

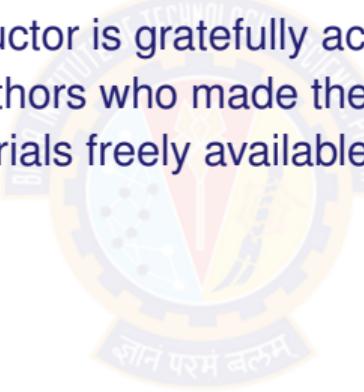


BITS Pilani
Pilani | Dubai | Goa | Hyderabad

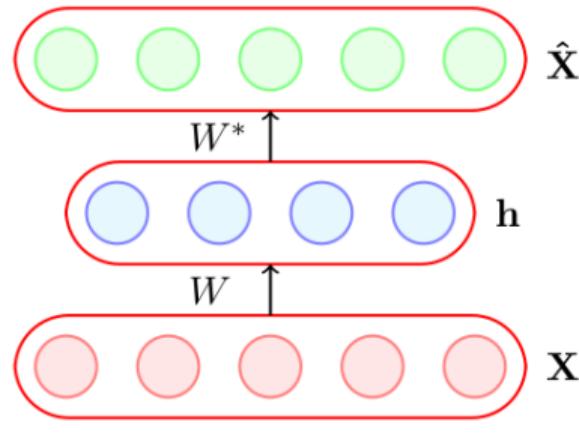
DEEP LEARNING MODULE 7 : VARIATIONAL AUTOENCODERS

Seetha Parameswaran
Asst Prof, BITS Pilani

The instructor is gratefully acknowledging
the authors who made their course
materials freely available online.



AUTOENCODER – SUMMARY



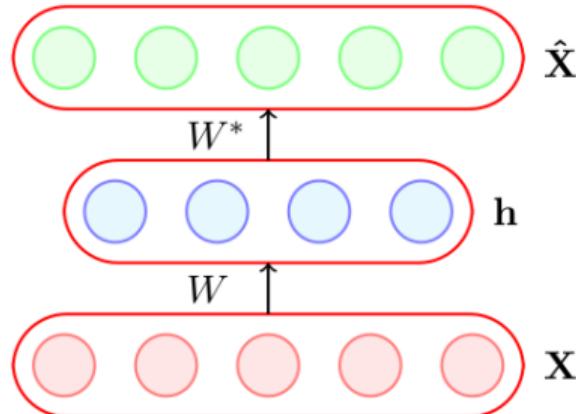
$$h = g(WX + b)$$
$$\hat{X} = f(W^*h + c)$$

- Before we start talking about VAEs, let us quickly revisit autoencoders
- An autoencoder contains an encoder which takes the input X and maps it to a hidden representation
- The decoder then takes this hidden representation and tries to reconstruct the input from it as \hat{X}
- The training happens using the following objective function

$$\min_{W, W^*, \mathbf{c}, \mathbf{b}} \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n (\hat{x}_{ij} - x_{ij})^2$$

- where m is the number of training instances, $\{x_i\}_{i=1}^m$ and each $x_i \in R^n$ (x_{ij} is thus the j -th dimension of the i -th training instance)

AUTOENCODER – SUMMARY

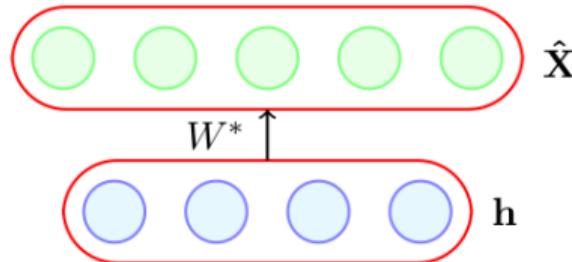


$$\mathbf{h} = g(\mathbf{W}\mathbf{X} + \mathbf{b})$$

$$\hat{\mathbf{X}} = f(\mathbf{W}^*\mathbf{h} + \mathbf{c})$$

- Can we do generation with autoencoders ?
- In other words, once the autoencoder is trained can I remove the encoder, feed a hidden representation h to the decoder and decode a \hat{X} from it ?
- In principle, yes! But in practice there is a problem with this approach
- h is a very high dimensional vector and only a few vectors in this space would actually correspond to meaningful latent representations of our input

AUTOENCODER – SUMMARY



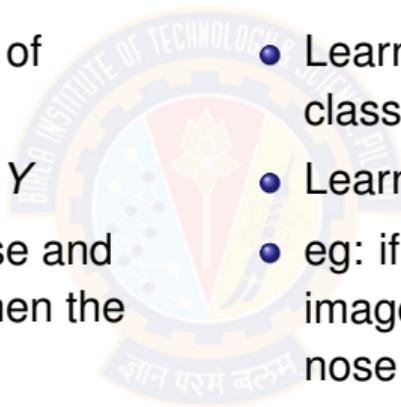
$$\hat{\mathbf{X}} = f(W^*\mathbf{h} + \mathbf{c})$$

- Ideally, we should only feed those values of h which are highly *likely*
- In other words, we are interested in sampling from $P(h|X)$ so that we pick only those h 's which have a high probability
- But unlike RBMs, autoencoders do not have such a probabilistic interpretation
- They learn a hidden representation h but not a distribution $P(h|X)$
- Similarly the decoder is also deterministic and does not learn a distribution over X (given a h we can get a X but not $P(X|h)$)

DISCRIMINATIVE MODELS VS GENERATIVE MODELS

DISCRIMINATIVE MODELS

- Set of features X and Set of categories Y
- Learn mapping from X to Y
- eg: if features are wet nose and tongue out and no purr; then the pet is a dog.
- Model $P(Y | X)$



GENERATIVE MODELS

- Learn realistic representation of a class
- Learn mapping from Y to X
- eg: if we have to generate a dog image, then it should have wet nose and tongue out.
- Model $P(X | Y)$

VARIATIONAL AUTOENCODER

- Generative model
- Variational autoencoder generates a latent representation and then uses this representation to generate new images.
- Data are assumed to be represented by a set of normally-distributed latent factors.
- Encoder generates parameters of these parameters namely μ and σ .
- Images are generated by sampling from these distributions.

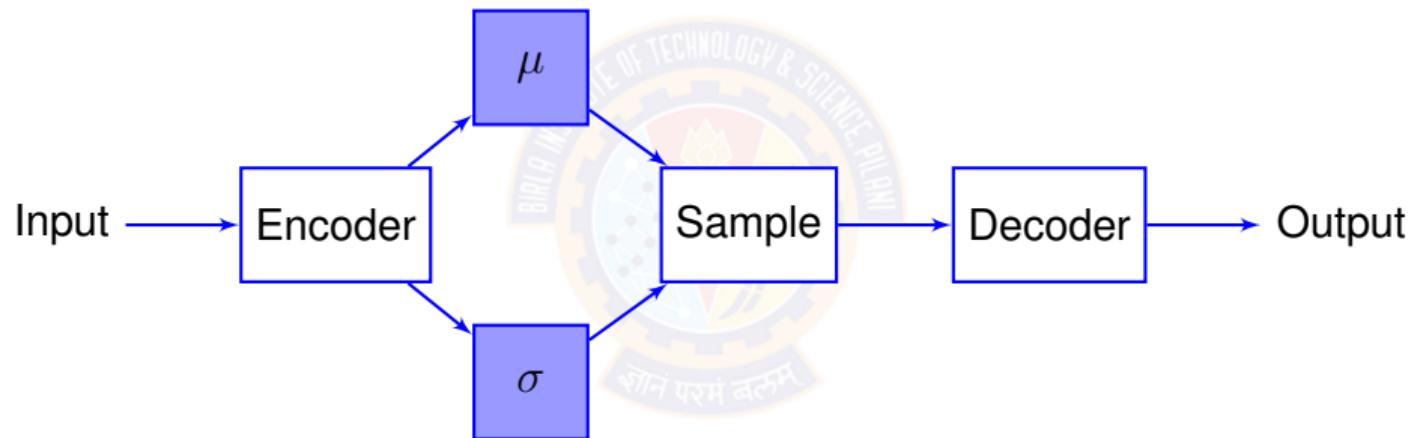
VARIATIONAL AUTOENCODER

- Primary Goal: Generate images using decoder.
- Latent vector: each element drawn from a normal distribution.
- Parameters are learned by the encoder.
- Secondary Goal: Similar images be close together in latent space.



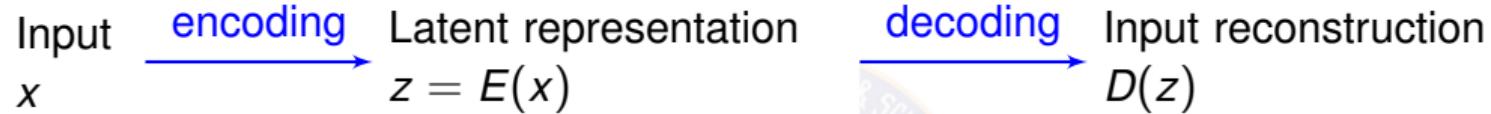
IBM

VARIATIONAL AUTOENCODER



AUTOENCODER VS VARIATIONAL AUTOENCODER

AUTOENCODER

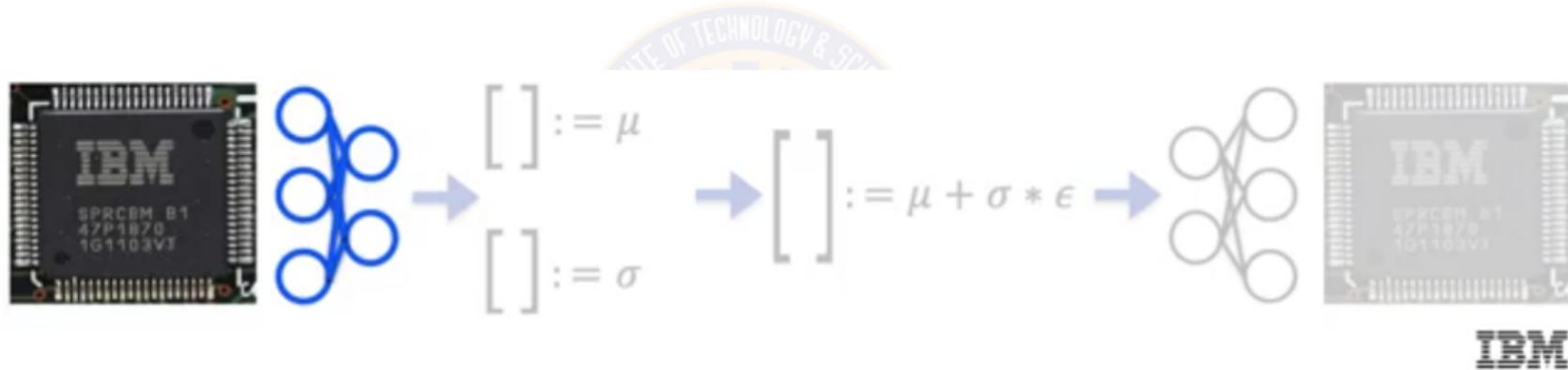


VARIATIONAL AUTOENCODER



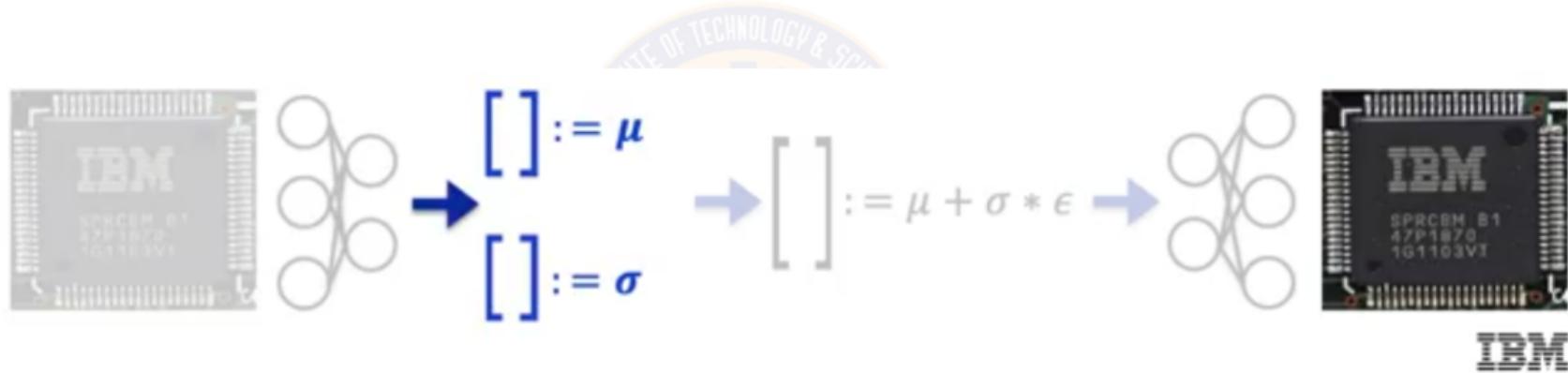
VARIATIONAL AUTOENCODER

- Step 1 : Image is fed through the encoder network.



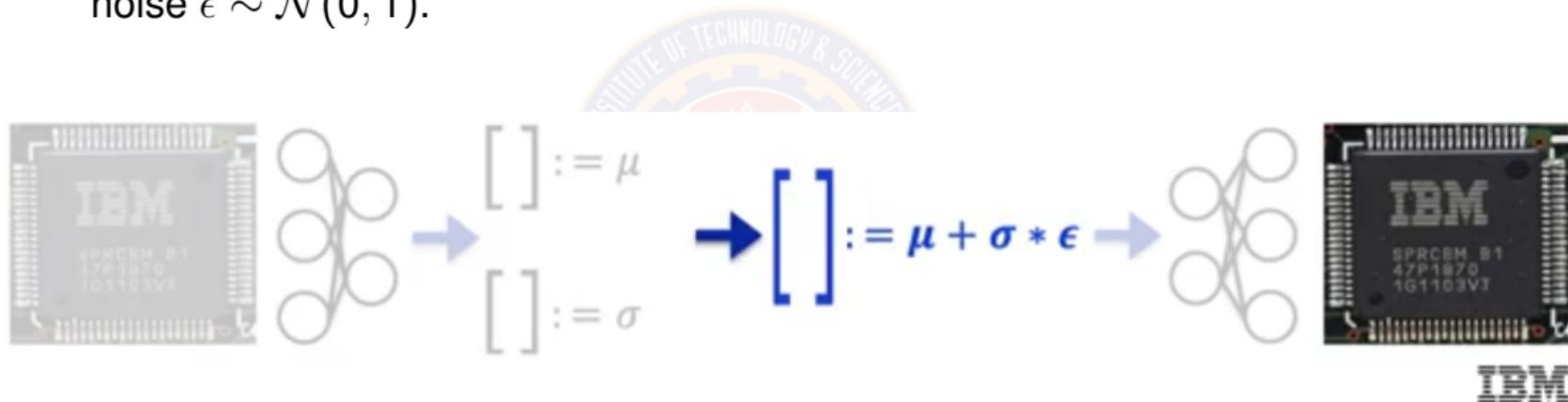
VARIATIONAL AUTOENCODER

- Step 2 : The encoder network learns the parameters μ and σ .



VARIATIONAL AUTOENCODER

- Step 3 : Combine the parameters μ and σ and normally distributed random noise $\epsilon \sim \mathcal{N}(0, 1)$.



VARIATIONAL AUTOENCODER EXAMPLE

Let $\mu = \begin{bmatrix} 0.7 \\ -0.6 \end{bmatrix}$

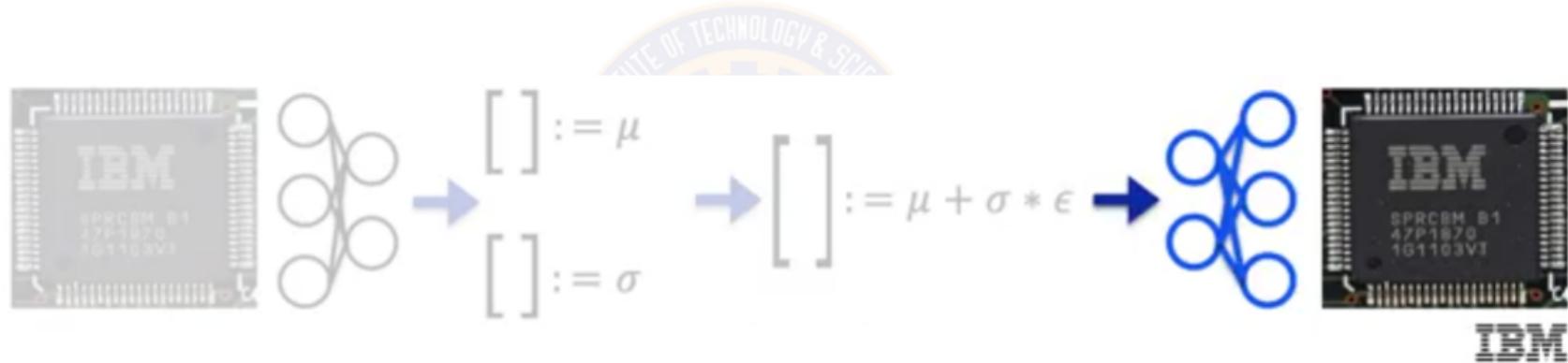
$$\sigma = \begin{bmatrix} 1.0 \\ 0.6 \end{bmatrix}$$

$$\epsilon = \begin{bmatrix} 1.4 \\ 1.9 \end{bmatrix}$$

$$\begin{aligned} Decoder_{in} &= \mu + \sigma * \epsilon \\ &= \begin{bmatrix} 0.7 \\ -0.6 \end{bmatrix} + \begin{bmatrix} 1.0 \\ 0.6 \end{bmatrix} * \begin{bmatrix} 1.4 \\ 1.9 \end{bmatrix} \\ &= \begin{bmatrix} 2.10 \\ 0.54 \end{bmatrix} \end{aligned}$$

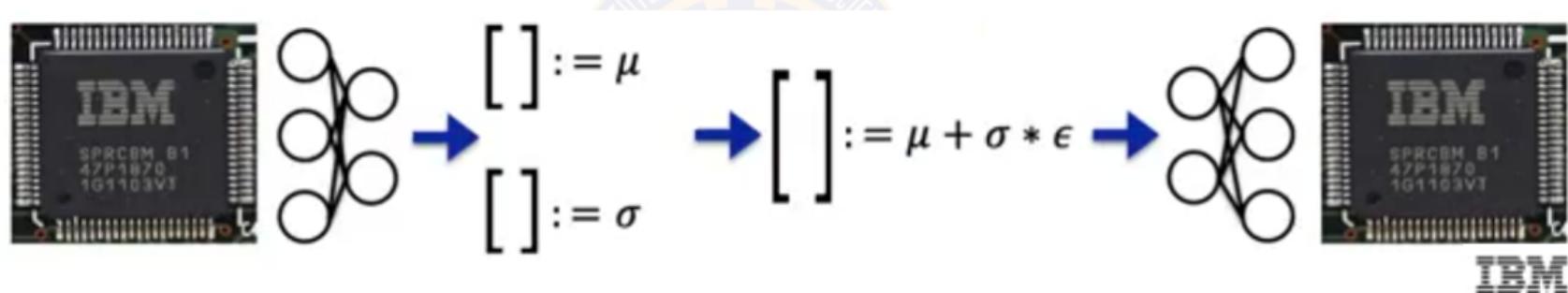
VARIATIONAL AUTOENCODER

- Step 4 : Feed the vector through the decoder.



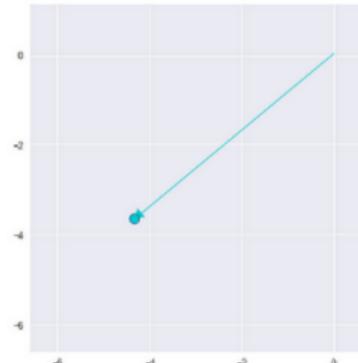
VARIATIONAL AUTOENCODER

- Step 5 : The reconstructed image is generated using the decoder.

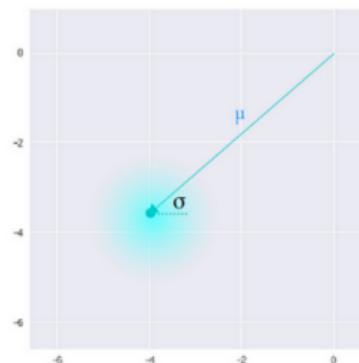


VARIATIONAL AUTOENCODER

- Encoder learns two encodings or outputs two vectors of size n
 - a vector of means μ – controls where the encoding of an input should be centered around.
 - a vector of standard deviations σ – controls (area) how much from the mean the encoding can vary.
- Decoder learns that a single point in latent space refers to a sample of that class. It also learns that all nearby points refer to the same class as well.



Standard Autoencoder
(direct encoding coordinates)



Variational Autoencoder
(μ and σ initialize a probability distribution)

VARIATIONAL AUTOENCODER – LOSS FUNCTION

- VAE should learn to reconstruct the original images.
- The loss function will have 2 components.
 - ① A penalty for NOT reconstructing the image correctly.
 - ★ Pixel-wise difference between the reconstructed image and the original image.
 - ★ Use Mean squared error (MSE)
 - ★ Use binary cross entropy
 - ② A penalty for generating parameter vectors μ and σ that are different than 0 and 1, respectively.
 - ★ The difference between vectors produced by the encoder and the standard normal distribution parameters $\mathcal{N}(0, 1)$.
 - ★ Use Kullback–Leibler divergence (KL divergence).

KULLBACK–LEIBLER DIVERGENCE

- The KL divergence between two probability distributions measures how much they diverge from each other.
- Minimizing the KL divergence means optimizing the probability distribution parameters (μ and σ) to closely resemble that of the target distribution.
- KL loss is equivalent to the sum of all the KL divergences between the component $X_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$ in X , and the standard normal $\mathcal{N}(\mu = 0, \sigma = 1)$

$$KL\ loss = \sum_{i=1}^n \sigma_i^2 + \mu_i^2 - \log(\sigma_i) - 1$$

KL DIVERGENCE EXAMPLE

Compare $\mu = \begin{bmatrix} 0.7 \\ -0.6 \end{bmatrix} \longleftrightarrow \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

Compare $\log(\sigma) = \begin{bmatrix} 1.0 \\ 0.6 \end{bmatrix} \longleftrightarrow \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

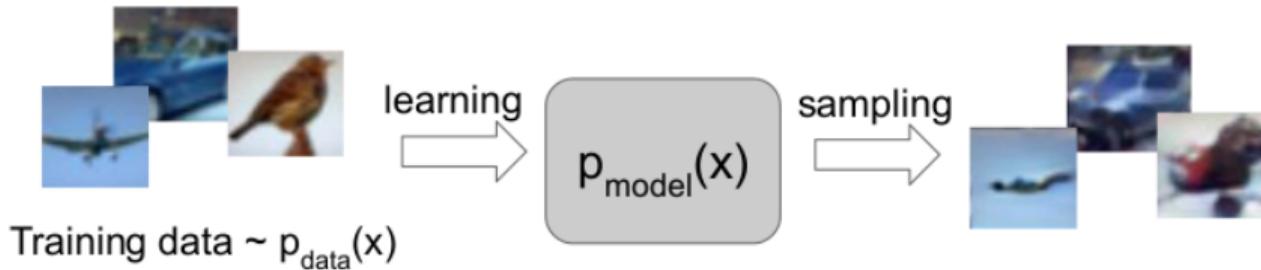
$\log(\sigma)$ compared to 0 since $\log(1) = 0$.

σ is interpreted as $\log(\sigma)$ since variance cannot be negative.

$$\text{KL Loss} = \frac{1}{2} * (e^{\log(\sigma)} - \log(\sigma) - 1 + \mu^2)$$

GENERATIVE MODEL

Given training data, generate new samples from same distribution

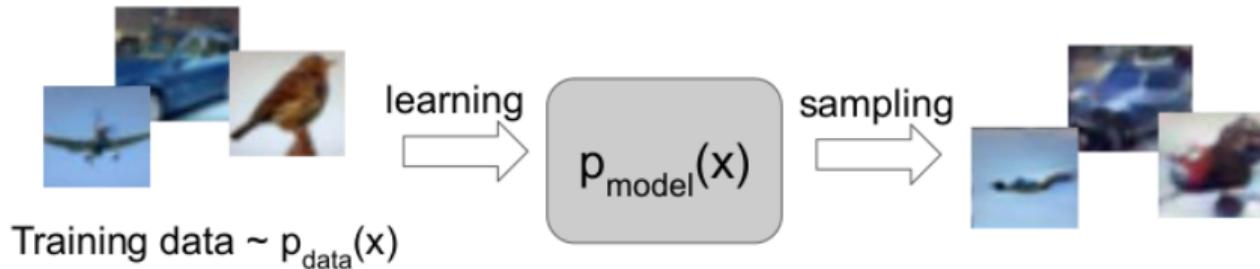


Objectives:

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

GENERATIVE MODEL

Given training data, generate new samples from same distribution



Formulate as density estimation problems:

- **Explicit density estimation:** explicitly define and solve for $p_{\text{model}}(x)$
- **Implicit density estimation:** learn model that can sample from $p_{\text{model}}(x)$ **without explicitly defining it.**

GENERATIVE MODEL

Taxonomy of Generative Models

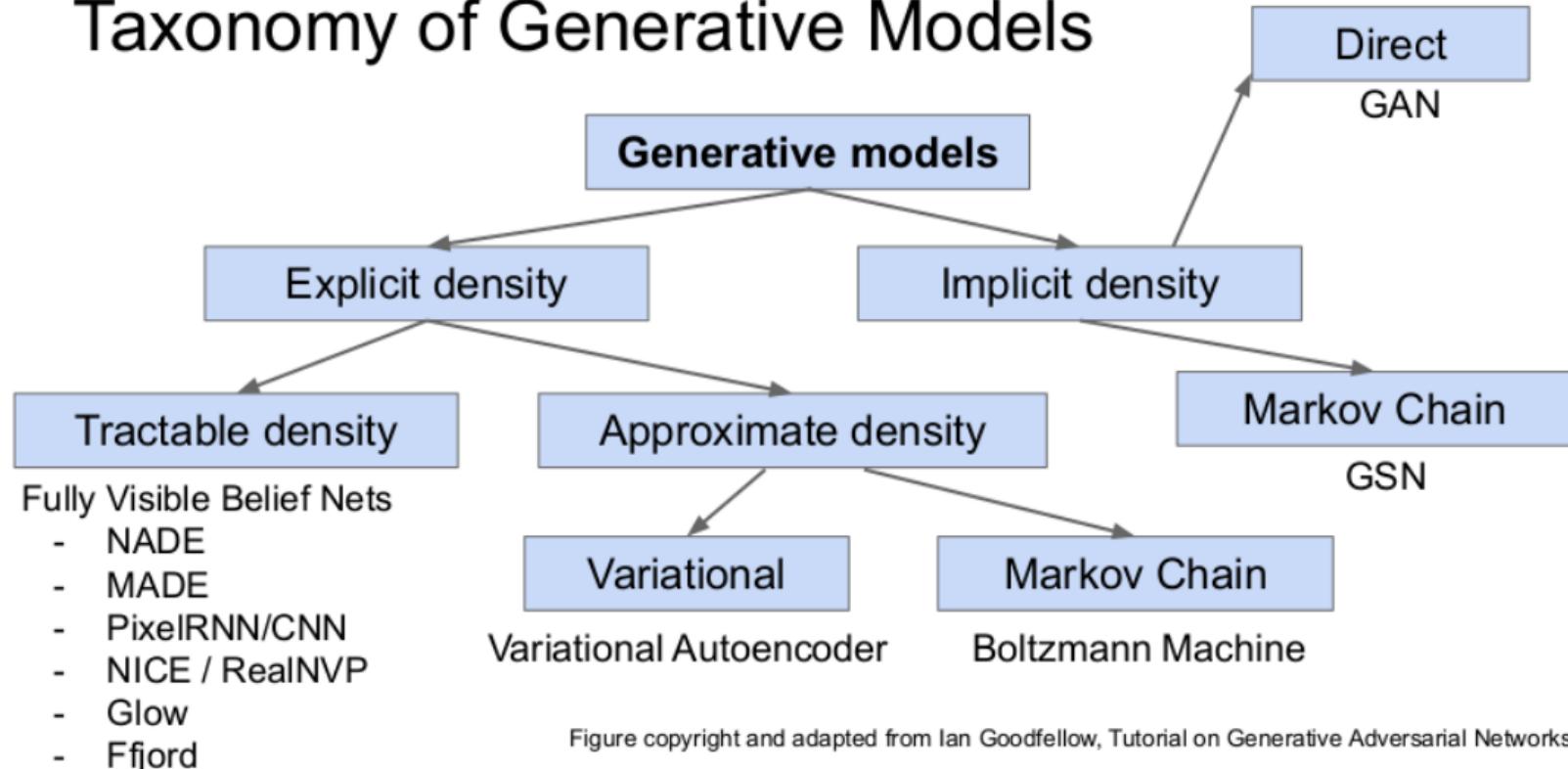


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

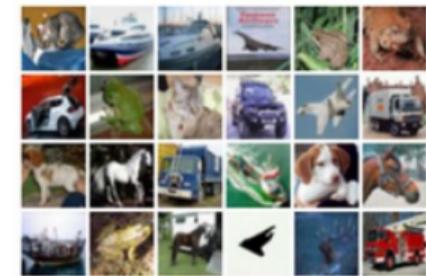
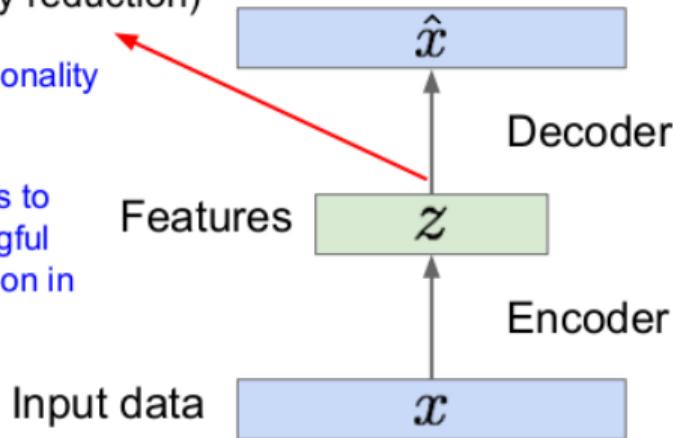
AUTOENCODERS – RECAP

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

z usually smaller than x
(dimensionality reduction)

Q: Why dimensionality reduction?

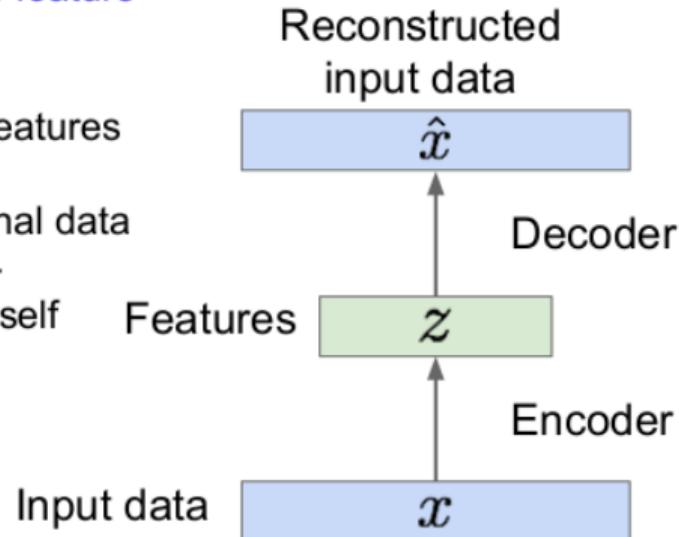
A: Want features to capture meaningful factors of variation in data



AUTOENCODERS – RECAP

How to learn this feature representation?

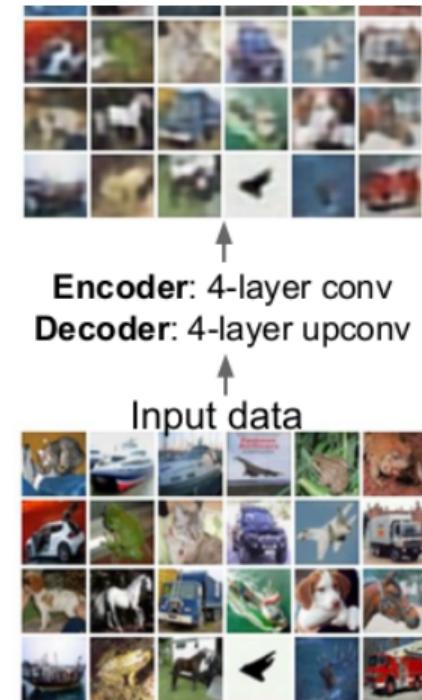
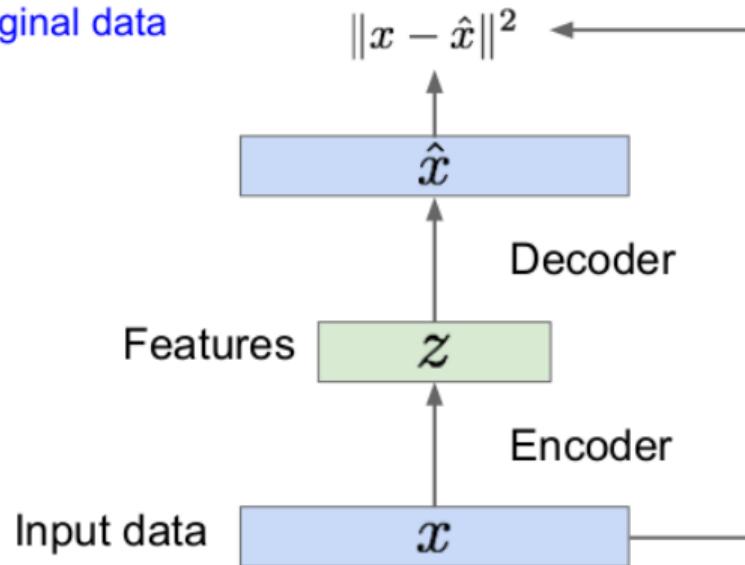
Train such that features can be used to reconstruct original data
“Autoencoding” - encoding input itself



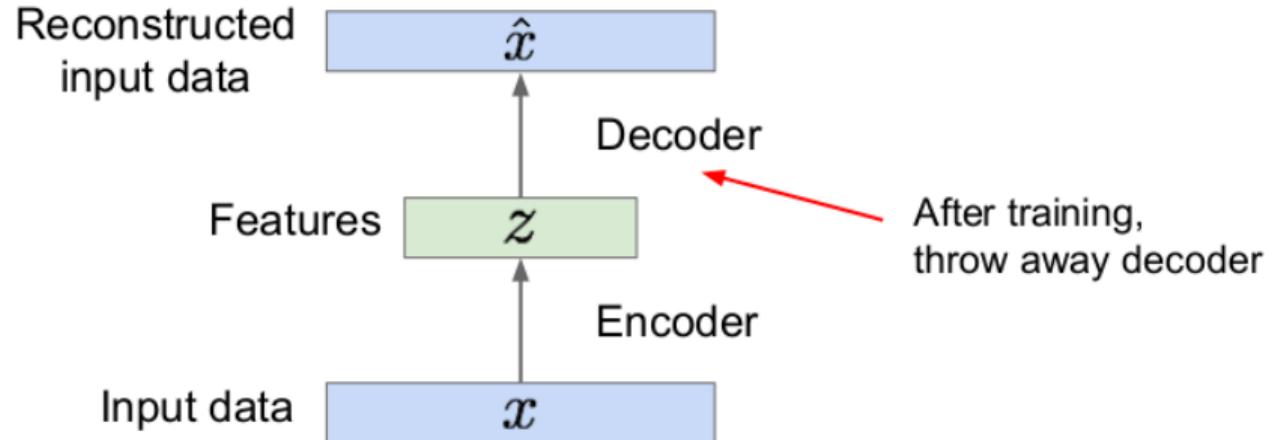
AUTOENCODERS – RECAP

Train such that features can be used to reconstruct original data

L2 Loss function:
Doesn't use labels!



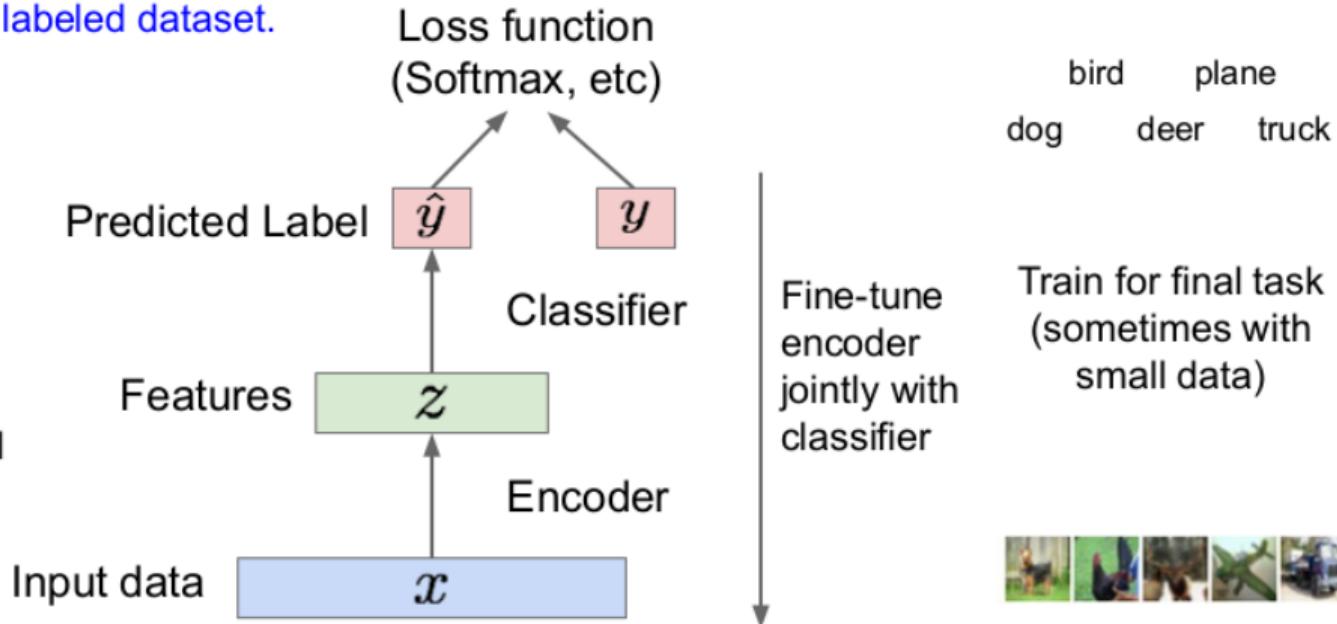
AUTOENCODERS – RECAP



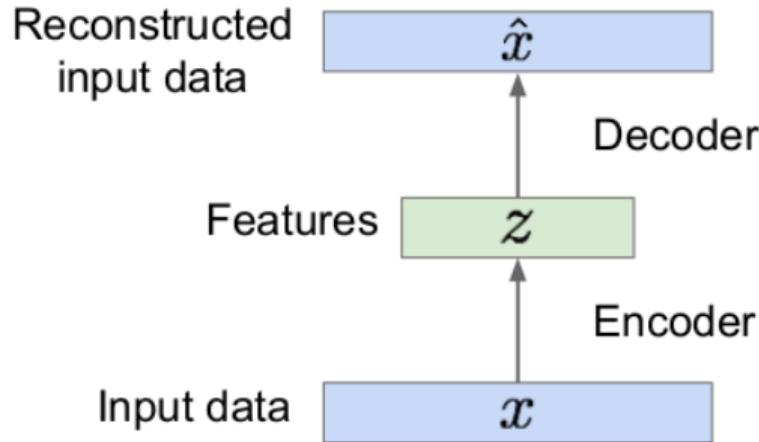
AUTOENCODERS – RECAP

Transfer from large, unlabeled dataset to small, labeled dataset.

Encoder can be used to initialize a **supervised** model



AUTOENCODERS – RECAP



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

Features capture factors of variation in training data.

But we can't generate new images from an autoencoder because we don't know the space of z .

How do we make autoencoder a **generative model**?

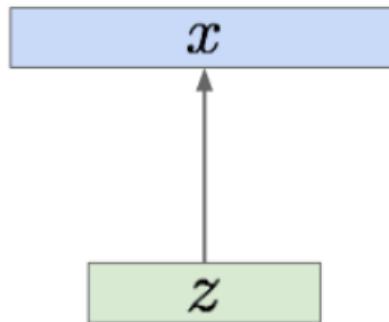
VARIATIONAL AUTOENCODERS

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 z

Sample from
true conditional

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



Sample from
true prior

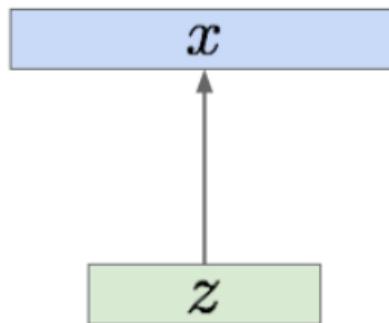
$$z^{(i)} \sim p_{\theta^*}(z)$$

VARIATIONAL AUTOENCODERS

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 \mathbf{z}

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



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

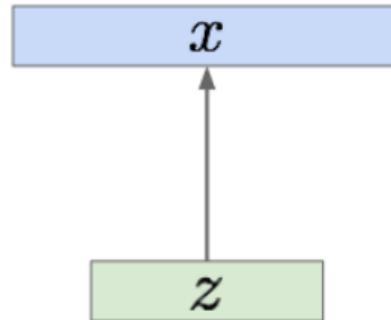
Intuition (remember from autoencoders!):
 \mathbf{x} is an image, \mathbf{z} is latent factors used to generate \mathbf{x} : attributes, orientation, etc.

VARIATIONAL AUTOENCODERS

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

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

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

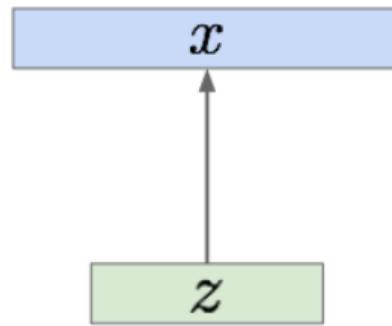


Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

VARIATIONAL AUTOENCODERS

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

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



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

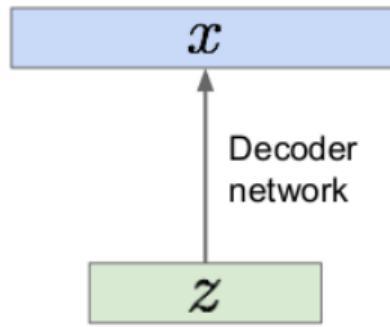
How should we represent this model?

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

VARIATIONAL AUTOENCODERS

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

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



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

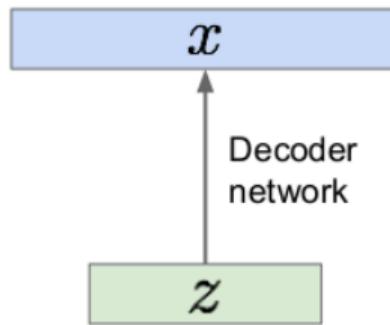
$$p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$$

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

VARIATIONAL AUTOENCODERS

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

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



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

VARIATIONAL AUTOENCODERS: INTRACTABILITY

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

↑
Simple Gaussian prior

VARIATIONAL AUTOENCODERS: INTRACTABILITY

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

The diagram illustrates the data likelihood equation for a variational autoencoder. The equation is $p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$. Two green checkmarks are placed above the integral sign. A blue arrow points from the text "Decoder neural network" to the term $p_{\theta}(x|z)$, indicating that this term is generated by the 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: $p_\theta(x) = \int p_\theta(z)p_\theta(x|z)dz$

Posterior density: $p_\theta(z|x) = p_\theta(x|z)p_\theta(z)/p_\theta(x)$

Intractable data likelihood

VARIATIONAL AUTOENCODERS

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

Posterior density also intractable: $p_{\theta}(z|x) = p_{\theta}(x|z)p_{\theta}(z)/p_{\theta}(x)$

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

Will see that the approximate posterior allows us to derive a lower bound on the data likelihood that is tractable, which we can optimize.

Variational inference is to approximate the unknown posterior distribution from only the observed data x

VARIATIONAL AUTOENCODERS

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



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



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



$p_\theta(z|x)$ intractable (saw earlier), can't compute this KL term :(But we know KL divergence always ≥ 0 .

VARIATIONAL AUTOENCODERS

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

We want to
maximize the
data
likelihood

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

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

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

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

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

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

$p_{\theta}(z|x)$ intractable (saw
earlier), can't compute this KL
term :(But we know KL
divergence always ≥ 0 .

VARIATIONAL AUTOENCODERS

We want to maximize the data likelihood

$$\begin{aligned}\log p_\theta(x^{(i)}) &= \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \underbrace{\mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))}_{\geq 0}\end{aligned}$$

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

VARIATIONAL AUTOENCODERS

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

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

Decoder:
reconstruct
the input data

Encoder:
make approximate
posterior distribution
close to prior

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

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

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

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

VARIATIONAL AUTOENCODERS

Putting it all together: maximizing the likelihood lower bound

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

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

Have analytical solution

Make approximate posterior distribution close to prior

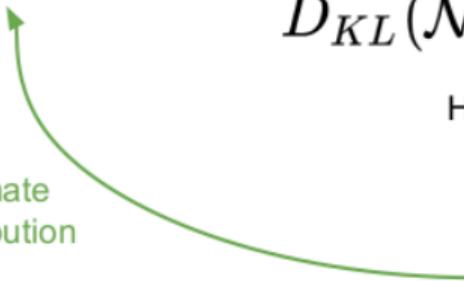
Encoder network

$$q_\phi(z|x)$$

Input Data

$$\mu_{z|x}$$

$$\Sigma_{z|x}$$



VARIATIONAL AUTOENCODERS

Putting it all together: maximizing the likelihood lower bound

$$\mathcal{L}(x^{(i)}, \theta, \phi) = \mathbb{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))$$

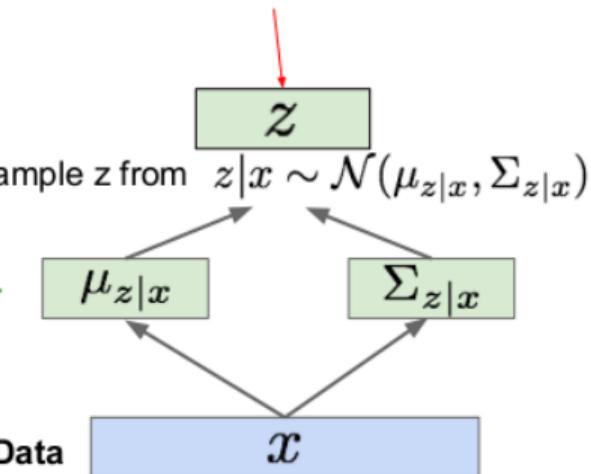
Make approximate posterior distribution close to prior

Encoder network

$$q_\phi(z|x)$$

Input Data

Not part of the computation graph!



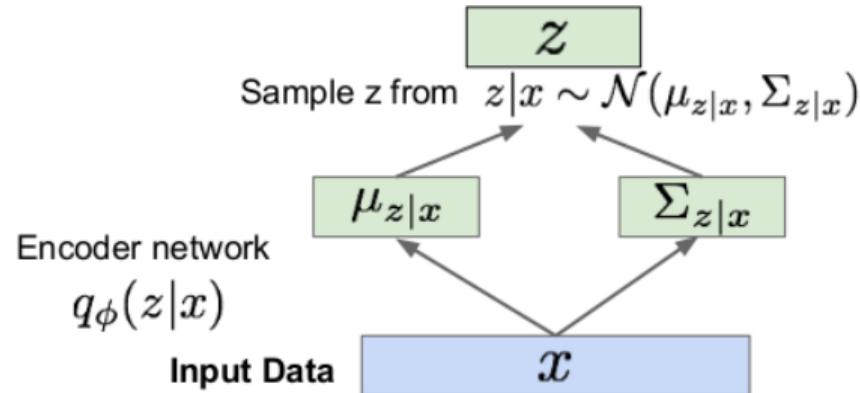
VARIATIONAL AUTOENCODERS

Putting it all together: maximizing the likelihood lower bound

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

$$\text{Sample } \epsilon \sim \mathcal{N}(0, I)$$

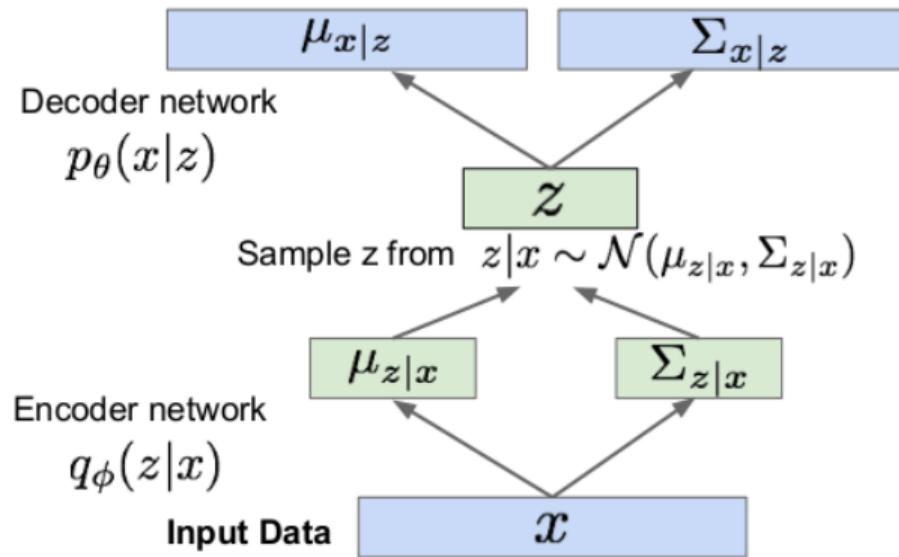
$$z = \mu_{z|x} + \epsilon \sigma_{z|x}$$



VARIATIONAL AUTOENCODERS

Putting it all together: maximizing the likelihood lower bound

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

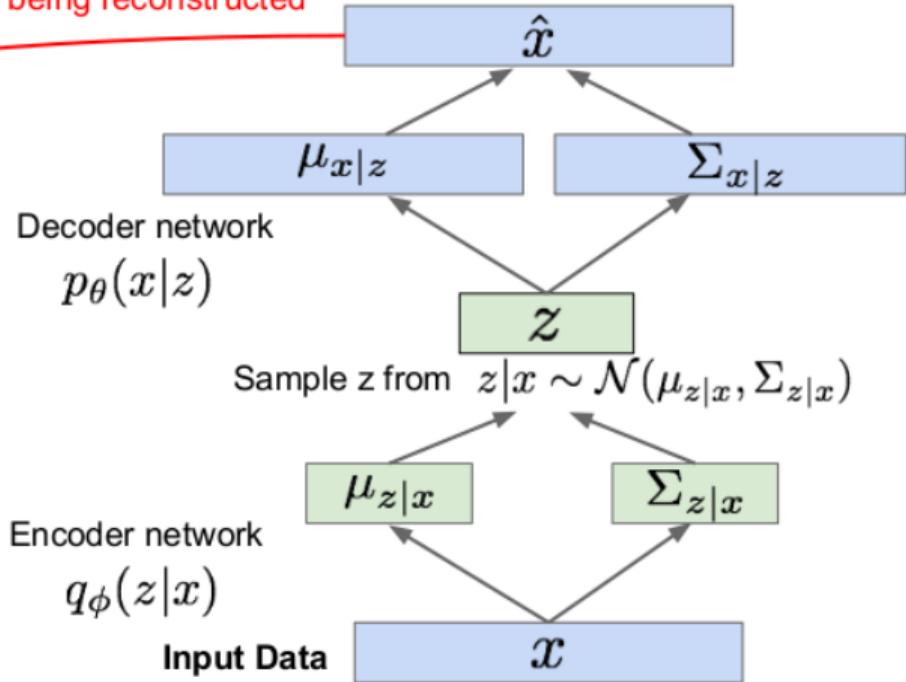


VARIATIONAL AUTOENCODERS

Putting it all together: maximizing the likelihood lower bound

$$\mathcal{L}(x^{(i)}, \theta, \phi) = \mathbb{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))$$

Maximize likelihood of original input being reconstructed



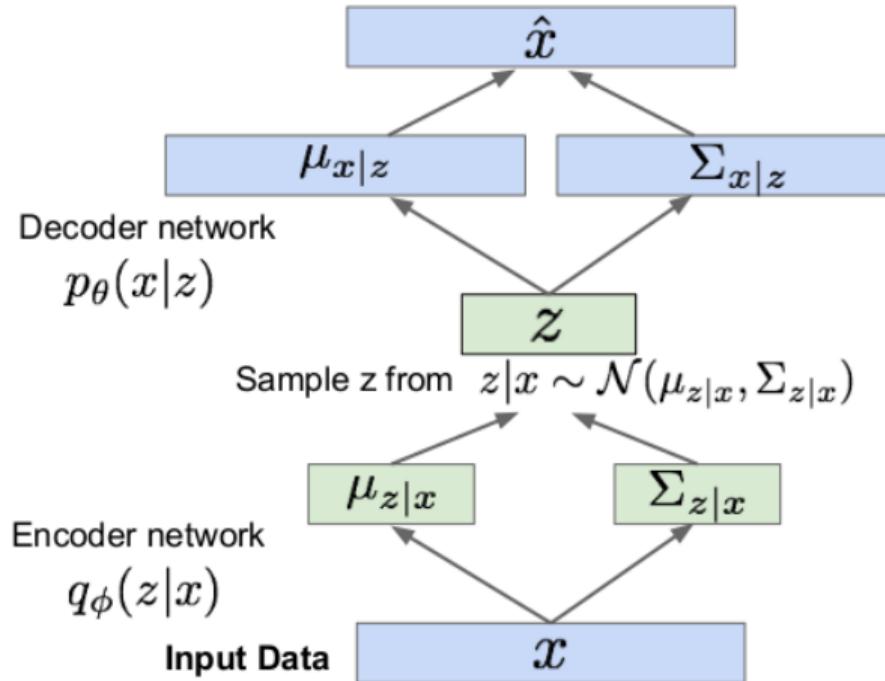
VARIATIONAL AUTOENCODERS

Putting it all together: maximizing the likelihood lower bound

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

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

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



VARIATIONAL AUTOENCODERS: GENERATING DATA

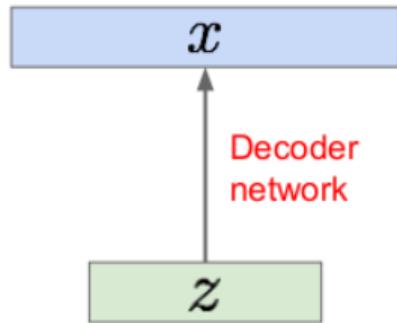
Our assumption about data generation process

Sample from true conditional

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

Sample from true prior

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



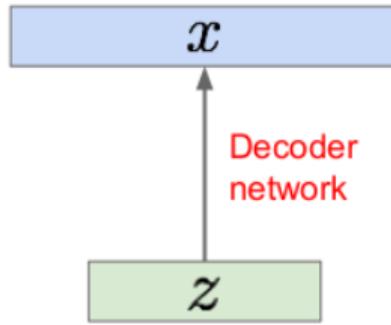
Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

VARIATIONAL AUTOENCODERS: GENERATING DATA

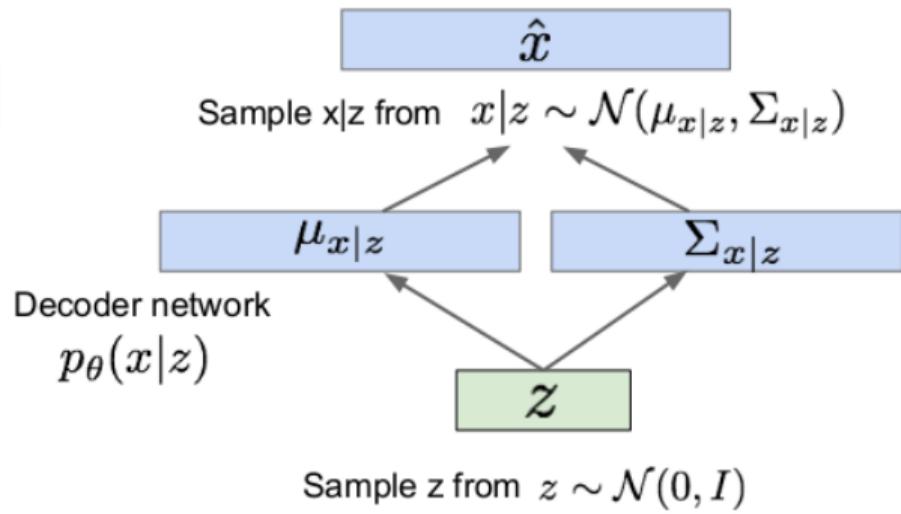
Our assumption about data generation process

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

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



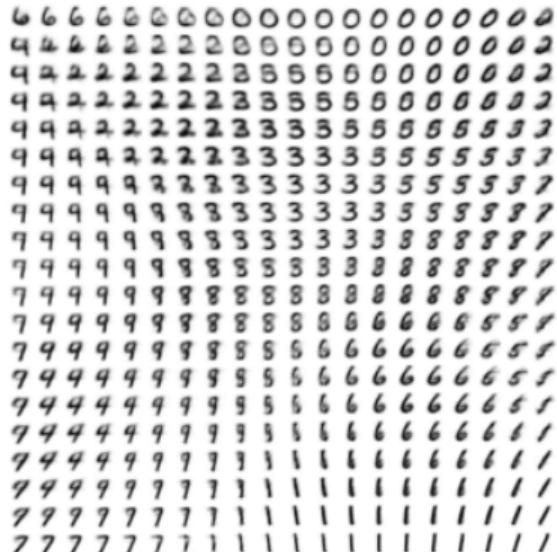
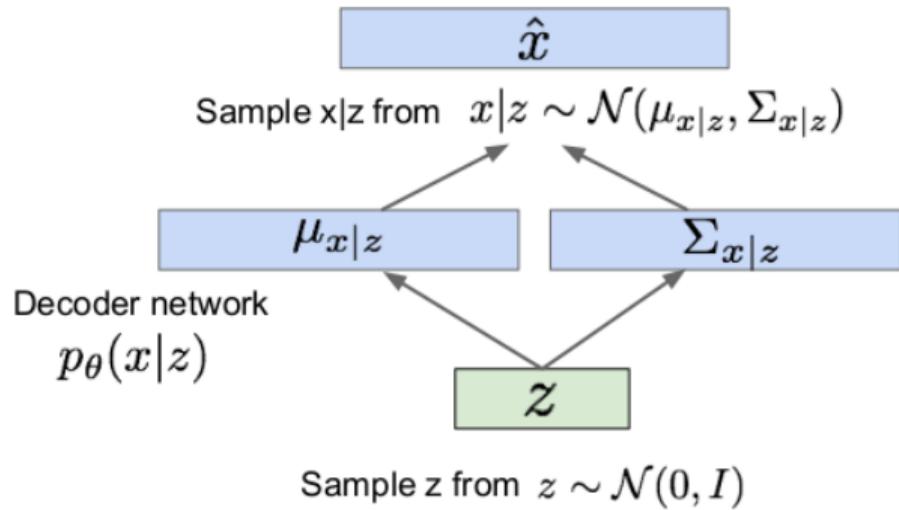
Now given a trained VAE:
use decoder network & sample z from prior!



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

VARIATIONAL AUTOENCODERS: GENERATING DATA

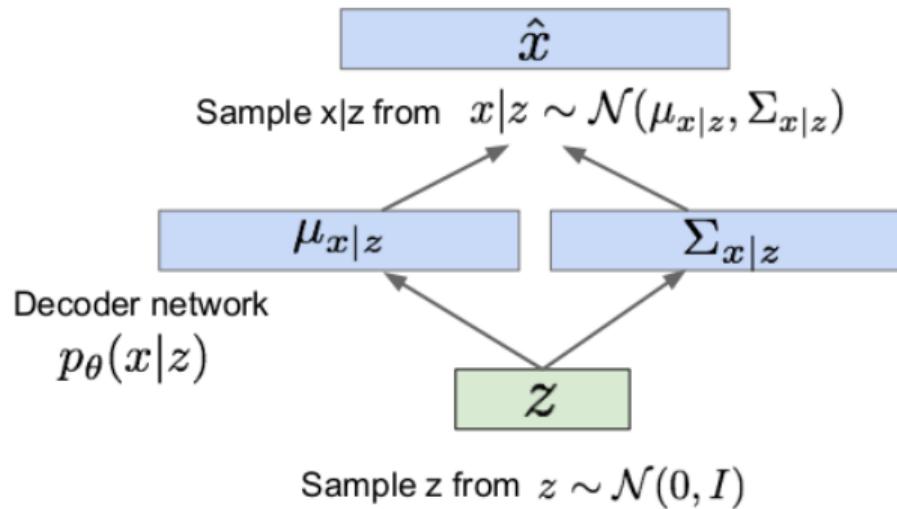
Use decoder network. Now sample z from prior!



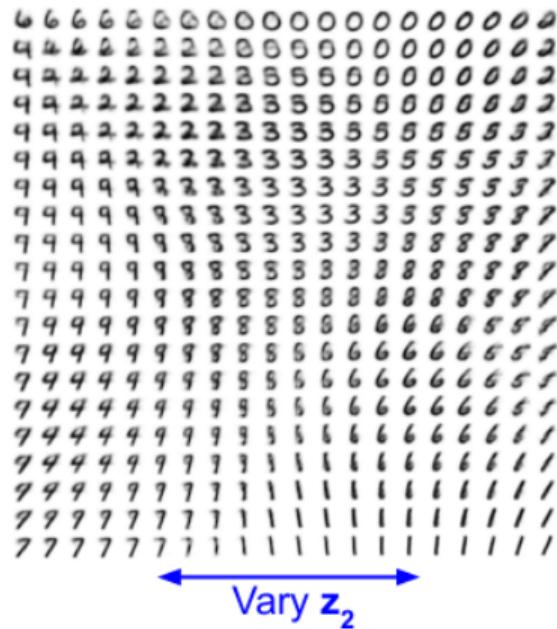
Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

VARIATIONAL AUTOENCODERS: GENERATING DATA

Use decoder network. Now sample z from prior!



Data manifold for 2-d \mathbf{z}



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

VARIATIONAL AUTOENCODERS: GENERATING DATA

Diagonal prior on \mathbf{z}
=> independent
latent variables

Different
dimensions of \mathbf{z}
encode
interpretable factors
of variation

Degree of smile
 \uparrow
Vary \mathbf{z}_1
 \downarrow



\leftarrow Vary \mathbf{z}_2 \rightarrow Head pose

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

VARIATIONAL AUTOENCODERS

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
- Interpretable latent space.
- Allows inference of $q(z|x)$, can be useful feature representation for other tasks

Cons:

- Maximizes lower bound of likelihood: okay, but not as good evaluation as PixelRNN/PixelCNN
- Samples blurrier and lower quality compared to state-of-the-art (GANs)

Active areas of research:

- More flexible approximations, e.g. richer approximate posterior instead of diagonal Gaussian, e.g., Gaussian Mixture Models (GMMs), Categorical Distributions.
- Learning disentangled representations.

References

- ① Deep Learning by Ian Goodfellow, Yoshua Bengio, Aaron Courville
<https://www.deeplearningbook.org/>
- ② Deep Learning with Python by Francois Chollet.
<https://livebook.manning.com/book/deep-learning-with-python/>

Thank You!