{W}_{\text{dec}}$  : encoder and decoder weights

$$\mathbf{z}(\mathbf{x}_i) = \mathbf{W}_{\text{enc}} \mathbf{x}_i \quad \tilde{\mathbf{x}}_i = \mathbf{W}_{\text{dec}} \mathbf{z}(\mathbf{x}_i)$$



We optimize the weight matrices

$$\mathbf{X} = \mathbf{U}\mathbf{\Psi}\mathbf{V}^T \text{ and } \mathbf{W}_{\text{dec}} = \mathbf{U}\mathbf{\Psi} \text{ and } \mathbf{W}_{\text{enc}} = \mathbf{\Psi}^{-1}\mathbf{U}^T \Rightarrow \tilde{\mathbf{x}}_i = \mathbf{u}_K \mathbf{u}_K^T \mathbf{x}_i$$

Relation to PCA (centered data)

- ★ PCA:  $\mathbf{\Sigma} = \frac{1}{N} \mathbf{X} \mathbf{X}^T = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$  so that  $\tilde{\mathbf{x}}_i = \mathbf{u}_K \mathbf{u}_K^T \mathbf{x}_i$
  - ★ AE:  $\mathbf{X} = \mathbf{U} \mathbf{\Psi} \mathbf{V}^T \Rightarrow \mathbf{\Sigma} = \frac{1}{N} \mathbf{X} \mathbf{X}^T = \frac{1}{N} \mathbf{U} \mathbf{\Psi}^2 \mathbf{U}^T$  and  $\tilde{\mathbf{x}}_i = \mathbf{u}_K \mathbf{u}_K^T \mathbf{x}_i$
- $$\Rightarrow \mathbf{\Lambda} = \frac{1}{N} \mathbf{\Psi}^2 = \left( \frac{1}{\sqrt{N}} \mathbf{\Psi} \right) \left( \frac{1}{\sqrt{N}} \mathbf{\Psi} \right) \text{ so that } \text{PCA}(\mathbf{X}) \leftrightarrow \text{AE} \left( \frac{1}{\sqrt{N}} \mathbf{X} \right)$$

$\Rightarrow$  A linear AE performs a PCA if the data is centered and scaled

Note:  $\mathbf{W}_{\text{dec}} = \mathbf{W}_{\text{enc}} = \text{Id}_D$  cannot be a solution if  $K < D$  and  $\mathbf{W}_{\text{enc}}$  is not the inverse of  $\mathbf{W}_{\text{dec}}$

# PCA and linear AutoEncoders

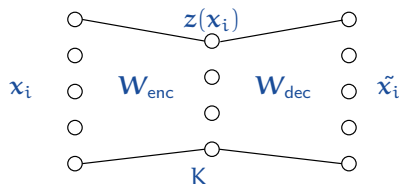
A linear AE is a simple structure

- ★  $\mathbf{W}_{\text{enc}}$  and  $\mathbf{W}_{\text{dec}}$  : encoder and decoder weights

$$\mathbf{z}(\mathbf{x}_i) = \mathbf{W}_{\text{enc}} \mathbf{x}_i \quad \tilde{\mathbf{x}}_i = \mathbf{W}_{\text{dec}} \mathbf{z}(\mathbf{x}_i)$$

We optimize the weight matrices

⇒ A linear AE performs a PCA if the data is centered and scaled



- ★ If adding hidden layers, the latent space becomes non-linear

⇒ non-linear AutoEncoders

- ★ Changing the loss function (ELBO) enables sparsity in the latent space

⇒ Variational AutoEncoders (VAE)

Note:  $\mathbf{W}_{\text{dec}} = \mathbf{W}_{\text{enc}} = \text{Id}_D$  cannot be a solution if  $K < D$  and  $\mathbf{W}_{\text{enc}}$  is not the inverse of  $\mathbf{W}_{\text{dec}}$

## Alternative formulations

- ★ **Probabilistic PCA** reformulates PCA with an explicit latent distribution  $\mathbf{P}(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$  and the conditional model is  $\mathbf{P}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \sigma^2 \mathbf{I}_D)$  so that
    - the model accomodates “measurement noise” (with variance  $\sigma^2$ )
    - the model can run in a generative mode
    - parameters can be estimated via Maximum Likelihood
    - an EM algorithm (see later) can be derived for saving computations
    - Bayesian PCA reverts the conditional so as to find  $K$  by training
  - ★ **Kernel PCA** embarks a **nonlinear mapping**  $\boldsymbol{\phi}(\mathbf{x}_i)$  via a kernel function  $k(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\phi}(\mathbf{x}_i)^T \boldsymbol{\phi}(\mathbf{x}_j)$  to perform PCA within a more favorable space
  - ★ **Local PCA** perform PCA on data neighborhoods to consider the **local intrinsic dimensionality** only
- ⚠ The **limitation of PCA** is often the decomposition of  $\boldsymbol{\Sigma}$  in  $O(D^3)$

# Summary

- ★ PCA is part of the linear latent models
- ★ PCA applies on centered data and uses variance as a criteria for decomposition
- ★ PCA is an exact complete decomposition into decorrelated components
- ★ PCA assumes a Normal distribution of the data
- ★ PCA can be equivalently formulated as a regression
- ★ PCA offers a sound decimation strategy based on a low rank approximation
- ★ PCA can be used for denoising via a Gaussian noise model
- ★ PCA is equivalent to a linear AutoEncoder
- ★ PCA may be given a stochastic formulation
- ★ PCA may be generalized to the non-linear case via kernels

## Example questions [mostly require formal – mathematical – answers]

- ★ Explain how PCA uses variance as a criterion
- ★ Show that variance maximization is equivalent to error minimization
- ★ Show how PCA uses the Eckart-Young theorem
- ★ Given some data, how do you apply PCA?
- ★ What information does it provide you with?
- ★ How do you reconstruct data with  $K < D$  components?
- ★ How do you select the components to keep?
- ★ Can you apply PCA on any data?
- ★ Is it relevant to apply PCA on any data?
- ★ How can I apply PCA over clustered data?
- ★ Show that PCA is equivalent to a linear AE

Note: Make sure you can explain in detail what is: linear transform, orthogonal matrix, coordinate, rank, mean, variance, projection, eigen decomposition, trace, Frobenius norm, Lagrange Multiplier



# References I

- [1] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, Berlin, Heidelberg, 2006. (available online).
- [2] Avrim Blum, John Hopcroft, and Ravindran Kannan. *Foundations of Data Science*. Cambridge University Press, 2020. (available online).
- [3] Richard O. Duda, Peter E. Hart, and David G. Stork. *Pattern Classification*. Wiley, New York, 2 edition, 2001.