

# Lecture 4 21/10/2024

## AUTOENCODERS

PCA ↗ Classical MDS  $\Delta_{i,l}^2 = \|x^i - x^l\|^2$   
MDS

↳ ISOMAP  $\Delta_{i,l} \rightarrow$  geodesic distance

Gram matrix  $\xrightarrow[\text{trick}]{\text{kernel}}$  PCA

Not classical MDS  $\begin{cases} \nearrow \text{Diffusion Maps} \\ \searrow \text{Sketch Map} \end{cases}$

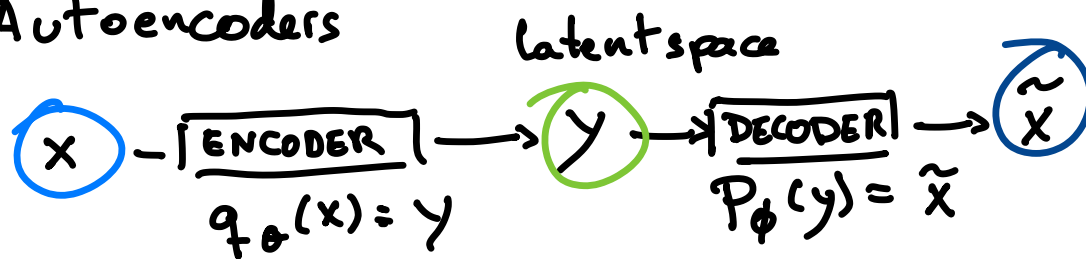
MDS  $L = f(\Delta, \Theta)$

Probabilistic Embeddings

$\Delta \rightarrow \mathcal{P}(N)$  Original space  
...  $\downarrow$  ...  $\downarrow$  ...  
 $\Theta \leftarrow Q(N)$  Projected Space  
(y)

*Note: A red arrow points from  $\Delta$  to  $Q(N)$ , and a blue arrow points from  $\Theta$  to  $\Delta$ .*

# Autoencoders



$D \dots \dots d \dots \dots D$

$q_\theta$  &  $p_\phi$  are linear functions

$$\left. \begin{array}{l} q_\theta = W \\ p_\phi = U \end{array} \right\} \begin{array}{l} y^i = W x^i \\ \tilde{x}^i = U y^i \end{array}$$

$$\tilde{x}^i = \underbrace{UW}_{\text{mapping function}} x^i$$

$$L = \sum_i \|x^i - \tilde{x}^i\|^2 \quad \text{Reconstruction error}$$

$$\min_{U \in \mathbb{R}^{D,d}, W \in \mathbb{R}^{d,D}} L = \sum_i \|x^i - UWx^i\|^2 \quad \text{with respect}$$

(A)  $\operatorname{argmin}_{U, W} \sum_i \|x^i - UWx^i\|^2$  implies

$$W = U^T; \quad U^T U = \mathbb{I}$$

(B) minimizing  $L$  is equivalent to maximize the  $\operatorname{Tr}$  of the covariance matrix on the latent space

$$\textcircled{A} \quad \hat{x}^i = U W x^i$$

$$\dots V \quad \boxed{U W x^i = V \alpha^i}$$

$$\ell = \|x - V\alpha\|^2 = \|x\|^2 + \alpha^T \underbrace{V^T V}_{\mathbb{I}} \alpha - 2\alpha^T V^T x$$

$$= \|x\|^2 + \underbrace{\alpha^T \alpha}_{\|\alpha\|^2} - 2\alpha^T (V^T x)$$

$$\frac{\partial \ell}{\partial \alpha} = 2\alpha - 2(V^T x) \xrightarrow{\text{min}} \underline{V^T x = \alpha}$$

$$\arg \min_{\alpha} \ell \quad \|x - V\alpha\|^2 = \|x - V(V^T x)\|^2 \quad \forall x \in \mathcal{X}^i$$

$$\boxed{V = U; V^T = W; U^T U = \mathbb{I}}$$

$$\textcircled{L} \quad \sum_i \|x^i - U W x^i\|^2 = \sum_i \|x^i - U U^T x^i\|^2$$

$$\textcircled{B} \quad \sum_i \|x^i\|^2 - 2x^{iT} U \underbrace{U^T U}_{\mathbb{I}} x^i + x^{iT} U U^T U U^T x^i$$

$$= \sum_i \|x^i\|^2 - 2x^{iT} U U^T x^i + x^{iT} U U^T x^i =$$

$$\sum_i \|x^i\|^2 - \underbrace{x^{iT} U}_{y^{iT}} \underbrace{U^T x^i}_{y^i} = \sum_i \|x^i\|^2 - \sum_i y^{iT} y^i$$

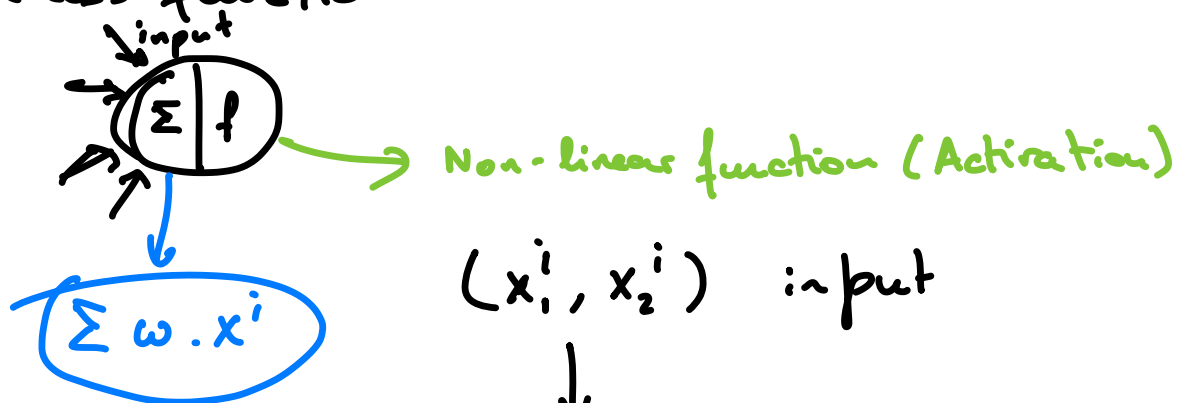
$$= \underbrace{\sum_i \|x^i\|^2}_{\text{constant}} - \underbrace{\text{Tr}(Y^T Y)}_{\hookrightarrow \text{PCA}}$$

Three ways of deriving PCA:

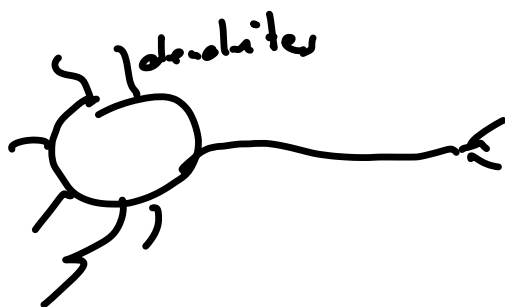
- ① Diagonalize Covariance matrix
- ② Reproducing the euclidean distances in the original space by minimizing the Frobenious Norm of the difference  $\|\Delta - \Theta\|^2$
- ③ Linear Autoencoders

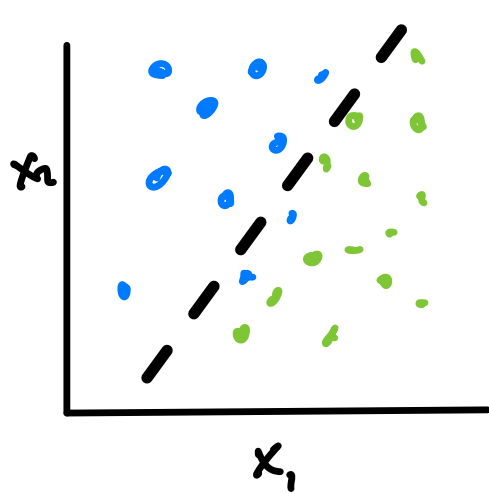
$P_\Phi; q_\Theta$  are not linear  $\rightarrow$  Neural Networks

Find functions (not-linear) that minimize a loss function



$$(x_1^i, x_2^i) \text{ input}$$
$$\downarrow$$
$$\tanh(\underbrace{w_1 x_1}_{\text{bias}} + \underbrace{w_2 x_2}_{\text{bias}} + w_3)$$
$$\downarrow$$
$$\text{output} \in [-1, 1]$$

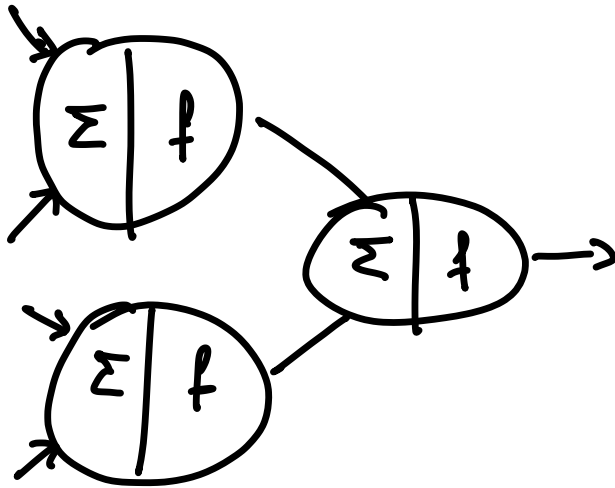




blue  $\rightarrow -1$

green  $\rightarrow 1$

$$L = \sum_i \| y^i - \tanh(\underbrace{\omega_0 x_0 + \omega_1 x_1 + \omega_2 x_2}_{\text{net input}}) \|^2$$



$$\begin{array}{c}
 L_1 \qquad \qquad L_2 \\
 \omega^1 \quad \omega^2 \quad \omega^3 \dots \omega^L \\
 \begin{array}{cccc}
 \bigcirc & \bigcirc & \bigcirc & \bigcirc \\
 \bigcirc & \bigcirc & \bigcirc & \bigcirc \\
 \bigcirc & \bigcirc & \bigcirc & \bigcirc \\
 \bigcirc & \bigcirc & \bigcirc & \bigcirc \\
 \bigcirc & \bigcirc & \bigcirc & \bigcirc
 \end{array}
 \end{array}$$

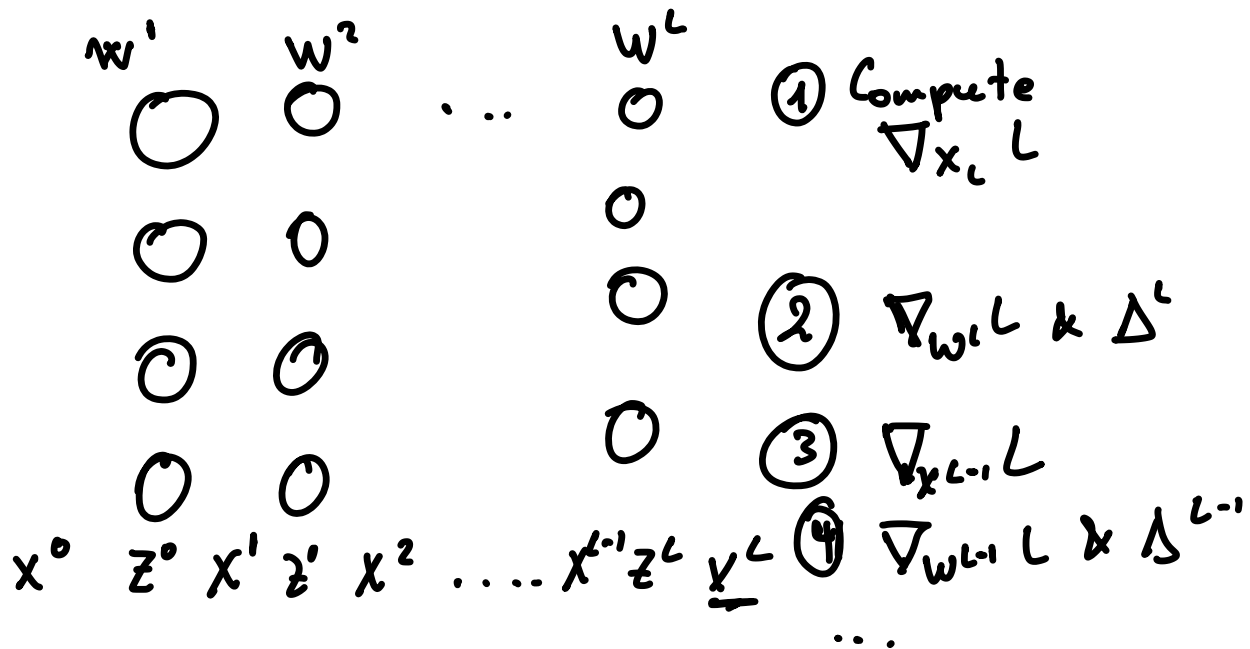
$$x^0 z^0 \rightarrow x'$$

$$[\omega]_{t+1} = [\omega]_t - \alpha \underbrace{\left[ \frac{\partial L}{\partial \omega} \right]_t}_{\text{gradient}}$$

$$\nabla_{\omega^L} L = (X^{L+1})^T \cdot \Delta^{L+1}$$

$$\Delta^{L+1} = \nabla_{x^{L+1}} L \odot f'^{L+1}(z^{L+1})$$

$$\nabla_{x^L} L = \Delta^{L+1} \cdot (W^{L+1})^T$$



## Stochastic Gradient Descent

$$\nabla = \frac{1}{N} \sum_{i=1}^N \delta_i \xrightarrow{n} \nabla \sim \tilde{\nabla} = \frac{1}{n} \sum_{i=1}^n \delta_i$$

GD

$t = \emptyset$

While not converged

$t++$

$\nabla_{w^L} L = \emptyset$

for  $i$  in  $N$

accumulate  $\nabla$

update weights

SGD

$t = \emptyset$

While not converged

$t++$

shuffle my data

for  $j$  in mini batches

$\nabla_{w^L} L = \emptyset$

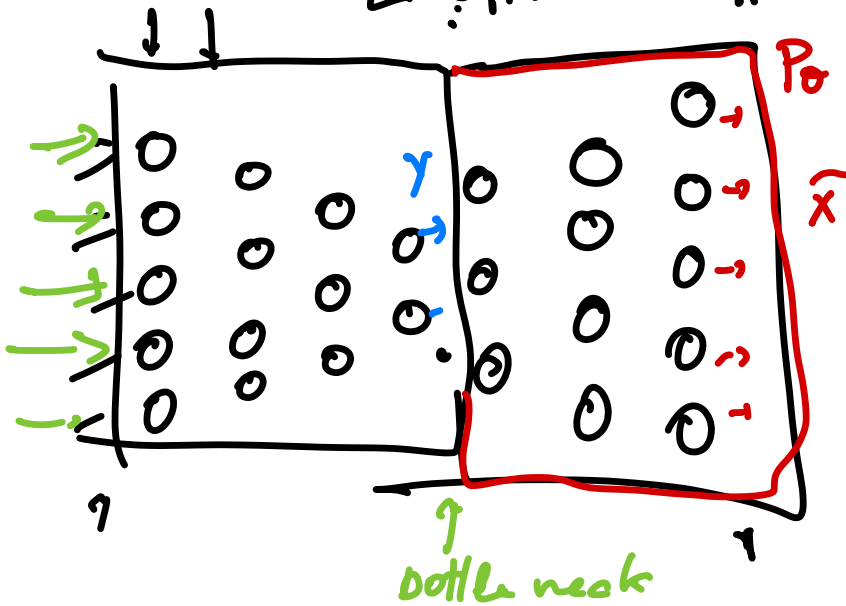
for  $i$  in  $N/m$

accumulate  $\nabla$

update weights

$$X \rightarrow \underbrace{[ENCO.]}_{q_\theta} \rightarrow Y \rightarrow \underbrace{[DECO.]}_{p_\phi} \rightarrow \tilde{X}$$

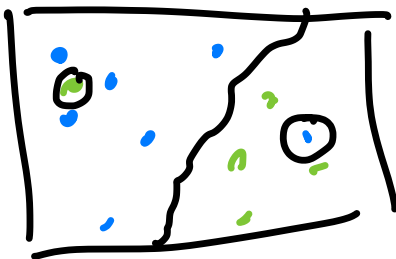
$$L = \sum_i \|x^i - \tilde{x}^i\|^2$$



$$L = \sum \|x^i - p_\phi(q_\theta(x^i))\|^2$$

$$y^i = q_\theta(x^i)$$

$\theta, \phi \Rightarrow$  parameters that minimize



$$L = \|x^i - \tilde{x}^i\|^2 + \lambda \sum_k \|\omega_k\|^p$$

$p = 1 \rightarrow$  LASSO reg.

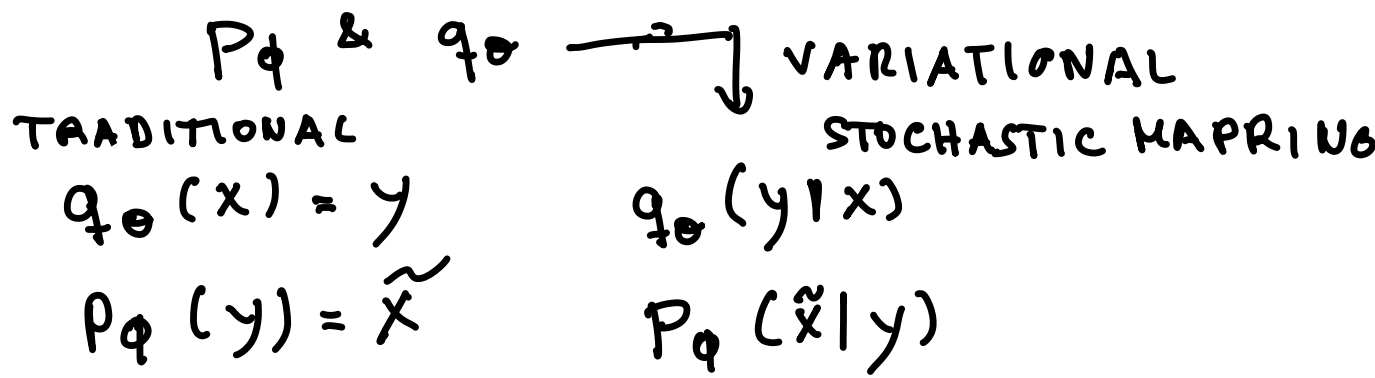
$p_1 \& p_2$

$p = 2 \rightarrow$  Ridge reg.

Elastic  
nets

$$+ \lambda (\alpha \|\omega_k\| + (1-\alpha) \|\omega_k\|^2)$$

# VARIATIONAL AUTOENCODERS



Max. Likelihood for the loss

prior to  $p(y) \sim N(\phi, 1)$

$$KL(p_\phi(x|y) \mid \underbrace{N(\phi, 1)})$$