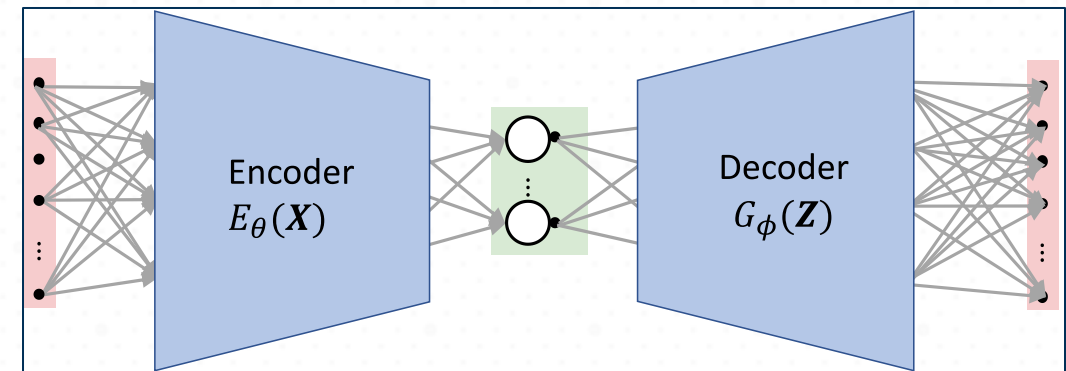


ECE 4252/8803: Fundamentals of Machine Learning (FunML) Fall 2024

Lecture 18: Autoencoders



Overview

In this Lecture..

AEs: Introduction and Motivation

- Unsupervised Learning
- Autoencoders

Fully-connected Autoencoders

Convolutional Autoencoders

Regularized Autoencoders

Variational Autoencoders

Introduction

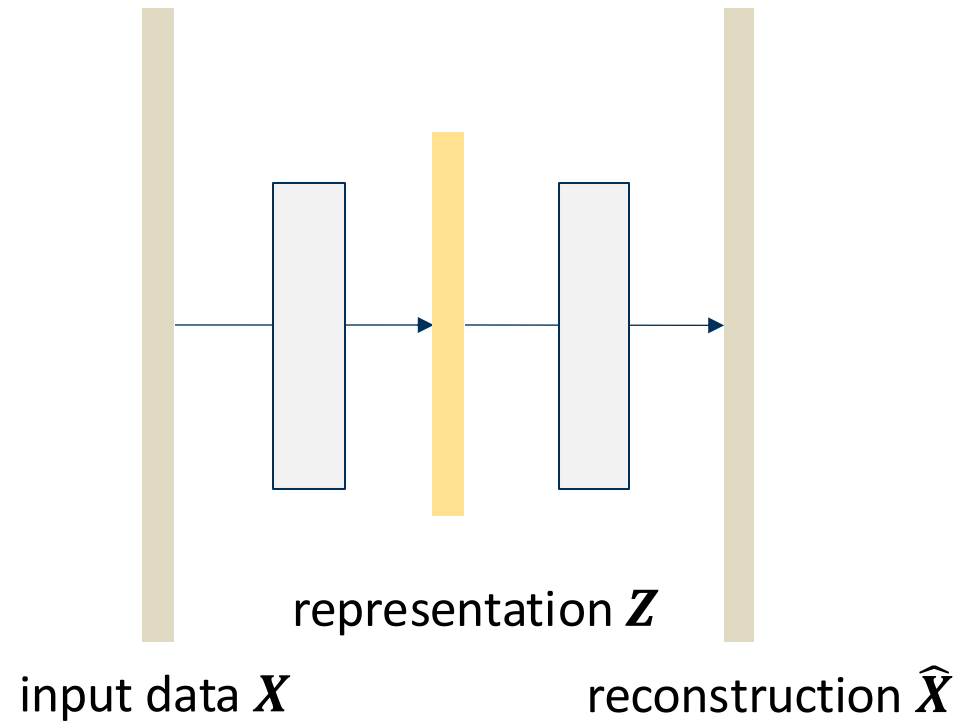
Types of Learning

- Supervised Learning
 - Requires labeled training data $(x_1, y_1), \dots, (x_N, y_N)$
 - Goal: Learn a *mapping function* $f_\theta: \mathcal{X} \rightarrow \mathcal{Y}$
 - Examples:
 - Classification, regression, etc.
- Unsupervised Learning
 - No labels required, only data x_1, \dots, x_N
 - Goal: Learn some underlying hidden **structure** of the data
 - Examples:
 - Clustering, dimensionality reduction, density estimation, etc.

Introduction

Autoencoders: Motivation

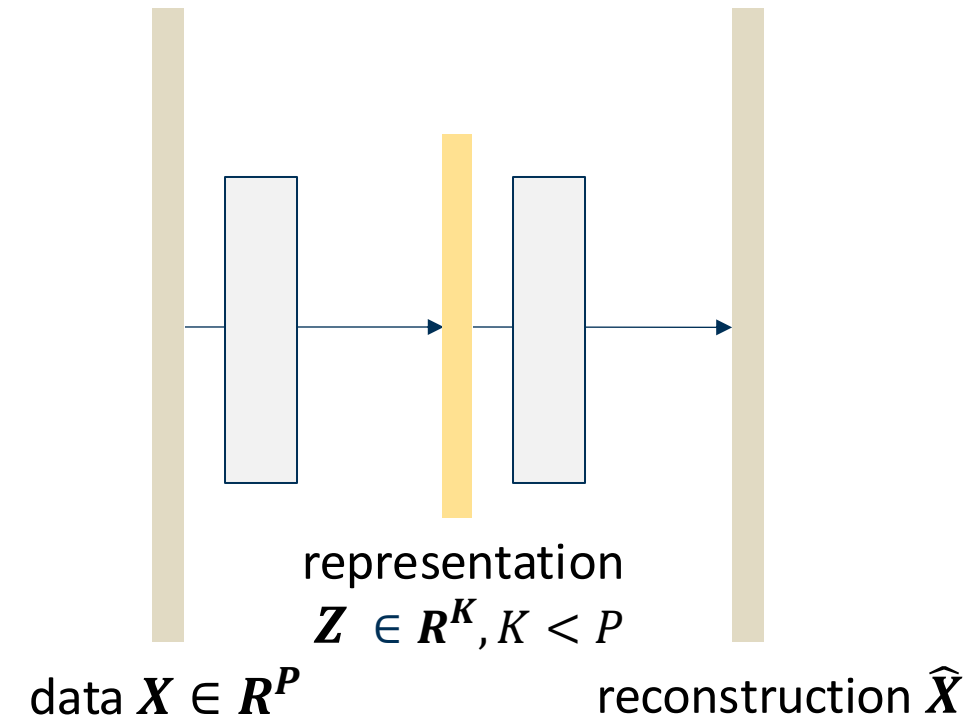
Unsupervised approach for learning a ***lower-dimensional latent representation*** from the **unlabeled data** and reconstructing the data from this representation



Introduction

Autoencoders: Motivation

- **Unsupervised learning**
 - Learn a **lower-dimensional** representation \mathbf{Z}
 - Reconstruct the data
- \mathbf{Z} describes the underlying **structure** of \mathbf{X} that can be used for:
 - Visualization
 - Map \mathbf{X} to $\mathbf{Z} \in \mathbb{R}^2$
 - **Dimensionality reduction**
 - Compress images, documents, text ... etc.
 - Subsequent supervised tasks
 - Use \mathbf{Z} to train/fine-tune for supervised tasks



Introduction

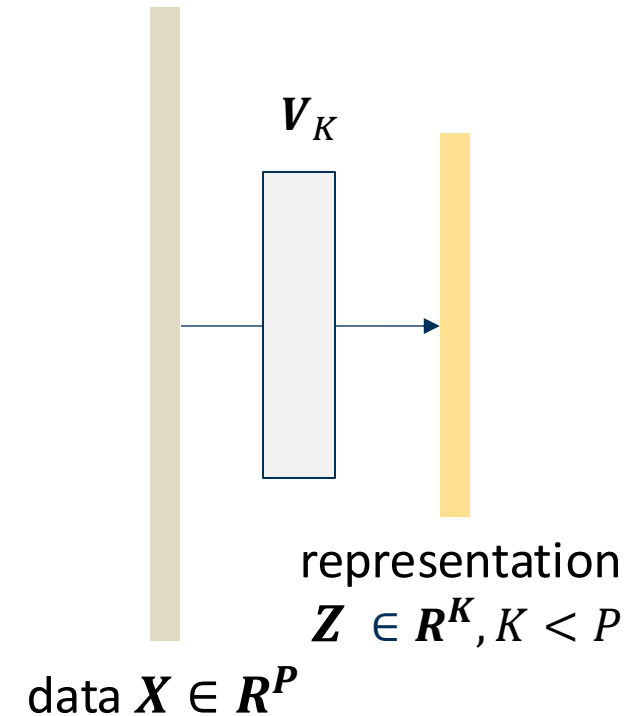
Autoencoders: Motivation

- Autoencoder: learning a lower-dimensional representation \mathbf{Z} from unlabeled input data \mathbf{X}

- **Linear** dimensionality reduction by PCA :

$$\mathbf{Z} = \mathbf{X}\mathbf{V}_K$$

where \mathbf{V}_K transforms \mathbf{X} to a K -dimensional subspace ($K < P$)



Introduction

Autoencoders: Motivation

- Autoencoder: learning a lower-dimensional representation \mathbf{Z} from unlabeled input data \mathbf{X}
- **Linear** dimensionality reduction by PCA :

$$\mathbf{Z} = \mathbf{X}\mathbf{V}_K$$

where \mathbf{V}_K transforms \mathbf{X} to a K -dimensional subspace ($K < P$)

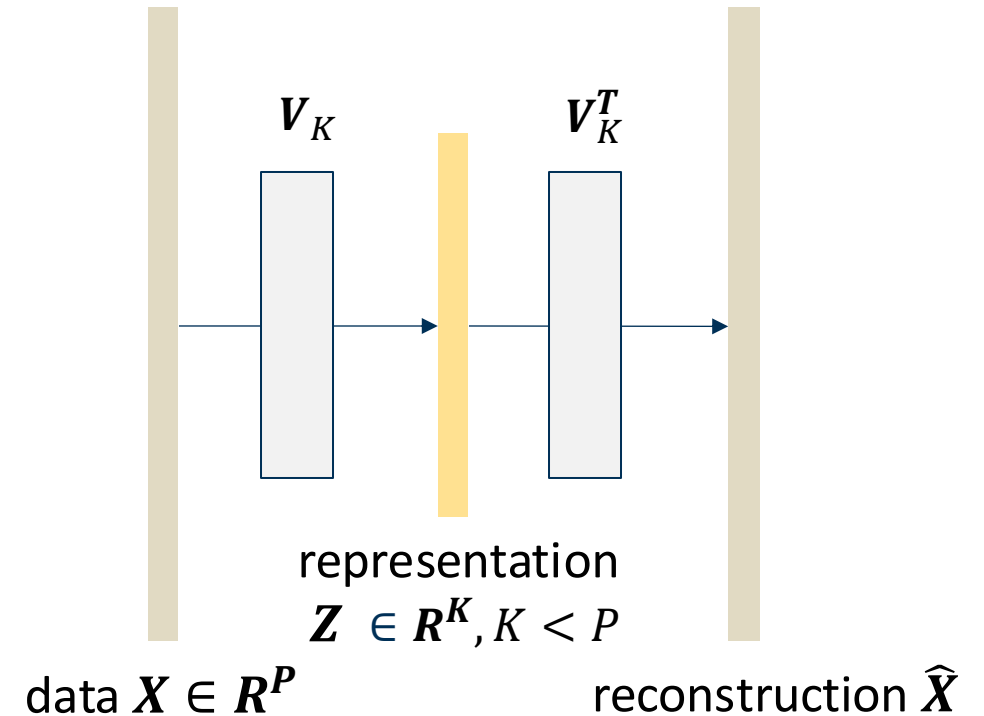
\mathbf{X} can be reconstructed from \mathbf{Z} :

$$\hat{\mathbf{X}} = \mathbf{X}\mathbf{V}_K\mathbf{V}_K^T$$

Reconstruction error:

$$\|\mathbf{X} - \mathbf{X}\mathbf{V}_K\mathbf{V}_K^T\|_F$$

\mathbf{V}_K maximizes the variance in \mathbf{X} that has been preserved, while minimizing the reconstruction error



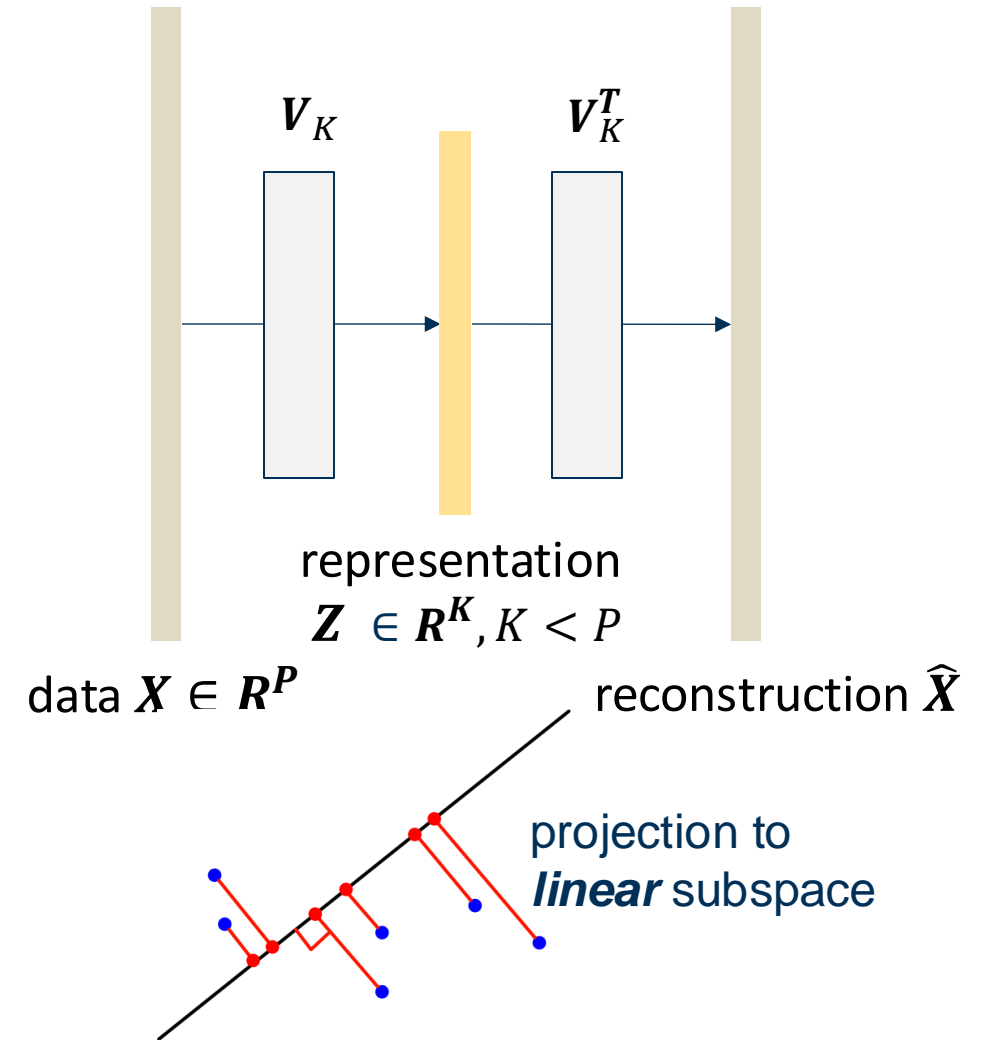
Introduction

Autoencoders: Motivation

- Autoencoder: learning a lower-dimensional representation \mathbf{Z} from unlabeled input data \mathbf{X}
- **Linear** dimensionality reduction by PCA :

$$\mathbf{Z} = \mathbf{X}\mathbf{V}_K$$
$$\hat{\mathbf{X}} = \mathbf{X}\mathbf{V}_K\mathbf{V}_K^T$$

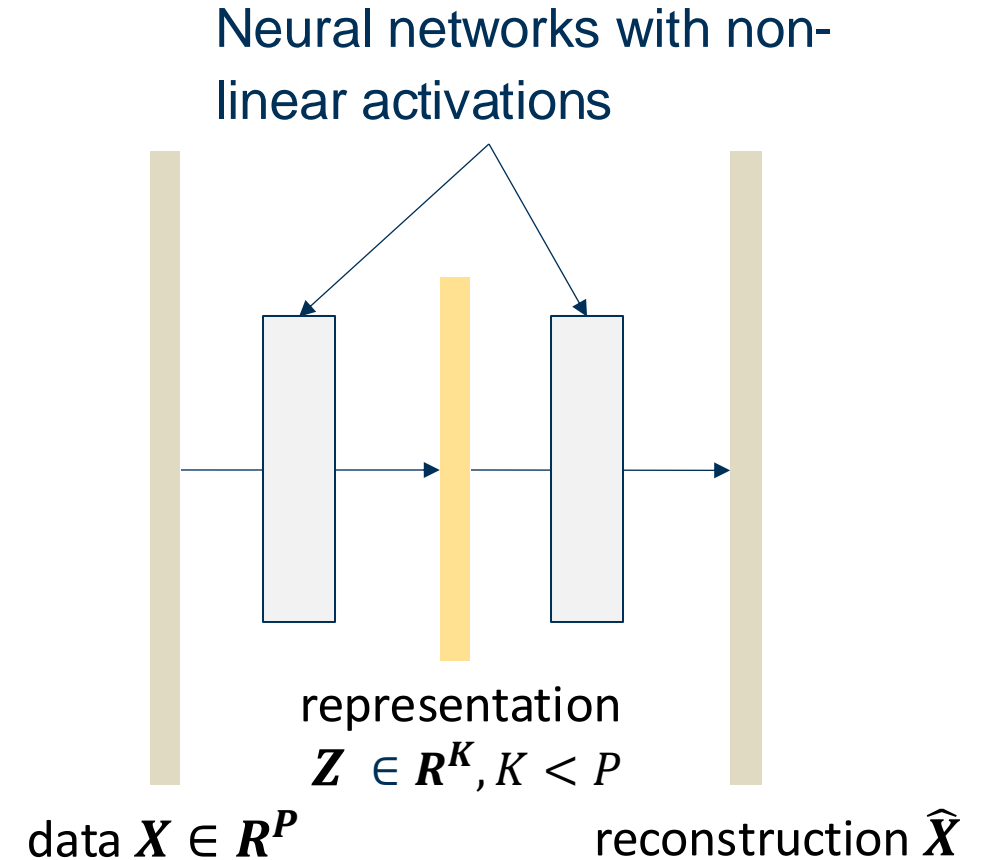
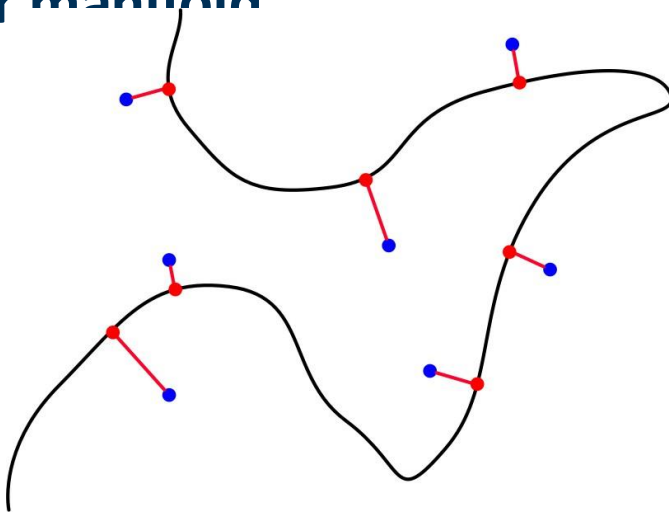
- Linear autoencoder learns PCA when:
 - one fc layer with weights \mathbf{V}_K and *linear* activation
 - one fc layer with weights \mathbf{V}_K^T and *linear* activation
 - loss function $L(\mathbf{X}, \hat{\mathbf{X}}) = \|\mathbf{X} - \hat{\mathbf{X}}\|_F$
 - **s.t.** \mathbf{V}_K is the matrix of K eigenvectors of $\mathbf{X}^T\mathbf{X}$



Introduction

Autoencoders: Beyond PCA

- Autoencoder: learning a lower-dimensional representation \mathbf{Z} from unlabeled input data \mathbf{X}
- Nonlinear autoencoders learn to project the data \mathbf{X} , not onto a linear subspace, but onto a **nonlinear manifold**



Overview

In this Lecture..

AEs: Introduction and Motivation

Fully-connected Autoencoders

- Architecture
- Visualization

Convolutional Autoencoders

Regularized Autoencoders

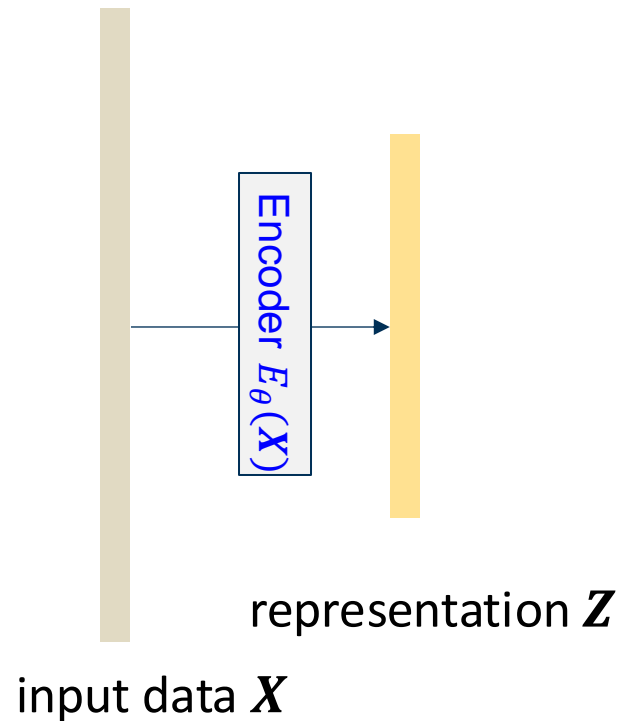
Variational Autoencoders

Autoencoders

Architecture

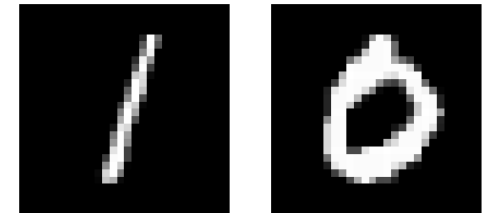
Learning a lower-dimensional representation \mathbf{Z} from unlabeled data \mathbf{X}

- an **encoding function** $E_{\theta}: \mathbf{X} \rightarrow \mathbf{Z}$



Want \mathbf{Z} to capture meaningful patterns in data \mathbf{X} .

e.g., the digit “1” usually has a straight line, the digit “0” is circular

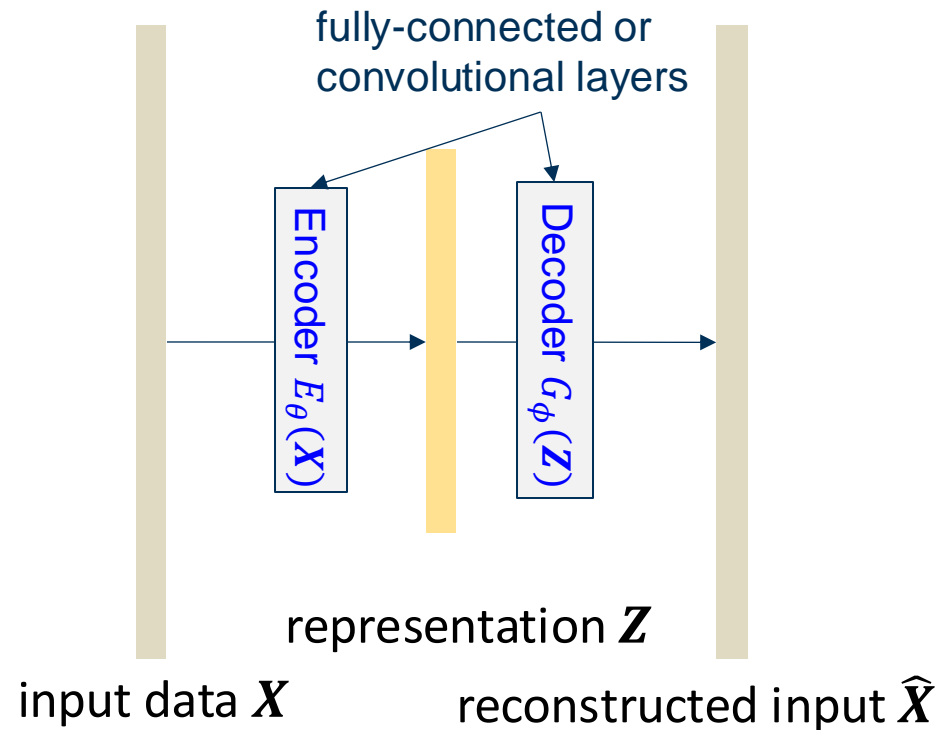


Autoencoders

Architecture

Learning a lower-dimensional representation \mathbf{Z} from unlabeled data \mathbf{X}

- an encoding function $E_\theta: \mathbf{X} \rightarrow \mathbf{Z}$
- a decoding function $G_\phi: \mathbf{Z} \rightarrow \hat{\mathbf{X}}$



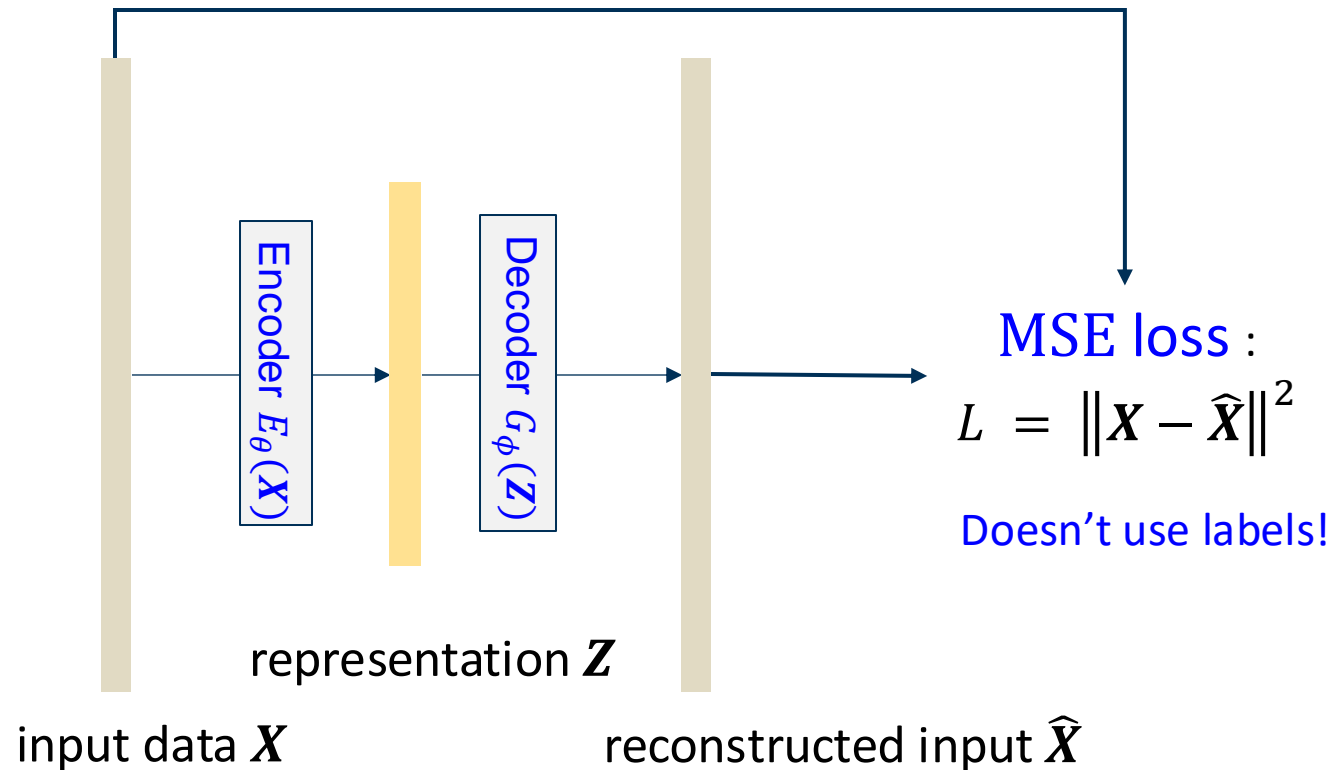
“Autoencoding” - encoding itself

Autoencoders

Architecture

How to learn this representation?

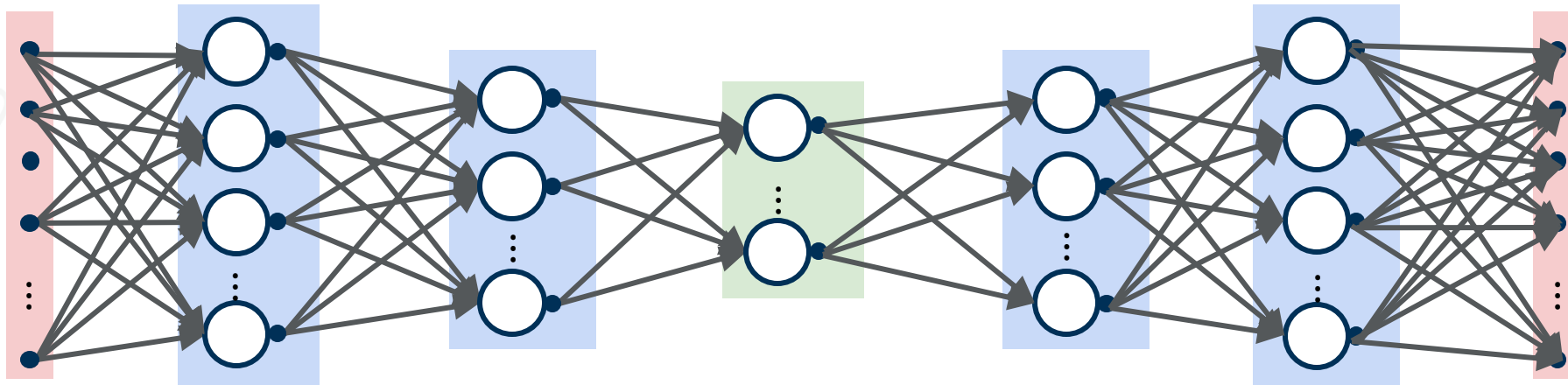
- Train such that representations can be used to reconstruct input data



Fully-connected Autoencoders

Architecture

- Encoder and decoder are fully-connected layers



Original
input X

Encoder
 $E_\theta(X)$

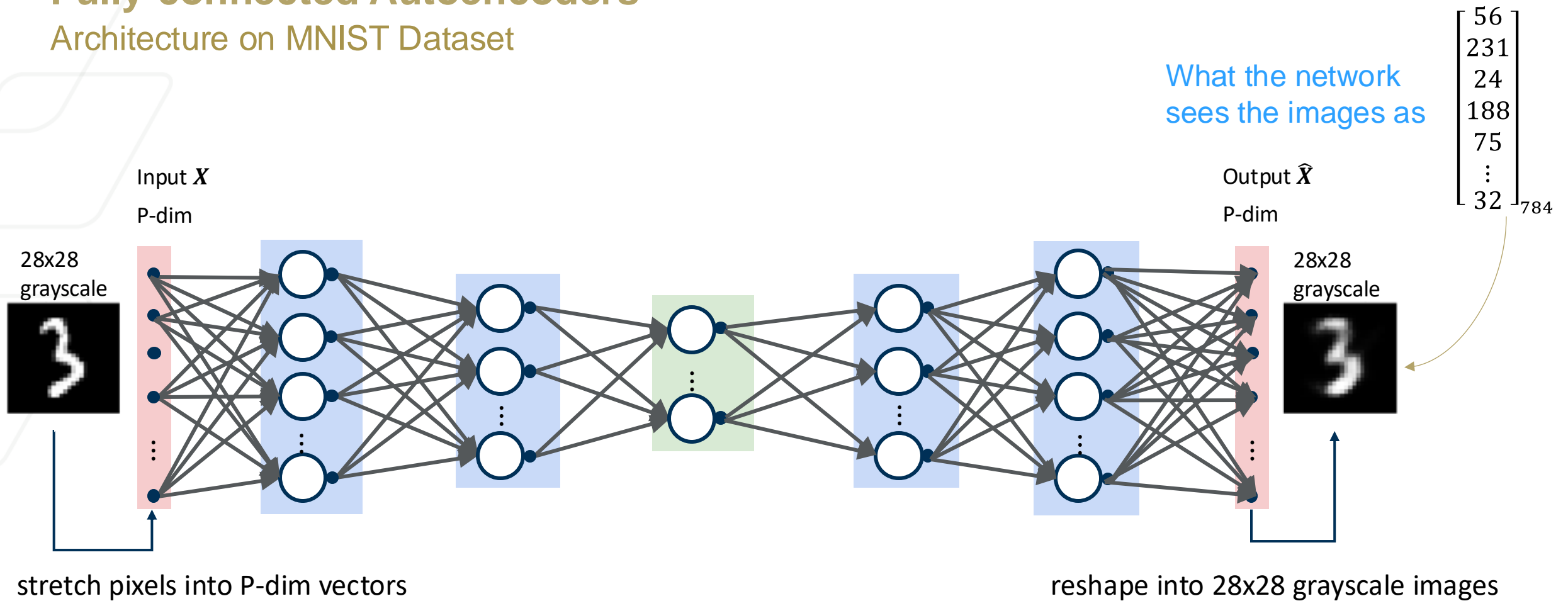
Latent
representation Z

Decoder
 $G_\phi(Z)$

Reconstructed
output \hat{X}

Fully-connected Autoencoders

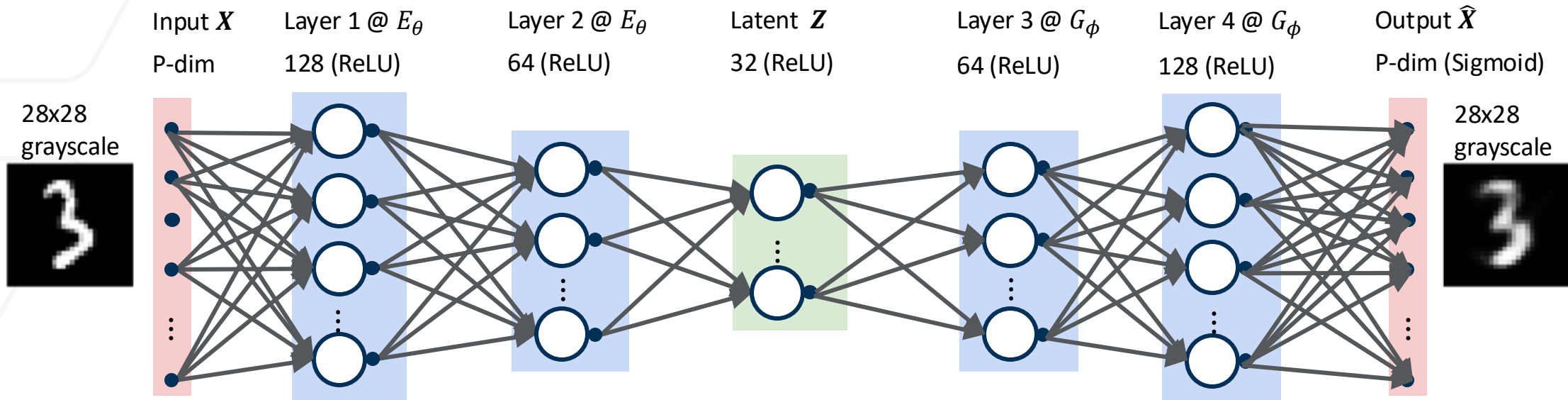
Architecture on MNIST Dataset



- $P = 784$ ($28 \times 28 \times 1$) for MNIST

Fully-connected Autoencoders

Architecture on MNIST Dataset



- $P = 784$ (28x28x1) for MNIST
- ReLU activation in the intermediate layers
- Sigmoid activation in the output layer

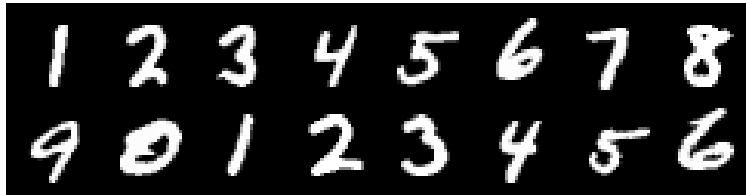
Fully-connected Autoencoders

Autoencoder Reconstructions on MNIST

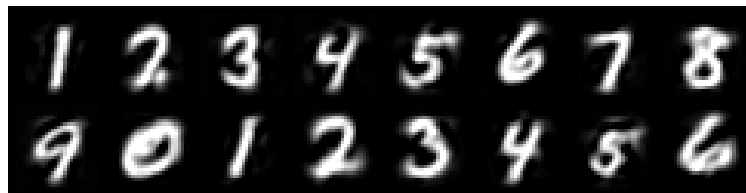
32-dimensional latent representation

Fully-connected AE (**linear**)

input



Reconstruction (MSE = 0.0677)

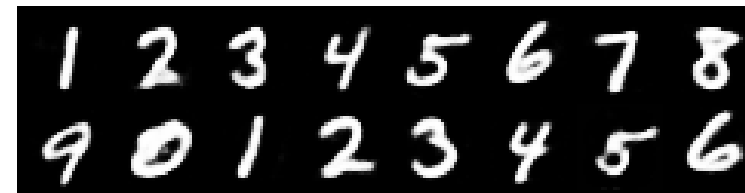


Fully-connected AE (**nonlinear**)

input



Reconstruction (MSE = 0.0306)



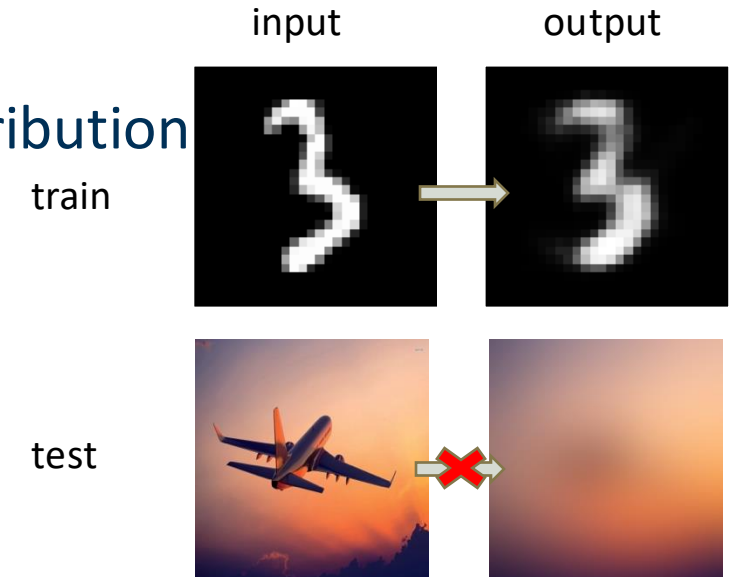
Linear AE: NOT using nonlinear activations; Nonlinear AE: using ReLU and Sigmoid

AE with nonlinearity learns more complex nonlinear latent representation to generate sharper reconstructions

Fully-connected Autoencoders

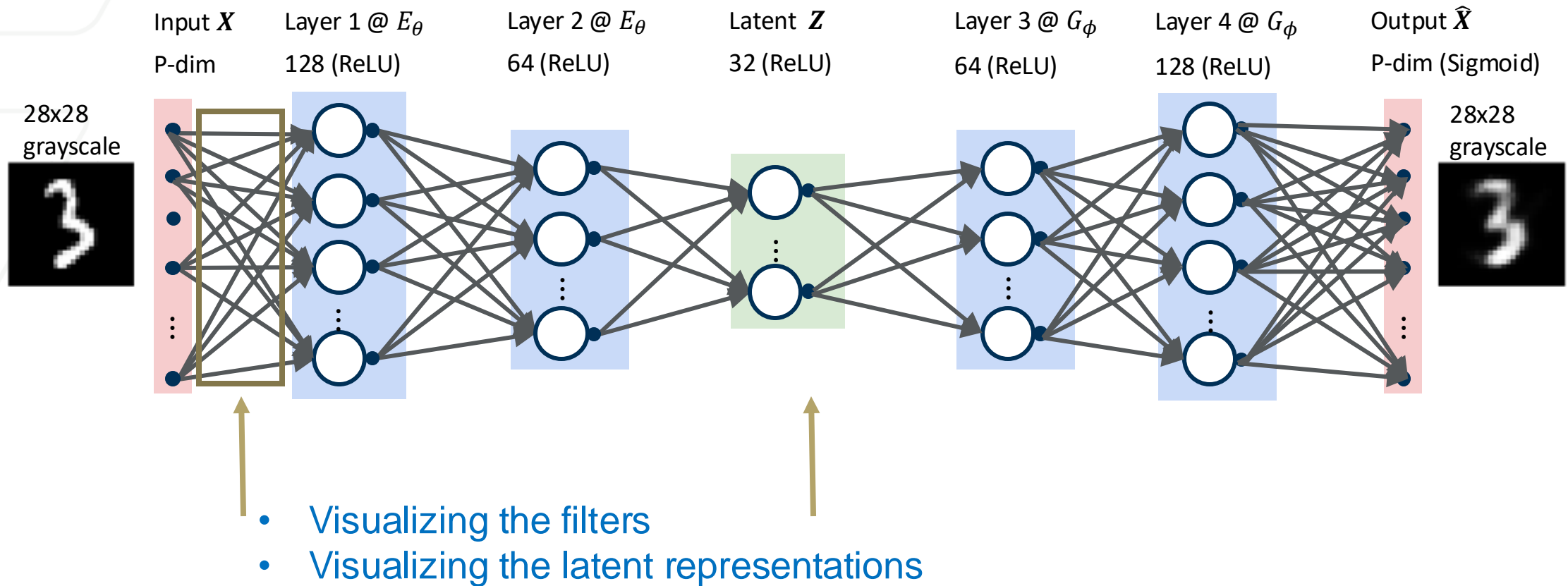
Dimensionality Reduction

- *Unsupervised*: training without image labels
- *Lossy*: the reconstruction output is more blurry
- *Data-specific*: test and training data follow same distribution
 - Can be used for anomaly detection



Fully-connected Autoencoders

Visualization



Fully-connected Autoencoders

Visualization of Learned Filters

Mathematical Notation:

- $\mathbf{X} \in \mathbb{R}^{N \times P}$: P -dimensional input vectors $\mathbf{x}_i \in \mathbb{R}^P, i = 1 \dots N$
- $\mathbf{W}^{(k)} \in \mathbb{R}^{P^{(k)} \times P^{(k-1)}}$: weights matrix of neurons in layer k
- $P^{(k)}$: the number of neurons in layer k
- $P^{(k-1)}$: the number of neurons in layer $k - 1$

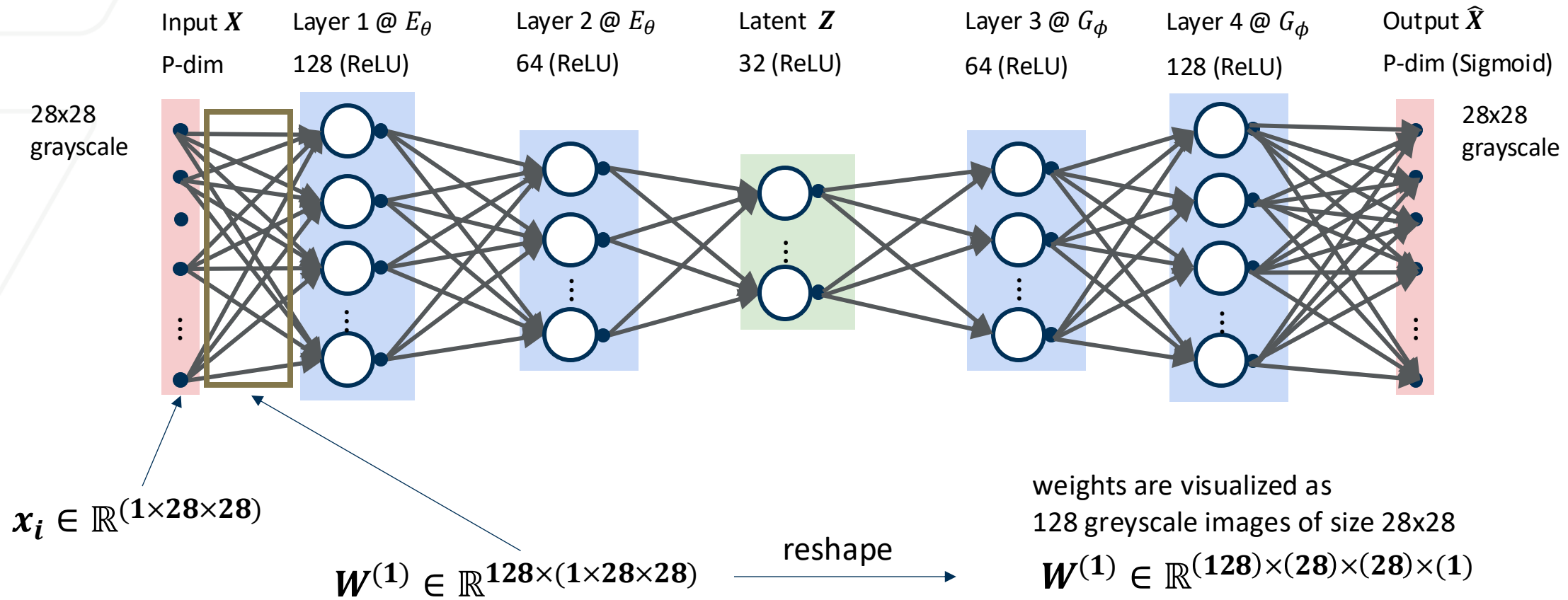
Visualization:

Extract the weights $\mathbf{W}^{(k)}$ (w/o bias) in layer k and reshape into an array of shape $(P^{(k)}, h, w, c)$, where $h \times w \times c = P^{(k-1)}$

Fully-connected Autoencoders

Visualization of Learned Filters

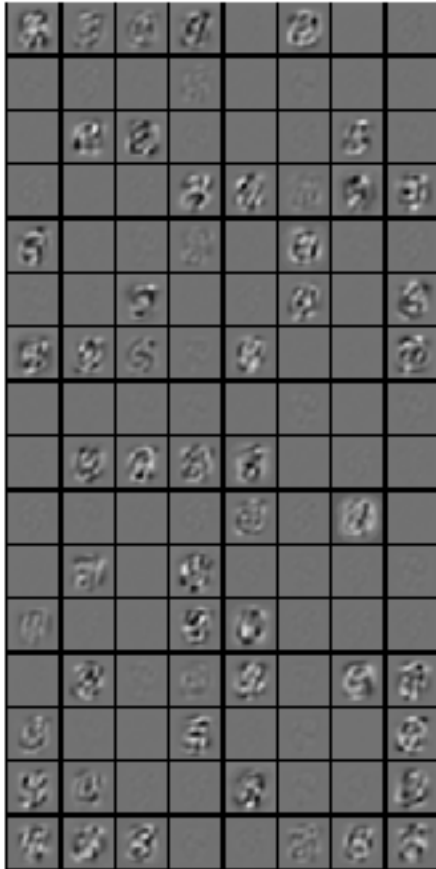
Extract $W^{(1)}$ in the first layer of the fully-connected autoencoder



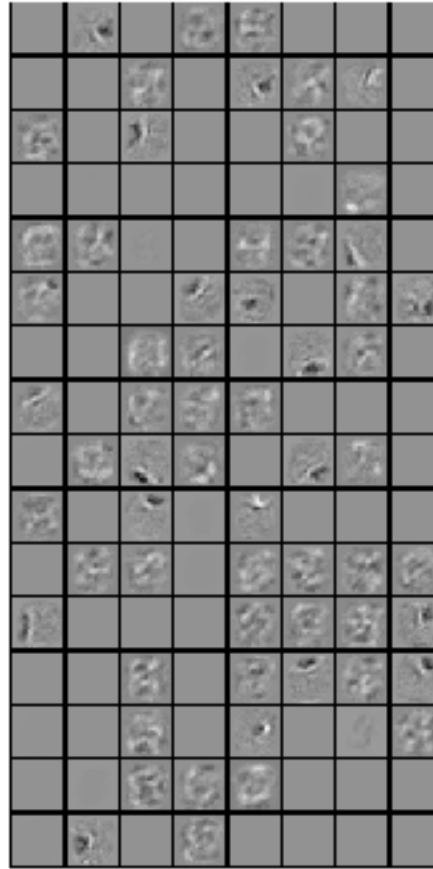
Fully-connected Autoencoders

Visualization of Learned Filters: First Layer

Fully-connected AE (**linear**)



Fully-connected AE (**nonlinear**)



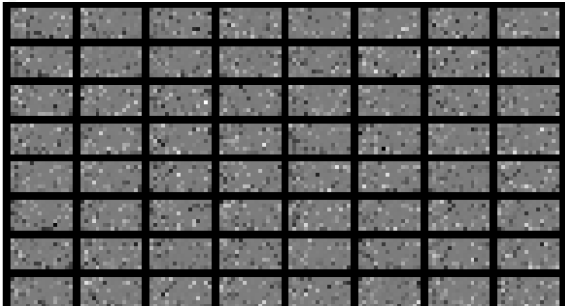
Extract the weights $\mathbf{W}^{(1)} \in \mathbb{R}^{128 \times (1 \times 28 \times 28)}$
in **layer 1** @ E_θ

Reshape $\mathbf{W}^{(1)}$ into 128 filters of size 28x28

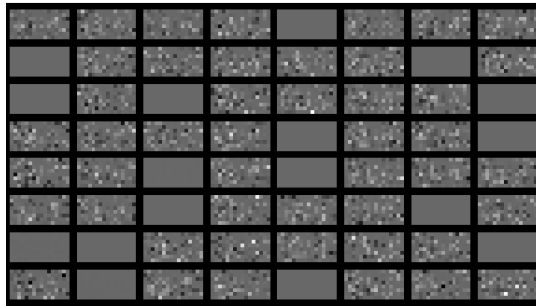
Fully-connected Autoencoders

Visualization of Learned Filters: Second and Third Layers

Fully-connected AE (**linear**)



Fully-connected AE (**nonlinear**)

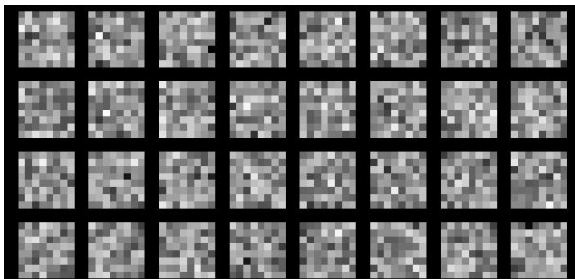


Extract the weights $\mathbf{W}^{(2)} \in \mathbb{R}^{64 \times (1 \times 8 \times 16)}$
in **layer 2** @ E_θ

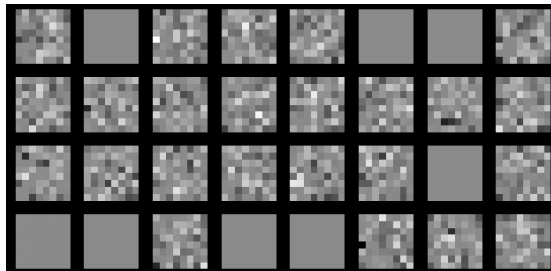
Reshape $\mathbf{W}^{(2)}$ into 64 filters of size 8x16

Takeaway: More neurons are utilized when we go deeper – not that interesting!

Fully-connected AE (**linear**)



Fully-connected AE (**nonlinear**)



Extract the weights $\mathbf{W}^{(3)} \in \mathbb{R}^{32 \times (1 \times 8 \times 8)}$
in **layer 3** @ E_θ

Reshape $\mathbf{W}^{(3)}$ into 32 filters of size 8x8

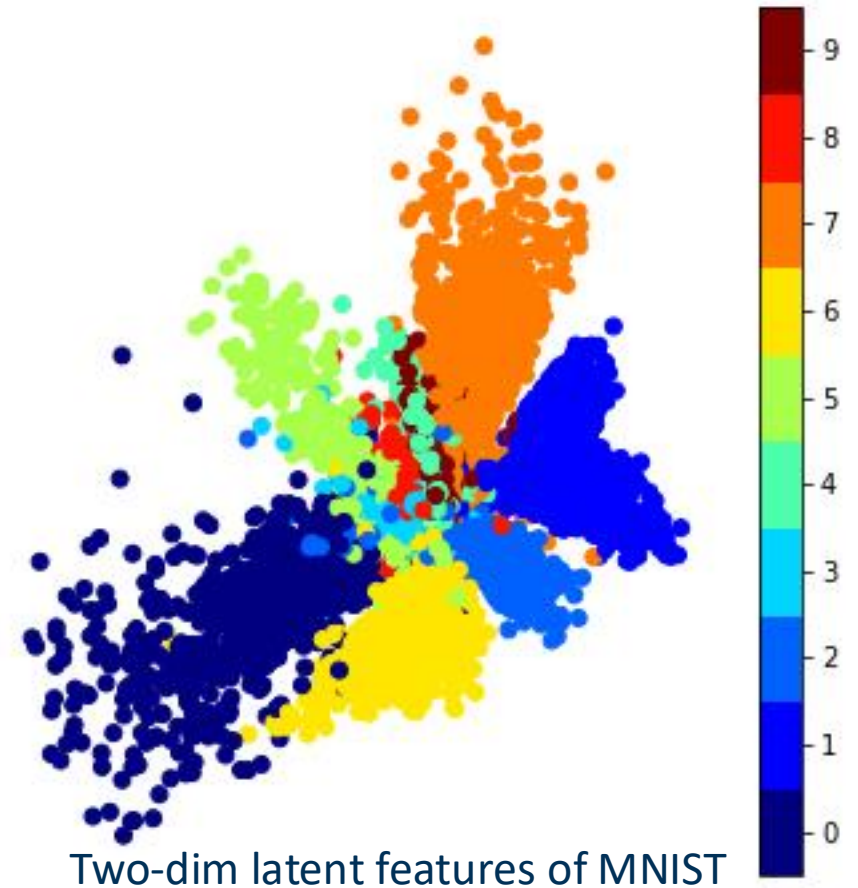
Fully-connected Autoencoders

Visualization of 2D Latent Space

We want to visualize the latent representations

Train a fc-autoencoder on MNIST:

- 2-dimensional latent representation
- visualize color-coded representations (according to their classes)



Two-dim latent features of MNIST

Fully-connected Autoencoders

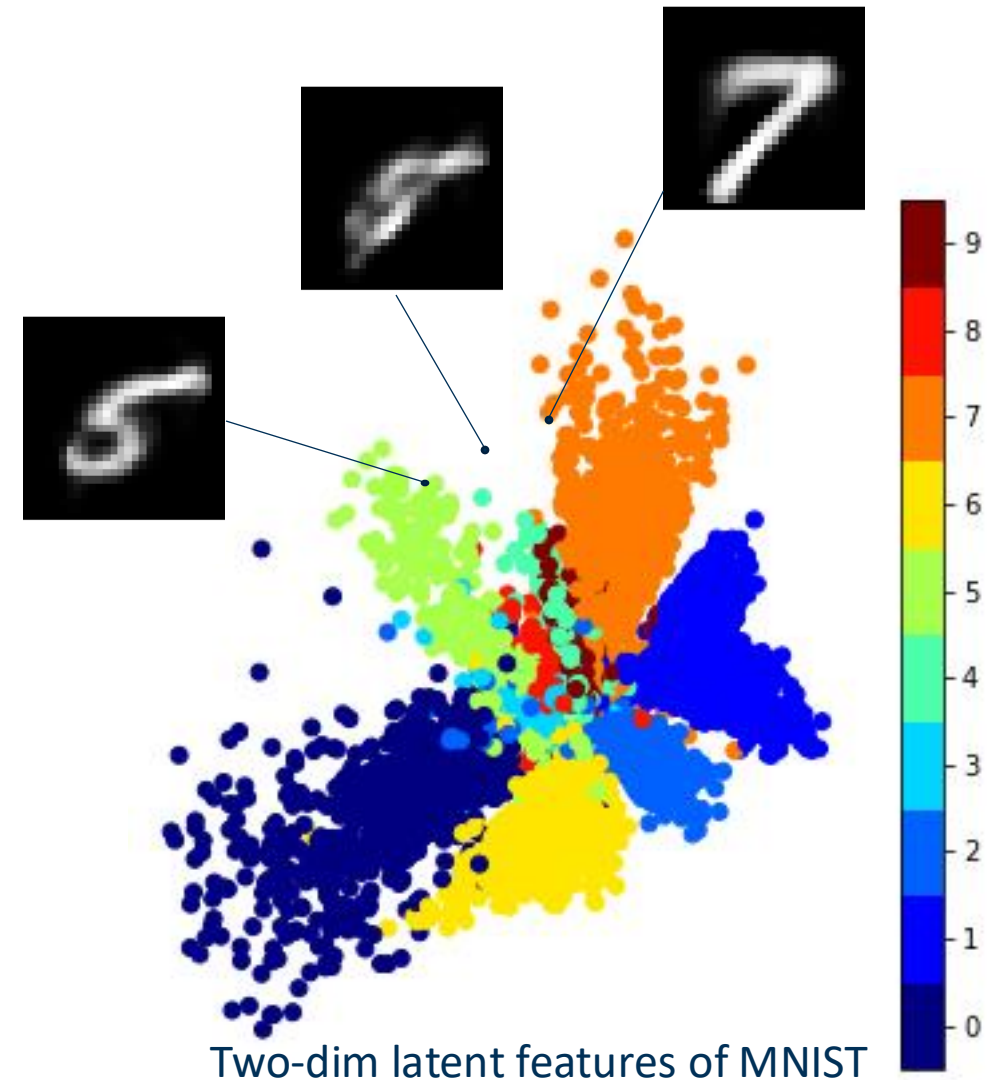
Visualization of 2D Latent Space

Train a fc-autoencoder on MNIST:

- 2-dimensional latent representation
- visualize color-coded representations (according to their classes)

The observations on latent features:

- Weak discriminability
 - Unsupervised learning discovers data structures but not aims to separate different classes
- Discontinuity:
 - The “digit-like” image decoded from gap regions look unrealistic

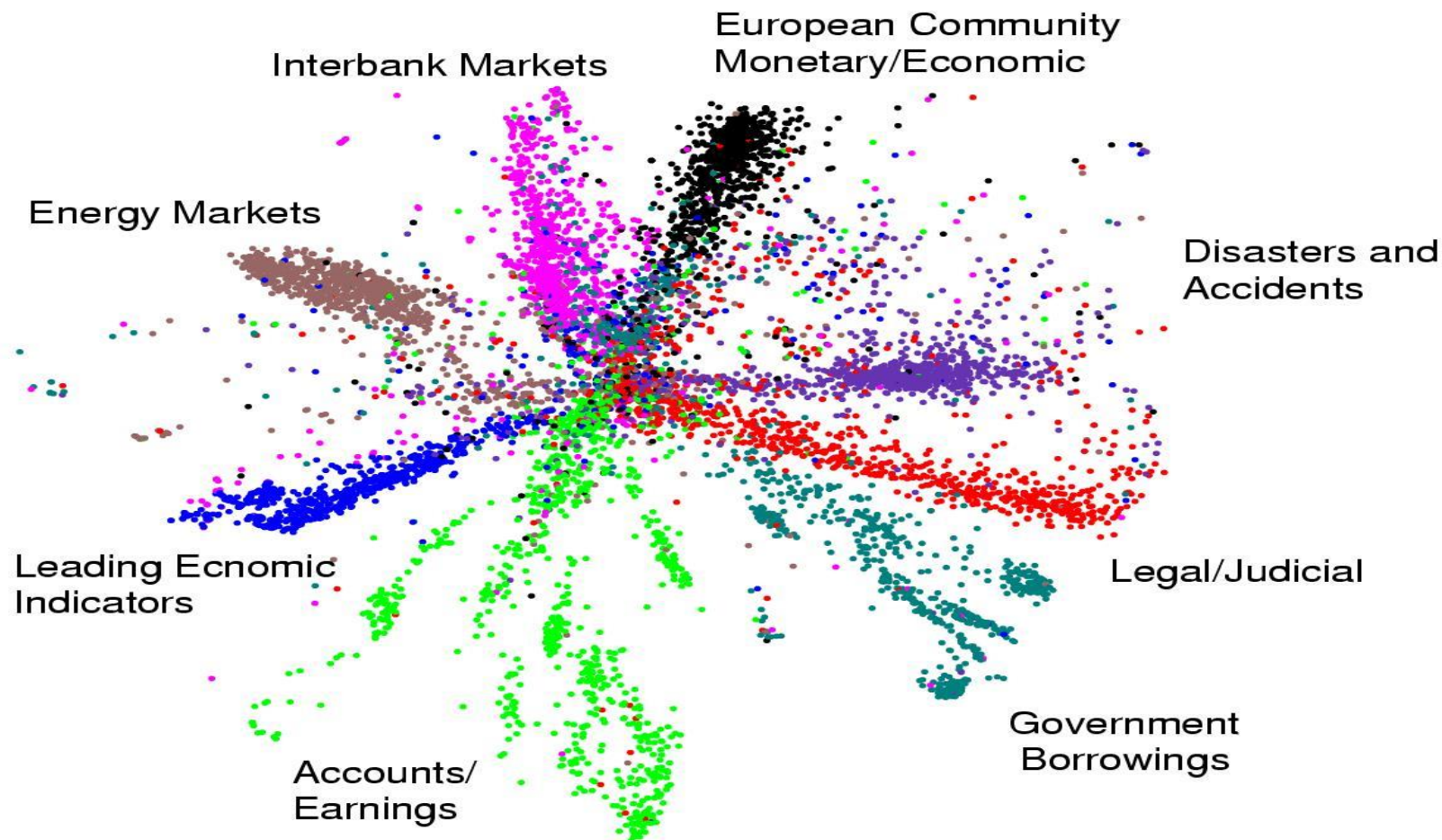


Fully-connected Autoencoders

Visualization of 2D Latent Space

- document categories

Autoencoder 2-D Topic Space



Overview

In this Lecture..

Introduction and Motivation

Fully-connected Autoencoders

Convolutional Autoencoders

- Convolutional Filters
- Unpooling

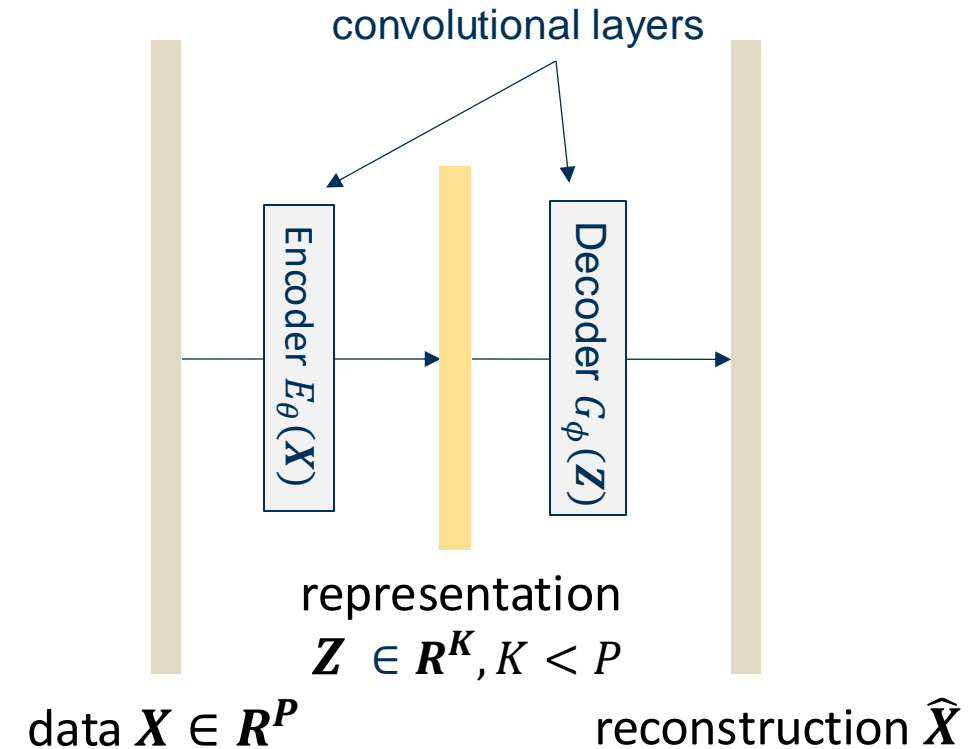
Regularized Autoencoders

Variational Autoencoders

Autoencoders

Convolutional Autoencoders

- Fully-connected autoencoders:
 - Does not leverage stationarity
 - Redundancy in the parameters
 - Learn from every sample as an independent feature
- Convolutional autoencoders:
 - Preserving **spatial locality**
 - **Sharing weights** across different locations.
 - Reconstruction based on the **local latent features**



Autoencoders

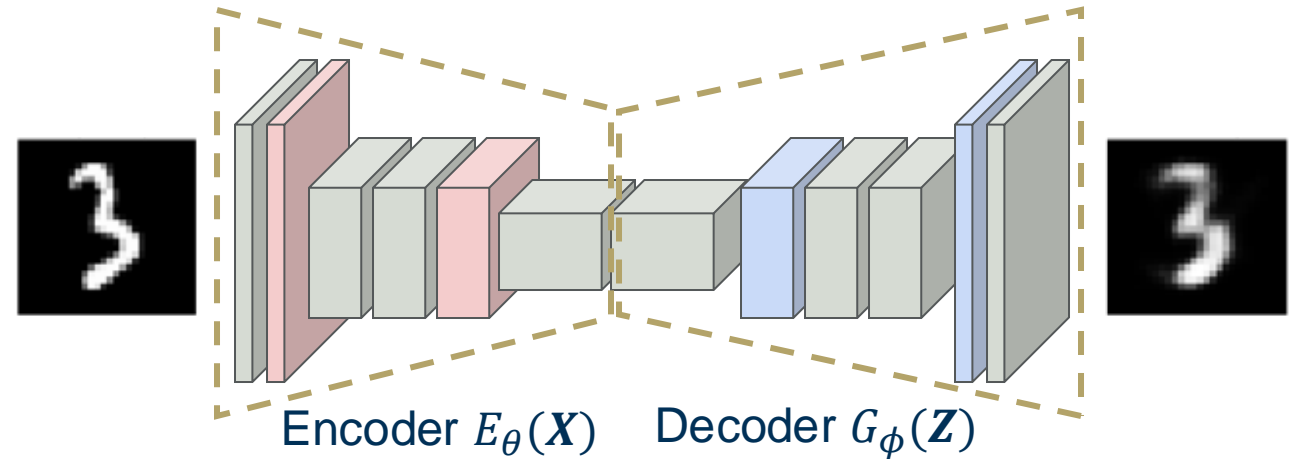
Convolutional Autoencoders

Encoder: (dimension decreases)

- Convolution (with stride)
- or Convolution + Pooling

Decoder: (dimension increases)

- Unpooling + Convolution
- or Transposed Convolution with stride
(learning the “Unpooling + Convolution”)



Convolutional Autoencoders

Unpooling

Nearest Neighbor

1	2
3	4

Input 2x2

The output feature map becomes blocky

1	1	2	2
1	1	2	2
3	3	4	4
3	3	4	4

Output 4x4

Bed of "Nails"

1	2
3	4

Input 2x2

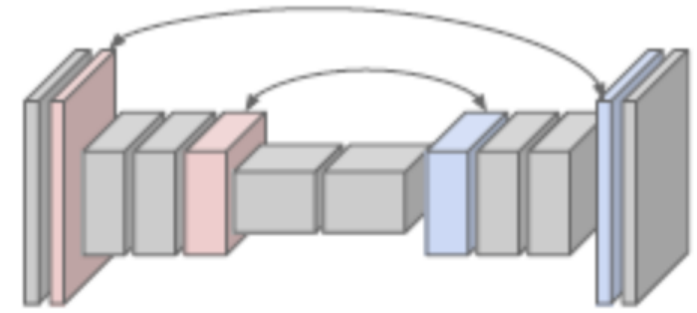
The upsampled elements always have a fixed location

1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

Output 4x4

Convolutional Autoencoders

Max-unpooling with indices



Max Pooling in $E_\theta(X)$

Remember which element was max!

1	2	6	3
3	5	2	1
1	2	2	1
7	3	4	8

Input 4x4



5	6
7	8

Output 2x2

Indices: $\begin{bmatrix} 5 & 2 \\ 12 & 15 \end{bmatrix}$

Max Unpooling in $G_\phi(Z)$

Use positions from pooling layer

5	6
7	8

Input 2x2



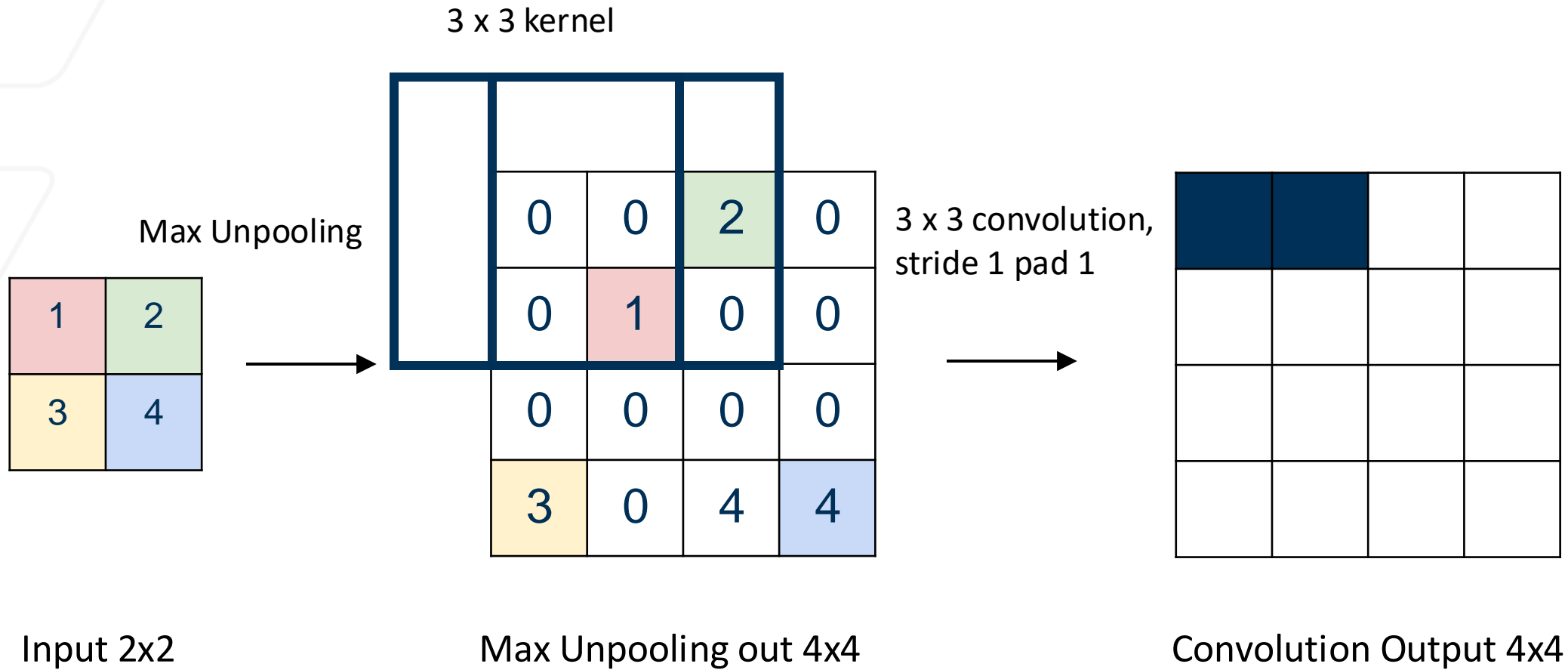
0	0	6	0
0	5	0	0
0	0	0	0
7	0	4	8

Output 4x4

Indices: $\begin{bmatrix} 5 & 2 \\ 12 & 15 \end{bmatrix}$

Decoder with convolutions

Max-unpooling with indices



Learnable Upsampling

Transposed Convolution with Stride

2 x 2 **transposed** convolution, stride 2

learn “unpooling + convolution”

1	1
1	1

filter 2x2

1	2
3	4

Input 2x2



The filter is weighted by every entry of the input, and is placed in the corresponding output location

1	1		
1	1		

Output 4x4

Learnable Upsampling

Transposed Convolution with Stride

2 x 2 **transposed** convolution, stride 2

1	1
1	1

filter 2x2

1	2
3	4

Input 2x2



Input gives
weight for filter

1	1	2	2
1	1	2	2

Output 4x4

stride 2: moves 2
pixels in the output
for every one pixel
in the input

Learnable Upsampling

Transposed Convolution with Stride

2 x 2 **transposed** convolution, stride 2

1	1
1	1

filter 2x2

1	2
3	4

Input 2x2



Input gives
weight for filter

1	1	2	2
1	1	2	2
3	3		
3	3		

Output 4x4

stride 2: moves 2
pixels in the output
for every one pixel
in the input

Learnable Upsampling

Transposed Convolution with Stride

2 x 2 **transposed** convolution, stride 2

1	1
1	1

filter 2x2

1	2
3	4

Input 2x2



Input gives
weight for filter

1	1	2	2
1	1	2	2
3	3	4	4
3	3	4	4

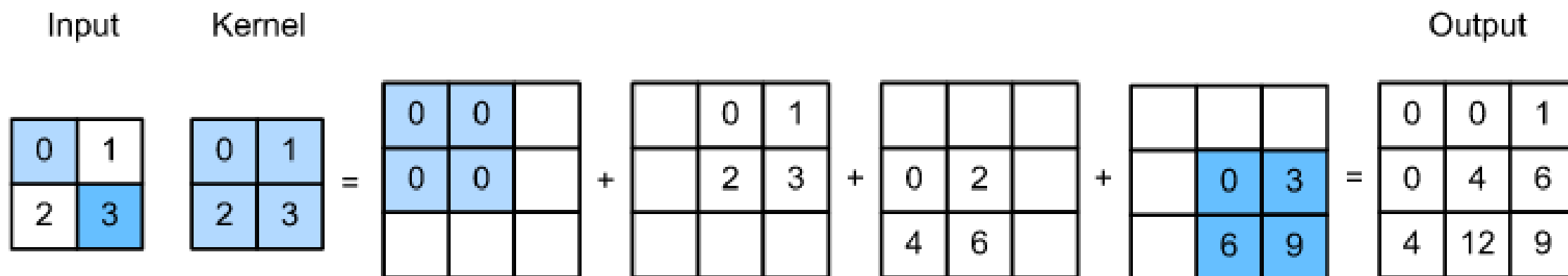
Output 4x4

stride 2: moves 2
pixels in the output
for every one pixel
in the input

Learnable Upsampling

Transposed Convolution with Stride

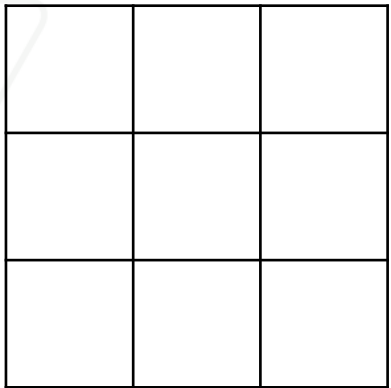
2 x 2 **transposed** convolution, stride 1



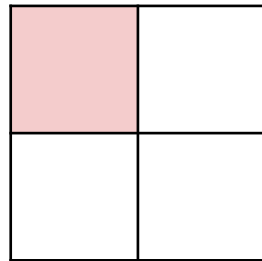
Learnable Upsampling

Transposed Convolution with Stride

3 x 3 **transposed** convolution, stride 2 pad 1



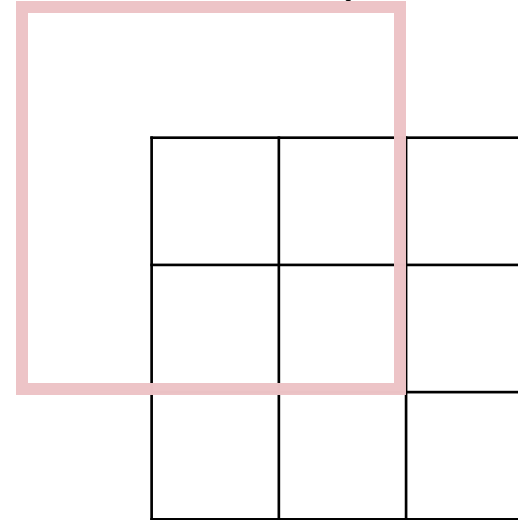
filter 3x3



Input 2x2



Input gives
weight for filter



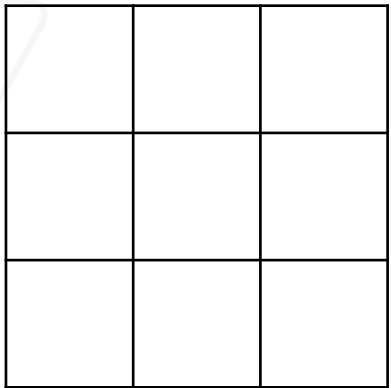
Output 3x3

“Implicit Padding” at
the output

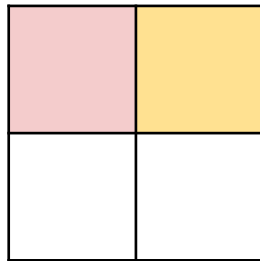
Learnable Upsampling

Transposed Convolution with Stride

3 x 3 **transposed** convolution, stride 2 pad 1



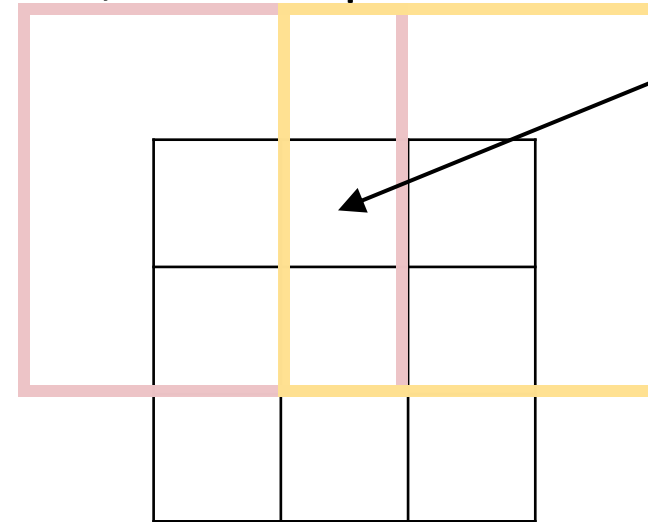
filter 3x3



Input 2x2



Input gives
weight for filter



Output 3x3

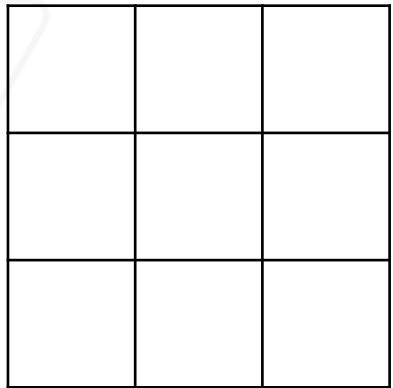
Sum where
output overlaps

moves 2 pixels in
the output for
every one pixel in
the input

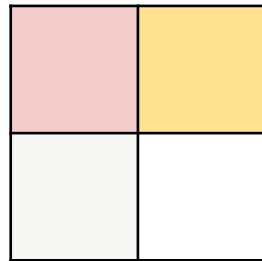
Learnable Upsampling

Transposed Convolution with Stride

3 x 3 **transposed** convolution, stride 2 pad 1



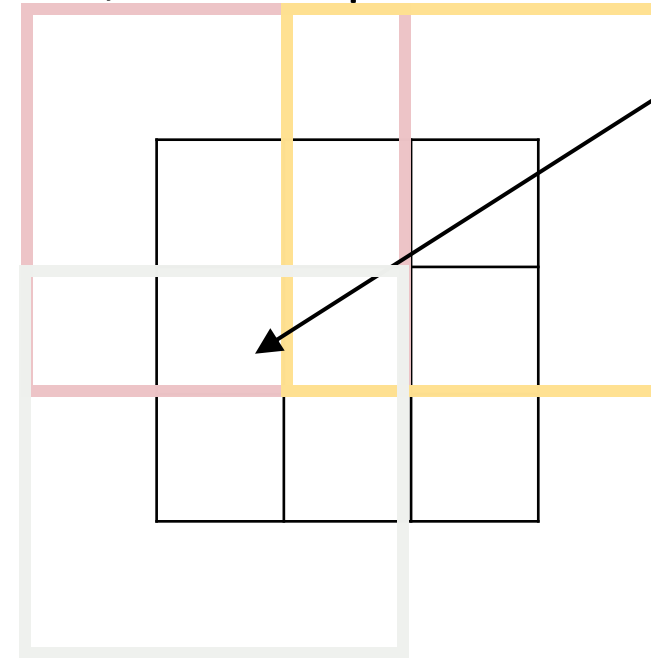
filter 3x3



Input 2x2



Input gives
weight for filter



Output 3x3

Sum where
output overlaps

moves 2 pixels in
the output for
every one pixel in
the input

Learnable Upsampling

Transposed Convolution with Stride

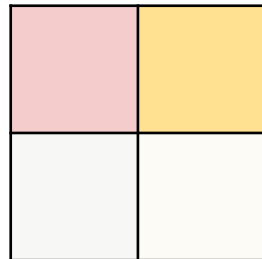
learn “unpooling + convolution”

Transposed convolution is NOT the inverse of a convolution

Other names:

- Upconvolution
- Fractionally strided convolution
- Backward strided convolution

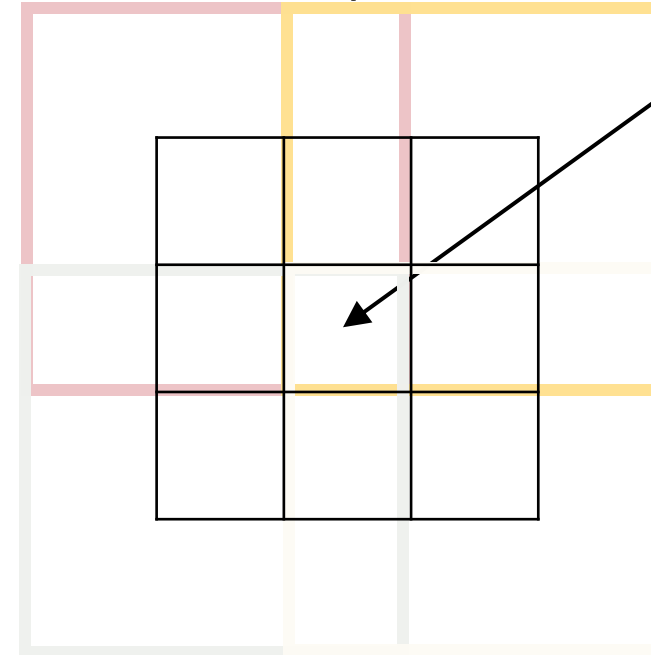
3 x 3 **transposed** convolution, stride 2 pad 1



Input 2x2



Input gives weight for filter



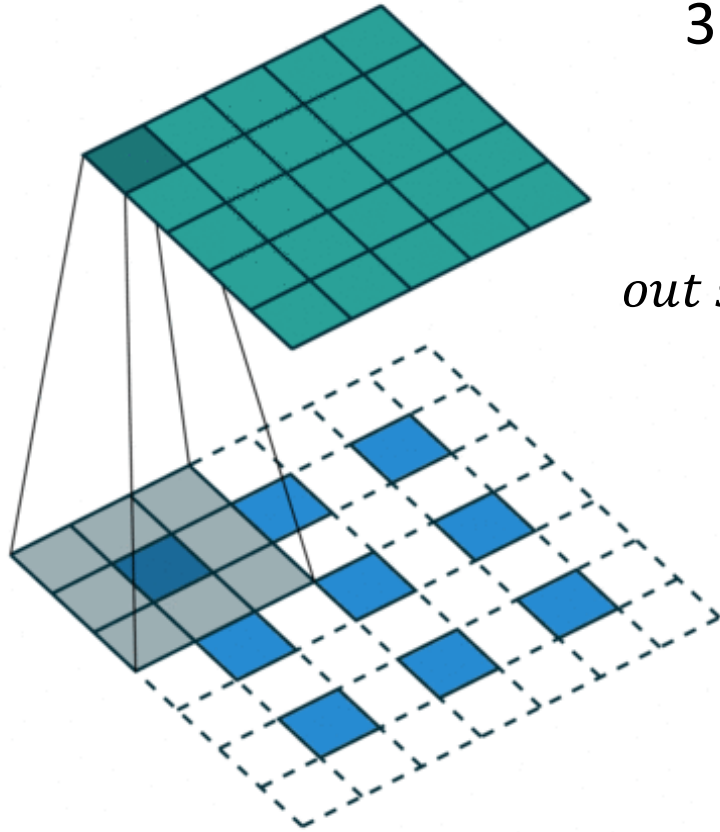
Output 3x3 (“crop” 1 pixel at the border)

Sum where output overlaps

moves 2 pixels in the output for every one pixel in the input

Transposed Convolution

Upsampling



3 x 3 **transposed** convolution, stride 2 pad 1

$$out\ size = stride \times (input\ size - 1) + filter\ size - 2 \times padding$$

$$5 = 2 \times (3 - 1) + 3 - 2 \times 1$$

Recall for convolutional operation:

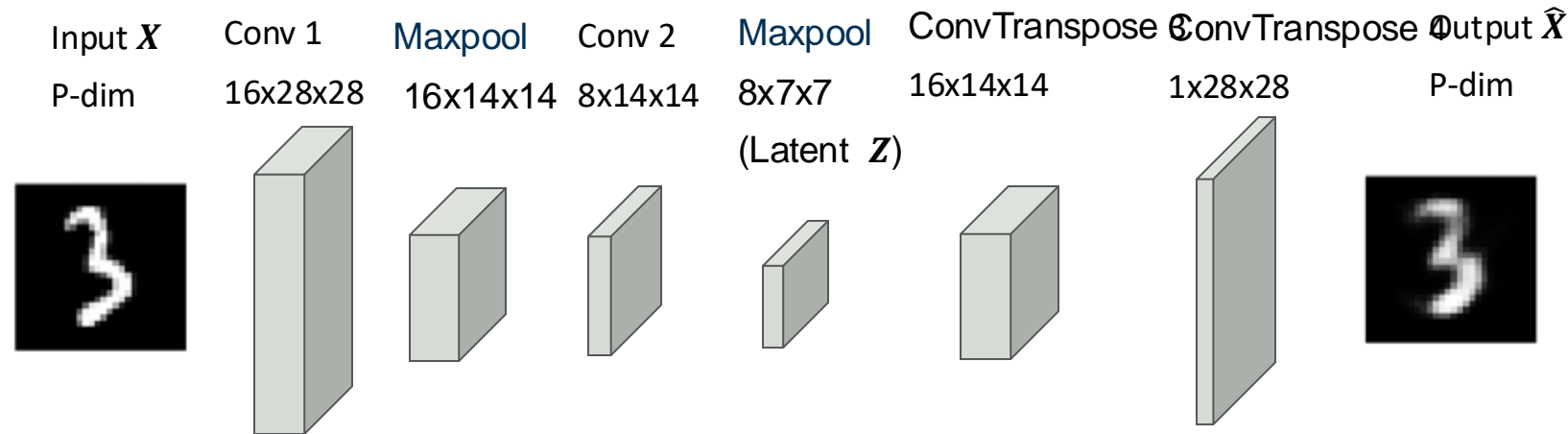
$$out\ size = \left\lfloor \frac{input\ size + 2 \times padding - filter\ size}{stride} + 1 \right\rfloor$$

Inputs: 3 x 3 blue maps, outputs: 5 x 5 green maps

Convolutional Autoencoder

Training on MNIST

Convolutional autoencoder:



Encoder:

- Conv 1: 16 kernels of 3x3, padding=1 (ReLU)
- Maxpooling: stride = 2
- Conv 2: 8 kernels of 3x3, padding=1 (ReLU)
- Maxpooling: with stride = 2

Decoder:

- Transposed conv 3: 16 kernels of 2x2, stride=2 (ReLU)
- Transposed conv 4: 1 kernels of 2x2, stride=2 (Sigmoid)

Experimental Results

FC-linear AE vs FC-nonlinear AE vs Convolutional AE

Fully-connected AE (linear)



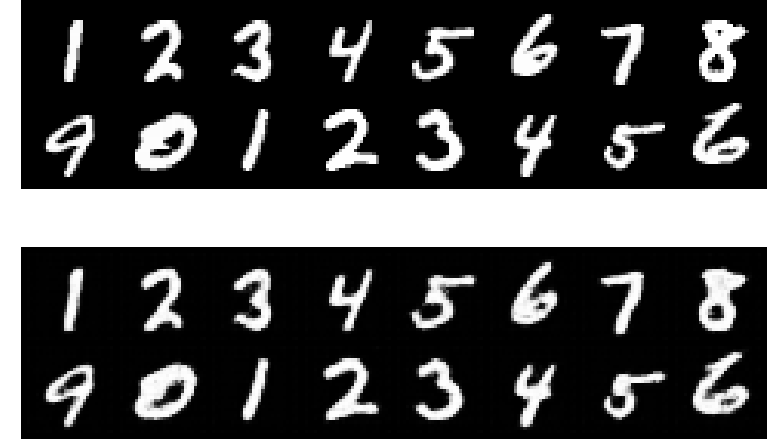
MSE = 0.0677

Fully-connected AE (nonlinear)



MSE = 0.0306

Convolutional AE



MSE = 0.0130

input (top row) and reconstruction (bottom row)

Terminology

- *Distribution*: (sample space) the set of all possible samples
- *Dataset*: a set of samples drawn from a distribution
- *Batch*: a subset of samples drawn from the dataset
- *Sample*: a single data object represented as a set of features
- *Feature*: value of a single attribute, property, in a sample. Could be numeric or categorical.

Appendix A: Notations

- x_i : a single feature
- \mathbf{x}_i : feature vector (a data sample)
- $\mathbf{x}_{:,i}$: feature vector of all data samples
- \mathbf{X} : matrix of feature vectors (dataset)
- \mathbf{W} : weight matrix
- \mathbf{Z} : latent representation
- E_θ : encoding function
- G_ϕ : decoding function
- $\hat{\mathbf{X}}$: reconstruction of data
- $\Omega(\mathbf{Z})$: sparsity constraint
- $\hat{\rho}_j$: average activation of neuron z_{ij}
- $\tilde{\mathbf{X}}$: corrupted input
- N : number of data samples
- P : number of features in a feature vector
- $P^{(k)}$: the number of neurons in layer k
- α : learning rate
- Bold letter/symbol: vector
- Bold capital letters/symbol: matrix