# Some Reminders: Spring Break !(?)

- What is Spring Break?
  - Undergrads have fun (take a trip, visit home, etc.)
  - $\circ$  Ph.D. students have fun (stop worrying about classes/grades and work on research  $\odot$  )
  - Faculty have fun (hide from the world, writing retreat, submit grant proposals)



JORGE CHAM OTHE STANFORD DAILY





# Some Reminders: Spring Break !(?)

- What is Spring Break?
  - Undergrads have fun (take a trip, visit home, etc.)
  - $\circ$  Ph.D. students have fun (stop worrying about classes/grades and work on research  $\odot$  )
  - Faculty have fun (hide from the world, writing retreat, submit grant proposals)



# Some Reminders: Spring Break !(?)

- What is Spring Break?
  - Undergrads have fun (take a trip, visit home, etc.)

  - $\circ$  Ph.D. students have fun (stop worrying about classes/grades and work on research  $\odot$  ) • Faculty have fun (hide from the world, writing retreat, submit grant proposals)
- Midterm Exam is on March 31
  - In this room (ITE 231)
  - During class hours (4:00 PM to 5:00 PM)
  - $\circ$  Syllabus: every lecture until and including 03/26





tejasgokhale.com

# Lecture 7: Representation Learning





CMSC 475/675 Neural Networks



Some VAE slides are borrowed from Ranjay Krishna





## Machine Learning Problems



## **SUPERVISED LEARNING**



 $\{\mathbf{x}^{(t)}, y^{(t)}\}$ 

 $\mathbf{x}^{(t)}, y^{(t)} \leftarrow p(\mathbf{x}, y)$ 

### Example

Input:  $x^{(t)}$  is an image

*Output*:  $y^{(t)}$  is an image category



## MULTITASK LEARNING



 $\{\mathbf{X}^{(t)}, y_{1}^{(t)}, \ldots, y_{\pi\pi}^{(t)}\}$ 

 $\mathbf{x}^{(t)}, y_1^{(t)}, \ldots, y_M^{(t)} \longleftarrow$  $p(\mathbf{x}, y_1, \ldots, y_M)$ 

### Example

Input:  $x^{(t)}$  is an image *Outputs*:

- $y_1^{(t)}$ : image category
- $y_1^{(t)}$ : object detection
- $y_1^{(t)}$ : depth estimation
- $y_1^{(t)}$ : semantic segmentation



# MULTITASK LEARNING

### **Topics:** multitask learning





# DOMAIN ADAPTATION



$$\{\mathbf{x}^{(t)}, y^{(t)}\}$$

$$\overline{\mathbf{x}}^{(t)} \leftarrow q(\mathbf{x})$$
$$y^{(t)} \leftarrow p(y|\overline{\mathbf{x}}^{(t)})$$

### • Example

- classify sentiment (positive vs negative)
- in reviews of different products
- training on Amazon Reviews but testing on Yelp Reviews





# **ONE-SHOT LEARNING**



$$\{\mathbf{x}^{(t)}, y^{(t)}\}$$

$$x^{(t)}, y^{(t)} \leftarrow p(\mathbf{x}, y)$$

• 
$$y^{(t)} 2 \{C + 1, \dots, C + M\}$$

### side information :

- a single labeled example from each of the M new classes

### • Example

 recognizing a person based on a single picture of him/her



# ZERO-SHOT LEARNING



$$\{\mathbf{x}^{(t)}, y^{(t)}\}$$

$$x^{(t)}, y^{(t)} \leftarrow p(\mathbf{x}, y)$$

• 
$$y^{(t)} 2 \{C + 1, \dots, C + M\}$$

### side information :

- description vector  $Z_c$  of each of the new M classes

### • Example

 recognizing an object based on a worded description of it





## **SEMI-SUPERVISED LEARNING**



 $\{\mathbf{x}^{(t)}, y^{(t)}\}$ 

 $\mathbf{x}^{(t)}, y^{(t)} \leftarrow p(\mathbf{x}, y)$ 

# UNSUPERVISED LEARNING



 $\{\mathbf{x}^{(t)}\}$  $\mathbf{x}^{(t)} \leftarrow p(\mathbf{x})$ 

## How can we train models "unsupervised"?

## How can we train models "unsupervised"?

### This is the focus of representation learning

### This is the focus of representation learning



### How can we train models "unsupervised"?

tions		
		Eng
cation	<u>h5-index</u>	<u>h5-r</u>
e	<u>488</u>	-
CVF Conference on Computer Vision and Pattern Recognition	<u>440</u>	(
lew England Journal of Medicine	<u>434</u>	8
ce	<u>409</u>	(
e Communications	<u>375</u>	
ancet	<u>368</u>	(
al Information Processing Systems	<u>337</u>	(
nced Materials	<u>327</u>	4
	<u>320</u>	4
ational Conference on Learning Representations	<u>304</u>	



### How

### International Conference

h5-index:304 h5-median:584

#2 Artificial Intelligence #4 Engineering & Computer Sci

Title / Author

An Image is Worth 16x A Dosovitskiy, L Beyer, A Kol ICLR

Decoupled Weight Dec I Loshchilov, F Hutter ICLR (Poster)

Measuring and Improving Y Hou, J Zhang, J Cheng, K ICLR

ALBERT: A Lite BERT f Z Lan, M Chen, S Goodman ICLR

Large Scale GAN Train A Brock, J Donahue, K Simo ICLR

DARTS: Differentiable A H Liu, K Simonyan, Y Yang ICLR (Poster)

LoRA: Low-Rank Adapt EJ Hu, Y Shen, P Wallis, Z A ICLR

Deformable DETR: Def X Zhu, W Su, L Lu, B Li, X W ICLR

BERTScore: Evaluating T Zhang, V Kishore, F Wu, K ICLR

ELECTRA: Pre-training K Clark, MT Luong, QV Le, C ICLR

# This is th

### There's an entire



ICLR has become one of the top CS and Engineering (not just AI) publication although it just started in 2013.

In fact top-10 in ALL OF SCIENCE

ence on Learning Representations	9
<u>ience</u>	
Cited by	
16 Words: Transformers for Image Recognition at Scale. Iesnikov, D Weissenborn, X Zhai,	
ay Regularization.	
ng the Use of Graph Information in Graph Neural Networks. Ma, RTB Ma, H Chen, MC Yang	
for Self-supervised Learning of Language Representations. , K Gimpel, P Sharma, R Soricut	E
ning for High Fidelity Natural Image Synthesis. Imyan	<u>ndex h</u> 8
Architecture Search.	<u>10</u> 34
tation of Large Language Models. Allen-Zhu, Y Li, S Wang, L Wang, W Chen	) <u>9</u> 7 <u>5</u>
formable Transformers for End-to-End Object Detection. Vang, J Dai	<u>88</u> 87
g Text Generation with BERT. (Q Weinberger, Y Artzi	2 <u>7</u> 20
Text Encoders as Discriminators Rather Than Generators. CD Manning	)4



## Warning

I might use the terms "latent", "embedding", "representation", "feature" interchangeably.

# Motivation (kind of): Compression

- The idea is similar to compression (signal processing) or hashing (data structures):
  - o encode an image into a smaller vector s.t. you can decode it back to its original form
    - Example: images, audio, video are stored in a compressed form on your computer using compression algorithms like JPEG, MP3, MPEG etc. The computer has software to decode it back so that you can view it (everytime you "open" a JPEG file to view an image, the decoder runs and converts code to RGB)





# Motivation (kind of): Compression

- The idea is similar to compression (signal processing) or hashing (information theory):
  - o encode an image into a smaller vector s.t. you can decode it back to its original form
    - Example: images, audio, video are stored in a compressed form on your computer using compression algorithms like JPEG, MP3, MPEG etc. The computer has software to decode it back so that you can view it (everytime you "open" a JPEG file to view an image, the decoder runs and converts code to RGB)
- Representation Learning
  - $\sim$  convert inputs automatically into "codes" (called representations/embeddings / features) s.t. the representations are:
    - Useful for downstream tasks (e.g. classification, regression, ...)
    - o "explain the data" and are "meaningful"
- Main difference: "meaningful" representation spaces to do "tasks" (The goal for compression is only efficient storage — not data classification/clustering etc.)



# **Representation Learning Paradigm**

### Raw Data (e.g. text, image, audio, video, ...)

Representation Learning (Encoding)

Do tasks / actions with these representations



Similar representations for similar concepts









A Semblance of "Context" should be encoded ...





### If you know the answer, don't share it with the class yet.

People from lands between Greece and India might know the answer ...

A Semblance of "Context" should be encoded ...

# l am very <u>hungry</u>, I will <u>eat</u> HALWA

A Semblance of "Context" should be encoded ...

## l am very <u>hungry</u>, I will <u>eat</u> HALWA

If you speak Marathi, this word has two meanings depending on context

Halwa (1): a food item Halwa (2): (an instruction to) move (something)



derived from: Farsi derived from: Sanskrit

A Semblance of "Context" should be encoded ...

## Is it a good idea to <u>eat HALWA after a meal</u>?

A Semblance of "Context" should be encoded ...



A Semblance of "Context" should be encoded ...

# sugary dessert **Oh no! I forgot to put sugar in the HALWA**



Parts, properties, attributes, ontology ?

- "bird" "wing", "beak", "feathers" has
- "bird" "fly" can
- "bird" "animal" is under category
- "bird" has subcategories

"eagle", "peacock", "sparrow", "seagull", "pigeon"

## **Representation Learning is a Philosophy for Learning**

### Key assumptions in this philosophy:

• You can convert a high-dimensional input space into a low-dimensional representation space

 $\circ$  Example: RGB images  $\rightarrow$  100 dim vectors

- A good representation space will have a "structure"
  - Example: Similarity, Symmetry, Relations will be easy to understand
  - Why? So that we can do arithmetic in representation space to do tasks
- Representations can be learned from data
- Representations can be leveraged for doing tasks

### Parallel Work in Cog.Sci.

### Trends in Cognitive Sciences 🥏



Volume 28, Issue 9, September 2024, Pages 844-856

Review

### Why concepts are (probably) vectors

Steven T. Piantadosi <sup>12</sup>  $\stackrel{>}{\sim}$   $\stackrel{\boxtimes}{\sim}$  , Dyana C.Y. Muller <sup>2</sup>, Joshua S. Rule <sup>1</sup>, Karthikeya Kaushik<sup>1</sup>, Mark Gorenstein<sup>2</sup>, Elena R. Leib<sup>1</sup>, Emily Sanford<sup>1</sup>

For decades, cognitive scientists have debated what kind of representation might characterize human concepts. Whatever the format of the representation, it must allow for the computation of varied properties, including similarities, features, categories, definitions, and relations. It must also support the development of theories, ad hoc categories, and knowledge of procedures. Here, we discuss why vectorbased representations provide a compelling account that can meet all these needs while being plausibly encoded into neural architectures. This view has become especially promising with recent advances in both large language models and vector symbolic architectures. These innovations show how vectors can handle many properties traditionally thought to be out of reach for neural models, including compositionality, definitions, structures, and symbolic computational processes.

### Highlights

Modern language models and vectorsymbolic architectures show that vector-based models are capable of handling the compositional, structured, and symbolic properties required for human concepts.

Vectors are also able to handle key phenomena from the psychology, including computation of features and similarities, reasoning about relations and analogies, and representation of theories.

Language models show how vector representation of word semantics and sentences can interface between concepts and language, as seen in definitional theories of concepts or ad hoc concepts.

The idea of Church encoding, from logic, allows us to understand how meaning can arise in vector-based or symbolic systems.

By combining these recent computational results with classic findings in psychology, vector-based models provide a compelling account of human conceptual representation.







## **Representation Learning is a Philosophy for Learning**

### Key assumptions in this philosophy:

• You can convert a high-dimensional input space into a low-dimensional representation space

 $\circ$  Example: RGB images  $\rightarrow$  100 dim vectors

- A good representation space will have a "structure"
  - Example: Similarity, Symmetry, Relations will be easy to understand
  - Why? So that we can do arithmetic in representation space to do tasks
- Representations can be learned from data
- Representations can be leveraged for doing tasks



## Ok whatever. Tell us how it works ...

## **Types of Modeling (Probabilistic Interpretation)**

the state



(probabilities of all inputs sum to 1)

### **Discriminative Model**

Learn Prob. Dist. P(y|x)



Tiger Horse Cat  $\forall x, \sum_{c} P(y = c | x) = 1$ 

## **Types of Modeling (Probabilistic Interpretation)**



(probabilities of all inputs sum to 1)

### **Generative Model**

Learn Marginal Prob. Dist. P(x)



### **Conditional Generative Model** Learn conditional probability P(x|y)

Discriminative Model (Unconditional) **Generative Model**  $P(v \mid$  $P(x \mid$ Conditional **Prior over labels Generative Model** 

## Types of Modeling (Probabilistic Interpretation)



(probabilities of all inputs sum to 1)

• **Discriminative Model** 

Learn Prob. Dist. P(y|x)

Generative Model

Learn Marginal Prob. Dist. P(x)

<u>Conditional Generative Model</u>

Learn conditional probability P(x|y)

## Types of Modeling (Probabilistic Interpretation) APPLICATIONS

Classification, Regression, Representation Learning (with labels)



• **Discriminative Model** 

Learn Prob. Dist. P(y|x)

Generative Model

Learn Marginal Prob. Dist. P(x)

• <u>Conditional Generative Model</u>

Learn conditional probability P(x|y)
### Types of Modeling (Probabilistic Interpretation) APPLICATIONS

Data Generation Outlier Detection Representation Learning (without labels)



• **Discriminative Model** 

Learn Prob. Dist. P(y|x)

Generative Model

Learn Marginal Prob. Dist. P(x)

<u>Conditional Generative Model</u>

Learn conditional probability P(x|y)

### Types of Modeling (Probabilistic Interpretation) APPLICATIONS

Machine Translation Text-to-image generation

(pretty much every "GenAl" product you see is a conditional generative model) • **Discriminative Model** 

Learn Prob. Dist. P(y|x)

Generative Model

Learn Marginal Prob. Dist. P(x)

<u>Conditional Generative Model</u>

Learn conditional probability P(x|y)

# **Generative Models**

- What's a Generative Model?
  - $\circ$  A model for the probability distribution of data x
  - A model that can be used to "generate" data



 Generative Models can be *learned* • You are given some observed data X • You choose a function (e.g. neural network) to model  $P(x; \theta)$  using parameters  $\theta$ • You estimate  $\theta$  s.t.  $P(x; \theta)$  best fits the observations X

#### (e.g. face images)



### **Generative Models**

- Generative Models can be *learned* 
  - You estimate  $\theta$  s.t.  $P(x; \theta)$  best fits the observations X
- 'Best fit'' in what sense?

• Maximum Likelihood

• How to model the distribution of high dimensional data?

$$P(x) = \int_{Z} P(x,z) dz = \int_{Z} P(z)P(x|z)$$

 $\circ P_{\theta}(z)$  and  $P_{\theta}(x|z)$  can be factorized

$$P_{\theta}(x|z) = P_{\theta}(x^{1} \dots, x^{D} | z) = \prod_{i} P_{\theta}(x^{i} | z)$$
$$P_{\theta}(x) = \underset{\theta}{\operatorname{argmax}} \prod_{i} P_{\theta}(x^{i} | z) = \underset{\theta}{\operatorname{argmax}} \log \sum_{i} P_{\theta}(x^{i} | z)$$

$$P_{\theta}(x|z) = P_{\theta}(x^{1} \dots, x^{D} | z) = \prod_{i} P_{\theta}(x^{i} | z)$$
$$\theta^{*} = \underset{\theta}{\operatorname{argmax}} P_{\theta}(x) = \underset{\theta}{\operatorname{argmax}} \prod_{i} P_{\theta}(x^{i} | z) = \underset{\theta}{\operatorname{argmax}} \log \sum_{i} P_{\theta}(x^{i} | z)$$

$$\theta^* = \operatorname*{argmax}_{\theta} P(x;\theta)$$

#### Ok whatever. Tell us how it works ...

Let's start simple ...

# The idea of an "Auto-encoder"

- NN trained to reproduce the input
- F() is a composition of two functions:
   o Embedding / Feature / Latent
  - $\circ$  Output



 $\hat{x} = F(x)$ 

encoder E() and decoder D() z = E(x) $\hat{x} = D(z) = D(E(x))$ 

# How would you train an autoencoder? Loss Function?



# Autoencoder: Loss Function

- The objective is to minimize the "distance" between x and  $\hat{x}$ • If  $d(x, \hat{x}) = 0$  then we get perfect reconstruction
- Mean squared error!
- Cross Entropy (for binary inputs)

$$l(f(\mathbf{x})) = -\sum_k$$

• For both cases, gradient is very simple:

$$l(f(\mathbf{x})) = \frac{1}{2} \sum_{k} (\widehat{x}_k - x_k)^2$$

 $(x_k \log(\hat{x}_k) + (1 - x_k) \log(1 - \hat{x}_k))$ 

 $\nabla_x L(f(x), x) = \hat{x} - x$ 





# Autoencoder: Simple Example



Source: Sebastian Raschka

# **Convolutional Autoencoder**







Source: Sebastian Raschka

## **Convolutional Autoencoder**





### **Convolutional Autoencoder: Expand Dimensions? Transposed Convolution!**

- The decoder needs to "expand dimensions"
   O Convert a small feature z into a large input x
- Use transposed convolution! A.K.A. fractionally stride convolution

   Often (incorrectly) called "de"convolution
   This is an incorrect term because mathematically "deconvolution" is "inverse of convolution"

# **Convolutional Autoencoder: Expand Dimensions?** Transposed Convolution!

#### **Regular Convolution:**





# **Transposed Conv in PyTorch**

```
: import torch
```

```
[-0.3054, -0.4644, -0.3286, -0.2444],
[-0.2332, -0.2557, -0.1876, -0.3970]]]],
grad_fn=<ThnnConvTranspose2DBackward>)
```

torch.Size([1, 1, 6, 6])

#### output = s(n-1) + k - 2p

torch.Size([1, 1, 10, 10])

# **Denoising Autoencoder**

- The input is "noisy"  $\tilde{x}$ . The expected output is a clean image (denoised image)
- Noise Examples:
  - $\tilde{x} = x + z; \quad z \sim N(0, \sigma^2 I)$ Gaussian:
  - Zero-out some of the components of x • Masking: (for images, make some pixels 0)
    - Can be random masks
    - Can be square masks
- Adding noise makes representations more robust

• Expect  $D(E((\tilde{x})) = x)$  for all z



Vincent, Larochelle, Bengio, Manzagol. ICML 2008





### **Example: Face Auto-Encoder**

#### Once trained, what can you do with this model?





# With Generative Models, there are 2 objectives:



Training data ~ p<sub>data</sub>(x)

# Objectives: 1. Learn p<sub>model</sub>(x) that approximates p<sub>data</sub>(x) 2. Sampling new x from p<sub>model</sub>(x)

figure adapted from Ranjay Krishna

5:

# An Auto-Encoder is a Generative Model **Probabilistic Interpretation:**

- Encoder E() estimates
- Decoder D() estimates  $P_{\theta_D}(x|z)$
- The marginal Bayes/Chain Rule ...
- Once the AE is "trained"
  - o you get a generative model that generates "x" given a latent code "z"
  - A conditional generative model takes an additional input "y"
    - E.g. generating images from text
    - More on this later ...

 $P_{\theta_F}(z|x)$ 

#### $P(x) = \int_{z} P(x,z) dz = \int_{z} P(z)P(x|z)$

P(x|z,y)

y=text, x = image

### **Example: Face Auto-Encoder**

#### Once trained, what can you do with this model?

Encode images into vectors (1)(throw away the decoder ...)





# **Example: Face Auto-Encoder**

Once trained, what can you do with this model?

(1) Encode images into vectors

(2) Generate new faces ...

(throw away the encoder)





#### Types of Autoencoders



#### So far, we have not enforced any "structure" on the latents z

But a structure is desirable

(Remember our motivations / goals for representation learning)



#### So far, we have not enforced any "structure" on the latents z

We can't generate *new images* from D() if we don't understand the z-space





- So far, we have not enforced any "structure" on the latents z
- We can't generate *new images* from D() if we don't understand the z-space
- For example, if I ask you to generate a "face with beard, glasses, brown hair" which z would you choose?





# VAE: Variational Autoencoder



- Force a "prior" distribution on the latent space • Example: Gaussian N(0, I)
- Gaussians are nice because they are perfectly symmetrical in every dimension
  - Isotropic (covariance matrix is identity I)
  - Dimensions are independent,
    - i.e.  $P(z_1|z_2) = P(z_1) = N(0,I) \forall z_1, z_2$
  - Property holds for any linear combination of z elements
    - i.e.  $P(z_1|az_2 + bz_3) = N(0,I)$



















Probabilistic spin on autoencoders - will let us sample from the model to generate data!

representation z

Assume training data  $\{x^{(i)}\}_{i=1}^{N}$  is generated from the distribution of unobserved (latent)

#### Sample from true conditional $p_{\theta^*}(x \mid z^{(i)})$

Sample from true prior  $z^{(i)} \sim p_{ heta^*}(z)$ 



We want to estimate the parameters  $\theta^*$ given training real data x.

How should we represent this model?

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014



#### Sample from true conditional $p_{\theta^*}(x \mid z^{(i)})$

Sample from true prior  $z^{(i)} \sim p_{ heta^*}(z)$ 





We want to estimate the parameters  $\theta^*$ given training real data x.

How should we represent this model?

Choose prior p(z) to be simple, e.g. Gaussian. Reasonable for latent attributes, e.g. pose, how much smile.

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014





Sample from true conditional  $p_{\theta^*}(x \mid z^{(i)})$ 

Sample from true prior $z^{(i)} \sim p_{ heta^*}(z)$ 







We want to estimate the parameters  $\theta^*$  given training real data x.

How should we represent this model?

Choose prior p(z) to be simple, e.g. Gaussian. Reasonable for latent attributes, e.g. pose, how much smile.

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014



Decoder must be probabilistic: Decoder inputs z, outputs mean  $\mu_{x|z}$  and (diagonal) covariance  $\sum_{x|z}$ 

Sample from true conditional  $p_{\theta^*}(x \mid z^{(i)})$ 

Sample from true prior

 $z^{(i)} \sim p_{ heta^*}(z)$ 





We want to estimate the parameters  $\theta^*$ given training real data x.

How should we represent this model?

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014



Decoder must be probabilistic: Decoder inputs z, outputs mean  $\mu_{xz}$  and (diagonal) covariance  $\sum_{x|z}$ 

Sample from true conditional  $p_{\theta^*}(x \mid z^{(i)})$ 

Sample from true prior

 $z^{(i)} \sim p_{ heta^*}(z)$ 





We want to estimate the parameters  $\theta^*$ given training real data x.

How should we represent this model?

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014



Sample from true conditional  $p_{\theta^*}(x \mid z^{(i)})$ 

Sample from true prior  $z^{(i)} \sim p_{ heta^*}(z)$ 



We want to estimate the parameters  $\theta^*$ given training real data x.

How to train the model?

Learn model parameters to maximize likelihood of training data

$$p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$$

Q: What is the problem with this?



Sample from true conditional  $p_{\theta^*}(x \mid z^{(i)})$ 

Sample from true prior  $z^{(i)} \sim p_{ heta^*}(z)$ 



We want to estimate the parameters  $\theta^*$ given training real data x.

How to train the model?

Learn model parameters to maximize likelihood of training data

$$p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$$

Q: What is the problem with this? Intractable! Impossible to iterate over all z



# Variational Autoencoders: Intractability

Data likelihood:  $p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$ 

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# Variational Autoencoders: Intractability Data likelihood: $p_{ heta}(x) = \int p_{ heta}(z) p_{ heta}(x|z) dz$ **Simple Gaussian prior**

# Variational Autoencoders: Intractability Data likelihood: $p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$

- **Decoder neural network**
# Variational Autoencoders: Intractability Data likelihood: $p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$ Intractable to compute p(x|z) for every z!

# Variational Autoencoders: Intractability Data likelihood: $p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$ Intractable to compute p(x|z) for every z!

 $\log p(x) \approx \log \frac{1}{k} \sum_{i=1}^{k} p(x|z^{(i)})$ , where  $z^{(i)} \sim p(z)$ 

- Monte Carlo estimation is too high variance

# **Variational Autoencoders: Intractability** Data likelihood: $p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$ Another idea: $p_{\theta}(x) = \frac{p_{\theta}(x \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x)}$ Use Bayes rule

# Variational Autoencoders: Intractability Data likelihood: $p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$ Another idea: $p_{\theta}(x) = \frac{p_{\theta}(x \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x)}$

We know how to calculate these



#### Variational Autoencoders: Intractability

Data likelihood:  $p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$ Another idea:  $p_{\theta}(x) = \frac{p_{\theta}(x \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x)}$ 

**Solution**: In addition to modeling  $p_{\theta}(x|z)$ , Learn  $q_{\phi}(z|x)$  that approximates the true posterior  $p_{\theta}(z|x)$ .

**Encoder Network**  $q_{\phi}(z \mid x) = N(\mu_{z|x}, \Sigma_{z|x})$  $\mu_{z|x}$   $\Sigma_{z|x}$ 

x

#### But how do you calculate this?

#### **Decoder Network**



## Variational Autoencoders: Intractability

Data likelihood:  $p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$ Another idea:  $p_{\theta}(x) = \frac{p_{\theta}(x \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x)}$ 



- **x**: 28x28 image = 784-dim vector z: 20-dim vector
- **Encoder Network**

**Decoder Network** 

z: 20



$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)})\right] \qquad (p$$

Using this approximation, we can derive a lower bound on the data likelihood p(x), making it tractable, therefore, possible to optimize.

 $p_{\theta}(x^{(i)})$  Does not depend on z)

 $\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)})\right] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z)$ 

Taking expectation wrt. z (using encoder network) will come in handy later

$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[ \log p_{\theta}(x^{(i)}) \right] \quad (p_{\theta}(z)) = \mathbf{E}_{z} \left[ \log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \right] \quad (B$$

- $p_{\theta}(x^{(i)})$  Does not depend on z)
- Bayes' Rule)

$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)})\right] \quad (p_{\theta}(z)) = \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})}\right] \quad (B)$$
$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \frac{q_{\phi}(z \mid z)}{q_{\phi}(z \mid z)}\right]$$

 $p_{\theta}(x^{(i)})$  Does not depend on z)

Bayes' Rule)

 $\frac{|x^{(i)})}{|x^{(i)})}$  (Multiply by constant)

$$\begin{split} \log p_{\theta}(x^{(i)}) &= \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[ \log p_{\theta}(x^{(i)}) \right] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_{z} \left[ \log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_{z} \left[ \log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \frac{q_{\phi}(z \mid x^{(i)})}{q_{\phi}(z \mid x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_{z} \left[ \log p_{\theta}(x^{(i)} \mid z) \right] - \mathbf{E}_{z} \left[ \log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z)} \right] + \mathbf{E}_{z} \left[ \log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})} \right] \quad (\text{Logarithms}) \end{split}$$

$$\begin{split} \log p_{\theta}(x^{(i)}) &= \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[ \log p_{\theta}(x^{(i)}) \right] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_{z} \left[ \log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_{z} \left[ \log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \frac{q_{\phi}(z \mid x^{(i)})}{q_{\phi}(z \mid x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_{z} \left[ \log p_{\theta}(x^{(i)} \mid z) \right] - \mathbf{E}_{z} \left[ \log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z)} \right] + \mathbf{E}_{z} \left[ \log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \mathbf{E}_{z} \left[ \log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z)) + D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z \mid x^{(i)})) \\ \end{split}$$

The expectation wrt. z (using encoder network) let us write nice KL terms

$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)})\right] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})}\right] \quad (\text{Bayes' Rule})$$

$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \frac{q_{\phi}(z \mid x^{(i)})}{q_{\phi}(z \mid x^{(i)})}\right] \quad (\text{Multiply by constant})$$

$$= \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z)}\right] + \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})}\right] \quad (\text{Logarithms})$$

$$= \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z)) + D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z \mid x^{(i)}))$$

Decoder network gives  $p_{\theta}(x|z)$ , can compute estimate of this term through sampling (need some trick to differentiate through sampling).

Dec compute estimate of this term through sampling (need some trick to differentiate through sampling).

Gaussians for encoder and z prior) has nice closed-form solution!

De CO sampling (need some trick to

prior) has nice closed-form

term : ( But we know KL divergence always >= 0.

$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)})\right] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})}\right] \quad (\text{Bayes' Rule})$$
We want to
maximize the
data
$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \frac{q_{\phi}(z \mid x^{(i)})}{q_{\phi}(z \mid x^{(i)})}\right] \quad (\text{Multiply by constant})$$

$$= \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z)}\right] + \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})}\right] \quad (\text{Logarithms})$$

$$= \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z)) + D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z \mid x^{(i)}))$$

$$= \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z)) + D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z \mid x^{(i)}))$$

$$= \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z)) + D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z \mid x^{(i)}))\right]$$

sampling.

prior) has nice closed-form solution!

term : ( But we know KL divergence always >= 0.

$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)})\right] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})}\right] \quad (\text{Bayes' Rule})$$
We want to
maximize the
data
$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \frac{q_{\phi}(z \mid x^{(i)})}{q_{\phi}(z \mid x^{(i)})}\right] \quad (\text{Multiply by constant})$$

$$= \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z)}\right] + \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})}\right] \quad (\text{Logarithms})$$

$$= \underbrace{\mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z \mid x^{(i)}))}_{\geq 0}\right]$$

Tractable lower bound which we can take gradient of and optimize! ( $p_{\theta}(x|z)$  differentiable, KL term is differentiable)

$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)})\right] \quad (p_{\theta}(z)) = \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)}|z)p_{\theta}(z)}{p_{\theta}(z|x^{(i)})}\right] \quad (B_{z})$$
Decoder:  
reconstruct =  $\mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)}|z)p_{\theta}(z)}{p_{\theta}(z|x^{(i)})} \frac{q_{\phi}(z|z)}{q_{\phi}(z|z)}\right] = \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)}|z)\right] - \mathbf{E}_{z} \left[\log \frac{q_{\theta}(z|z)}{\mathcal{L}(x^{(i)},\theta,\phi)}\right]$ 

Tractable lower bound which we can take gradient of and optimize! ( $p_{\theta}(x|z)$  differentiable, KL term differentiable)

 $\theta_{\theta}(x^{(i)})$  Does not depend on z)





Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_{z}\left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid\mid p_{\theta}(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_{z}\left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid\mid p_{\theta}(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Let's look at computing the KL divergence between the estimated posterior and the prior given some data

#### **Input Data**

x

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_{z}\left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid\mid p_{\theta}(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$



Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_{z}\left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid\mid p_{\theta}(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$
Make approximate

posterior distribution close to prior

#### $D_{KL}(\mathcal{N}(\mu_{z|x},\Sigma_{z|x})||\mathcal{N}(0,I))$

#### This equation has an analytical solution



Putting it all together: maximizing the likelihood lower bound

$$\mathbf{E}_{z} \left[ \log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid\mid p_{\theta}(z))$$

$$\mathcal{L}(x^{(i)}, \theta, \phi)$$
Make approximate posterior distribution

close to prior



Putting it all together: maximizing the likelihood lower bound

$$\mathbf{E}_{z} \left[ \log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid\mid p_{\theta}(z))$$
$$\mathcal{L}(x^{(i)}, \theta, \phi)$$



Reparameterization trick to make sampling differentiable:

Sample 
$$\epsilon \sim \mathcal{N}(0,I)$$
  
 $z = \mu_{z|x} + \epsilon \sigma_{z|x}$ 



Putting it all together: maximizing the likelihood lower bound

$$\mathbf{E}_{z} \left[ \log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid\mid p_{\theta}(z))$$
$$\mathcal{L}(x^{(i)}, \theta, \phi)$$







Putting it all together: maximizing the likelihood lower bound

$$\mathbf{E}_{z} \left[ \log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid\mid p_{\theta}(z))$$
$$\mathcal{L}(x^{(i)}, \theta, \phi)$$



# Variational Autoencoders Maximize likelihood of original

Putting it all together: maximizing the likelihood lower bound

$$\mathbf{E}_{z}\left[\log p_{\theta}(x^{(i)} \mid z)\right] \blacktriangleleft D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid\mid p_{\theta}(z))$$

 $\mathcal{L}(x^{(i)}, \theta, \phi)$ 



Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_{z}\left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid\mid p_{\theta}(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

For every minibatch of input data: compute this forward pass, and then backprop!



Our assumption about data generation process

Sample from true conditional  $p_{\theta^*}(x \mid z^{(i)})$ 

Sample from true prior

 $z^{(i)} \sim p_{ heta^*}(z)$ 



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Our assumption about data generation process

Sample from true conditional  $p_{\theta^*}(x \mid z^{(i)})$ 

Sample from true prior

 $z^{(i)} \sim p_{ heta^*}(z)$ 



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

- Now given a trained VAE:
- use decoder network & sample z from prior!

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

Use decoder network. Now sample z from prior!



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

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Use decoder network. Now sample z from prior!



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

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

#### Data manifold for 2-d z

6600 Vary **z**<sub>1</sub> Vary z,

Diagonal prior on z => independent latent variables

Different dimensions of **z** encode interpretable factors of variation

Degree of smile Vary **z**<sub>1</sub>

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014



Vary **z**<sub>1</sub>

Diagonal prior on z => independent latent variables

Different dimensions of z encode interpretable factors of variation

Degree of smile

Also good feature representation that can be computed using  $q_{\phi}(z|x)!$ 

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014





#### 32x32 CIFAR-10



#### Labeled Faces in the Wild

## Editing images with VAEs

 Run input data through encoder to get a distribution over latent codes


- 1. Run input data through encoder to get a distribution over latent codes
- 2. Sample code z from encoder output



- 1. Run input data through encoder to get a distribution over latent codes
- 2. Sample code z from encoder output
- Modify some dimensions of 3. sampled code



- 1. Run input data through encoder to get a distribution over latent codes
- 2. Sample code z from encoder output
- Modify some dimensions of 3. sampled code
- 4. Run modified z through decoder to get a distribution over data sample



- 1. Run input data through encoder to get a distribution over latent codes
- 2. Sample code z from encoder output
- 3. Modify some dimensions of sampled code
- 4. Run modified z through decoder to get a distribution over data sample
- Sample new data from (4) 5.







