

# Autoencoders

University of Padova, A.A. 2022/23

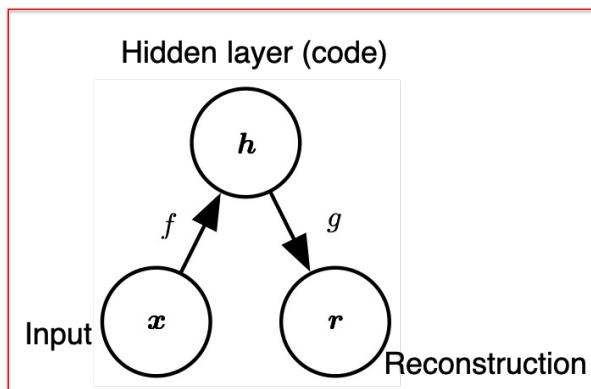
## Autoencoder: General Idea

**Learn a (compressed and/or sparse) code  $h$  for the input  $x$**

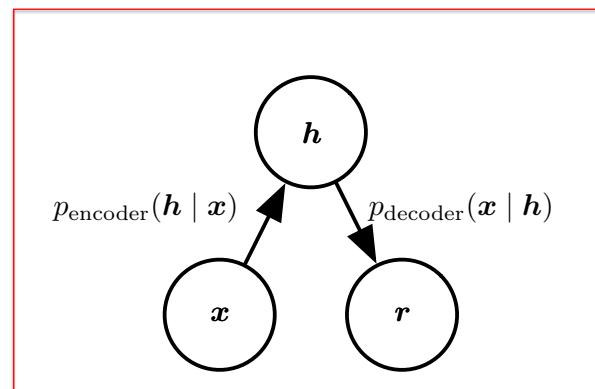
How ? By learning the identity function  $x = g(f(x))$  with constraints:

- on the architecture of the network (**undercomplete autoencoder**)
- adding a regularizing term to the loss (**overcomplete autoencoder**)

Deterministic



Stochastic



# Autoencoder: General Idea

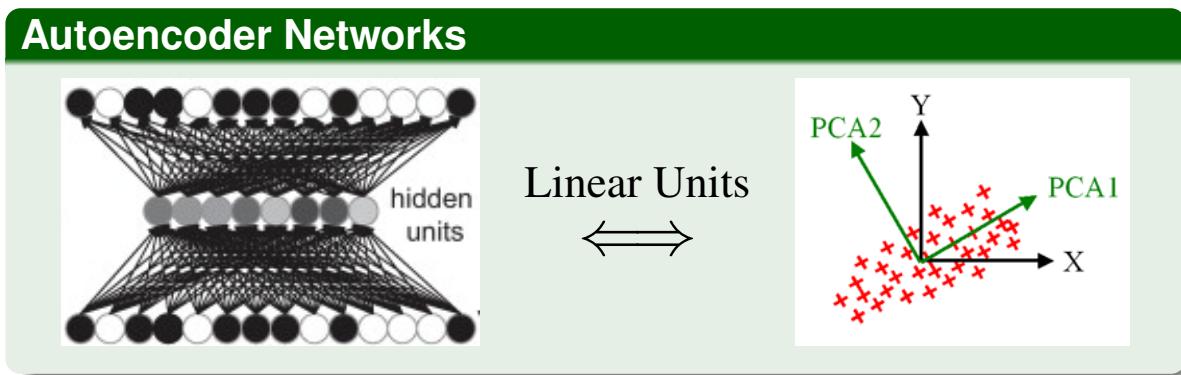
**Not useful to learn identity function on the whole input domain!**

Identity function only on  $x_i \in \mathcal{T} \dots$  and hopefully, similar points  $\Rightarrow$  manifold!

- Undercomplete autoencoders:
  - $h$  has lower dimension than  $x$
  - $f(\cdot)$  or  $g(\cdot)$  has low capacity (e.g., linear  $g$ )
  - must discard some information in  $h$
- Overcomplete autoencoders:
  - $h$  has higher dimension than  $x$
  - must be regularized

A.A. 2022/23: Deep Learning

## Undercomplete: Linear Autoencoder Networks

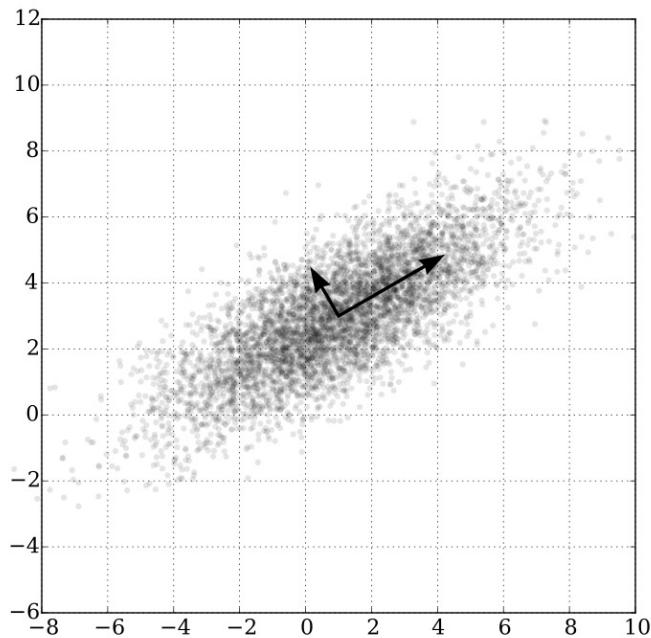


- Linear autoencoder networks are actually closely linked to PCA [Bourlard and Kamp, 1988; Baldi and Hornik, 1989]
- Principal Components: (orthogonal) directions of maximum variance of data
- Compression of information: define the optimal sub-space where to project data
- Used in Pattern Recognition to reduce the dimension of data and/or noise (see image processing)

A.A. 2022/23: Deep Learning

# How to Compute Linear Autoencoders

Compute an orthogonal basis (principal directions) minimizing the average squared distance from a point  $x$  in the dataset to the spanned subspace



A.A. 2022/23: Deep Learning

# How to Compute Linear Autoencoders

Singular Value Decomposition (SVD)

- Collect all the data  $x_i$  as rows (columns) in a matrix  $\Xi \in \mathbb{R}^{n \times m}$
- Compute the SVD decomposition of  $\Xi$

$$\Xi = V S U^\top$$

where

- $V \in \mathbb{R}^{n \times n}$  is an orthonormal matrix (columns are a basis for the rows of  $\Xi$ )
- $S \in \mathbb{R}^{n \times m}$  is a diagonal matrix with non-negative real numbers on the diagonal (*singular values*)
- $U \in \mathbb{R}^{m \times m}$  is an orthonormal matrix (columns are a basis for the columns of  $\Xi$ )

Example:

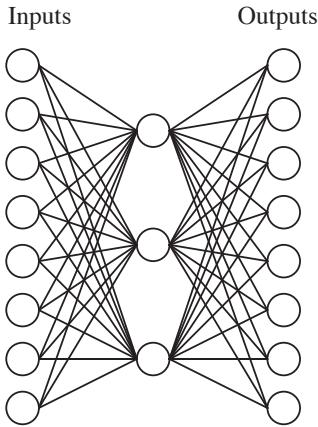
$$\underbrace{\begin{bmatrix} 1 & 0 & 0 & 0 & 2 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \end{bmatrix}}_{\Xi} = \underbrace{\begin{bmatrix} 0 & -1 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & -1 & 0 \end{bmatrix}}_{V} \underbrace{\begin{bmatrix} 3 & 0 & 0 & 0 & 0 \\ 0 & \sqrt{5} & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & \textcolor{red}{0} & 0 \end{bmatrix}}_{S} \underbrace{\begin{bmatrix} 0 & 0 & -1 & 0 & 0 \\ -\sqrt{0.2} & 0 & 0 & 0 & -\sqrt{0.8} \\ 0 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ -\sqrt{0.8} & 0 & 0 & 0 & \sqrt{0.2} \end{bmatrix}}_{U^\top}$$

$$x_i \text{ as rows: reduced representations} \Rightarrow \underbrace{V S U^\top}_{= I} U = V S$$

A.A. 2022/23: Deep Learning

# Undercomplete: Shallow Autoencoder Networks

Learning the identity function:  
(lossy) compression of data



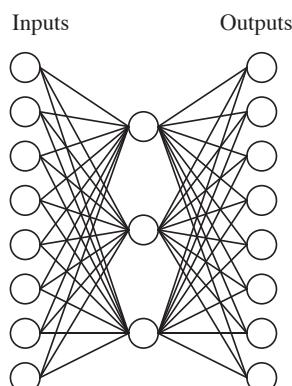
Supervised learning for..  
Unsupervised learning!!

Input	Output
00000001	00000001
00000010	00000010
00000100	00000100
00001000	00001000
00010000	00010000
00100000	00100000
01000000	01000000
10000000	10000000

A.A. 2022/23: Deep Learning

# Undercomplete: Shallow Autoencoders

After learning...

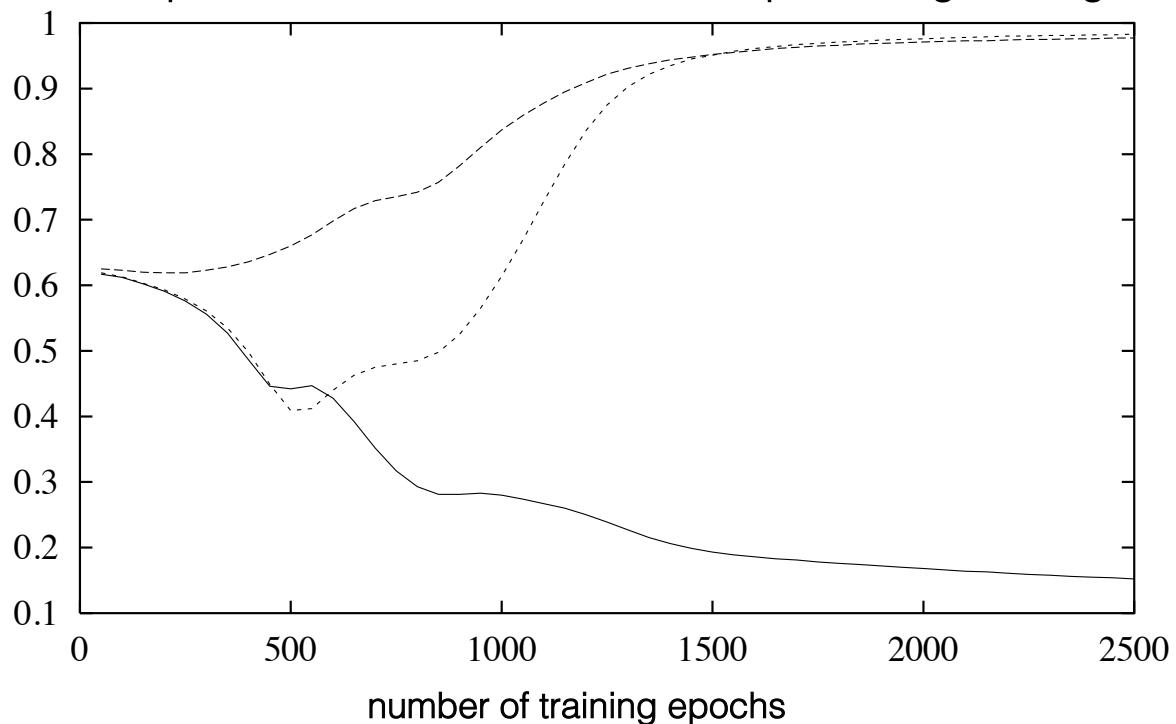


Input	Hidden values	Output
10000000	→ 0.89 0.04 0.08	→ 10000000
01000000	→ 0.01 0.11 0.88	→ 01000000
00100000	→ 0.01 0.97 0.27	→ 00100000
00010000	→ 0.99 0.97 0.71	→ 00010000
00001000	→ 0.03 0.05 0.02	→ 00001000
00000100	→ 0.22 0.99 0.99	→ 00000100
00000010	→ 0.80 0.01 0.98	→ 00000010
00000001	→ 0.60 0.94 0.01	→ 00000001

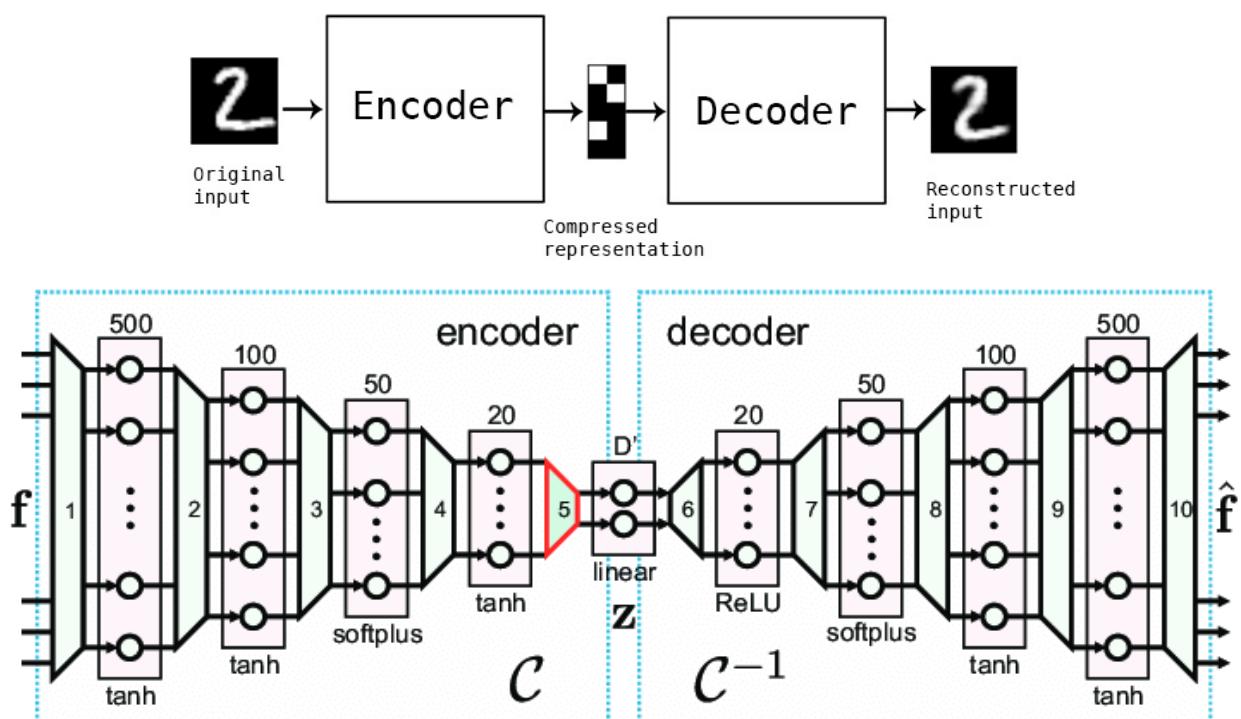
A.A. 2022/23: Deep Learning

# Undercomplete: Shallow Autoencoders

output of the hidden units for one input during training



# Undercomplete: Deep Autoencoders



# Overcomplete Autoencoders

## Regularized autoencoders

- Sparse autoencoders
- Denoising autoencoders
- Contractive autoencoders
- Autoencoders with Dropout on the hidden layer

A.A. 2022/23: Deep Learning

## Sparse Autoencoders

Limit the capacity of autoencoder

- By adding a term to the cost function penalizing the entries of the code  $\mathbf{h} = f(\mathbf{x})$  for being larger

$$L(\mathbf{x}, g(f(\mathbf{x}))) + \Omega(\mathbf{h})$$

- Special case of variational autoencoder (will present in the future)

- Probabilistic model  $p_{\text{model}}(\mathbf{x}, \mathbf{h}) = p_{\text{model}}(\mathbf{h})p_{\text{model}}(\mathbf{x}|\mathbf{h})$

- $\mathbf{h}$  latent variables *generating*  $\mathbf{x}$

- $\log p_{\text{model}}(\mathbf{x}, \mathbf{h}) = \log p_{\text{model}}(\mathbf{h}) + \log p_{\text{model}}(\mathbf{x}|\mathbf{h})$

- Laplace prior corresponds to L1 sparsity penalty

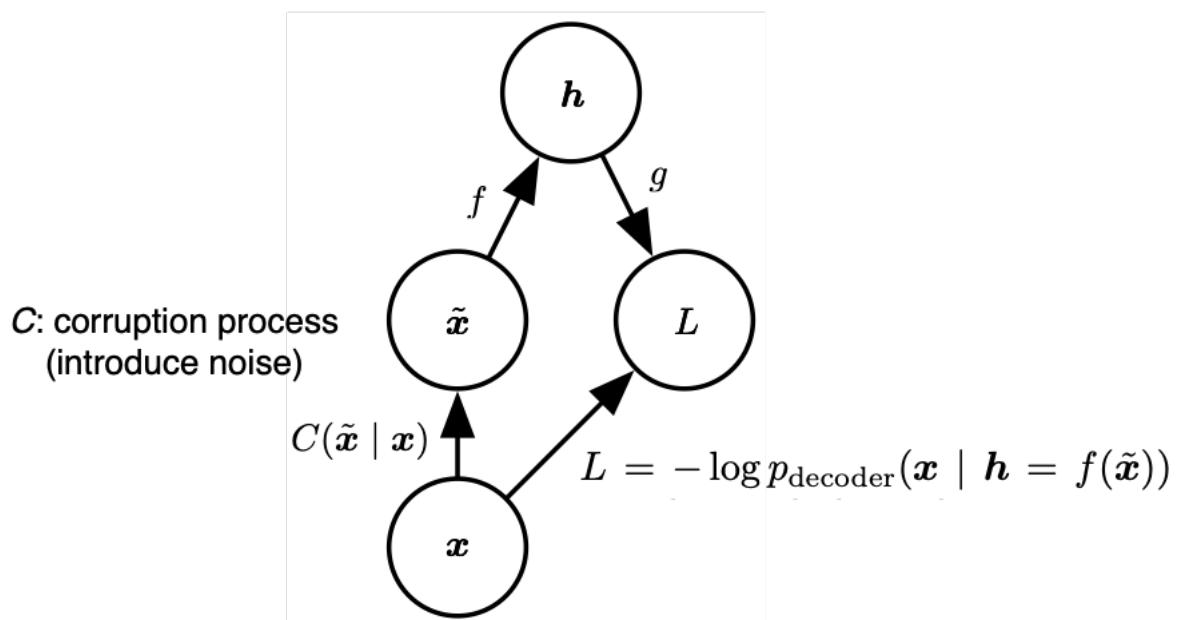
- $p_{\text{model}}(h_i) = \frac{\lambda}{2} \exp^{-\lambda|h_i|}$

$$-\log p_{\text{model}}(\mathbf{h}) = \sum_i \left( \lambda|h_i| - \log \frac{\lambda}{2} \right) = \lambda \sum_i |h_i| + \text{const}$$

- Other sparsity penalty terms have been proposed in the literature

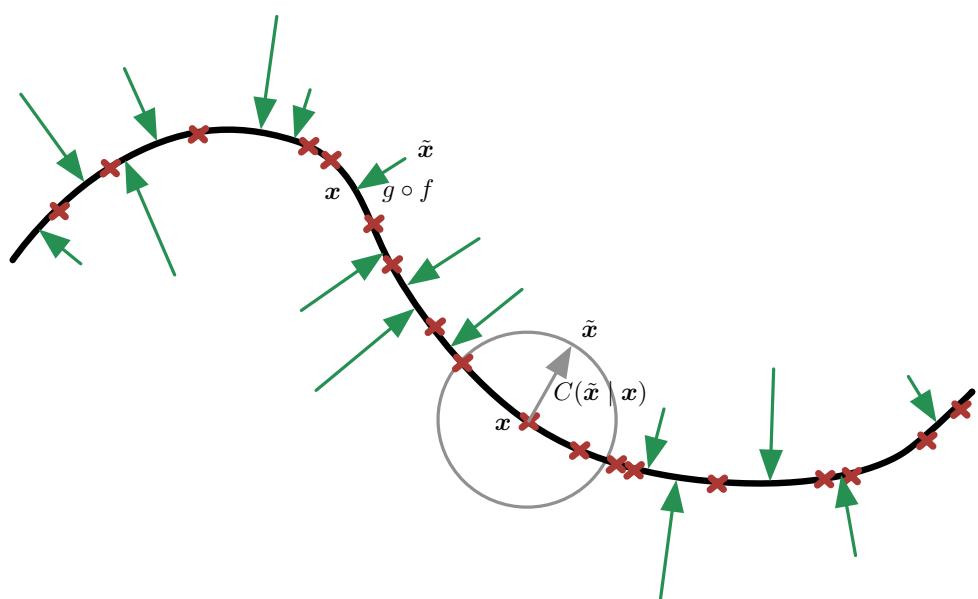
A.A. 2022/23: Deep Learning

# Denoising Autoencoder



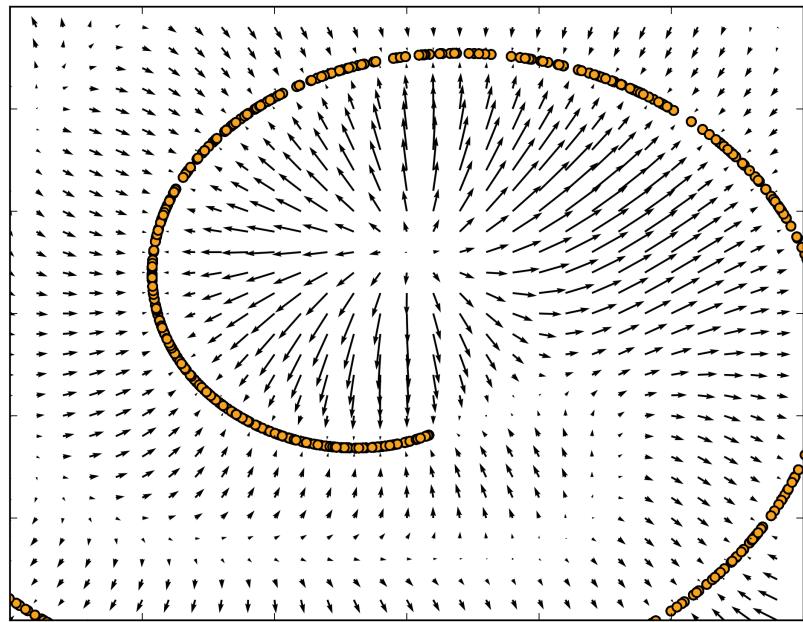
A.A. 2022/23: Deep Learning

## Denoising Autoencoder: Learn a Manifold



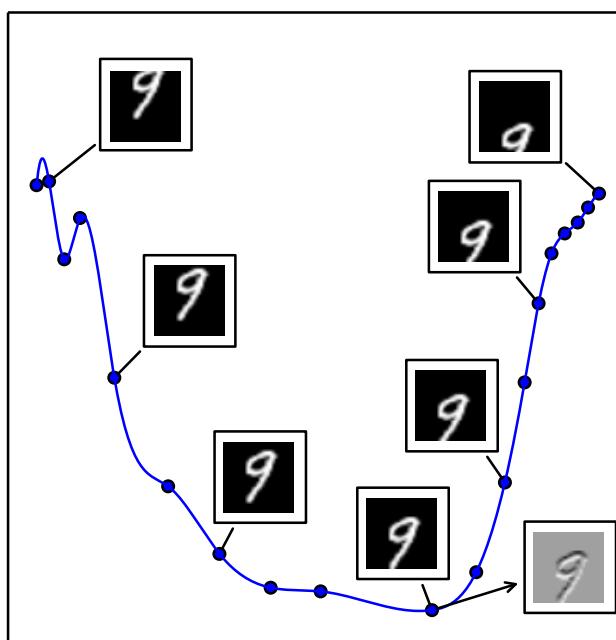
A.A. 2022/23: Deep Learning

# Vector Field Learned by a Denoising Autoencoder



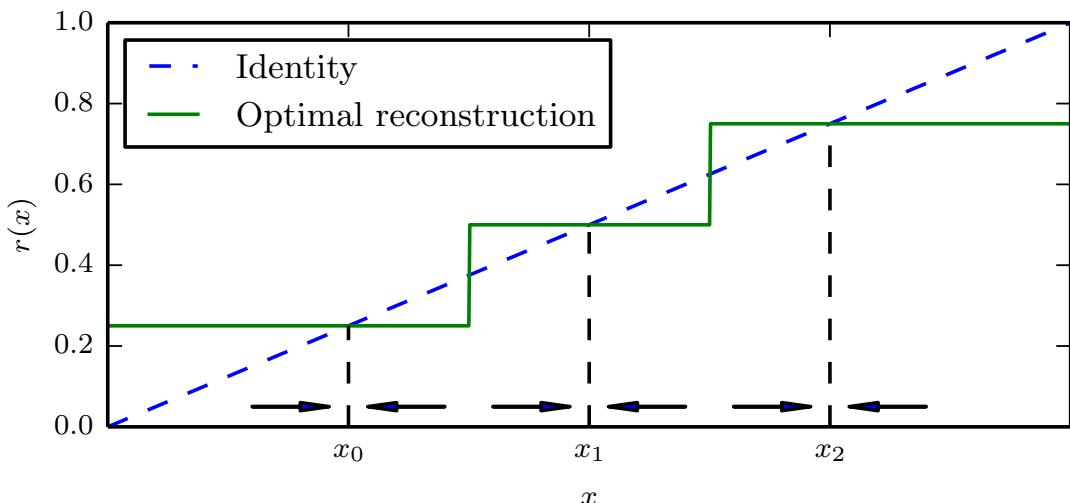
A.A. 2022/23: Deep Learning

## Tangent Hyperplane of a Manifold



A.A. 2022/23: Deep Learning

# A Simple Example



A.A. 2022/23: Deep Learning

## Contractive Autoencoders

$$\Omega(\mathbf{h}) = \lambda \left\| \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right\|_F^2 \quad \text{← Frobenius norm (sum of squared elements)}$$

Practical issues

- heavy computation of Jacobian for deep networks: build one layer at time!
- scale of decoder should be the same as encoder: impose weight matrix of decoder to be the transpose of weight matrix of encoder!

Input point	Tangent vectors
	Local PCA (no sharing across regions)
	Contractive autoencoder

A.A. 2022/23: Deep Learning