CS 4644-DL / 7643-A: LECTURE 17 DANFEI XU

Generative Models:

PixelCNN / PixelRNN

Variational AutoEncoders (VAEs)

Administrative

- Milestone Report is due EOD 11/3 NO GRACE PERIOD
- HW3 now due EOD 10/22 + 2 late days
- HW4 is out, due 11/11

Supervised Learning

- Train Input: {X, Y}
- Learning output: $f: X \to Y$, e.g. P(y|x)

Unsupervised Learning

- Input: {X}
- Learning output: P(x)
- Example: Clustering, density estimation, etc.

Reinforcement Learning

- Supervision in form of reward
- No supervision on what action to take

Very often combined, sometimes within the same model!

Supervised Learning

- Train Input: {X, Y}
- Learning output: $f: X \to Y$, e.g. P(y|x)

Unsupervised Learning

- Input: {X}
- Learning output: P(x)
- Example: Clustering, density estimation, etc.

Reinforcement Learning

- Supervision in form of reward
- No supervision on what action to take

Very often combined, sometimes within the same model!

What if all we have are data without label?



We have lots of *raw* data (e.g., Internet)! Can we still learn useful things without labels?

Generative Models

Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn a *function* to map x -> y

Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.

Supervised Learning

Data: (x, y) x is data, y is label

Goal: Learn a *function* to map x -> y

Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.



Classification

Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn a *function* to map x -> y

Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.



A cat sitting on a suitcase on the floor

Image captioning

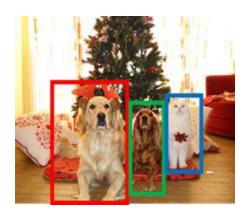
Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn a *function* to map x -> y

Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.



DOG, DOG, CAT

Object Detection

Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn a *function* to map x -> y

Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.



Semantic Segmentation

Unsupervised Learning

Data: x Just data, **no labels!**

Goal: Learn some underlying hidden *structure* of the data

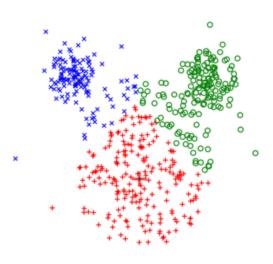
Examples: Clustering, dimensionality reduction, density estimation, etc.

Unsupervised Learning

Data: x Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, density estimation, etc.



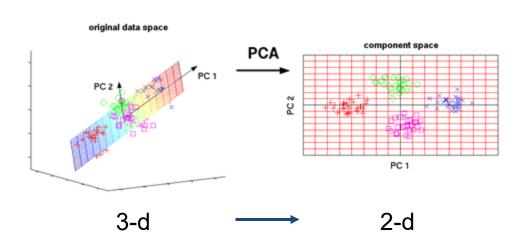
K-means clustering

Unsupervised Learning

Data: x Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, density estimation, etc.



Principal Component Analysis (Dimensionality reduction)

Unsupervised Learning

Data: x
Just data, no labels!

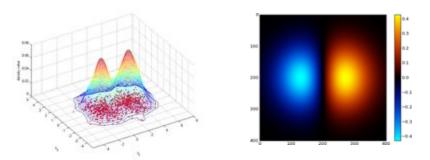
Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, density estimation, etc.



Figure copyright Ian Goodfellow, 2016. Reproduced with permission.

1-d density estimation



2-d density estimation

Modeling p(x)

Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn a *function* to map x -> y

Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.

Unsupervised Learning

Data: x

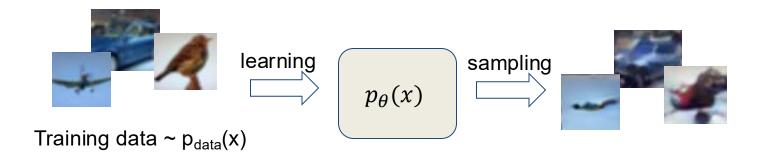
Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, density estimation, etc.

Generative Modeling

Given training data, generate new samples from same distribution

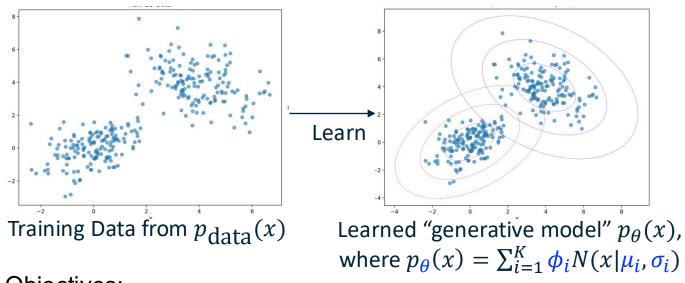


Objectives:

- 1. Learn $p_{\theta}(x)$ that approximates an unknown $p_{\text{data}}(x)$
- 2. Sampling new x from $p_{model}(x)$

Generative Modeling

Gaussian Mixture Model (GMM) as a generative model

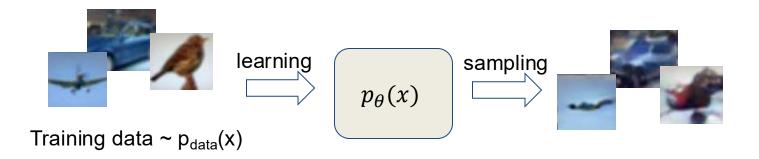


Objectives:

- 1. Learn $p_{\theta}(x)$ that approximates an unknown $p_{\text{data}}(x)$
- 2. Sampling new x from $p_{model}(x)$

Generative Modeling

Given training data, generate new samples from same distribution

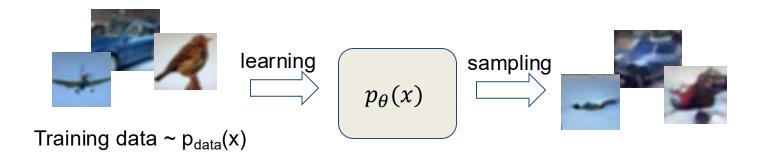


Formulate as density estimation problems:

- **Explicit density estimation**: explicitly define and solve for $p_{\theta}(x)$, e.g., a high-dimensional Gaussian Mixture Model (GMM)
- **Implicit density estimation**: learn model that can sample from $p_{\theta}(x)$ without explicitly defining it.

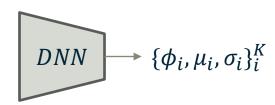
Deep Generative Models

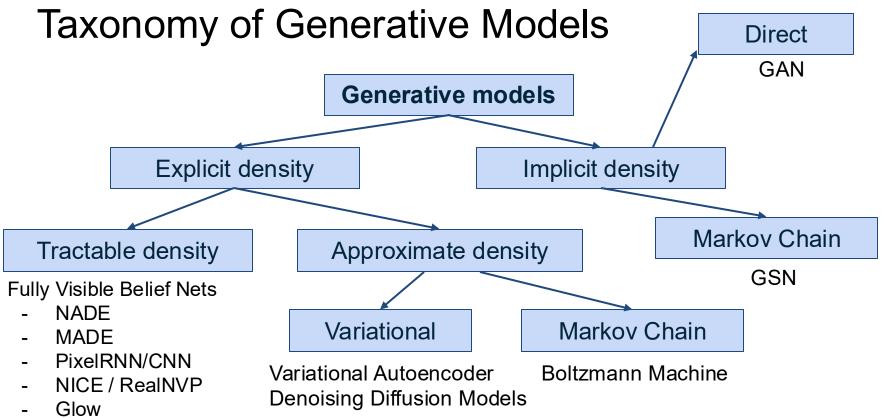
Given training data, generate new samples from same distribution



Use deep neural networks to represent $p_{\theta}(x)$!

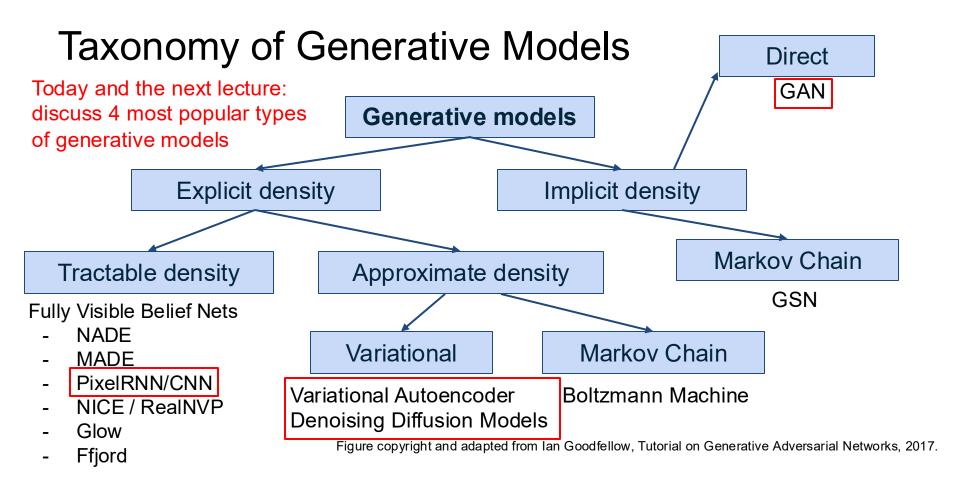
Example: a DNN with GMM output





Ffjord

Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

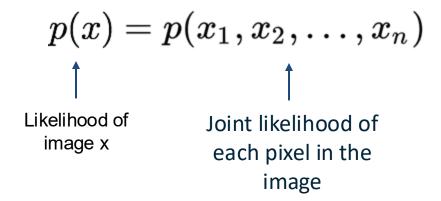


PixelRNN and PixelCNN

(Autoregressive Generative Model)

Fully visible belief network (FVBN)

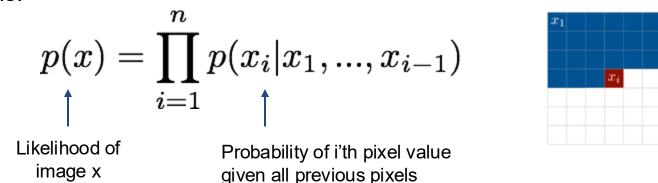
Explicit density model



Fully visible belief network (FVBN)

Explicit density model

Use probability chain rule to decompose likelihood of an image x into product of 1-d distributions:



Then maximize likelihood of training data

Fully visible belief network (FVBN)

Explicit density model

Use probability chain rule to decompose likelihood of an image x into product of 1-d distributions:

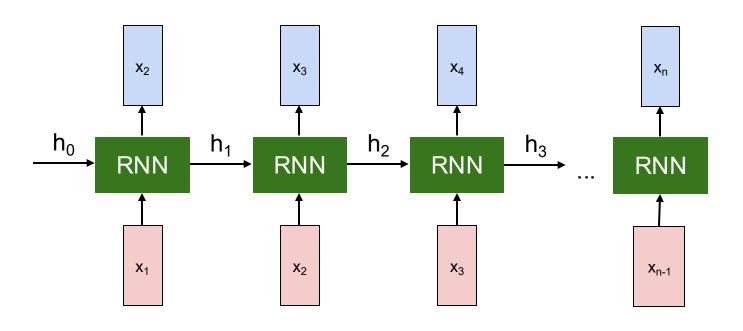
$$p(x) = \prod_{i=1}^{n} p(x_i|x_1,...,x_{i-1})$$

$$\uparrow \qquad \downarrow \qquad \uparrow$$
Likelihood of image x
$$\downarrow p(x_i|x_1,...,x_{i-1})$$
Probability of i'th pixel value given all previous pixels

Then maximize likelihood of training data

Complex distribution over pixel values => Express using a neural network!

Recurrent Neural Network



$$p(x_i|x_1,...,x_{i-1})$$

Generate image pixels starting from corner

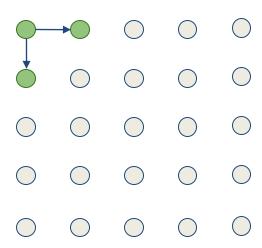
Dependency on previous pixels modeled using an RNN (LSTM)

| \bigcup | \bigcirc | \bigcup | \subset |
|-----------|------------|-----------|-----------|

$$\circ$$
 \circ \circ \circ

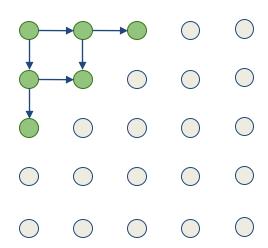
Generate image pixels starting from corner

Dependency on previous pixels modeled using an RNN (LSTM)



Generate image pixels starting from corner

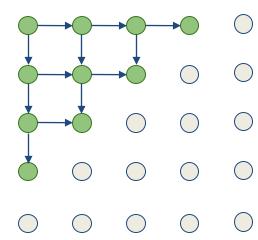
Dependency on previous pixels modeled using an RNN (LSTM)



Generate image pixels starting from corner

Dependency on previous pixels modeled using an RNN (LSTM)

Drawback: sequential generation is slow in both training and inference!



Still generate image pixels starting from corner

Dependency on previous pixels now modeled using a CNN over context region (masked convolution)

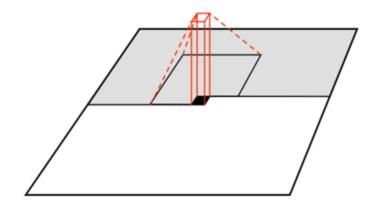


Figure copyright van der Oord et al., 2016. Reproduced with permission.

Still generate image pixels starting from corner

Dependency on previous pixels now modeled using a CNN over context region (masked convolution)

Training is faster than PixelRNN (can parallelize convolutions since context region values known from training images)

Generation is still slow:

For a 32x32 image, we need to do forward passes of the network 1024 times for a single image

Figure copyright van der Oord et al., 2016. Reproduced with permission.

Generation Samples



32x32 CIFAR-10



32x32 ImageNet

Figures copyright Aaron van der Oord et al., 2016. Reproduced with permission.

PixelRNN and PixelCNN

$$\rightarrow$$
 P(x) = 0.12

$$\rightarrow$$
 P(x) = 0.00003

Pros:

- Can explicitly compute likelihood p(x)
- Easy to optimize
- Good samples

Con:

Sequential generation => slow

Improving PixelCNN performance

- Gated convolutional layers
- Short-cut connections
- Discretized logistic loss
- Multi-scale
- Training tricks
- Etc...

See

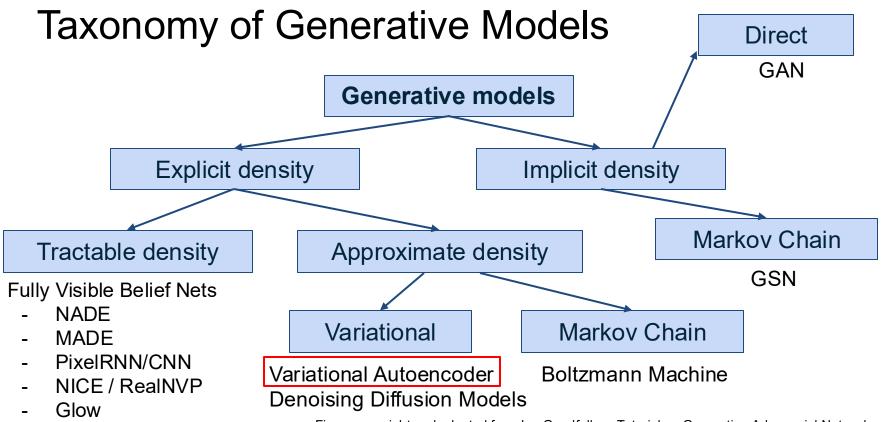
- Van der Oord et al. NIPS 2016
- Salimans et al. 2017 (PixelCNN++)

Aside: Why are LLMs "Generative"?

$$p(x) = \prod_{i=1}^n p(x_i|x_1,...,x_{i-1})$$
 Transformer/RNN

Language models, especially built on RNN or Transformers with proper causal masking, can be thought of as **autoregressive generative models (predict future based on the past)**, similar to PixelRNN and PixelCNN.

Sample an entire sentence by taking sequence of word samples following the probability chain rule decomposition.



Ffjord

Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Variational

Autoencoders (VAE)

PixelR/CNNs define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^{n} p_{\theta}(x_i|x_1, ..., x_{i-1})$$

PixelR/CNNs define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^{n} p_{\theta}(x_i|x_1, ..., x_{i-1})$$

Variational Autoencoders (VAEs) define intractable density function with latent z:

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

No dependencies among pixels, can generate all pixels at the same time! Latent variable z that captures important factors of variations in dataset

PixelR/CNNs define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^{n} p_{\theta}(x_i|x_1, ..., x_{i-1})$$

Variational Autoencoders (VAEs) define intractable density function with latent z:

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

No dependencies among pixels, can generate all pixels at the same time! Latent variable z that captures important *factors of variations* in dataset

Cannot optimize (maximum likelihood estimation) directly, derive and optimize lower bound on likelihood instead

PixelR/CNNs define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^{n} p_{\theta}(x_i|x_1, ..., x_{i-1})$$

Variational Autoencoders (VAEs) define intractable density function with latent z:

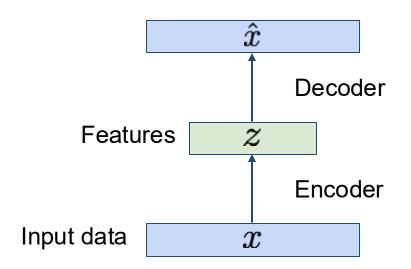
$$p_{ heta}(x) = \int p_{ heta}(z) p_{ heta}(x|z) dz$$

No dependencies among pixels, can generate all pixels at the same time!

Latent variable z that captures important factors of variations in dataset

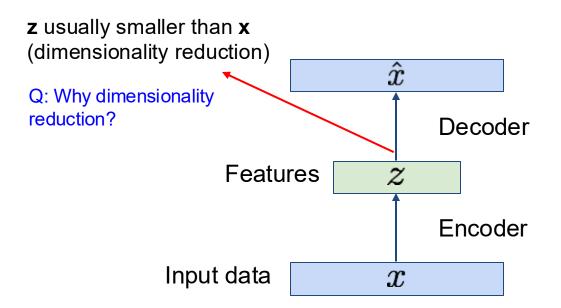
Cannot optimize (maximum likelihood estimation) directly, derive and optimize lower bound on likelihood instead

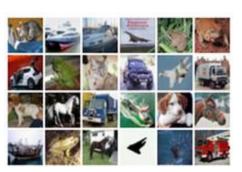
Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data



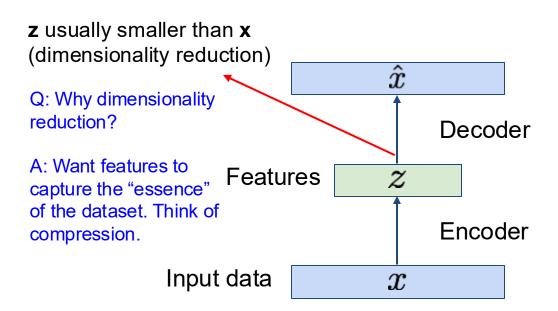


Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data





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





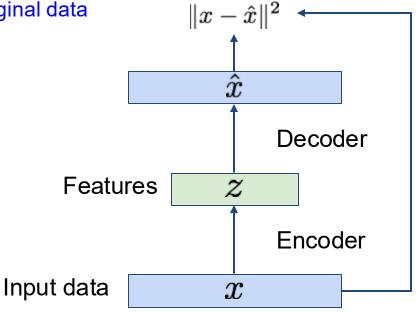
How to learn this feature Reconstructed representation? input data Train such that features \hat{x} can be used to reconstruct original data Decoder "Autoencoding" encoding input itself **Features** Encoder Input data x



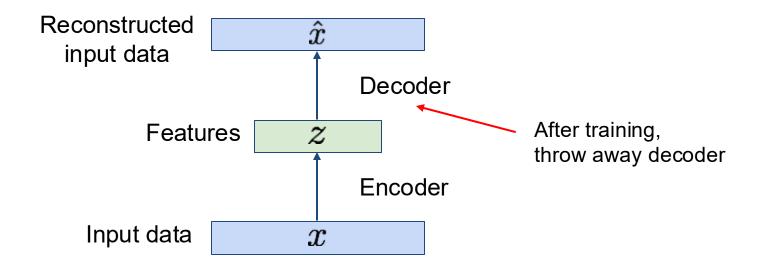
Train such that features can be used to reconstruct original data

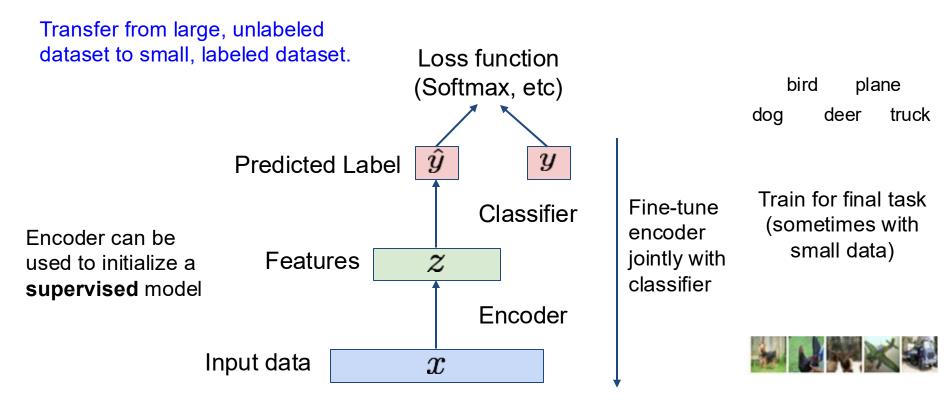
Doesn't use labels!

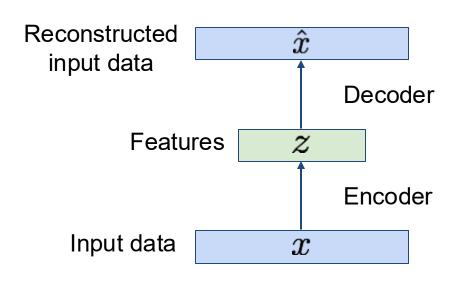
L2 Loss function:



Reconstructed data **Encoder**: 4-layer conv **Decoder**: 4-layer upconv Input data



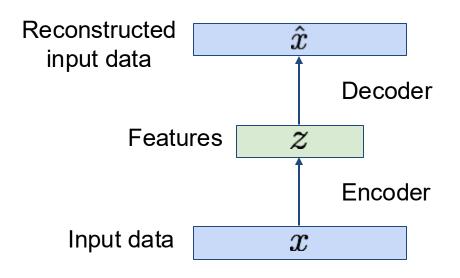




Autoencoders can reconstruct data, and can learn features to initialize a supervised model

Features capture factors of variation in training data.

Ideally, knowing the space of Z is sufficient to recover the *entire training set* through the decoder.



Autoencoders can reconstruct data, and can learn features to initialize a supervised model

Features capture factors of variation in training data.

Ideally, knowing the space of Z is sufficient to recover the *entire training set* through the decoder.

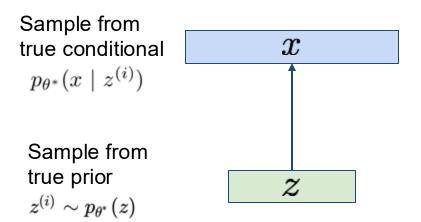
VAE: Model data distribution p(x) through a probabilistic latent space p(z) and a probabilistic decoder p(x|z).

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

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

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

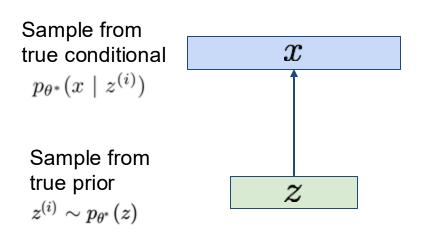
Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from the distribution of unobserved (latent) representation ${\bf z}$



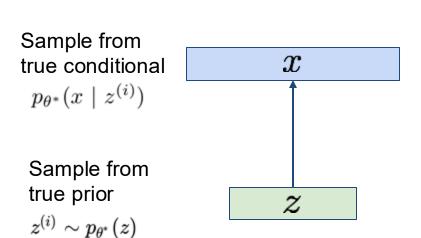
Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

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

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



Intuition (remember from autoencoders!): **x** is an image, **z** is latent code used to generate **x**.



Goal of Variational Autoencoder:

We want to estimate the true parameters θ^* of this generative model given training data x.

 θ^* includes both the decoder neural network parameters and the prior distribution

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

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

We want to estimate the true parameters θ^* of this generative model given training data x.

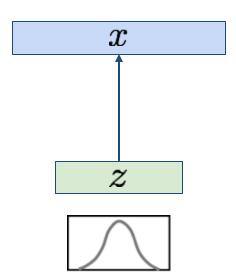
How should we represent this model?

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 true parameters θ^* of this generative model given training data x.

How should we represent this model?

Assume p(z) is *known* and *simple*, e.g. isotropic Gaussian. Reasonable for latent attributes, e.g. pose, how much smile.

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

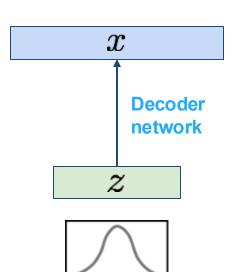


Sample from true conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

Sample from true prior

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



We want to estimate the true parameters θ^* of this generative model given training data x.

How should we represent this model?

Assume p(z) is *known* and *simple*, e.g. isotropic Gaussian. Reasonable for latent attributes, e.g. pose, how much smile.

Conditional p(x|z) is *complex* (generates image) => represent with neural network

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

Sample from true conditional $m{x}$ Decoder network Sample from true prior $m{z}^{(i)} \sim p_{ heta^*}(z)$

We want to estimate the true parameters θ^* of this generative model given training data x.

How to train the model?

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

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

We want to estimate the true parameters θ^* of this generative model given training data x.

How to train the model?

Learn model parameters to maximize likelihood of training data

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

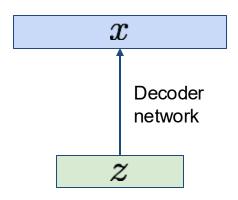
Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

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 true parameters θ^* of this generative model given training 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!

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

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

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

Simple Gaussian prior

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

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

Intractable to compute p(x|z) for every z!

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

Intractable to compute p(x|z) for every z!

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

Can we do Monte Carlo sampling?

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

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

Can we do Monte Carlo sampling?

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

We don't know which z corresponds to a sample (x)! Most z's will be sampled from where p(x|z) is zero.

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

Can we do Monte Carlo sampling?

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

Can we estimate posterior density?

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

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

Can we do Monte Carlo sampling?

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

Can we estimate posterior density? Not quite, but ...

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

Intractable data likelihood

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

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

Can we do Monte Carlo sampling?

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

Can we estimate posterior density? Not quite, but ...

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

VAE: We can use an approximate posterior $q_{\theta}(z|x)$ (variational distribution) to form a *tractable lower bound* of the data likelihood p(x).

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

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

Let's assume we can sample from some approximate posterior for now ...

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

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

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

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

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

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

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

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

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

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

Recall:
$$D_{KL}(q||p) = \mathbf{E}_q[\log \frac{q}{p}]$$

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

approximate log[p(x)] by raising ELBO. Higher ELBO -> lower KL(q(z|x)|p(z|x))

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

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

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

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

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

Maximize the likelihood of the sample $x^{(i)}$ (e.g., an image) given a latent prior sample.

Minimize KL -> Make the approximate posterior more like the prior!
Use NN to model the approximate posterior.

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

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

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

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

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

Maximize the likelihood of the sample $x^{(i)}$ (e.g., an image) given a latent prior sample. Can be thought of a decoder model

This KL term (between Gaussians for encoder and z prior) has nice closed-form solution!

```
\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)}) || p_{\theta}(z))
```

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

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

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

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

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

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

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

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

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

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

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

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

Sample z from the learned posterior (encoder) to train the decoder to reconstruct!

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

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

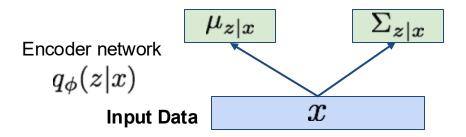
Putting it all together: maximizing the likelihood lower bound

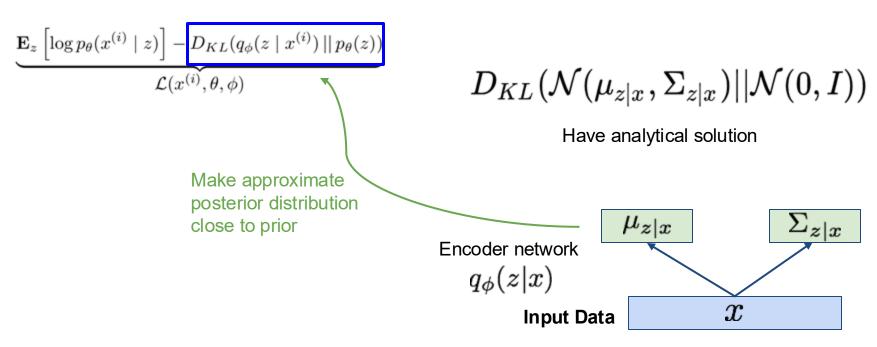
$$\underbrace{\mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid\mid p_{\theta}(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

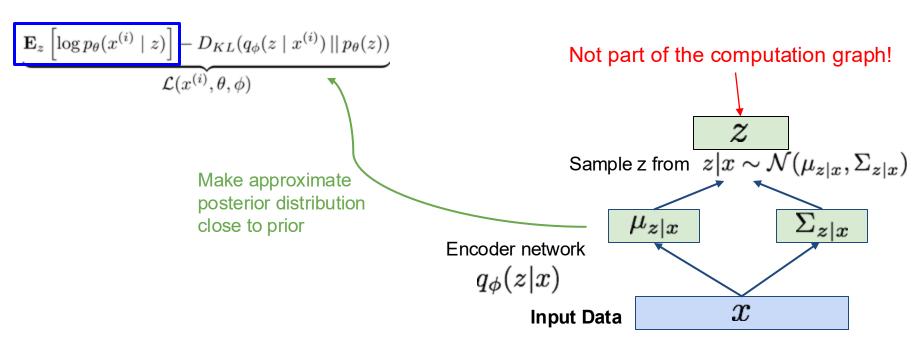
Let's look at computing the KL divergence between the estimated posterior and the prior given some data

Input Data

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





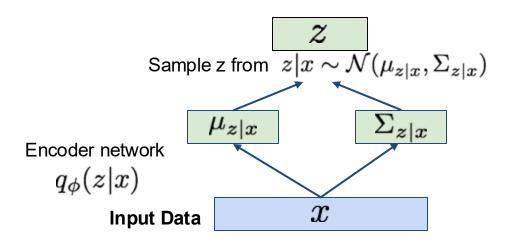


Putting it all together: maximizing the likelihood lower bound

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

Reparameterization trick to make sampling differentiable:

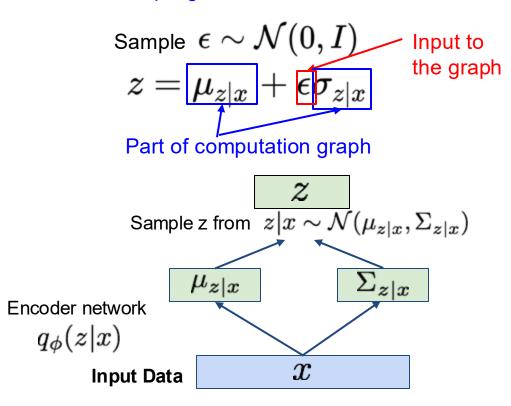
Sample
$$\epsilon \sim \mathcal{N}(0,I)$$
 $z = \mu_{z|x} + \epsilon \sigma_{z|x}$



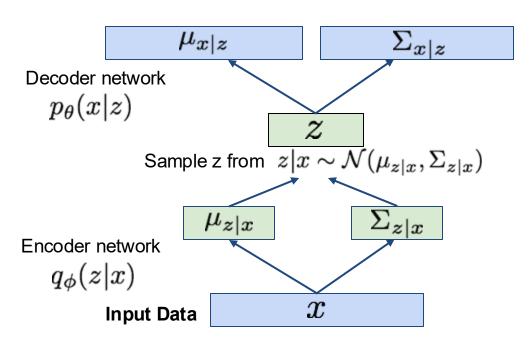
Putting it all together: maximizing the likelihood lower bound

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

Reparameterization trick to make sampling differentiable:



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



Putting it all together: maximizing the likelihood lower bound

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

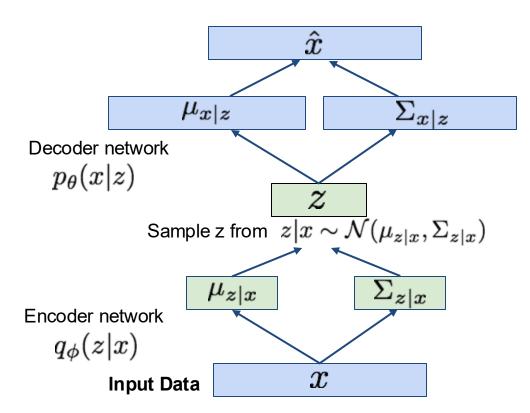
$$\mathcal{L}(x^{(i)}, \theta, \phi)$$

Maximize likelihood of original input being reconstructed \hat{x} $\mu_{x|z}$ $\sum_{x|z}$ Decoder network $p_{\theta}(x|z)$ Sample z from $\overline{z|x} \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$ $\Sigma_{z|x}$ $\mu_{z|x}$ Encoder network $q_{\phi}(z|x)$ x**Input Data**

Putting it all together: maximizing the likelihood lower bound

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

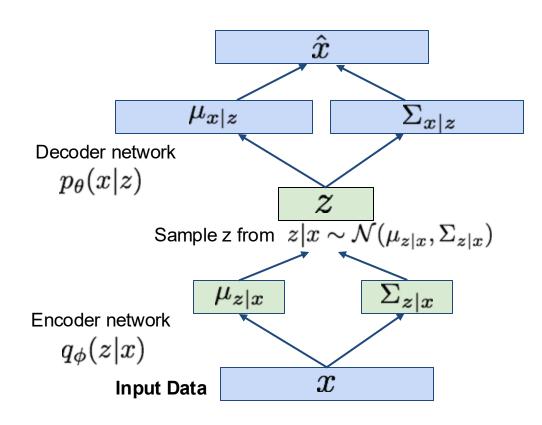
Hyperparameter to weigh the strength of the prior matching objective



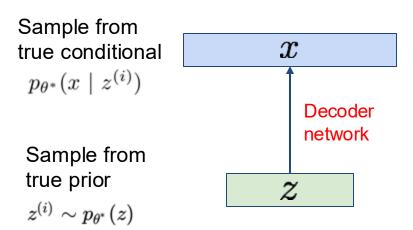
Putting it all together: maximizing the likelihood lower bound

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

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

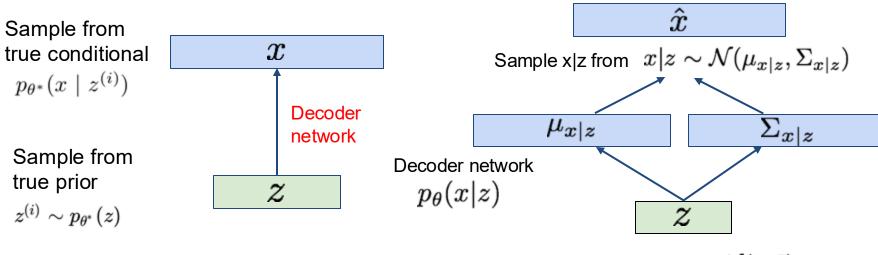


Our assumption about data generation process



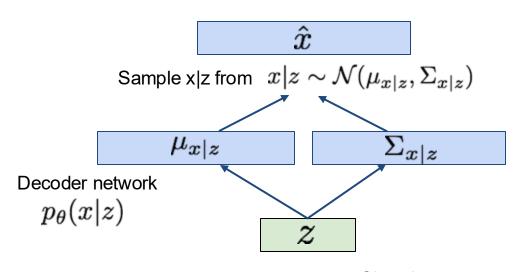
Our assumption about data generation process

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



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

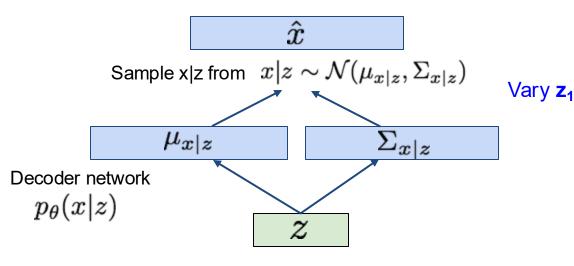
Use decoder network. Now sample z from prior!



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

```
66666666666666
```

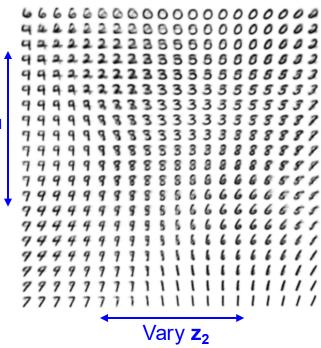
Use decoder network. Now sample z from prior!



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

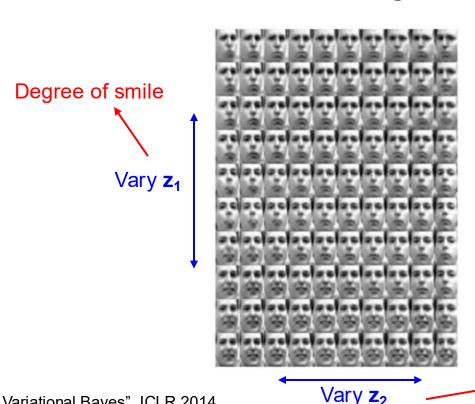
Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Data manifold for 2-d z



Diagonal prior on **z** => independent latent variables

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



Head pose

Diagonal prior on **z** => independent Degree of smile latent variables Different dimensions of **z** Vary **z**₁ encode interpretable factors elaciaciaciaciaciacia of variation Also good feature representation that can be computed using $q_{\phi}(z|x)!$ Head pose Vary **z**₂ Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014



32x32 CIFAR-10



Labeled Faces in the Wild

Figures copyright (L) Dirk Kingma et al. 2016; (R) Anders Larsen et al. 2017. Reproduced with permission.

Probabilistic spin to traditional autoencoders => allows generating data

Defines an intractable density => derive and optimize a (variational) lower bound

Pros:

- Principled approach to generative models
- Latent space z is interpretable and may be useful for other downstream tasks.

Cons:

- Samples are blurry
- KL weights are hard to tune
- Latent distributions are aggressive representation bottlenecks that may limit the expressiveness of the model.

Can be made more powerful by making VAE hierarchical (multiple layers of latents). **Diffusion model (denoising diffusion) can be thought of a type of hierarchical VAE!**

Next Time: Denoising Diffusion