

ADVANCED ARTIFICIAL INTELLIGENCE: COMPUTER VISION

Generative models: latent spaces and variational autoencoders

Schedule for today

- Lecture: VAEs—variational autoencoders for image generation
- Lab: add batch normalization to CNN for CIFAR10
- Lab: finetune ResNet18 with (augmented) MNIST data
- Lab: generate MNIST-like images
- Lab: generate CIFAR-like images

Generative modeling

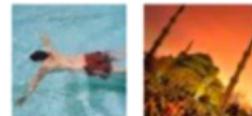
- take as input training samples from some distribution and learn a model that represents that distribution
- how to learn a distribution similar to that of the training data?

Sample Generation



Input samples

Training data $\sim P_{data}(x)$

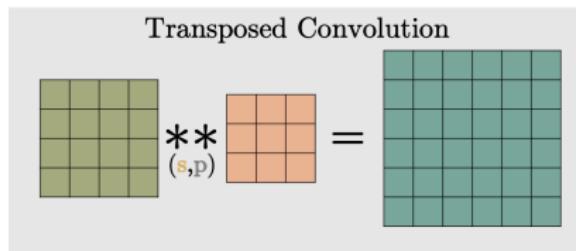


Generated samples

Generated $\sim P_{model}(x)$

Transposed convolution

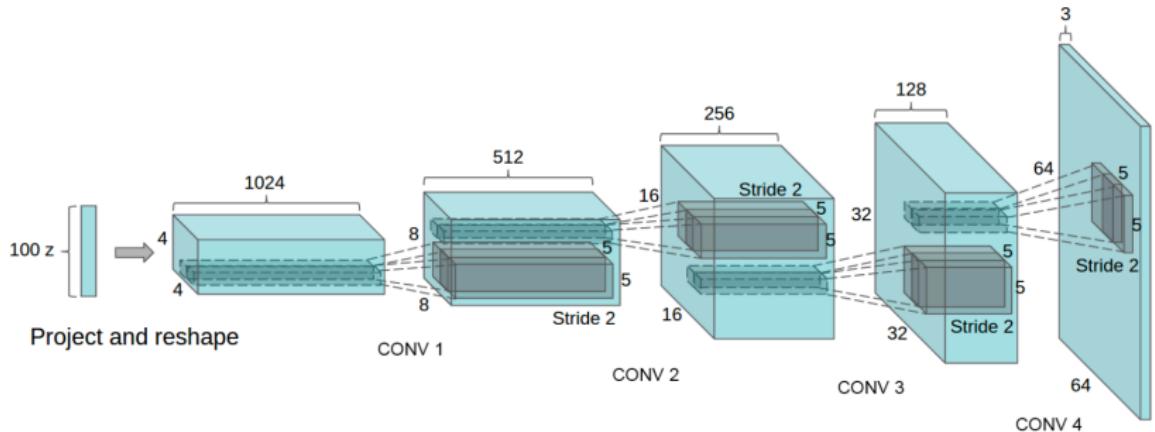
- operation that expands the spatial dimensions



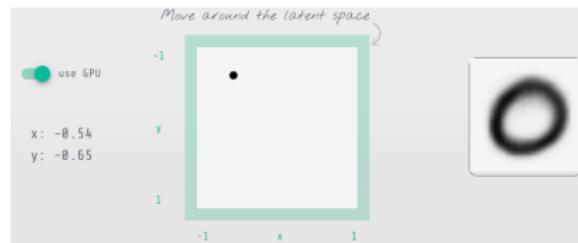
Input	Kernel	Output									
$\begin{matrix} 0 & 1 \\ 2 & 3 \end{matrix}$	$\begin{matrix} 0 & 1 \\ 2 & 3 \end{matrix}$	$=$	$\begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix}$	$+$	$\begin{matrix} 0 & 1 \\ 2 & 3 \end{matrix}$	$+$	$\begin{matrix} 0 & 2 \\ 4 & 6 \end{matrix}$	$+$	$\begin{matrix} 0 & 3 \\ 6 & 9 \end{matrix}$	$=$	$\begin{matrix} 0 & 0 & 1 \\ 0 & 4 & 6 \\ 4 & 12 & 9 \end{matrix}$

Transposed convolution

- used for trainable image generation

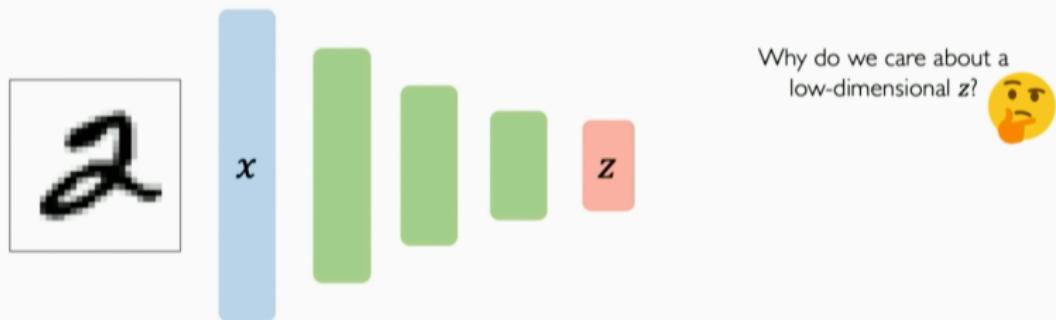


Demo



Autoencoders

Unsupervised approach for learning a **lower-dimensional** feature representation from unlabeled training data

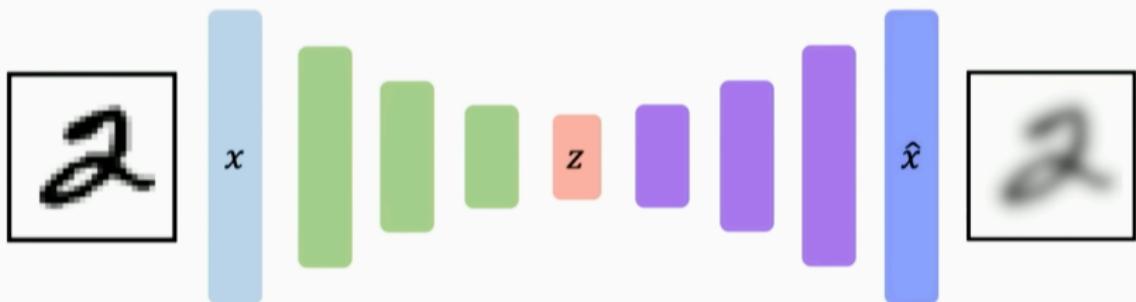


"Encoder" learns mapping from the data, x , to a low-dimensional latent space, z

Autoencoders

How can we learn this latent space?

Train the model to use these features to **reconstruct the original data**

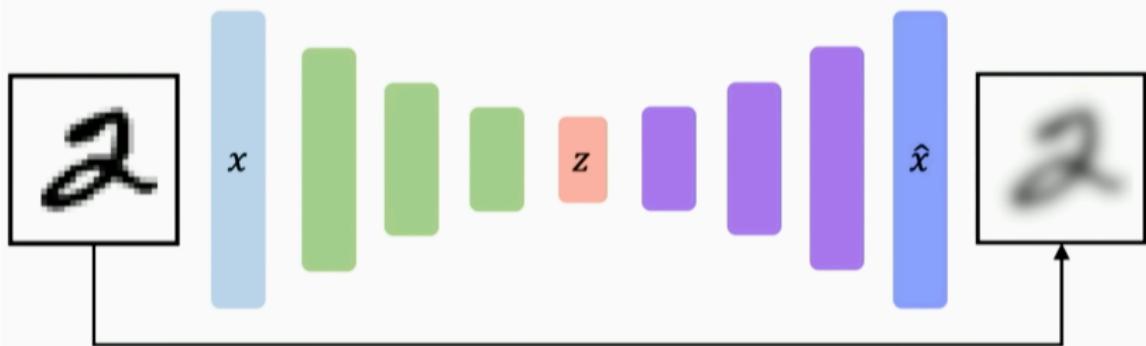


"Decoder" learns mapping back from latent space, z ,
to a reconstructed observation, \hat{x}

Autoencoders

How can we learn this latent space?

Train the model to use these features to **reconstruct the original data**



$$\mathcal{L}(x, \hat{x}) = \|x - \hat{x}\|^2$$

Dimensionality of latent space → reconstruction quality

Autoencoding is a form of compression!

Smaller latent space will force a larger training bottleneck

2D latent space

7	2	/	0	4	/	9	9	8	9
0	6	9	0	1	5	9	7	8	9
9	6	6	5	4	0	7	4	0	1
3	1	3	0	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	9	3	7	4	6	4	3	0
7	0	2	9	1	7	3	2	9	7
9	6	2	7	8	4	7	3	6	1
3	6	9	3	1	4	1	7	6	9

5D latent space

7	2	/	0	4	/	4	9	9	9
0	6	9	0	1	5	9	7	3	4
9	6	6	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	9	3	7	4	6	4	3	0
7	0	2	9	1	7	3	2	9	7
9	6	2	7	8	4	7	3	6	1
3	6	9	3	1	4	1	7	6	9

Ground Truth

7	2	/	0	4	/	4	9	9	5
0	6	9	0	1	5	9	7	8	4
9	6	6	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	9	3	7	4	6	4	3	0
7	0	2	9	1	7	3	2	9	7
9	6	2	7	8	4	7	3	6	1
3	6	9	3	1	4	1	7	6	9

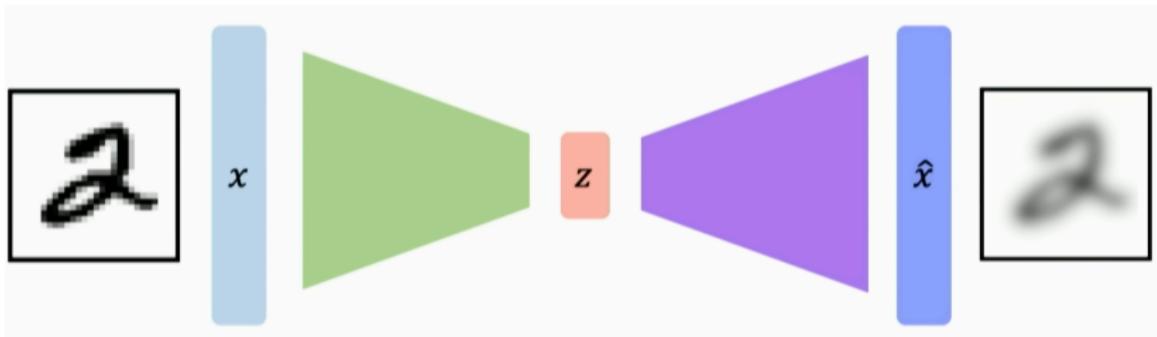
Autoencoders for representation learning

Bottleneck hidden layer forces network to learn a compressed latent representation

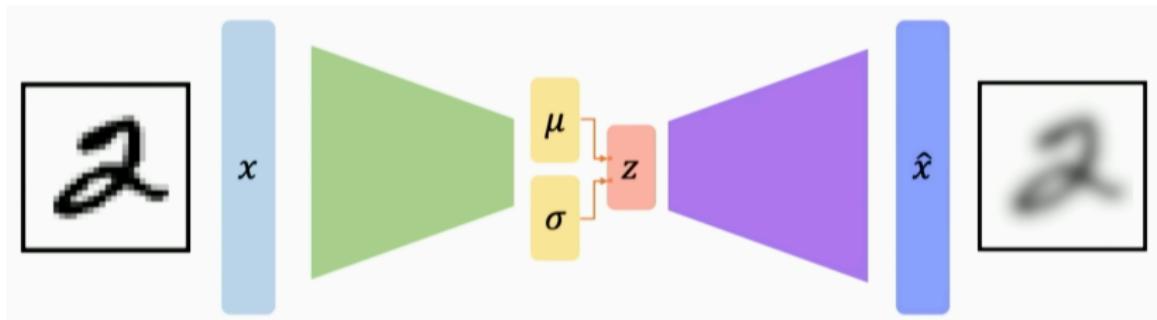
Reconstruction loss forces the latent representation to capture (or encode) as much "information" about the data as possible

Autoencoding = **A**utomatically **e**ncoding data; "Auto" = **s**elf-encoding

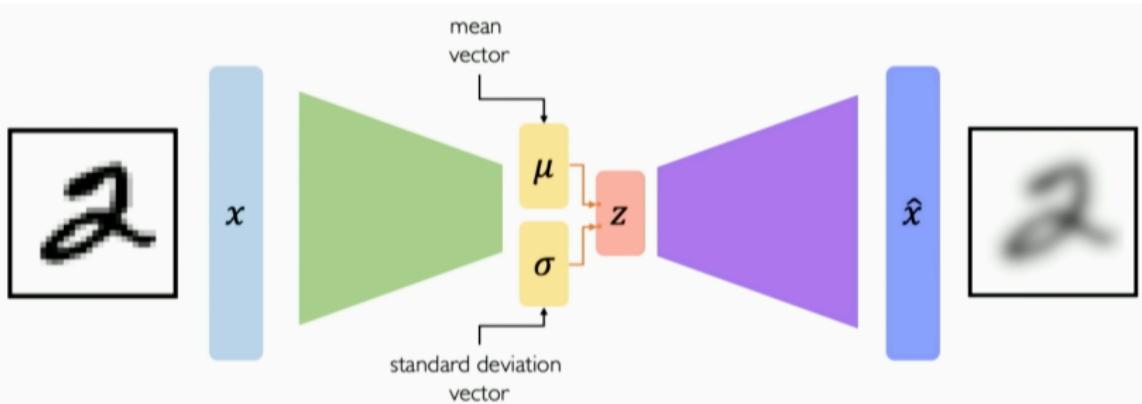
Traditional autoencoders



VAEs key difference from traditional autoencoders



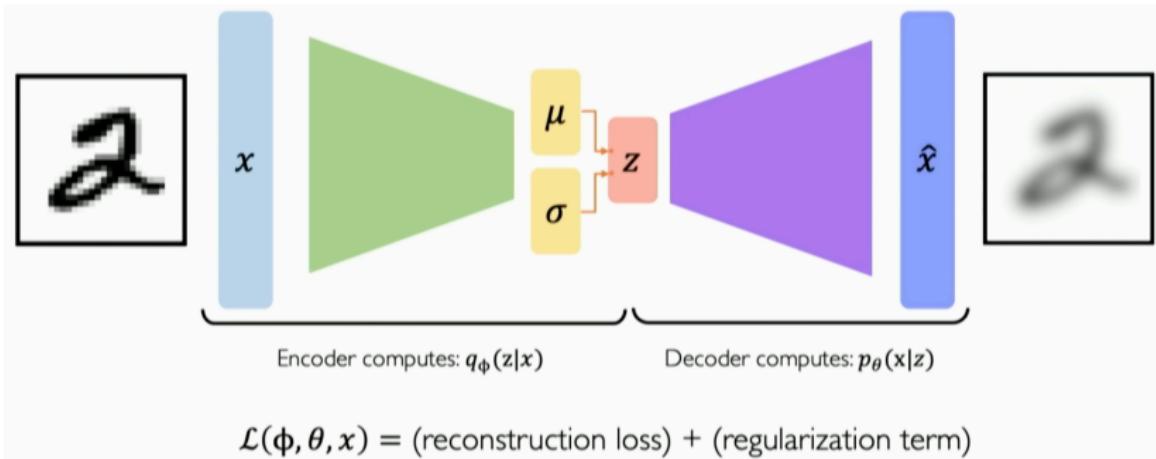
VAEs key difference from traditional autoencoders



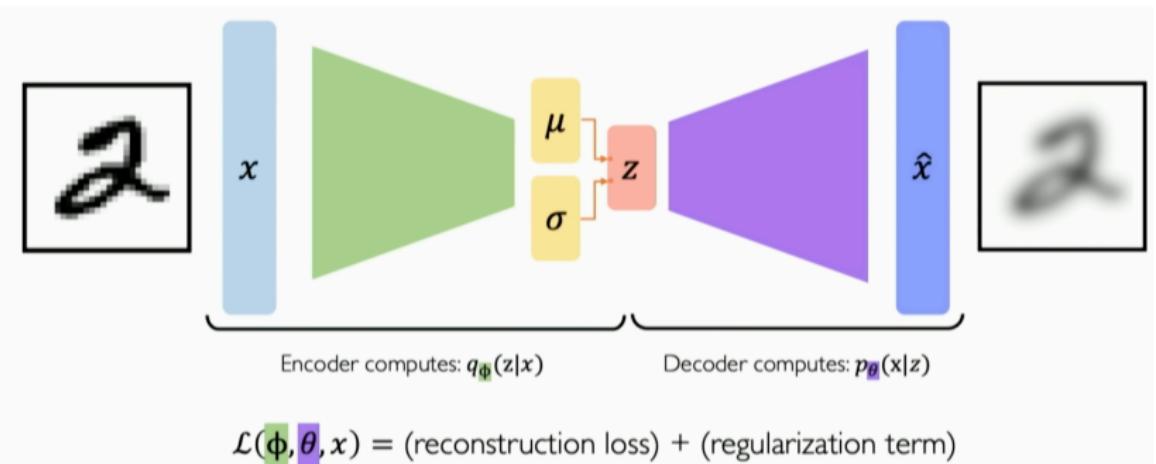
Variational autoencoders are a probabilistic twist on autoencoders!

Sample from the mean and standard deviation to compute latent sample

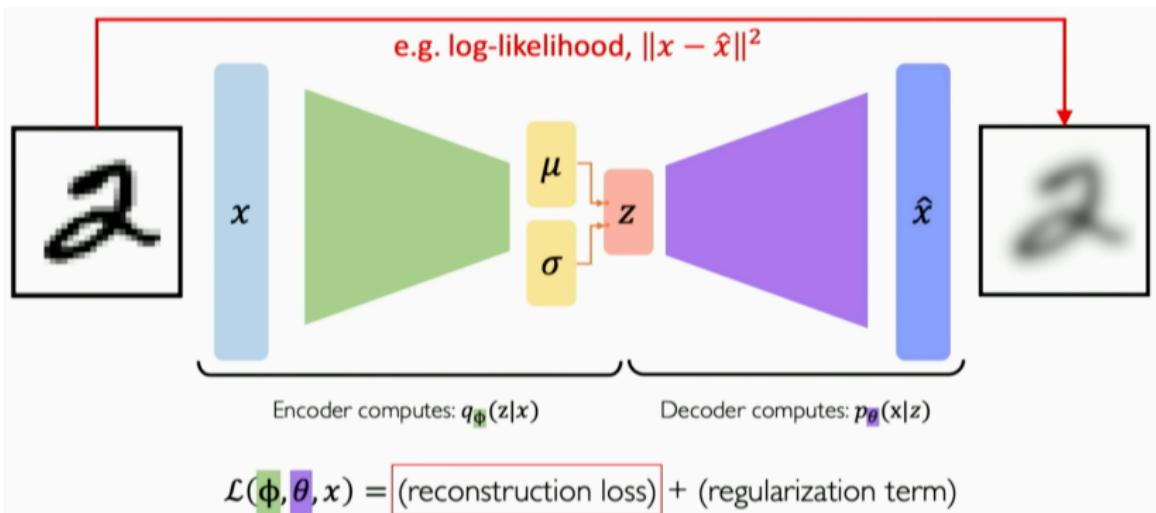
VAE optimization



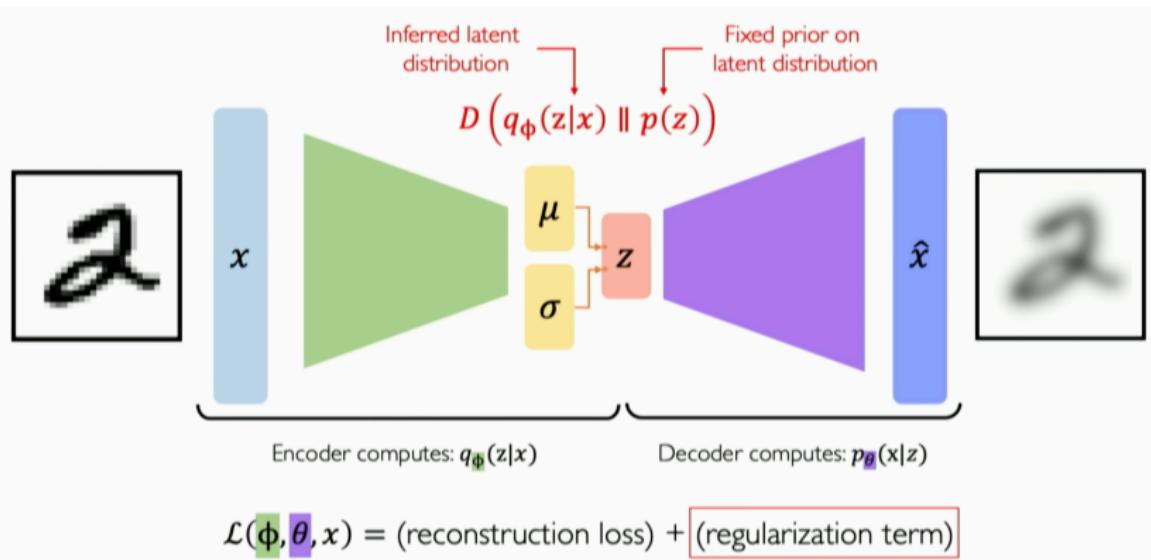
VAE optimization



VAE optimization



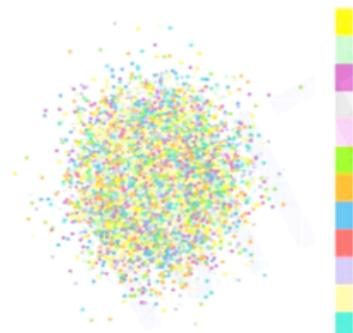
VAE optimization



Priors on the latent distribution

$$D \left(q_{\phi}(z|x) \parallel p(z) \right)$$

↑ ↑
Inferred latent Fixed prior on
distribution latent distribution



Common choice of prior – Normal Gaussian:

$$p(z) = \mathcal{N}(\mu = 0, \sigma^2 = 1)$$

- Encourages encodings to distribute encodings evenly around the center of the latent space
- Penalize the network when it tries to “cheat” by clustering points in specific regions (i.e., by memorizing the data)

Priors on the latent distribution

$$D\left(q_{\phi}(z|x) \parallel p(z)\right)$$

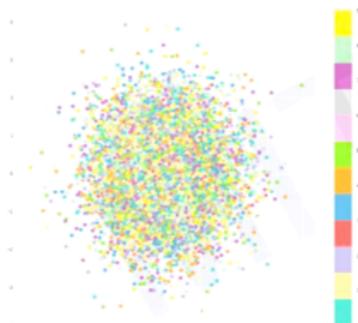
$$= -\frac{1}{2} \sum_{j=0}^{k-1} (\sigma_j + \mu_j^2 - 1 - \log \sigma_j)$$

KL-divergence
between the two
distributions

Common choice of prior – Normal Gaussian:

$$p(z) = \mathcal{N}(\mu = 0, \sigma^2 = 1)$$

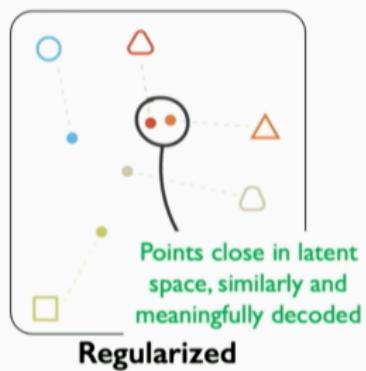
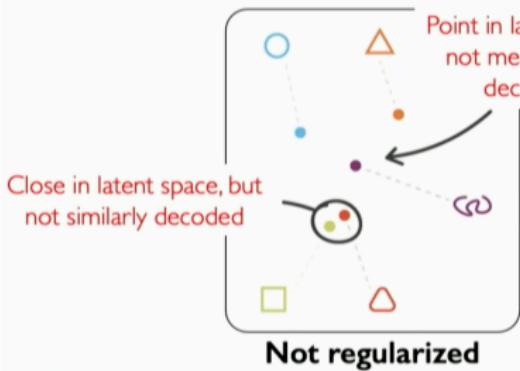
- Encourages encodings to distribute encodings evenly around the center of the latent space
- Penalize the network when it tries to "cheat" by clustering points in specific regions (i.e., by memorizing the data)



Intuition on regularization

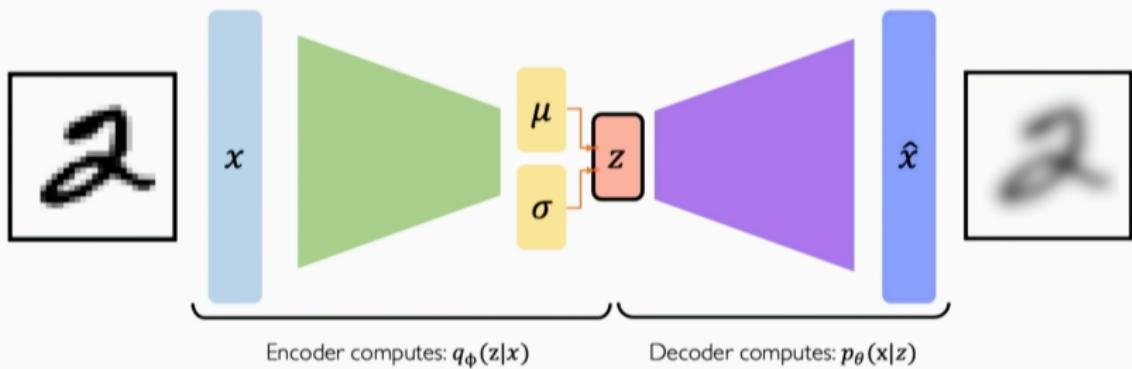
What properties do we want to achieve from regularization? 🤔

- Continuity:** points that are close in latent space → similar content after decoding
- Completeness:** sampling from latent space → “meaningful” content after decoding



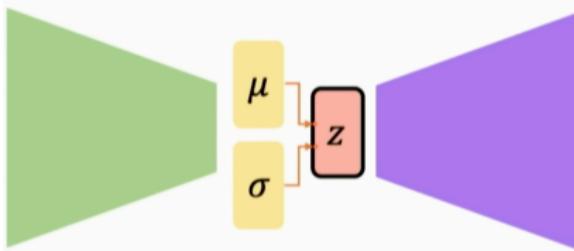
VAE computation graph

Problem: We cannot backpropagate gradients through sampling layers!



$$\mathcal{L}(\phi, \theta, x) = (\text{reconstruction loss}) + (\text{regularization term})$$

Reparameterizing the sampling layer



Key Idea:

$$z \sim \mathcal{N}(\mu, \sigma^2)$$

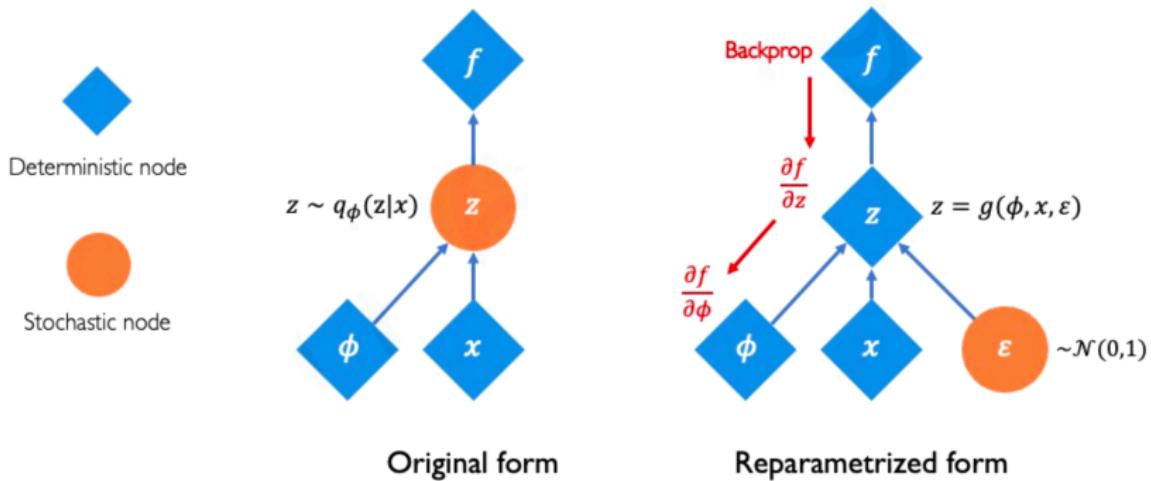
Consider the sampled latent vector z as a sum of

- a fixed μ vector;
- and fixed σ vector, scaled by random constants drawn from the prior distribution

$$\Rightarrow z = \mu + \sigma \odot \varepsilon$$

where $\varepsilon \sim \mathcal{N}(0,1)$

Reparameterizing the sampling layer



VAEs: latent perturbation

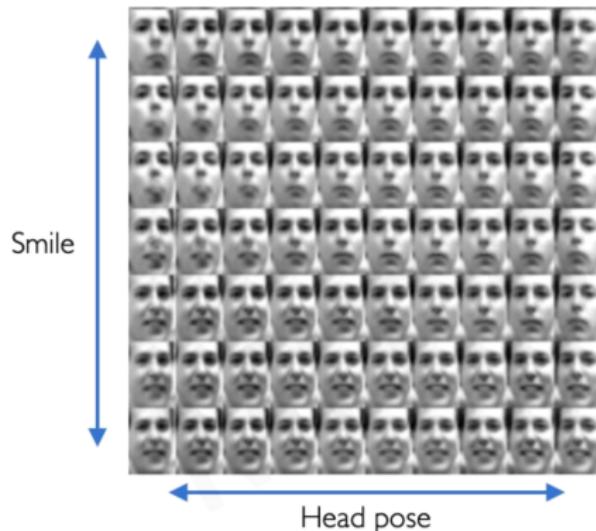
Slowly increase or decrease a **single latent variable**
Keep all other variables fixed



Head pose

Different dimensions of z encodes **different interpretable latent features**

VAEs: latent perturbation



Ideally, we want latent variables that are uncorrelated with each other

Enforce diagonal prior on the latent variables to encourage independence

Disentanglement

Increasing regularization

Head rotation (azimuth)



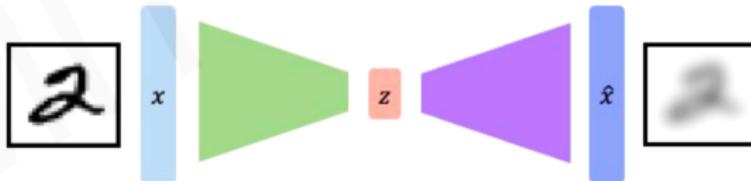
Standard VAE ($\beta = 1$)



β -VAE ($\beta = 250$)

VAE summary

1. Compress representation of world to something we can use to learn
2. Reconstruction allows for unsupervised learning (no labels!)
3. Reparameterization trick to train end-to-end
4. Interpret hidden latent variables using perturbation
5. Generating new examples



Latent variable models

- VAEs belong to a more general class of *latent variable models*

Autoencoders and Variational Autoencoders (VAEs)

Learn lower-dimensional **latent space** and **sample** to generate input reconstructions



Generative Adversarial Networks (GANs)

Competing **generator** and **discriminator** networks

