Introduction and Motivation

Fully-connected Autoencoders

Convolutional Autoencoders

Regularized Autoencoders

- Sparse Autoencoders
- Denoising Autoencoders

Variational Autoencoders

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

# Regularized Autoencoders
## Motivation

- Standard autoencoders learn useful representations by restricting the dimension to be small

- A higher dimensional representation helps model more complex data distributions

- We want to **regularize** representations to learn important and complex features without restricting the latent dimension

Input data $X$

10-dim latent $Z$

reconstruction $\widehat{X}$

Input data $X$

2-dim latent $Z$

reconstruction $\widehat{X}$



Encoder $E_\theta(X)$

Decoder $G_\phi(Z)$

representation $Z$

data $X$

reconstruction $\widehat{X}$
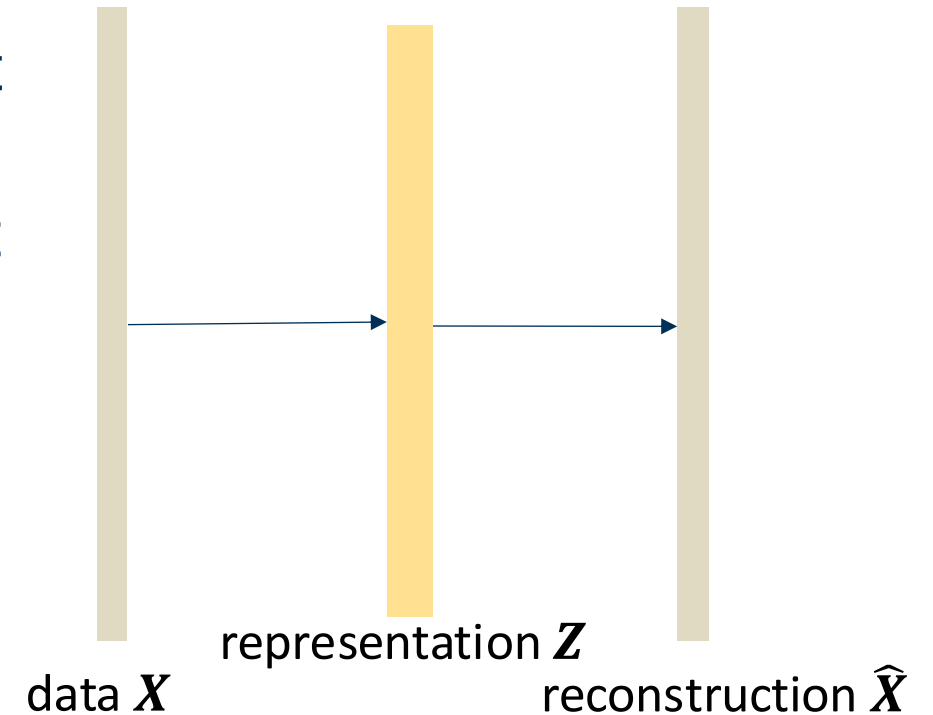
OLIVES
@GeorgiaTech

Georgia Tech

# Sparse Autoencoders
## Motivation

- We want to avoid identity mapping without limiting the representation capability to extract features from the complex data

- a high-dimensional latent space can represent more complex data
  - $\dim(\boldsymbol{Z})$ can be larger than $\dim(\boldsymbol{X})$

- reduce the redundancy to improve generalization



Redundant information learned by high-dim latent Z

representation $\boldsymbol{Z}$

data $\boldsymbol{X}$

reconstruction $\widehat{\boldsymbol{X}}$

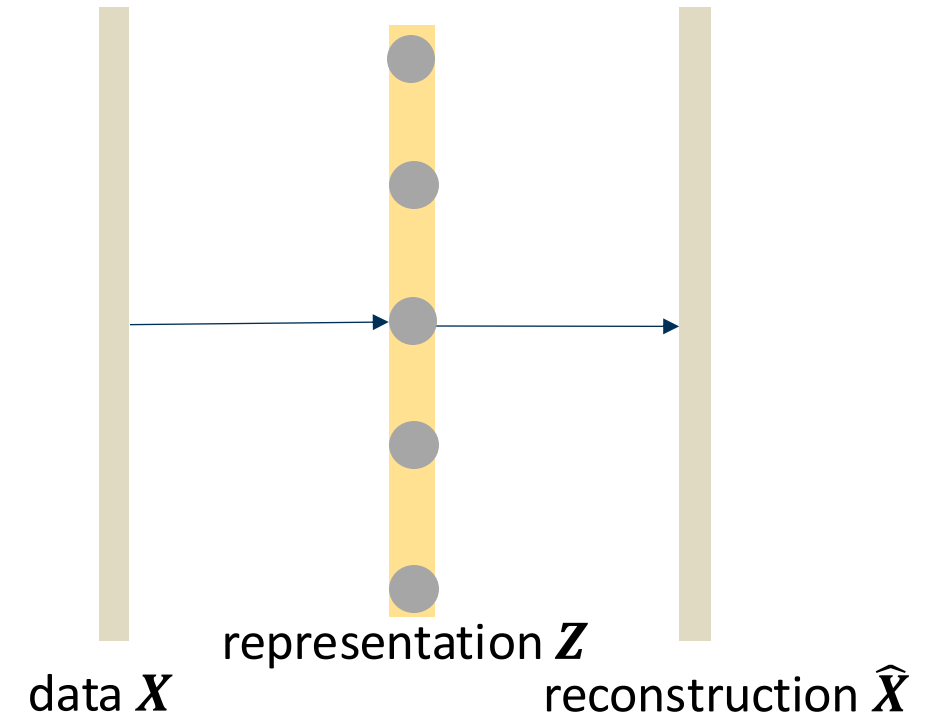[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

# Sparse Autoencoders
## Motivation

Sparse representations:

- Represent data using a subset of latent neurons being active at the same time

- Force the model to learn the unique statistical features that can be used for other tasks such as classification.



data $X$      representation $Z$      reconstruction $\widehat{X}$

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

Georgia Tech

# Sparse Autoencoders
## Motivation

Learning sparse representations:

- Loss $L = \left\| X - \widehat{X} \right\|^2 + \boxed{\lambda \mathbf{\Omega}(\mathbf{Z})}$

where

- $\mathbf{\Omega}(\mathbf{Z})$: *sparsity* constraint
- $\lambda$: constraint coefficient

- Two methods to enforce sparsity:
    - $L_1$ regularization
    - Kullback-Leibler (KL) divergence



representation $\mathbf{Z}$

data $\mathbf{X}$    reconstruction $\widehat{\mathbf{X}}$

OLIVES
@GeorgiaTech

Georgia Tech.

Exploiting the $L_1$ norm as sparsity constraint:

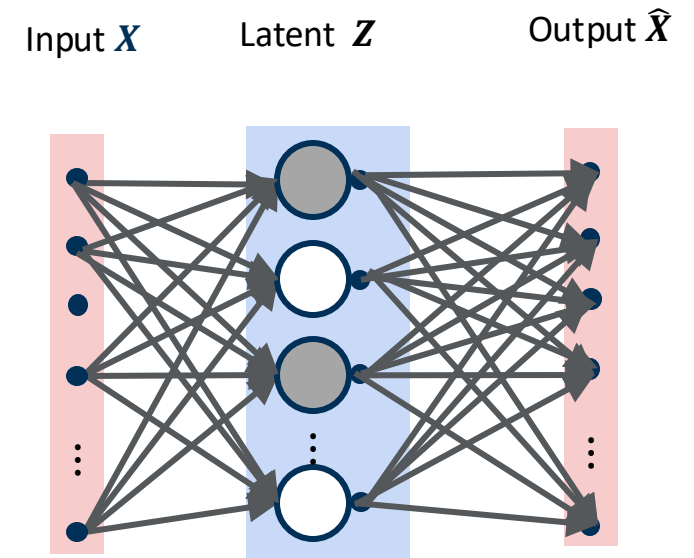L1 norm on activation $z_i$: $\|z_i\|_1 = \sum_j |z_{ij}|$

- $z_{ij}$: activation of $j$-th latent neuron given the input $x_i$

Thus,

$$\Omega(Z) = \sum_i \|z_i\|_1$$

suppresses activations towards 0 and achieve sparsity

Input $X$     Latent $Z$     Output $\widehat{X}$

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

Georgia Tech

# Sparse Autoencoders
## Sparsity via KL Divergence

Goal: The latent neuron activations must be **near** 0
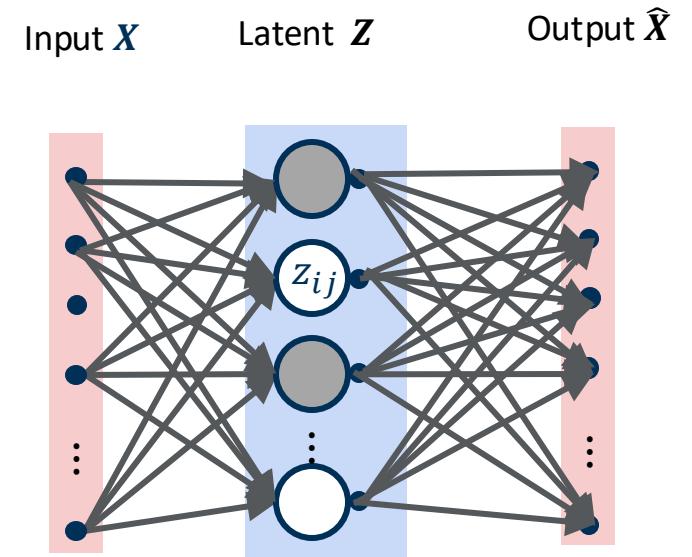
Let

$$\hat{\rho}_j = \frac{1}{N} \sum_{i=1}^{N} z_{ij}$$

be the average activation of neuron $z_{ij}$

where

- $z_{ij}$: activation of $j$-th latent neuron given the input $x_i$
- $N$: the number of training samples

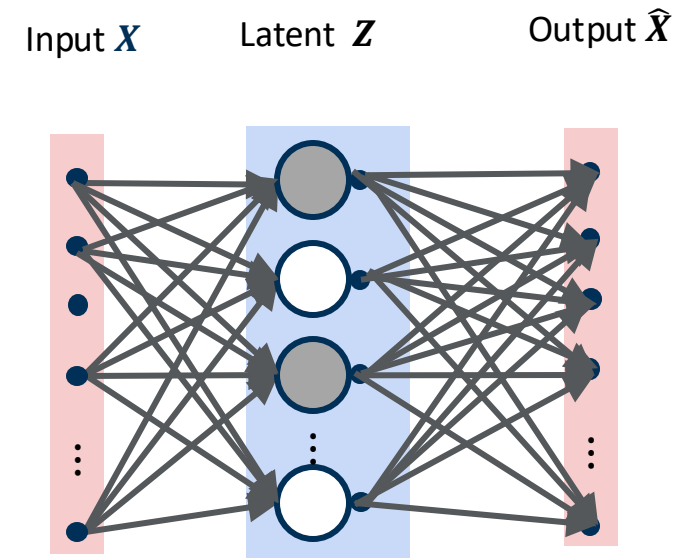Thus, we want to constrain $\hat{\rho}_j$ to be close to zero



Input $X$    Latent $Z$    Output $\hat{X}$

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

Georgia Tech

# Sparse Autoencoders
## Sparsity via KL Divergence

- Goal: The latent neuron activations must be **near** 0

- Assume:

  - using sigmoid

  - $z_{ij}$ is a Bernoulli random variable

  - $\hat{\rho}_j = \frac{1}{N}\sum_{i=1}^{N} z_{ij}$ measures the *observed* expectation

- We want $z_{ij}$ to be drawn from a Bernoulli distribution with mean $\rho$ to enforce sparsity

  - Penalize $\hat{\rho}_j$ for deviating from $\rho$

  - Smaller $\rho$ forces sparser $\boldsymbol{Z}$

- Kullback-Leibler (KL) divergence measures the difference between two distributions

Input $\boldsymbol{X}$      Latent $\boldsymbol{Z}$      Output $\hat{\boldsymbol{X}}$

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

Georgia Tech

# Sparse Autoencoders
Sparsity via KL Divergence

Exploiting the Kullback-Leibler (KL) divergence:

$$D_{\mathrm{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$$

$$\mathbf{\Omega(Z)} = \sum_j KL\left(\rho \| \hat{\rho}_j\right) = \sum_j \left( \rho \log \frac{\rho}{\hat{\rho}_j} + (1-\rho) \log \frac{1-\rho}{1-\hat{\rho}_j} \right)$$

where $\hat{\rho}_j = \frac{1}{N} \sum_{i=1}^{N} z_{ij}(x_i)$

- $KL\left(\rho \| \hat{\rho}_j\right) = 0$ if $\rho = \hat{\rho}_j$, and increases monotonically as $\hat{\rho}_j$ diverges from $\rho$

- $KL\left(\rho \| \hat{\rho}_j\right)$ has **control over sparsity** via $\rho$
  - Smaller $\rho$ forces sparser $\mathbf{Z}$



$\rho = 0.05$, plot the KL divergence for $\hat{\rho} \in (0, 1)$

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]
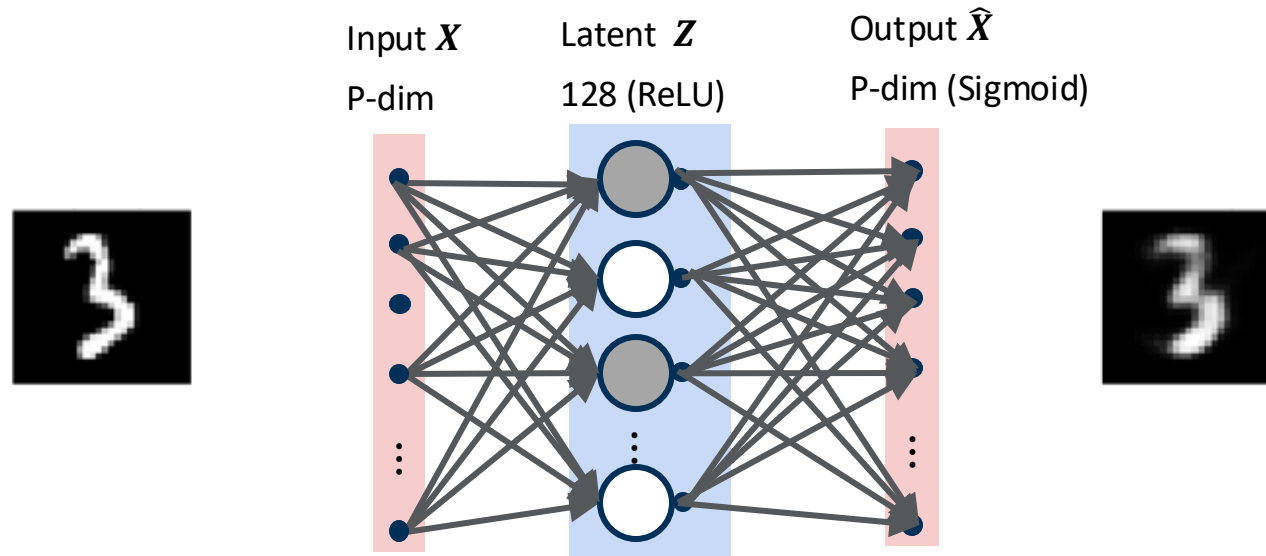
# Sparse Autoencoders
## Differences in Sparsity Constraints

- Kullback-Leibler (KL) divergence
  - control over sparsity via $\rho$
  - Approximate sparsity, suppress activations close to small value
  - differentiable at 0

- $L_1$ regularization
  - no explicit control over sparsity
  - True sparsity, drive activations towards 0
  - not differentiable at 0

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

Georgia Tech

# Sparse Autoencoders
## Training

- **Sparsity** constraint on the latent activation

- Loss $L = \left\| X - \widehat{X} \right\|^2 + \boxed{\lambda \Omega(Z),}$ where $\lambda$ is the constraint coefficient



Input $X$
P-dim

Latent $Z$
128 (ReLU)

Output $\widehat{X}$
P-dim (Sigmoid)

- For simpler implementation, we consider having one hidden layer besides input and output

OLIVES
@GeorgiaTech

Georgia Tech
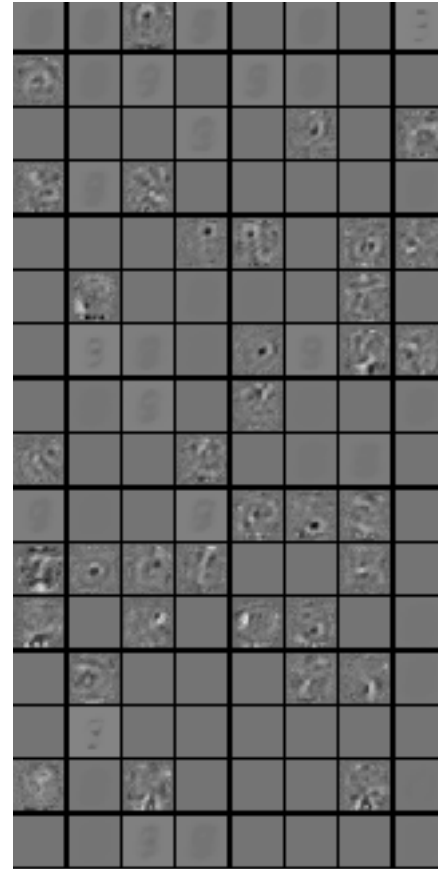
# Sparse Autoencoders
## Visualization of Trained Filters

Visualize the filters of the first layer of a fully-connected AE and a sparse AE with the same architecture components.

Fully-connected AE (nonlinear)    Sparse AE (**KL**, nonlinear), $\lambda = 0.001$    Sparse AE (**KL**, nonlinear), $\lambda = 0.003$



Extract the weights $W^{(1)} \in \mathbb{R}^{128 \times (1 \times 28 \times 28)}$ in **layer 1** @ $E_\theta$
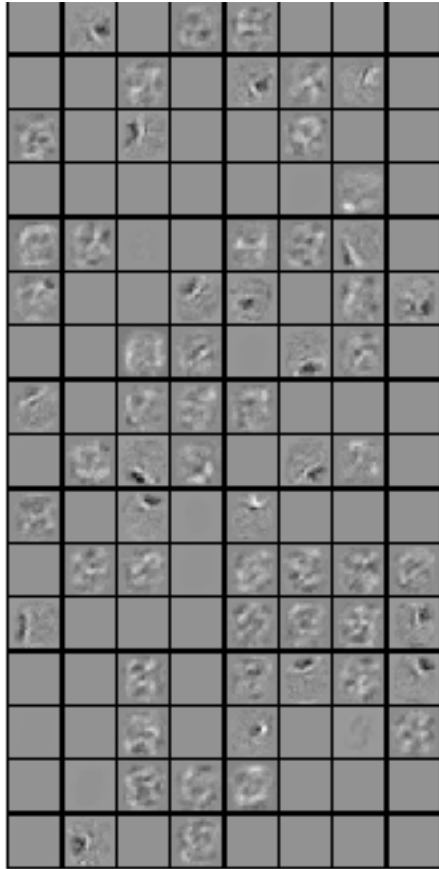
Reshape $W^{(1)}$ into 128 filters of size 28x28

Sparse AE uses **KL divergence** with $\rho = 0.05$

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

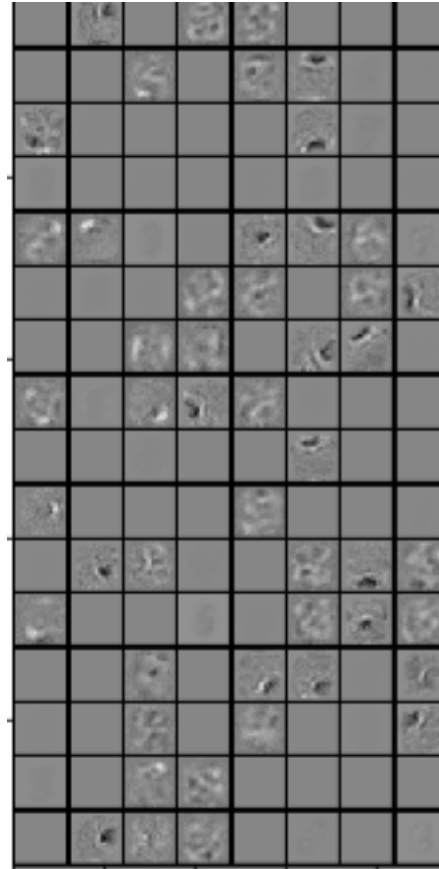Georgia Tech.

# Sparse Autoencoders
## Visualization of Trained Filters

Visualize the filters of the first layer of a fully-connected AE and a sparse AE with the same architecture components.
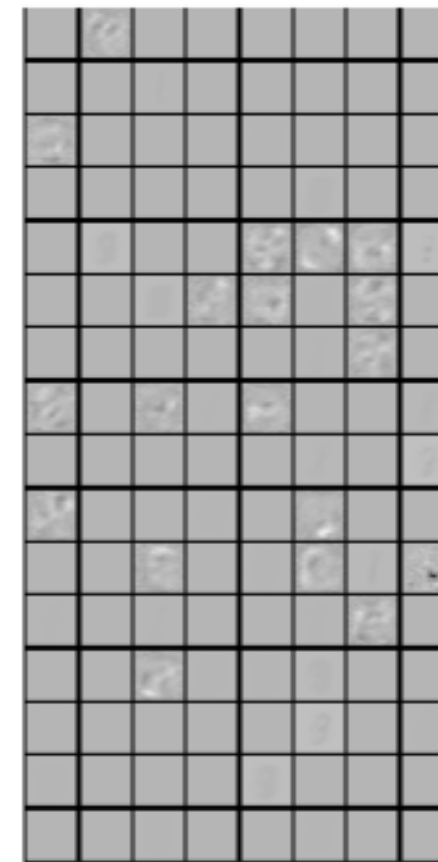
Fully-connected AE (nonlinear)     Sparse AE (**L1**, nonlinear), $\lambda = 0.001$     Sparse AE (**L1**, nonlinear), $\lambda = 0.005$



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

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

Sparse AE uses **L1 regularization**

OLIVES @GeorgiaTech

Georgia Tech

Introduction and Motivation

Fully-connected Autoencoders

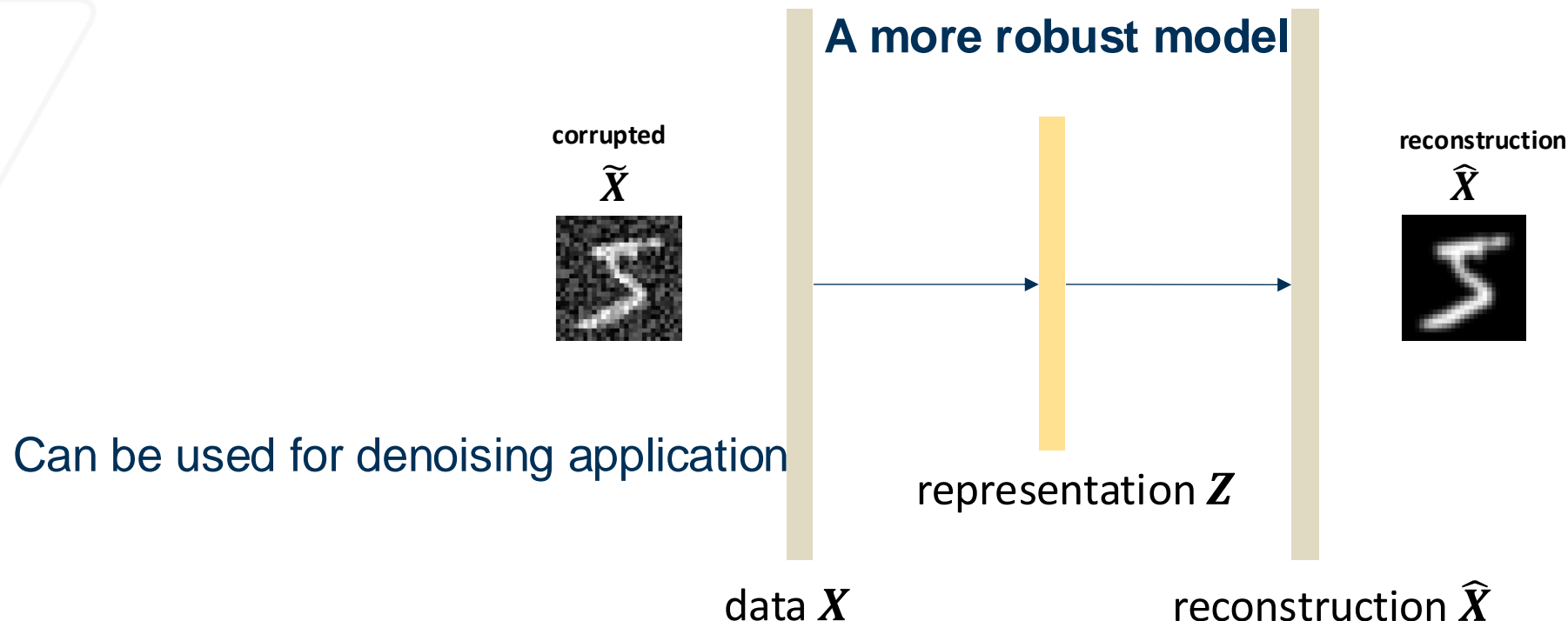Convolutional Autoencoders

Regularized Autoencoders

- Sparse Autoencoders

- Denoising Autoencoders

Variational Autoencoders

OLIVES
@GeorgiaTech

Georgia Tech

# Denoising Autoencoders
## Motivation

- We still aim to encode the input and avoid identity mapping
- We try to learn a **robust representation** that undo the effect of *corruption* applied to the input
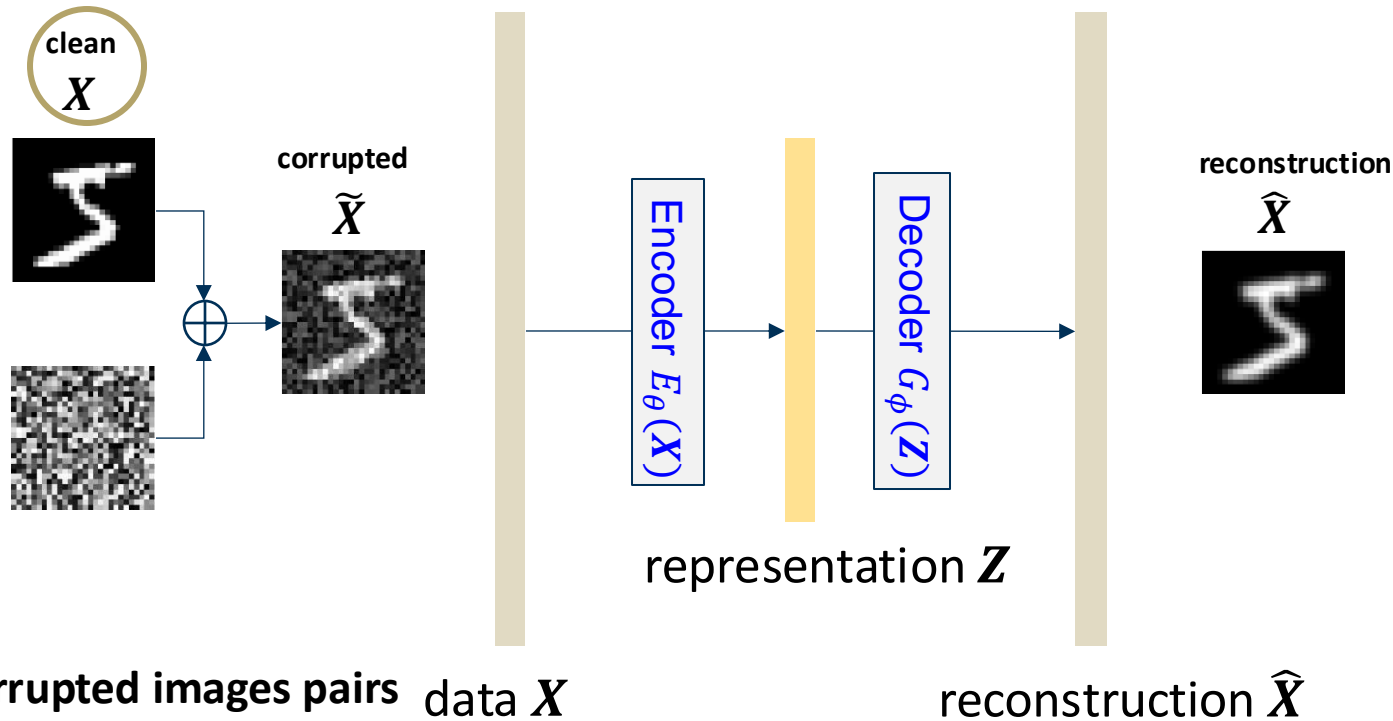
**A more robust model**

corrupted

$\widetilde{X}$



representation $\mathbf{Z}$

reconstruction

$\widehat{X}$



Can be used for denoising application

data $\mathbf{X}$      reconstruction $\widehat{\mathbf{X}}$

Learn a **robust representation** from **corrupted input** $\widetilde{X}$

Minimize a loss $L = \|X - \widehat{X}\|^2, \widehat{X} = G_\phi(E_\theta(\widetilde{X}))$



need **clean** and **corrupted images pairs**   data $X$      reconstruction $\widehat{X}$

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

A 'not prone to changes' model

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

Learn a **robust representation** from **corrupted input** $\widetilde{X}$

Minimize a loss $L = \left\| X - \widehat{X} \right\|^2, \widehat{X} = G_\phi(E_\theta(\widetilde{X}))$



need **clean** and **corrupted images pairs**

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]
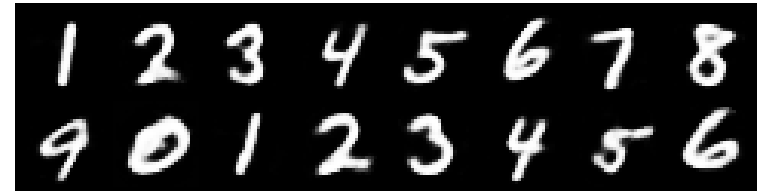
# Denoising Autoencoders
## Results

Denoising Fully-connected AE (linear)

Denoising Fully-connected AE (nonlinear)



MSE = 0.0694

MSE = 0.0328

input (upper) and reconstruction (bottom)

OLIVES
@GeorgiaTech

Georgia Tech

# Overview
## In this Lecture..

**Introduction and Motivation**

**Fully-connected Autoencoders**
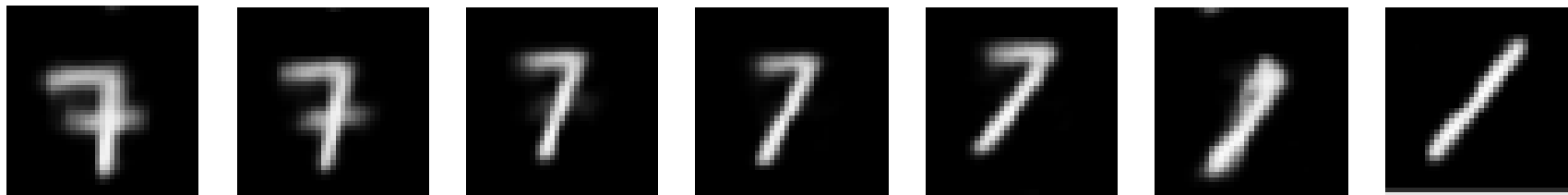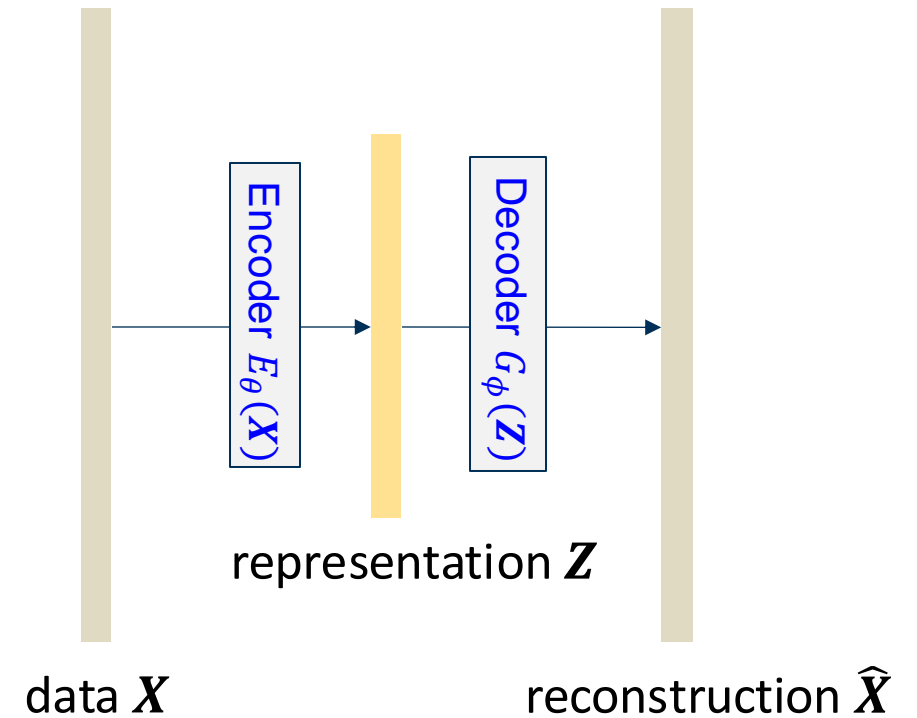
**Convolutional Autoencoders**

**Regularized Autoencoders**

**Variational Autoencoders**

- Motivation
- Variational Lower Bound
- Reparameterization Trick
- Generating Novel Samples
- Visualization

OLIVES
@GeorgiaTech

Georgia Tech

Motivation

- Autoencoders can **reconstruct** the input data via deterministic processes $Z = E_\theta(X)$, $\widehat{X} = G_\phi(Z)$

- Can we generate novel samples?
  - If $Z$ is a random variable drawn from some distribution, we can *sample* from the latent space and **generate novel variations** of inputs

representation $Z$

data $X$          reconstruction $\widehat{X}$

varying digits from 7 to 1

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

# Variational Autoencoders
## Motivation



Latent space of standard autoencoder trained on MNIST

Recall the latent space of regular autoencoders:

- Discontinuity
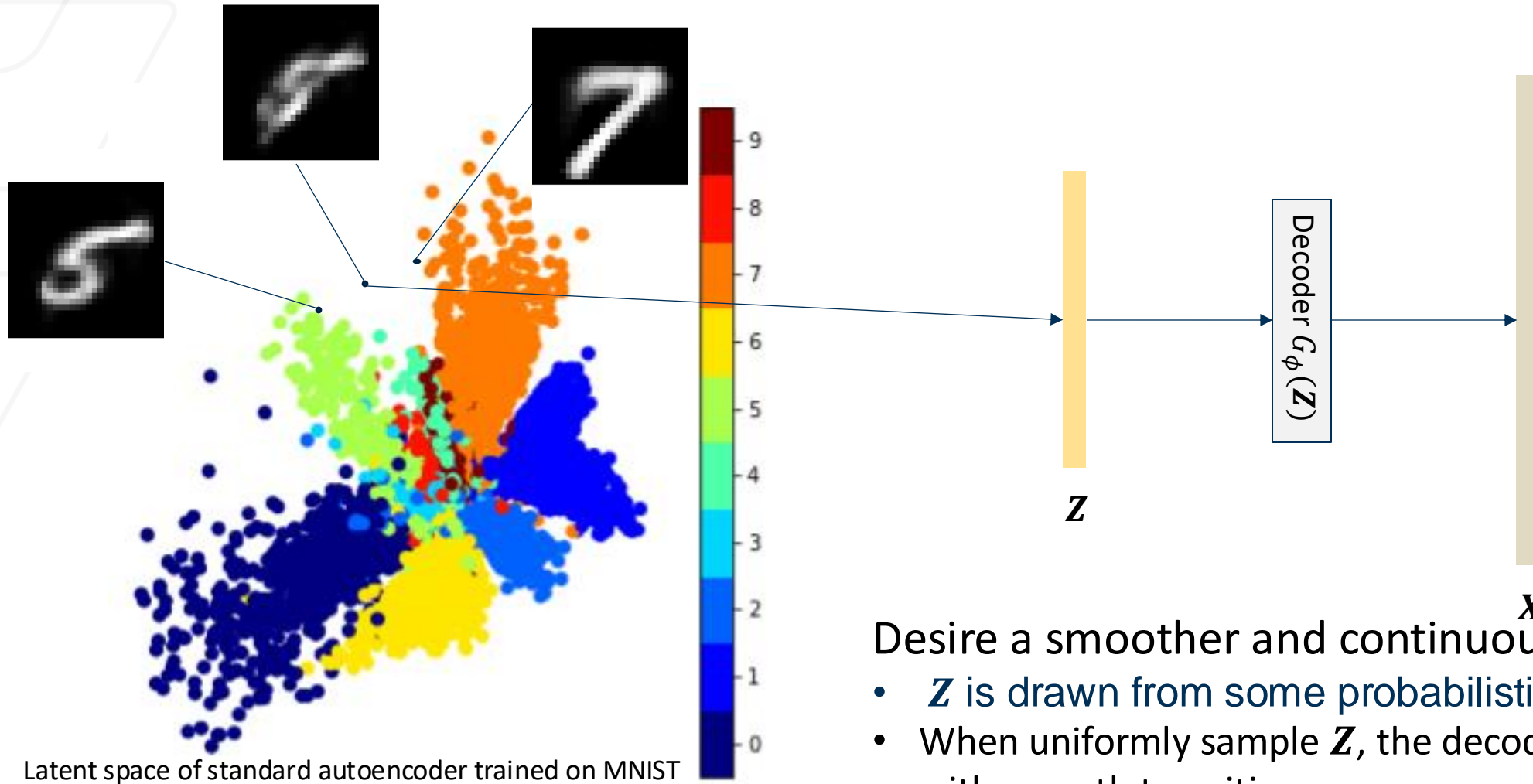  - decoder generates an unrealistic output if sample/generate a variation from there

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

Georgia Tech

# Variational Autoencoders
## Motivation



Latent space of standard autoencoder trained on MNIST

Decoder $G_\phi(\mathbf{Z})$

$\mathbf{Z}$

$X$

Desire a smoother and continuous latent space:
- $\mathbf{Z}$ is drawn from some probabilistic distribution
- When uniformly sample $\mathbf{Z}$, the decoded images vary with smooth transition

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

# Variational Autoencoders
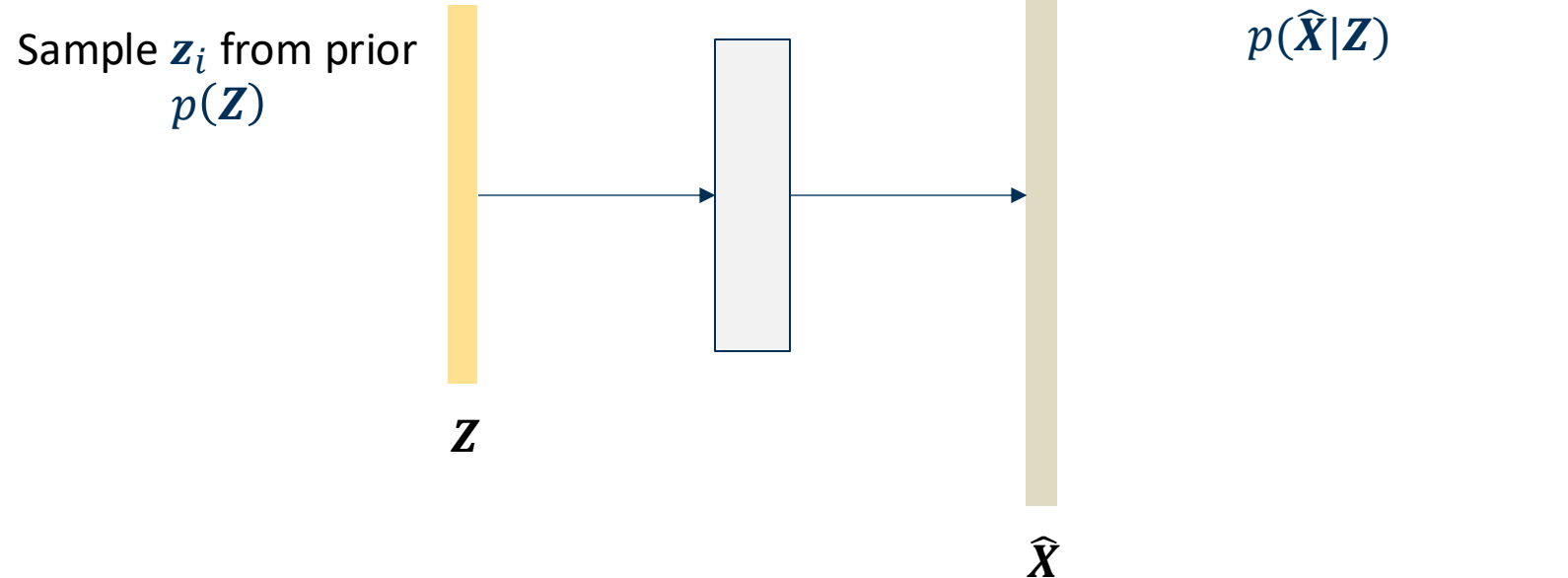Probabilistic Perspective

To generate datapoint $\hat{x}_i$:

- Assume we sample latent variables $z_i \sim p(Z)$
  - where $p(Z)$ is a prior distribution

- Given sampled $z_i \sim p(Z)$, draw $\hat{x}_i \sim p(\hat{X}|Z)$

Conditional distribution
$p(\hat{X}|Z)$

Sample $z_i$ from prior
$p(Z)$

$Z$

$\hat{X}$

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

We want to find a probabilistic modeling $p_\phi(X)$ to estimate true distribution of $X$

$$p_\phi(X) = \int p_\phi(Z) p_\phi(X|Z) dZ$$

$\phi$ are the parameters for the generation model

Sample $z_i$ from prior
$p(Z)$

Sample Conditional distribution
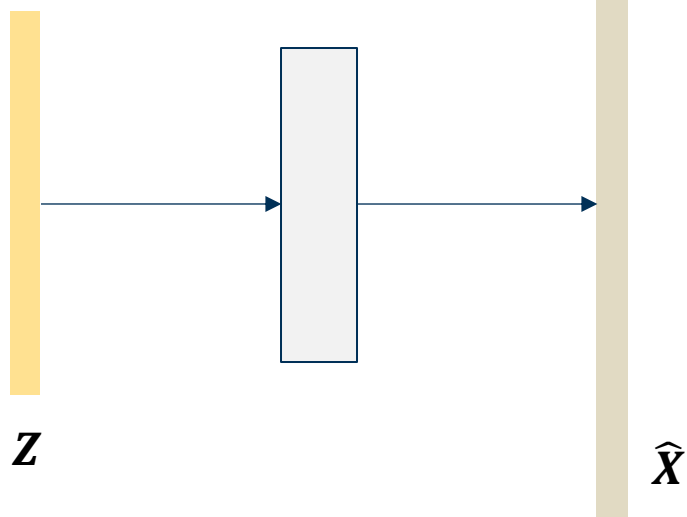$p(\widehat{X}|Z)$

$Z$

$\widehat{X}$

# Variational Autoencoders
Probabilistic Perspective

We want to find a probabilistic modeling $p_\phi(X)$ to estimate true distribution of $X$

$$p_\phi(X) = \int p_\phi(Z) p_\phi(X|Z) dZ$$

$\phi$ are the parameters for the generation model

Sample $z_i$ from prior
$p(Z)$

Sample Conditional distribution
$p(\widehat{X}|Z)$

$Z$

$\widehat{X}$

How should we represent this model?
- Choose **prior** $p(Z)$ to be **simple and tractable**, e.g., *Gaussian* $\mathcal{N}(0, \mathbf{I})$.
- **Conditional** $p(\widehat{X}|Z)$ is **complex** (generates image) → modeled by **decoder** neural network

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

Georgia Tech

# Variational Autoencoders
## Probabilistic Modeling via MLE
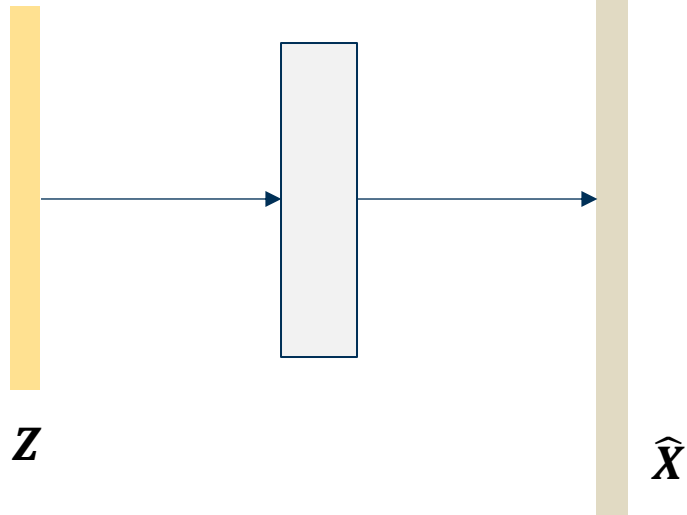
We want to find a probabilistic modeling $p_\phi(X)$ to estimate true distribution of $X$

$$p_\phi(X) = \int p_\phi(Z) p_\phi(X|Z) dZ$$

**Simple Gaussian prior**

$\phi$ are the parameters for the generation model

Sample $z_i$ from prior
$p(Z)$

Sample Conditional distribution
$p(\widehat{X}|Z)$

Find the optimal $\phi$ via **maximizing likelihood** of training data

$Z$

$\widehat{X}$

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

Georgia Tech

Data likelihood: $p_\phi(X) = \int p_\phi(Z) \boxed{p_\phi(X|Z)} dZ$

Intractable to compute $p(X|Z)$ and integrate for every dimension of $Z$!

Decoder $G_\phi(Z)$

$Z$

$\widehat{X}$

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

Posterior is Intractable

Data likelihood: $p_\phi(X) = \int p_\phi(Z) p_\phi(X|Z) dZ$

Posterior density also intractable: $p_\phi(Z|X) = p_\phi(X|Z) p_\phi(Z) / p_\phi(X)$

use an encoder to model $q_\theta(Z|X)$ that approximates $p_\phi(Z|X)$

$q_\theta(Z|X)$

$p_\phi(X|Z)$

Encoder $E_\theta(X)$

Decoder $G_\phi(Z)$

$Z$

$Z$

Approximate Posterior

$X$

$\widehat{X}$

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

Georgia
Tech

# Variational Autoencoders
## Variational Lower Bound

Reconstruct the input data

$$\mathbf{E_z}\big[\log p_\phi(\boldsymbol{x_i}|\mathbf{z})\big] - D_{KL}(q_\theta(\mathbf{z}|\boldsymbol{x_i})||p_\phi(\boldsymbol{z}))$$

$L(x_i, \theta, \phi)$

VAE maximizes variational lower bound ("ELBO")

Make approximate posterior distribution $q_\theta(\mathbf{z}|\boldsymbol{x_i})$ close to a tractable prior $p_\phi(\boldsymbol{z})$, e.g., Gaussian $\mathcal{N}(0, \mathbf{I})$

Encoder $E_\theta(X)$

$q_\theta(\boldsymbol{z}|\boldsymbol{x})$

$\boldsymbol{\mu_{z|x}}$

$\boldsymbol{\Sigma_{z|x}}$

$X$

Encoder estimates two quantities:
- $\boldsymbol{\mu_{z|x}}$: **mean** of $Z|X$
- $\boldsymbol{\Sigma_{z|x}}$: (diagonal)**covariance** of $Z|X$

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

Georgia Tech

$$\mathbf{E_z}\big[\log p_\phi(\boldsymbol{x_i}|\mathbf{z})\big] - D_{KL}(q_\theta\,(\mathbf{z}|\boldsymbol{x_i})||p_\phi(\boldsymbol{z}))$$

$L(x_i, \theta, \phi)$

Make approximate posterior distribution close to prior

$\boldsymbol{\mu}_{z|x}$

Encoder $E_\theta(\boldsymbol{X})$

$q_\theta\,(\boldsymbol{z}|\boldsymbol{x})$

$\boldsymbol{\Sigma}_{z|x}$

$\boldsymbol{Z}$

$\boldsymbol{X}$

Assume tractable prior $p_\phi(\boldsymbol{z})$, e.g., Gaussian.

We can **sample** $\boldsymbol{z}$ from $\boldsymbol{z}|\boldsymbol{x}{\sim}\mathcal{N}\big(\boldsymbol{\mu}_{z|x}, \boldsymbol{\Sigma}_{z|x}\big)$

(Z is not directly computed by VAE )

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

Georgia Tech

## Variational Lower Bound

$$\mathbf{E_z}\left[\log p_\phi(\boldsymbol{x_i}|\mathbf{z})\right] - D_{KL}(q_\theta\,(\mathbf{z}|\boldsymbol{x_i})||p_\phi(\mathbf{z}))$$



Encoder $E_\theta(\boldsymbol{X})$

$q_\theta\,(\boldsymbol{z}|\boldsymbol{x})$

$\boldsymbol{\mu_{z|x}}$

$\boldsymbol{\Sigma_{z|x}}$

Sample $z$
from $\boldsymbol{z}|\boldsymbol{x}\sim\mathcal{N}(\boldsymbol{\mu_{z|x}},\boldsymbol{\Sigma_{z|x}})$

$\boldsymbol{X}$

$\boldsymbol{Z}$

Decoder $G_\phi(\boldsymbol{Z})$

$p_\phi\,(\boldsymbol{x}|\boldsymbol{z})$

$\boldsymbol{\mu_{x|z}}$

$\boldsymbol{\Sigma_{x|z}}$

OLIVES
@GeorgiaTech

Georgia Tech

# Variational Autoencoders
## Variational Lower Bound

$$\mathbf{E_z}\big[\log p_\phi(\pmb{x_i}|\mathbf{z})\big] - D_{KL}\big(q_\theta\,(\mathbf{z}|\pmb{x_i})||p_\phi(\pmb{z})\big)$$

Maximize the likelihood of input being reconstructed (Minimize Binary Cross-entropy between $\pmb{x_i}$ and $\widehat{\pmb{x_i}}$)

$L(x_i, \theta, \phi)$

Make approximate posterior distribution close to prior

(Minimizing KL between $\mathcal{N}(0, \mathbf{I})$ and $\mathcal{N}(\pmb{\mu_{z|x}}, \pmb{\Sigma_{z|x}})$ )

Encoder $E_\theta(\pmb{X})$

$q_\theta\,(\pmb{z}|\pmb{x})$

$\pmb{\mu_{z|x}}$

$\pmb{\Sigma_{z|x}}$

$\pmb{Z}$

Sample $z$ from $\pmb{z}|\pmb{x} \sim \mathcal{N}(\pmb{\mu_{z|x}}, \pmb{\Sigma_{z|x}})$

Decoder $G_\phi(\pmb{Z})$

$p_\phi\,(\pmb{x}|\pmb{z})$

$\pmb{\mu_{x|z}}$

$\pmb{\Sigma_{x|z}}$

Sample $x$ from $\pmb{x}|\pmb{z} \sim \mathcal{N}(\pmb{\mu_{x|z}}, \pmb{\Sigma_{x|z}})$

$\pmb{X}$

$\widehat{\pmb{X}}$

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

Georgia Tech

Backpropagation is not feasible through random sampling

$$z|x \sim \mathcal{N}\left(\mu_{z|x}, \Sigma_{z|x}\right)$$

Cannot back propagate
through a random variable Z

Original form



◇ : Deterministic node

● : Random node

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

Georgia
Tech

Image from: NIPS Workshop 2015 (Kingma & Welling)

Backpropagation is not feasible through random sampling

$$z|x \sim \mathcal{N}\left(\boldsymbol{\mu}_{z|x}, \boldsymbol{\Sigma}_{z|x}\right)$$

Cannot back propagate through a random variable Z

$$z|x = \boldsymbol{\mu}_{z|x} + \boldsymbol{\Sigma}_{z|x} \odot \varepsilon_i$$
where $\varepsilon_i \sim \mathcal{N}(0, \mathbf{I})$

$\odot$: element-wise product

Can back propagate through a deterministic part $\mu_{z|x}$,

Original form

Reparameterised form



◇ : Deterministic node

● : Random node

[Kingma, 2013]
[Bengio, 2013]
[Kingma and Welling 2014]
[Rezende et al 2014]

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

Georgia Tech

Image from: NIPS Workshop 2015 (Kingma & Welling)

# Variational Autoencoders
## FC-VAE on MNIST



Input $\boldsymbol{X}$

P-dim

28x28 grayscale

Layer 1 @ $E_\theta$

400 (ReLU)

Layer $2_1$ @ $E_\theta$

***mean*** of $\boldsymbol{Z}$

Layer $2_2$ @ $E_\theta$

***variance*** of $\boldsymbol{Z}$

Latent $\boldsymbol{Z}$

20

Layer 3 @ $G_\phi$

400 (ReLU)

Output $\widehat{X}$

P-dim (Sigmoid)

28x28 grayscale

sample $\varepsilon \sim \mathcal{N}(0, \mathbf{I})$

- P = 784 (28x28x1) for MNIST
- ReLU activation in the intermediate layers
- Sigmoid activation in the output layer
  - Input $\boldsymbol{X}$ is normalized to [0,1]

OLIVES
@GeorgiaTech

Georgia Tech

Input MNIST digits

Reconstruction fc-VAE (**latent dim: 20**)

After training:

- Sample $Z$ from Gaussian prior

- Pass $Z$ through decoder network

- 20-dim latent dimension

These novel samples do not visually resemble training data.
**There is no encoder at inference!**
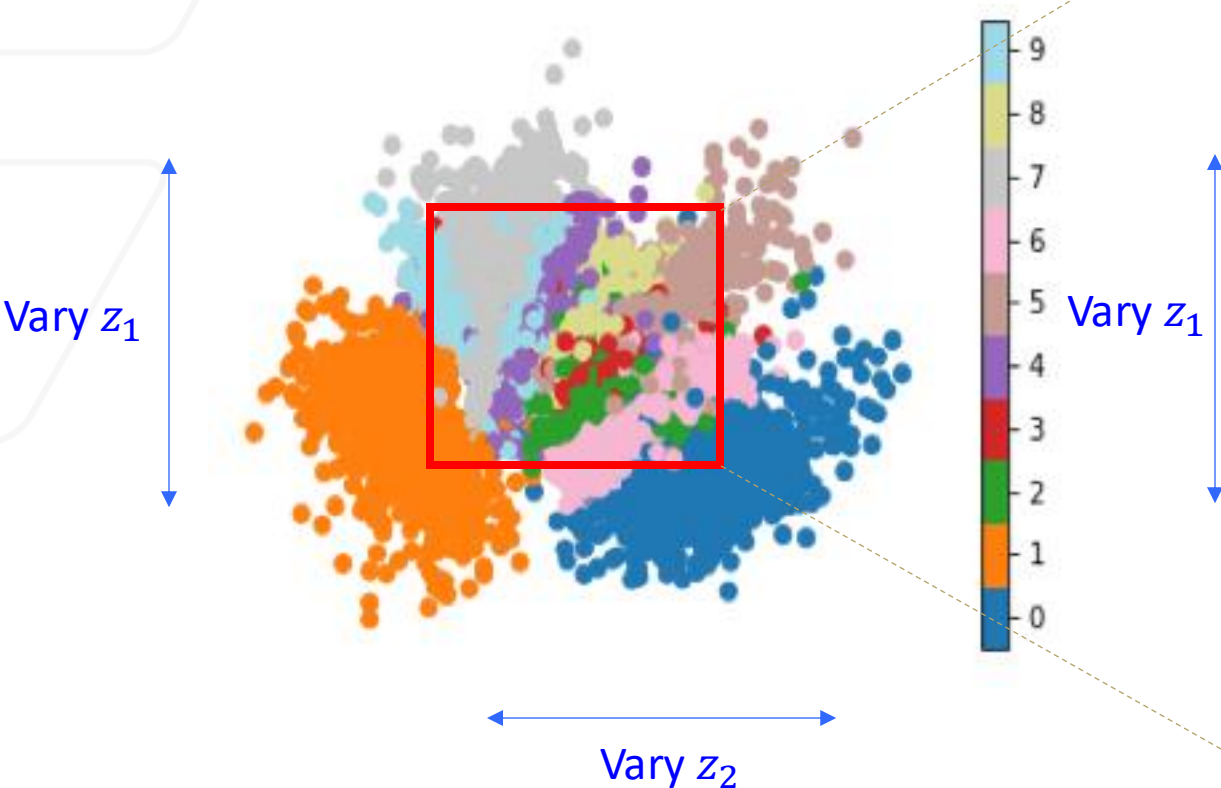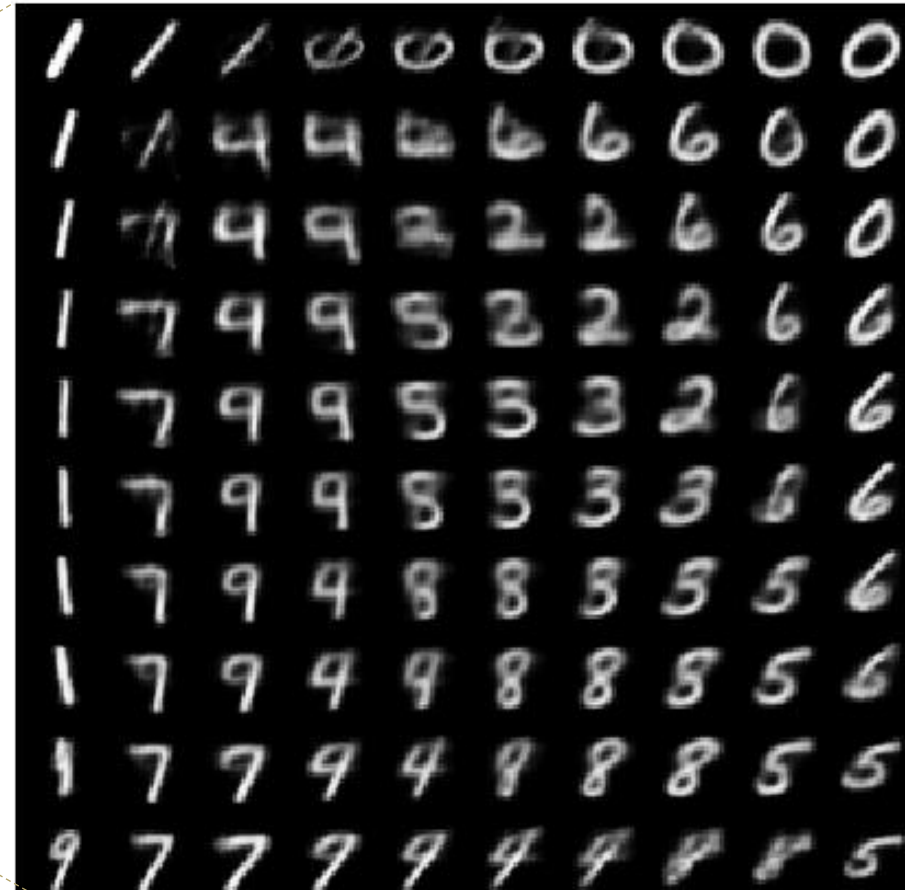
Randomly generated novel samples

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

# Variational Autoencoders
## FC-VAE on MNIST

MNIST Manifold (2D) visualization

MNIST Manifold (2D) visualization

Training a 2-dim FC VAE



Vary $z_1$

Vary $z_2$

Vary $z_1$

Vary $z_2$

Varying a single dimension ($z_1$ or $z_2$), it is possible to generate variation of digits

Varying $z_2$ (left to right): straight digits to digits with curves

OLIVES
@GeorgiaTech

Georgia Tech

- *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.

OLIVES
@GeorgiaTech

Georgia Tech

- $x_i$: a single feature
- $\boldsymbol{x}_i$: feature vector (a data sample)
- $\boldsymbol{x}_{:,i}$: feature vector of all data samples
- $X$: matrix of feature vectors (dataset)
- $\boldsymbol{W}$: weight matrix
- $\boldsymbol{Z}$: latent representation
- $E_\theta$: encoding function
- $G_\phi$: decoding function
- $\widehat{X}$ : reconstruction of data
- $\Omega(\boldsymbol{Z})$: sparsity constraint

- $\hat{\rho}_{\boldsymbol{j}}$: average activation of neuron $z_{ij}$
- $\widetilde{X}$ : corrupted input
- $N$: number of data samples
- $P$: number of features in a feature vector
- $P^{(k)}$: the number of neurons in layer $k$
- $\alpha$: learning rate
- Bold letter/symbol: vector
- Bold capital letters/symbol: matrix

[FunML L19: VAEs] | [Ghassan AlRegib and Mohit Prabhushankar] | [Oct 30, 2024]

OLIVES
@GeorgiaTech

Georgia Tech.

- Started going through the Progress reports and came across the following upload!



- Please send me the correct PDF!!