

Cornell Bowers CIS

## Discriminative Models

typically supervised

Goal: model  $p(Y|X)$   
from samples of  $p(X,Y)$   
(\* so that we can list  
most likely labels)

### Questions:

- Does one reduce to the other?
- Which is more difficult?

## Generative Models

typically unsupervised

Goal: model  $p(X)$   
from samples  
(\* so that we can  
draw new random samples)

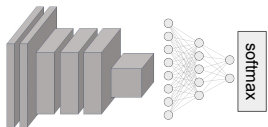
no labels Y!

Cornell Bowers CIS

## Discriminative Models

typically supervised

Goal: model  $p(Y|X)$   
from samples of  $p(X,Y)$   
(\* so that we can list  
most likely labels)



## Generative Models

typically unsupervised

Goal: model  $p(X)$   
from samples  
(\* so that we can  
draw new random samples)

no labels Y!

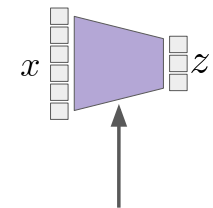
### Examples:

- GANs + variants
- Normalizing Flow Models
- Variational Autoencoders
  - Diffusion Models

Cornell Bowers CIS

## Dimensionality Reduction

Want to compress image  $x \in \mathbb{R}^D$   
to code  $z \in \mathbb{R}^d$



for the purposes of

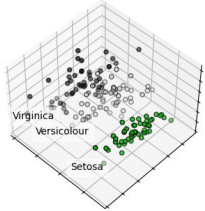
- visualization
- extracting important features  
(for downstream tasks)
- a more useful space, where  
geometry has semantic meaning

What properties should  
this mapping have?

## Principal Component Analysis (PCA)

- a linear transformation
- that capture as much variance as possible
- the components of  $z$  are independent

$$\begin{matrix} \boxed{X} & \boxed{W} & = & \boxed{Z} \\ n \times D & D \times d & & n \times d \end{matrix}$$

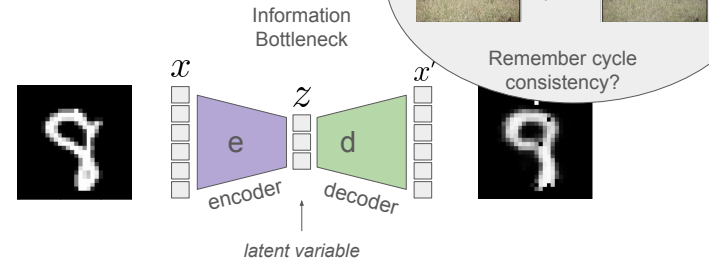


sklearn demo: Iris dataset (4 features, 3 classes)

Can be computed directly with linear algebra: take leading eigenvectors of  $X^T X$

$$\text{SVD: } \begin{matrix} \boxed{X} \\ n \times D \end{matrix} = \begin{matrix} \boxed{U} \\ n \times n \end{matrix} \begin{matrix} \boxed{\Sigma} \\ N \times D \end{matrix} \begin{matrix} \boxed{W^T} \\ D \times D \end{matrix}$$

## Autoencoders

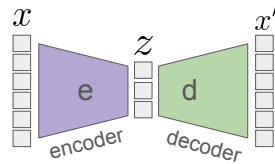


Question: What loss function should we use?

## Reconstruction Loss, first attempt

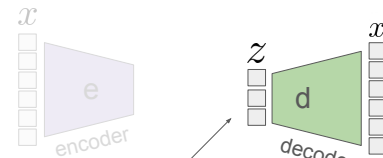
- "the obvious loss"

$$\sum_{x \in \mathcal{D}} (x - x')^2 \quad \text{where } x' = e(d(x))$$



The Result: an Autoencoder.  
[Kramer, 1991]

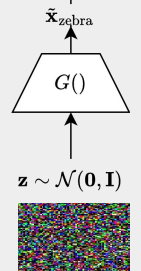
## Sampling from an Autoencoder



$$z \sim \mathcal{N}(0, I)$$

feed decoder  
(Gaussian) noise?

Recall how we could sample with GANs...



## Autoencoder trained on MNIST: latent space

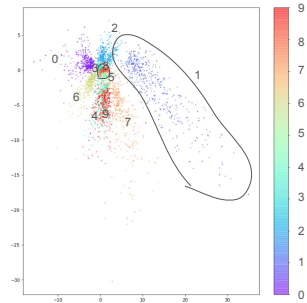


Figure 3-8. Plot of the latent space, colored by digit

Not a very nice representation...

- no symmetries between digit representations
- lots of empty space

### Question:

What does this mean for sampling?

## What's needed is some kind of "regularization"

to "encourage" the encoder to have "nice properties"...

- Contractive Autoencoders [2011]
- Sparse Autoencoder [2013]
- Variational Autoencoders [2014]

## A Probabilistic Perspective

Building Blocks:

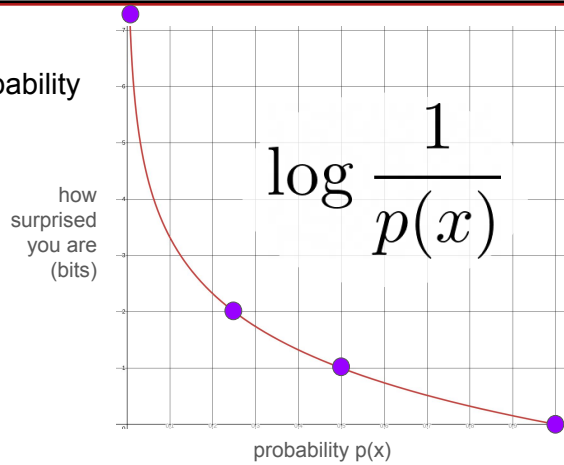
- Conditional and marginal probabilities
- Surprisal / Negative Log Likelihood
- Relative Entropy / KL Divergence

## Conditional and Marginal Probabilities

$$p(X, Y) = p(Y|X)p(X)$$

$$p(X) = \int p(X, y) dy$$

## Negative Log Probability = "Surprisal"



## KL Divergence (a.k.a. relative entropy)

$$D(p \parallel q) := \mathbb{E}_{x \sim p} \left[ \log \frac{p(x)}{q(x)} \right]$$

↑ reality (e.g., dataset)
↑ model

$$= \mathbb{E}_{x \sim p} \left[ \log \frac{1}{q(x)} \right] - \log \frac{1}{p(x)}$$

Cross Entropy!
(constant; does not depend on model q)

- non-negative  $D(p \parallel q) \geq 0$
- zero means same  $D(p \parallel q) = 0 \iff p = q$
- not symmetric
- has many other, uniquely nice properties ...

## KL Divergence

Justin's Coin



Varsha's Coin



### Question:

Is it just as easy to mistake the output of Justin's coin for that of Varsha's coin, as vice versa?

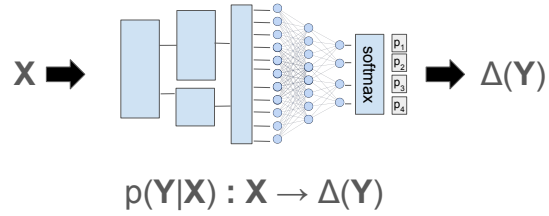
[Link to visualization](#)

### Building Blocks:

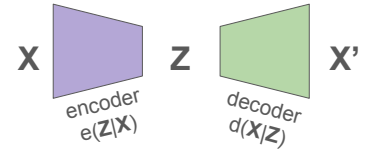
- ✓ Conditional and marginal probabilities
- ✓ Surprisal / Negative Log Likelihood
- ✓ Relative Entropy / KL Divergence

## Neural Networks as Conditional Probabilities

A network with a softmax encodes a conditional probability distribution



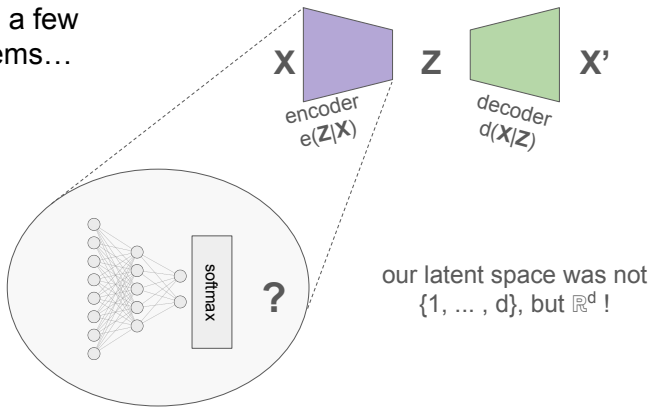
## Reconstruction Loss, using surprisal



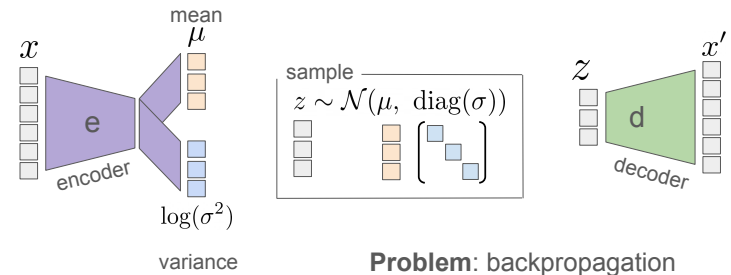
How surprising would it be to encode  $x$ , decode the result, and recover  $x$ ?

$$\mathbb{E}_{z \sim e(x)} \left[ \log \frac{1}{d(x|z)} \right]$$

## Fixing a few problems...



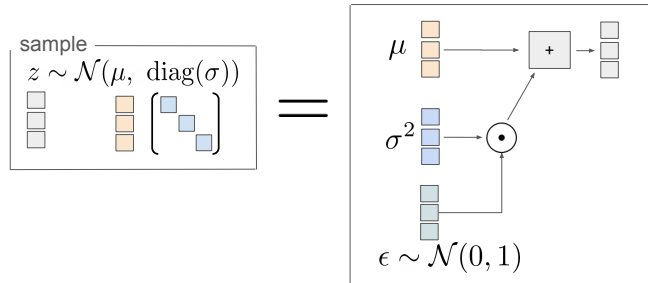
## An Architecture for Gaussians



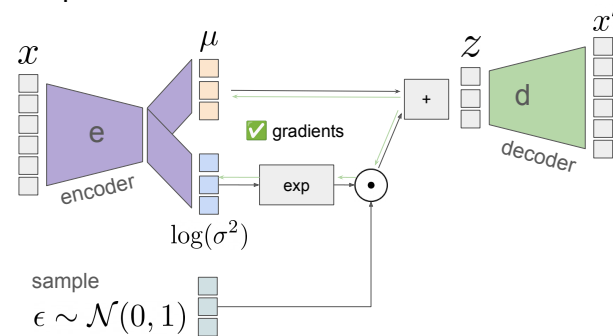
**Problem:** backpropagation through sampling process?

### The Reparameterization Trick

$$\mathcal{N}(\mu, \text{diag}(\sigma)) = \mu + \sigma^2 \odot \mathcal{N}(0, I)$$

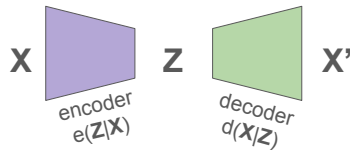


### The Reparameterization Trick



### Reconstruction Loss, using surprisal

How surprising would it be to encode  $x$ , decode it, and recover the same sample?



$$\mathbb{E}_{z \sim e(x)} \left[ \log \frac{1}{d(x|z)} \right] \quad \text{Essentially MSE, again!}$$

$$\propto \exp \left( -\frac{1}{2} (x - f(z))^2 \right)$$

### We're back at an autoencoder, but probabilistic

The upshot: we can now add a regularization term

$$D(e(Z|x) \parallel p(Z))$$

Want each encoding ... to match a prior (e.g., a standard Gaussian)

#### Questions:

Does this have a connection to PCA?  
Is there a conceptual problem with this regularization?

# Variational Inference

## Motivating VAEs

- Have joint model  $p(X, Z)$
- observe  $x$  (but not  $z$ );
- want to calculate posterior  $p(Z|x) = \frac{p(x, Z)}{p(x)}$ ,
- which requires  $p(x) = \int p(x, z) dz$ 
  - i.e., the "evidence".
  - But the integral is often intractable!
- so, instead ...

Optimize!

$$\min_{\phi} D(q_{\phi}(Z) \parallel p(Z|x))$$

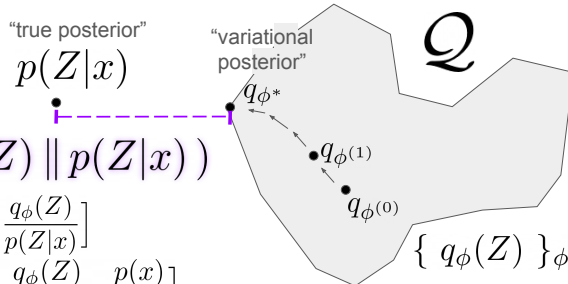
$$= \mathbb{E}_{Z \sim q_{\phi}} \left[ \log \frac{q_{\phi}(Z)}{p(Z|x)} \right]$$

$$= \mathbb{E}_{Z \sim q_{\phi}} \left[ \log \frac{q_{\phi}(Z)}{p(Z|x)} \cdot \frac{p(x)}{p(x)} \right]$$

$$= \mathbb{E}_{Z \sim q_{\phi}} \left[ \log \frac{q_{\phi}(Z)}{p(x, Z)} \right] + \log p(x)$$

(constant; does not depend on  $\phi$ )

Tractable;  $-\text{ELBO}(x, \phi)$



## Variational Bound

$$D(q_{\phi}(Z) \parallel p(Z|x)) = \underbrace{\mathbb{E}_{Z \sim q_{\phi}} \left[ \log \frac{q_{\phi}(Z)}{p(x, Z)} \right]}_{-\text{ELBO}(x)} + \log p(x)$$

$$\text{ELBO}_{p, \phi}(x) + \underbrace{D(q_{\phi}(Z) \parallel p(Z|x))}_{\text{non-negative}} = \log p(x)$$

$$\text{ELBO}_{p, \phi}(x) \leq \log p(x)$$

“(log) evidence”

### What does this have to do with autoencoders?

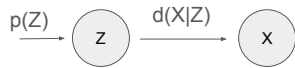
Replace:

$$q(Z) \rightsquigarrow e(Z|X)$$

$$p(X, Z) \rightsquigarrow p(Z)d(Z|X)$$

↑  
prior / regularizer

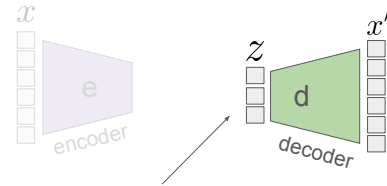
$$-\text{ELBO}_{p,d,e}(x) = \mathbb{E}_{Z \sim e_\phi(x)} \left[ \log \frac{e_\phi(Z|x)}{p(Z)d(x|Z)} \right]$$



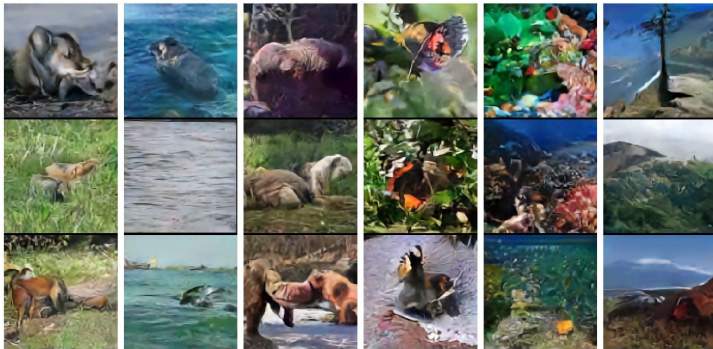
**Questions:**

1. What is the evidence, in this case?
2. Is there something strange about it?

### Sampling from a VAE

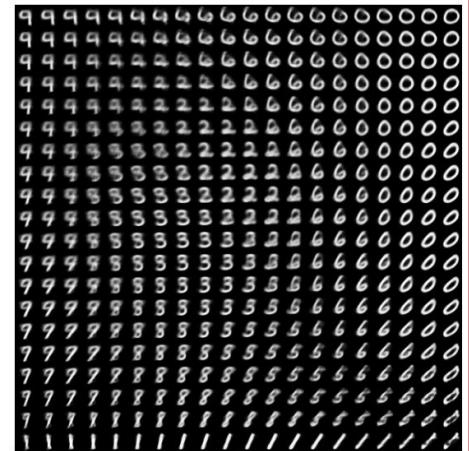


$$z \sim p_{\text{prior}} = \mathcal{N}(0, 1)$$



### a much nicer space...

can smoothly interpolate digits in a meaningful, digit-y kind of way





Cornell Bowers CIS

## a much nicer space

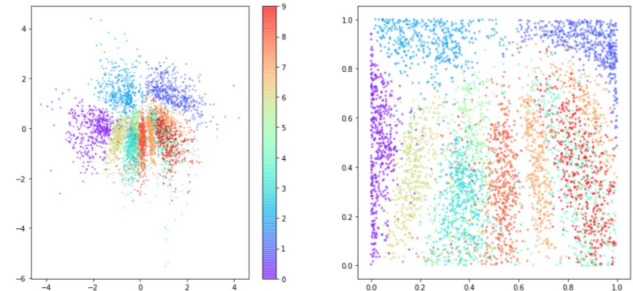
dimensions in latent space correspond to meaningful concepts, like sentiment and orientation



Cornell Bowers CIS

## Back to MNIST: Visualizing latent space again

VAE Latent space, note the distribution is centered, and each digit has an equal portion



Cornell Bowers CIS

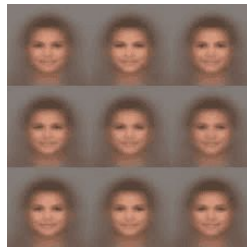
## The Biggest Drawback of VAEs

- Out of the box, generated images can be blurry.

Question: Why?



[VAE v. GAN](#)



<https://borisburkov.net/2022-12-31-1/>

Cornell Bowers CIS

## Hierarchical VAEs

The generative process is modeled as a Markov chain, where each latent  $z_i$  is generated only from the previous latent  $z_{i-1}$

