

Generative Models. VAE

Evgeny Burnaev

Skoltech, Moscow, Russia



- 1 Unsupervised Learning
- 2 Autoencoders
- 3 Variational Autoencoders

1 Unsupervised Learning

2 Autoencoders

3 Variational Autoencoders

- **Supervised Learning**
- **Data:** (x, y) , x is a feature vector, y is a label
- **Goal:** Learn a function to map $x \rightarrow y$
- **Examples:** Classification, regression, object detection, semantic segmentation, image captioning, etc.



→ Cat

Classification

- **Supervised Learning**
- **Data:** (x, y) , x is a feature vector, y is a label
- **Goal:** Learn a function to map $x \rightarrow y$
- **Examples:** Classification, regression, object detection, semantic segmentation, image captioning, etc.



DOG, DOG, CAT

Object Detection

- **Supervised Learning**
- **Data:** (x, y) , x is a feature vector, y is a label
- **Goal:** Learn a function to map $x \rightarrow y$
- **Examples:** Classification, regression, object detection, semantic segmentation, image captioning, etc.



GRASS, CAT,
TREE, SKY

Semantic Segmentation

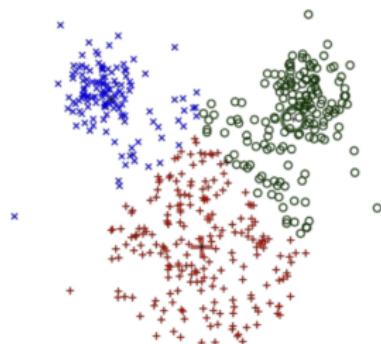
- **Supervised Learning**
- **Data:** (x, y) , x is a feature vector, y is a label
- **Goal:** Learn a function to map $x \rightarrow y$
- **Examples:** Classification, regression, object detection, semantic segmentation, image captioning, etc.



A cat sitting on a suitcase on the floor

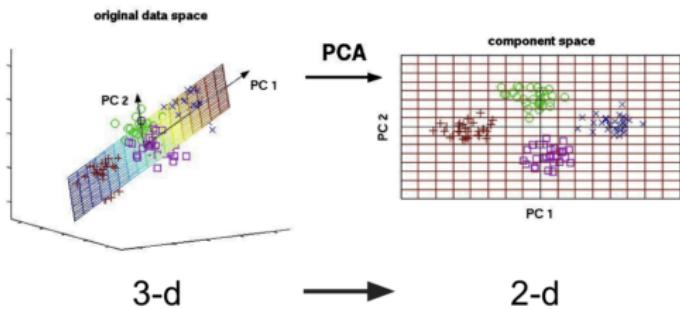
Image captioning

- **Unsupervised Learning**
- **Data:** x , x is a feature vector
- **Goal:** Learn some underlying hidden structure of the data
- **Examples:** Clustering, dimensionality reduction, feature learning, density estimation, etc.



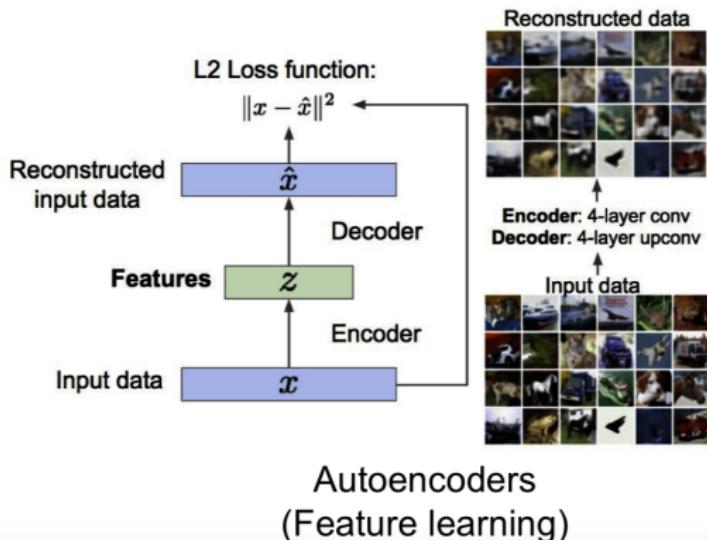
K-means clustering

- **Unsupervised Learning**
- **Data:** x , x is a feature vector
- **Goal:** Learn some underlying hidden structure of the data
- **Examples:** Clustering, dimensionality reduction, feature learning, density estimation, etc.



Principal Component Analysis
(Dimensionality reduction)

- **Unsupervised Learning**
- **Data:** x , x is a feature vector
- **Goal:** Learn some underlying hidden structure of the data
- **Examples:** Clustering, dimensionality reduction, feature learning, density estimation, etc.

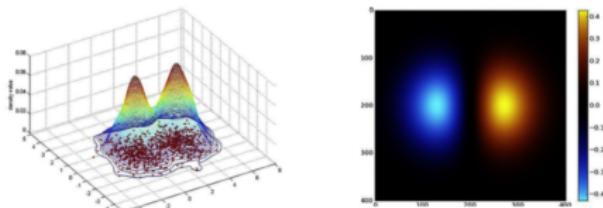


- **Unsupervised Learning**
- **Data:** x , x is a feature vector
- **Goal:** Learn some underlying hidden structure of the data
- **Examples:** Clustering, dimensionality reduction, feature learning, density estimation, etc.



Figure copyright Ian Goodfellow, 2016. Reproduced with permission.

1-d density estimation



2-d density estimation

- **Supervised Learning**

- **Data:** (x, y) , x is a feature vector, y is a label
- **Goal:** Learn a function to map $x \rightarrow y$
- **Examples:** Classification, regression, object detection, semantic segmentation, image captioning, etc.

Feature training data is cheap!

Holy grail: solve unsupervised learning \Rightarrow understand structure of visual world

- **Unsupervised Learning**

- **Data:** x , x is a feature vector
- **Goal:** Learn some underlying hidden structure of the data
- **Examples:** Clustering, dimensionality reduction, feature learning, density estimation, etc.

Given training data, generate new samples from same distribution



Training data $\sim p_{data}(\mathbf{x})$

Generated data $\sim p_{model}(\mathbf{x})$

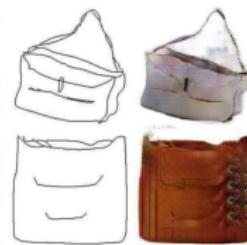
Want to learn $p_{model}(\mathbf{x})$ similar to $p_{data}(\mathbf{x})$

Addresses density estimation, a core problem in unsupervised learning.

Issues:

- Explicit density estimation: explicitly define and solve for $p_{model}(\mathbf{x})$
- Implicit density estimation: learn model that can sample from $p_{model}(\mathbf{x})$ w/o explicitly defining it

- Realistic samples for artwork, super-resolution, colorization, etc.



- Generative models of time-series data can be used for simulation and planning (reinforcement learning applications!)
- Training generative models can also enable inference of latent representations that can be useful as general features

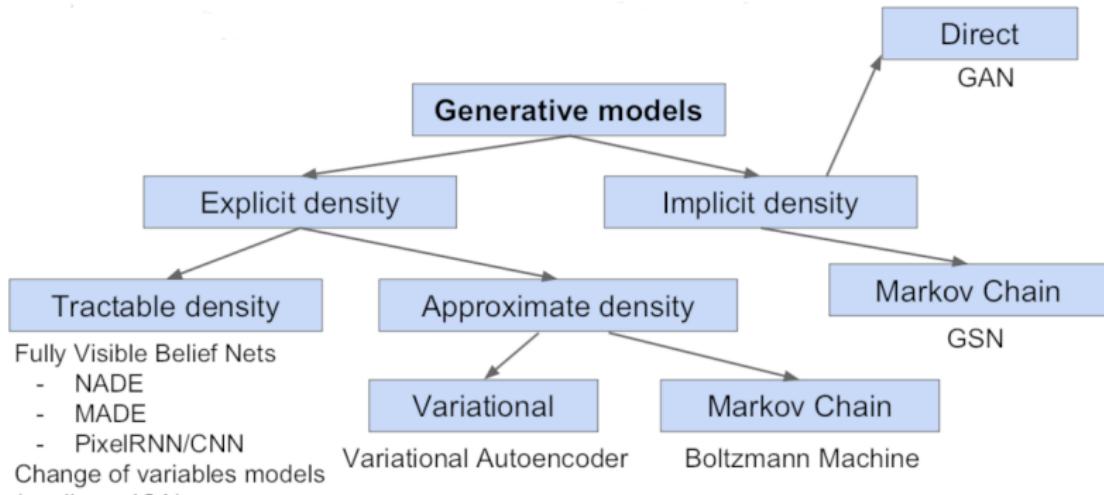


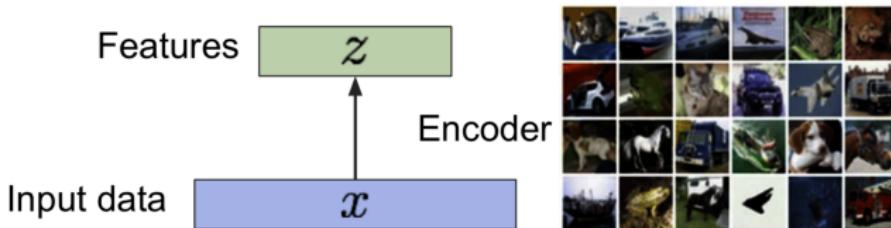
Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

1 Unsupervised Learning

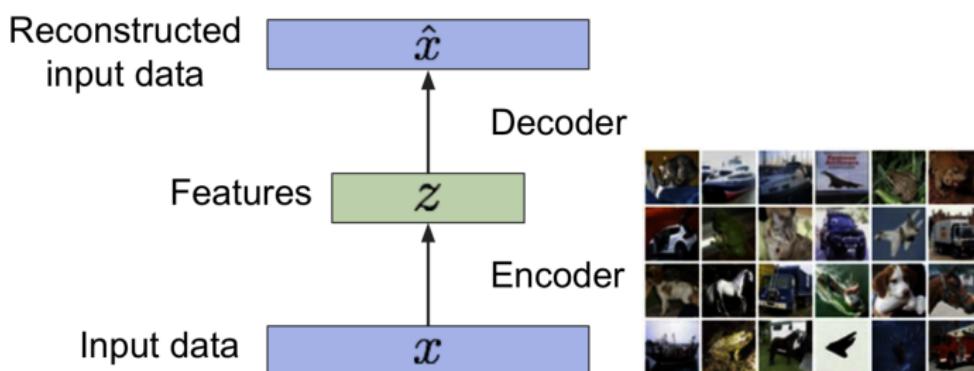
2 Autoencoders

3 Variational Autoencoders

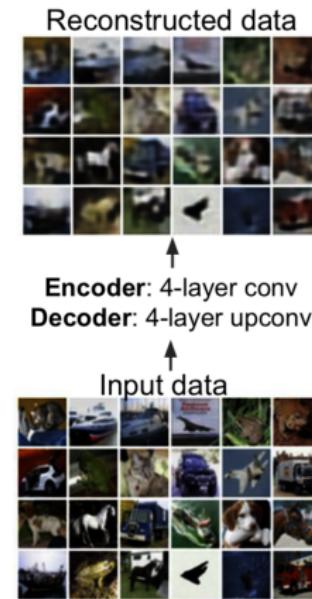
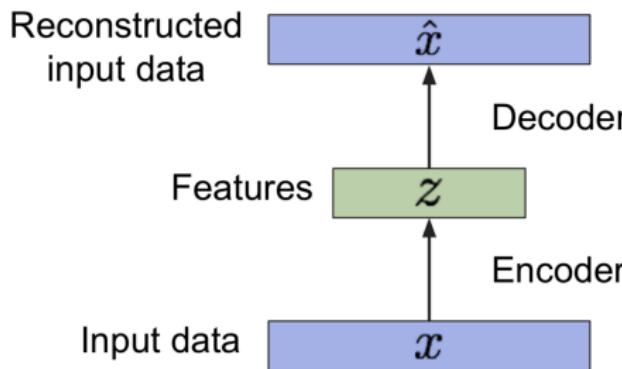
- Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data
- **Encoder:** Linear + nonlinearity (sigmoid) → Deep fully-connected → ReLU CNN
- z usually smaller than x (dimensionality reduction)
- Q: Why dimensionality reduction?
- A: Want features to capture meaningful factors of variation in data



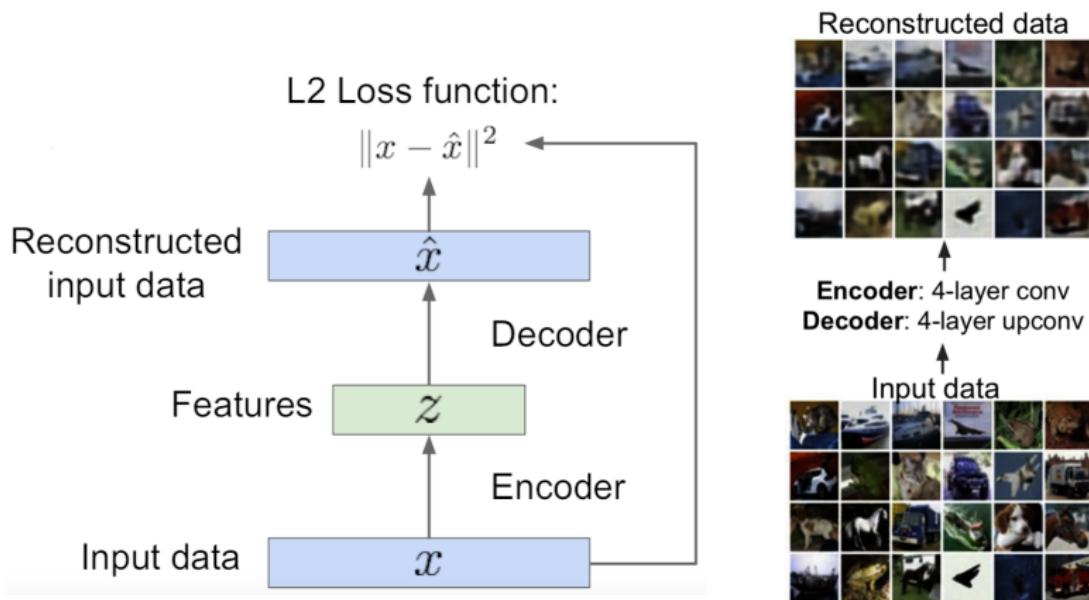
- How to learn this feature representation?
- Train such that features can be used to reconstruct original data
“Autoencoding” – encoding itself
- **Decoder:** Linear + nonlinearity (sigmoid) → Deep fully-connected
→ ReLU CNN



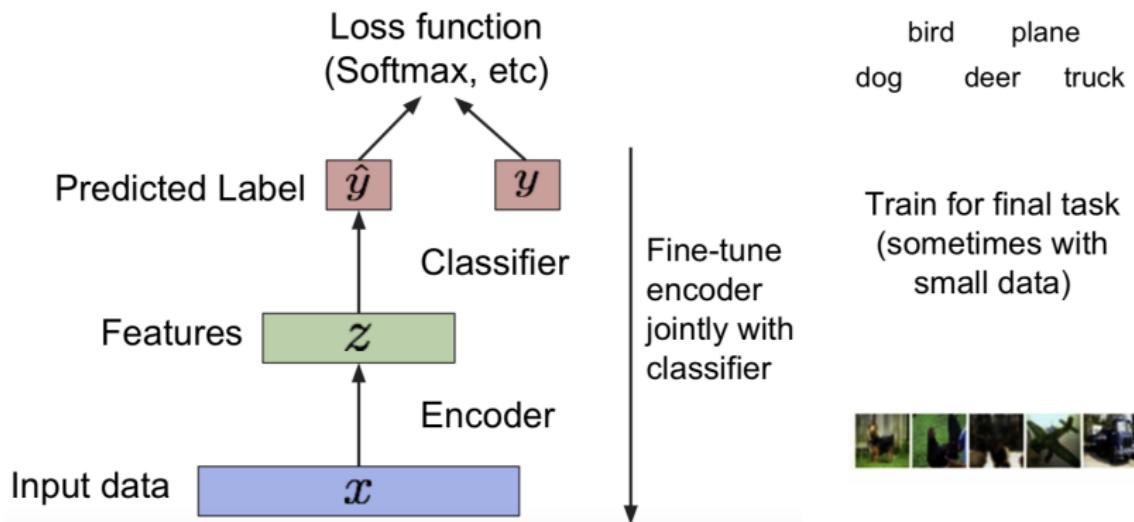
- How to learn this feature representation?
- Train such that features can be used to reconstruct original data
“Autoencoding” – encoding itself
- **Decoder:** Linear + nonlinearity (sigmoid) → Deep fully-connected
→ ReLU CNN



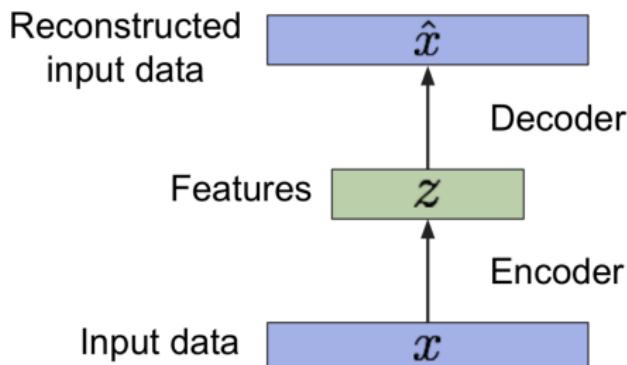
- Train such that features can be used to reconstruct original data
- Doesn't use labels!
- After training we can throw away decoder and ...



- Encoder can be used to initialize a supervised model



- Autoencoders can reconstruct data, and can learn features to initialize a supervised model
- Features capture factors of variation in training data. Can we generate new images from an autoencoder?



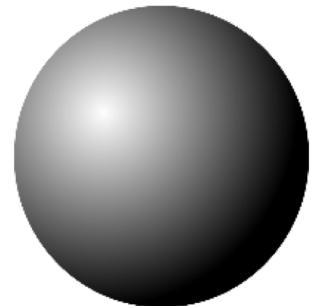
1 Unsupervised Learning

2 Autoencoders

3 Variational Autoencoders

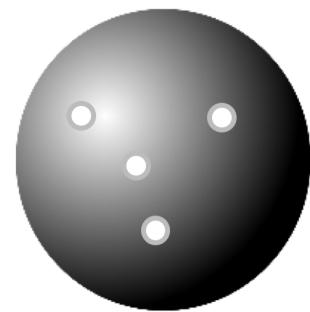
Probabilistic model for Data on manifolds

$$z \sim p(z)$$



Probabilistic model for Data on manifolds

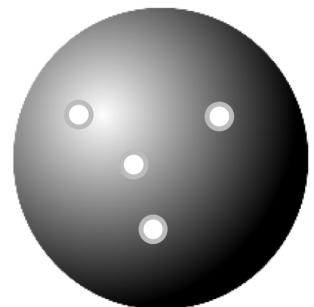
$$z \sim p(z)$$



$$z_1, z_2, \dots, z_n$$

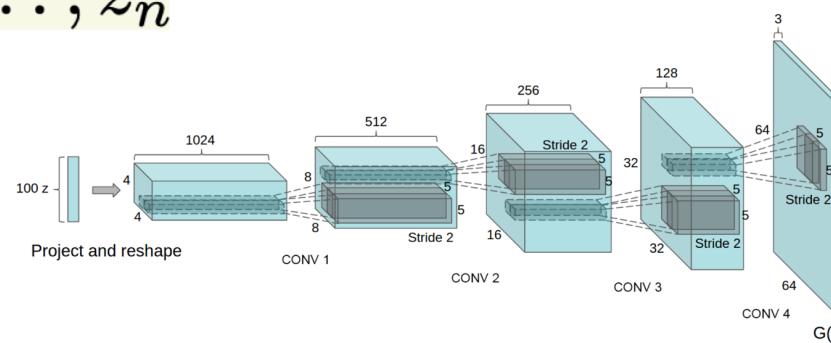
Probabilistic model for Data on manifolds

$$z \sim p(z)$$



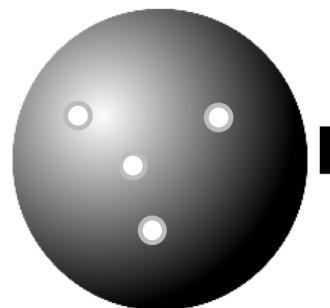
$$g_{\theta}$$

$$z_1, z_2, \dots, z_n$$



Probabilistic model for Data on manifolds

$$z \sim p(z)$$

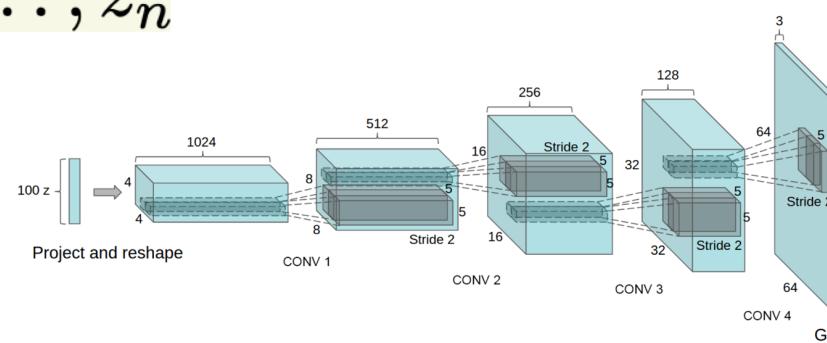


$$g_{\theta}$$



$$z_1, z_2, \dots, z_n$$

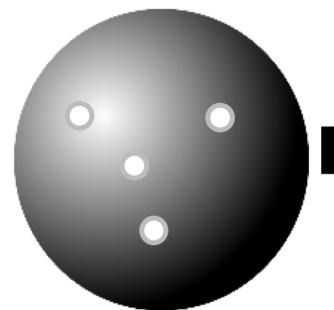
$$x_1, x_2, \dots, x_n$$



Probabilistic model for Data on manifolds

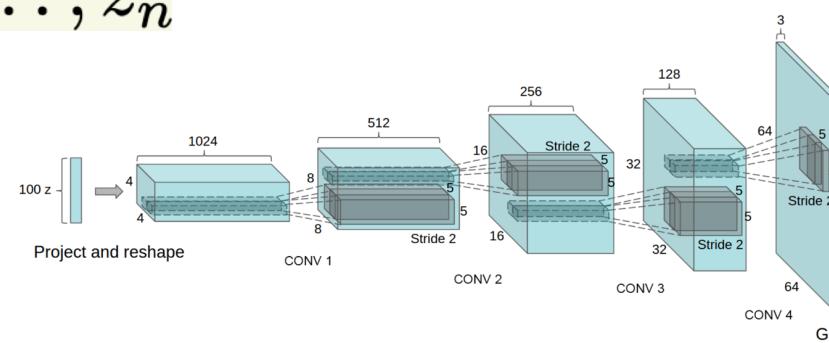
$$z \sim p(z)$$

$$x \sim p(x|g_\theta(z)) \cdot p(z)$$


$$g_\theta$$

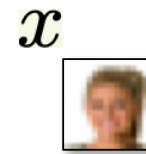

$$z_1, z_2, \dots, z_n$$

$$x_1, x_2, \dots, x_n$$



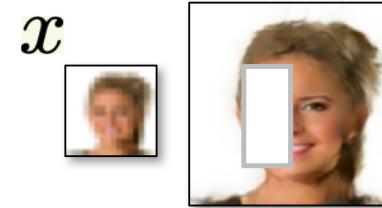
Why do I really need a latent model?

Answer: image restoration/editing/enhancement



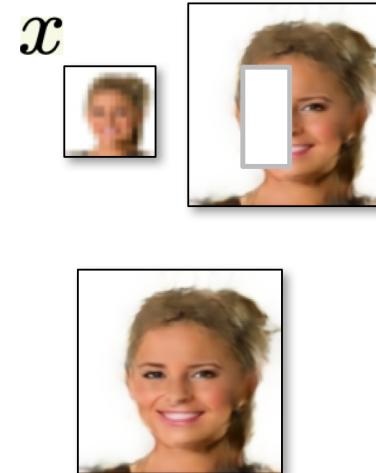
Why do I really need a latent model?

Answer: image restoration/editing/enhancement



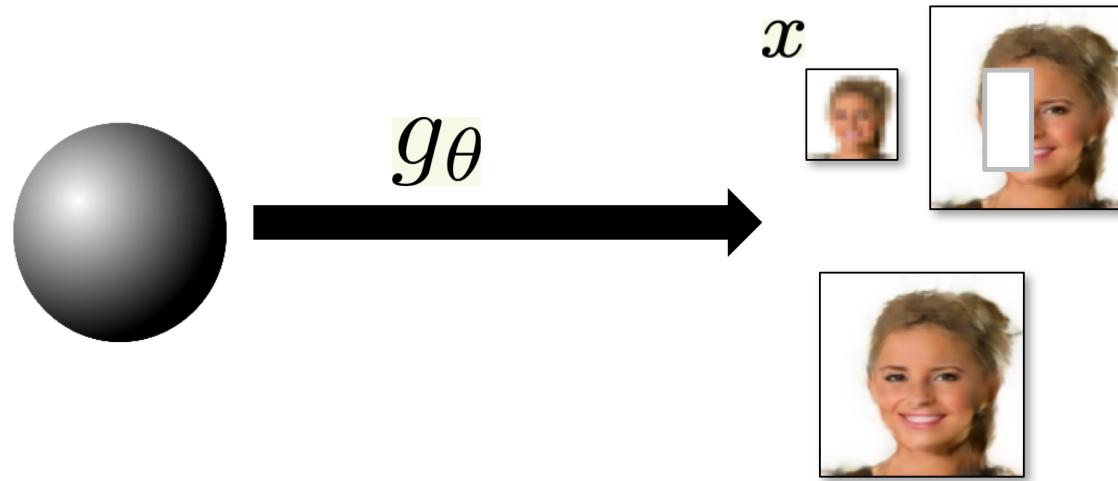
Why do I really need a latent model?

Answer: image restoration/editing/enhancement



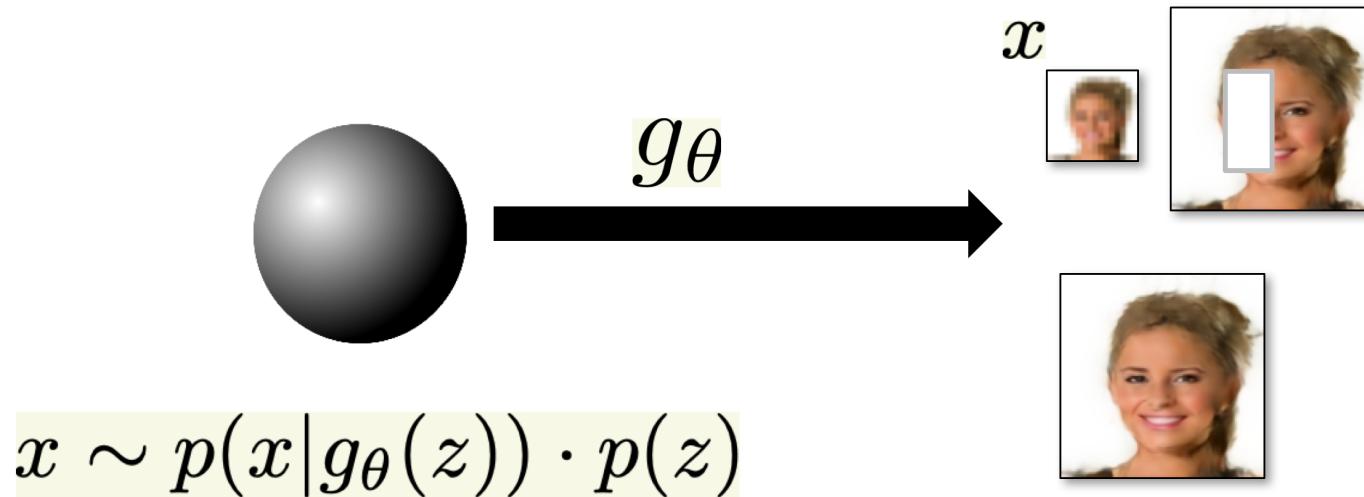
Why do I really need a latent model?

Answer: image restoration/editing/enhancement



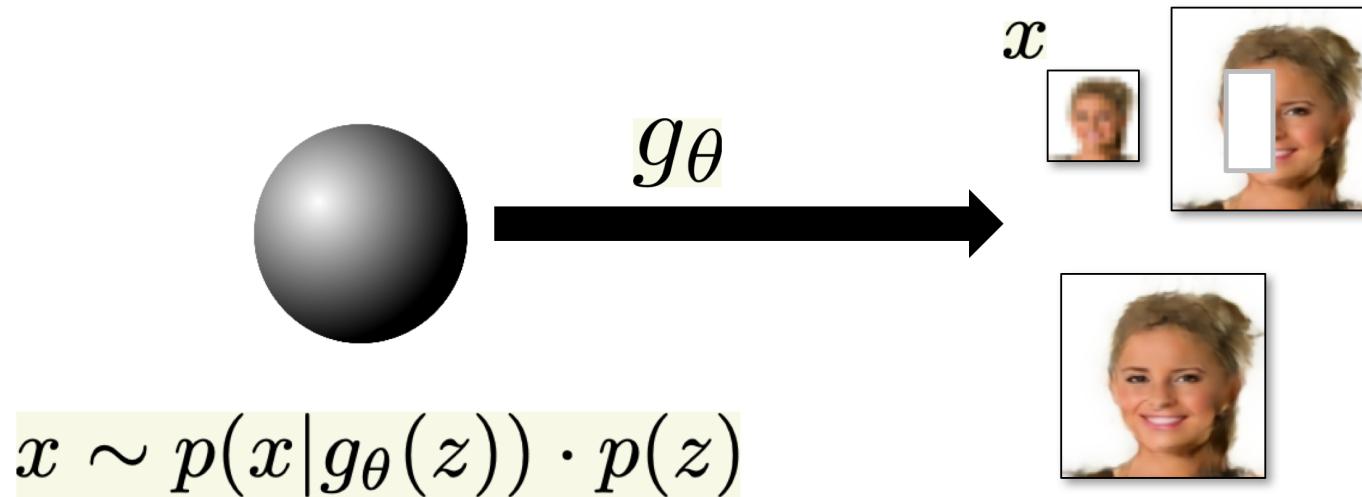
Why do I really need a latent model?

Answer: image restoration/editing/enhancement



Why do I really need a latent model?

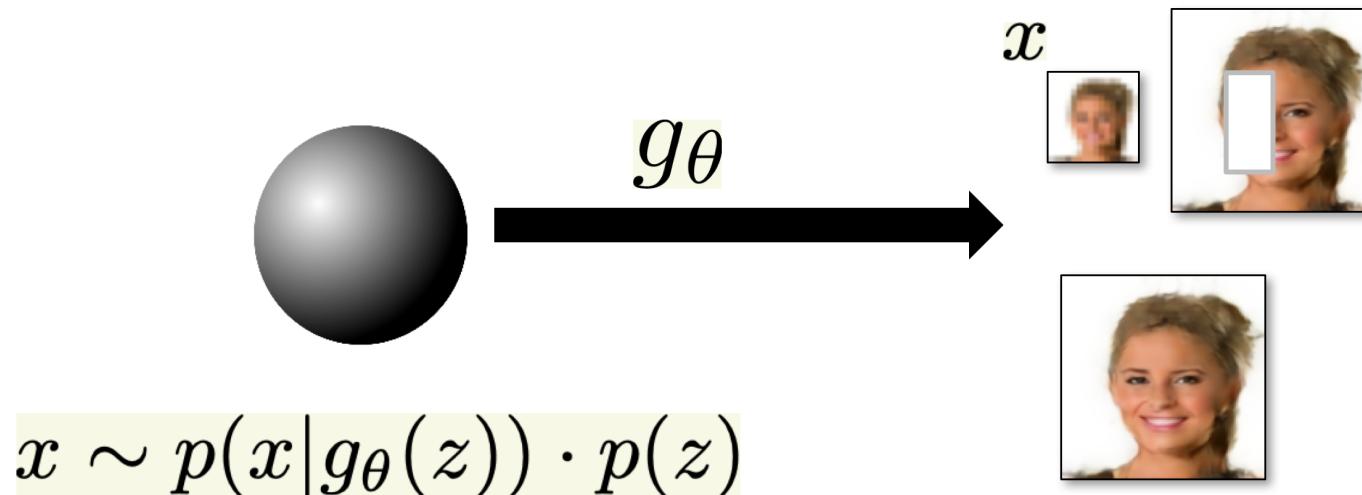
Answer: image restoration/editing/enhancement



$$\hat{z} = \arg \max_z [\log p(x|g_\theta(z)) + \log p(z)]$$

Why do I really need a latent model?

Answer: image restoration/editing/enhancement



$$\hat{z} = \arg \max_z [\log p(x|g_\theta(z)) + \log p(z)]$$

$$\hat{x} = g_\theta(\hat{z})$$



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



Variational AutoEncoders: Intractability

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

Simple Gaussian prior

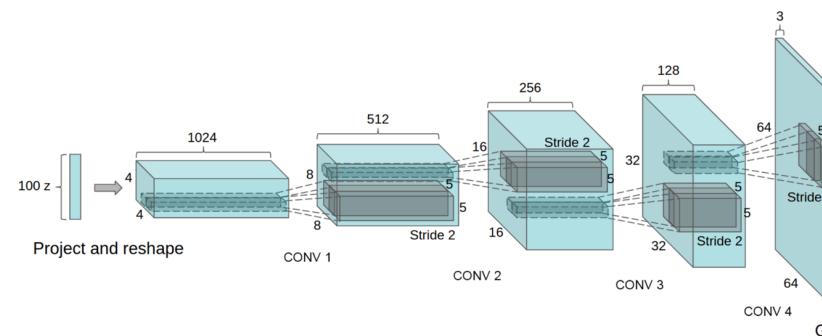


Data likelihood: $p_\theta(x) = \int p_\theta(z)p_\theta(x|z)dz$

↙
Decoder neural network

Decoder neural network:

$$p_\theta(x|z) = \mathcal{N}(x|\mu_\theta(z), \sigma_\theta^2(z) \cdot I)$$



Variational AutoEncoders: Intractability

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

↑
Intractible to compute!



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

Posterior density is

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



Variational AutoEncoders: Intractability

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

Posterior density is

also intractable: $p_\theta(z|x) = p_\theta(x|z) p_\theta(z) / p_\theta(x)$

Intractable data likelihood

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

Posterior density is

also intractable: $p_\theta(z|x) = p_\theta(x|z) p_\theta(z) / p_\theta(x)$

Solution: construct an encoder network $q_\phi(z|x)$

to approximate $p_\theta(z|x)$

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

Posterior density is

also intractable: $p_\theta(z|x) = p_\theta(x|z)p_\theta(z)/p_\theta(x)$

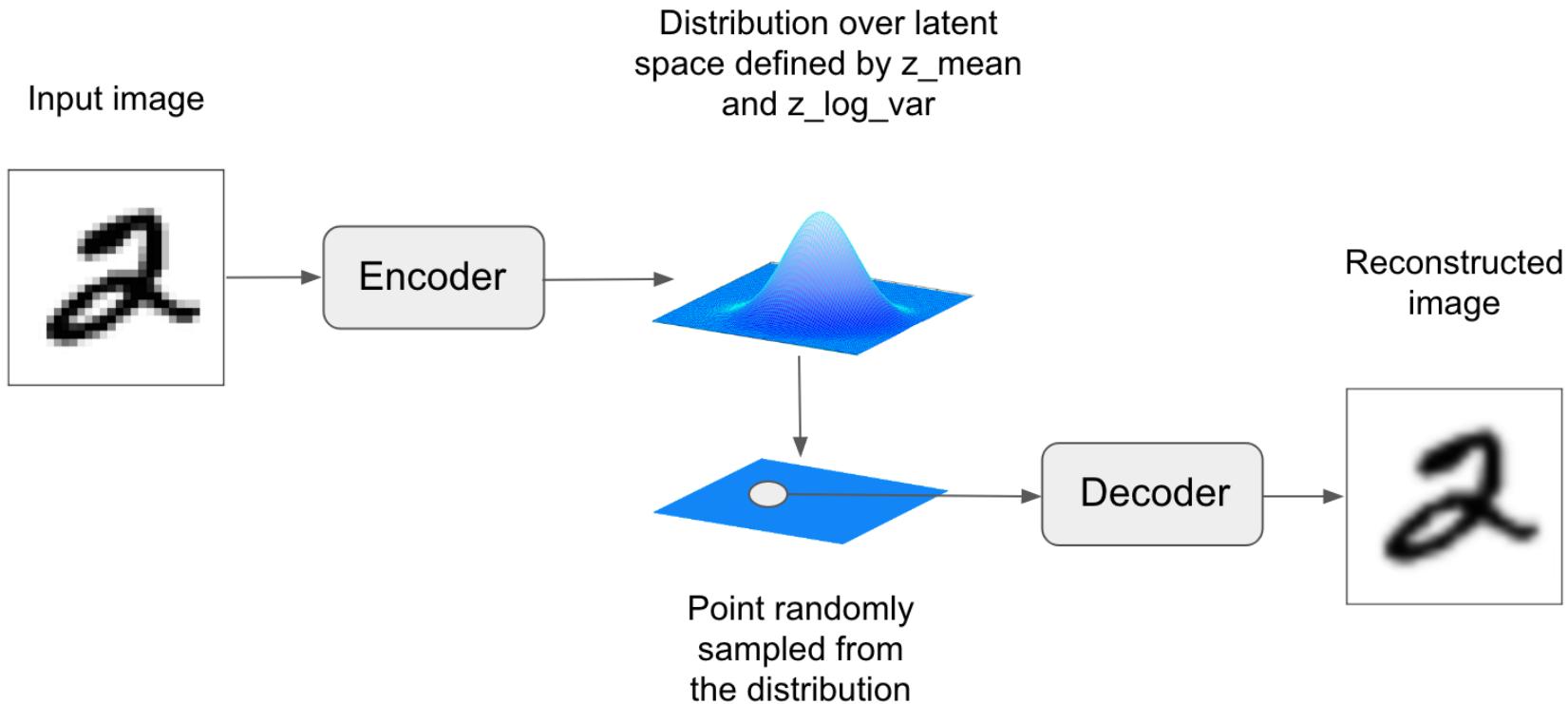
Solution: construct an encoder network $q_\phi(z|x)$
to approximate $p_\theta(z|x)$

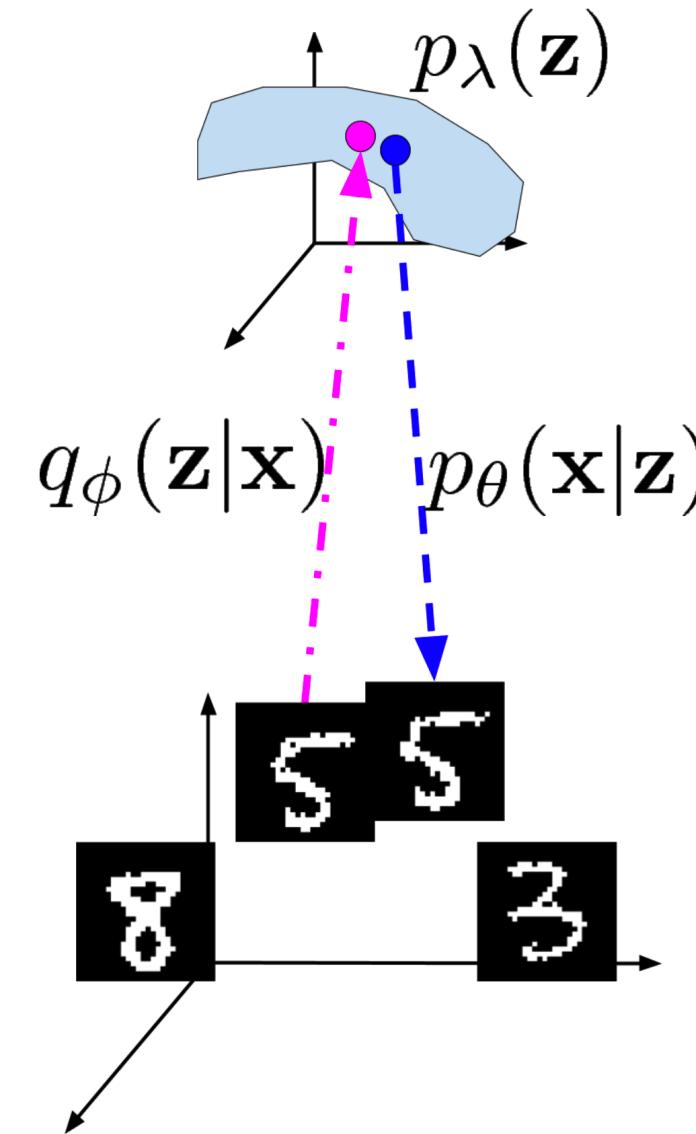
Encoder neural network:

$$q_\phi(z|x) = \mathcal{N}(z|\mu_\phi(x), \sigma_\phi^2(x) \cdot I)$$



Variational AutoEncoders





Marginal log-likelihood:

$$\log p_{\theta}(x_i) =$$

Marginal log-likelihood:

$$\begin{aligned}\log p_{\theta}(x_i) &= \\ &= KL[q_{\phi}(z|x)\|p_{\theta}(z|x)] + L(\theta, \phi, x_i)\end{aligned}$$

↑
always ≥ 0



Marginal log-likelihood:

$$\begin{aligned}\log p_{\theta}(x_i) &= \\ &= KL[q_{\phi}(z|x) \| p_{\theta}(z|x)] + L(\theta, \phi, x_i) \\ &\geq L(\theta, \phi, x_i)\end{aligned}$$

↑
always ≥ 0

Marginal log-likelihood:

$$\begin{aligned}\log p_{\theta}(x_i) &= \\ &= KL[q_{\phi}(z|x)||p_{\theta}(z|x)] + L(\theta, \phi, x_i) \\ &\geq L(\theta, \phi, x_i)\end{aligned}$$

↑
always ≥ 0

Variational lower-bound:

$$\begin{aligned}L(\theta, \phi, x_i) &= \\ &= -KL[q_{\phi}(z|x_i)||p(z)] + \mathbb{E}_{q_{\phi}(z|x_i)}[\log p_{\theta}(x_i|z)]\end{aligned}$$



Marginal log-likelihood:

$$\begin{aligned}\log p_{\theta}(x_i) &= \\ &= KL[q_{\phi}(z|x)||p_{\theta}(z|x)] + L(\theta, \phi, x_i) \\ &\geq L(\theta, \phi, x_i)\end{aligned}$$

↑
always ≥ 0

Variational lower-bound:

$$\begin{aligned}L(\theta, \phi, x_i) &= \\ &= -KL[q_{\phi}(z|x_i)||p(z)] + \mathbb{E}_{q_{\phi}(z|x_i)}[\log p_{\theta}(x_i|z)]\end{aligned}$$

↑
Regularization!



Marginal log-likelihood:

$$\begin{aligned}\log p_{\theta}(x_i) &= \\ &= KL[q_{\phi}(z|x)\|p_{\theta}(z|x)] + L(\theta, \phi, x_i) \\ &\geq L(\theta, \phi, x_i)\end{aligned}$$

↑
always ≥ 0

Variational lower-bound:

$$\begin{aligned}L(\theta, \phi, x_i) &= \\ &= -KL[q_{\phi}(z|x_i)\|p(z)] + \mathbb{E}_{q_{\phi}(z|x_i)}[\log p_{\theta}(x_i|z)]\end{aligned}$$

↑
Regularization! ↑
Latent representation



Marginal log-likelihood:

$$\begin{aligned}\log p_{\theta}(x_i) &= \\ &= KL[q_{\phi}(z|x)\|p_{\theta}(z|x)] + L(\theta, \phi, x_i) \\ &\geq L(\theta, \phi, x_i)\end{aligned}$$

↑
always ≥ 0

Variational lower-bound:

$$\begin{aligned}L(\theta, \phi, x_i) &= \\ &= -KL[q_{\phi}(z|x_i)\|p(z)] + \mathbb{E}_{q_{\phi}(z|x_i)}[\log p_{\theta}(x_i|z)]\end{aligned}$$

↑
Regularization!

↑
Latent representation

↓
Reconstruction error

ELBO:

ELBO:

$$\begin{aligned} \sum_{i=1}^n \log p_{\theta}(x_i) &\geq \\ &\geq \sum_{i=1}^n L(\theta, \phi, x_i) \end{aligned}$$

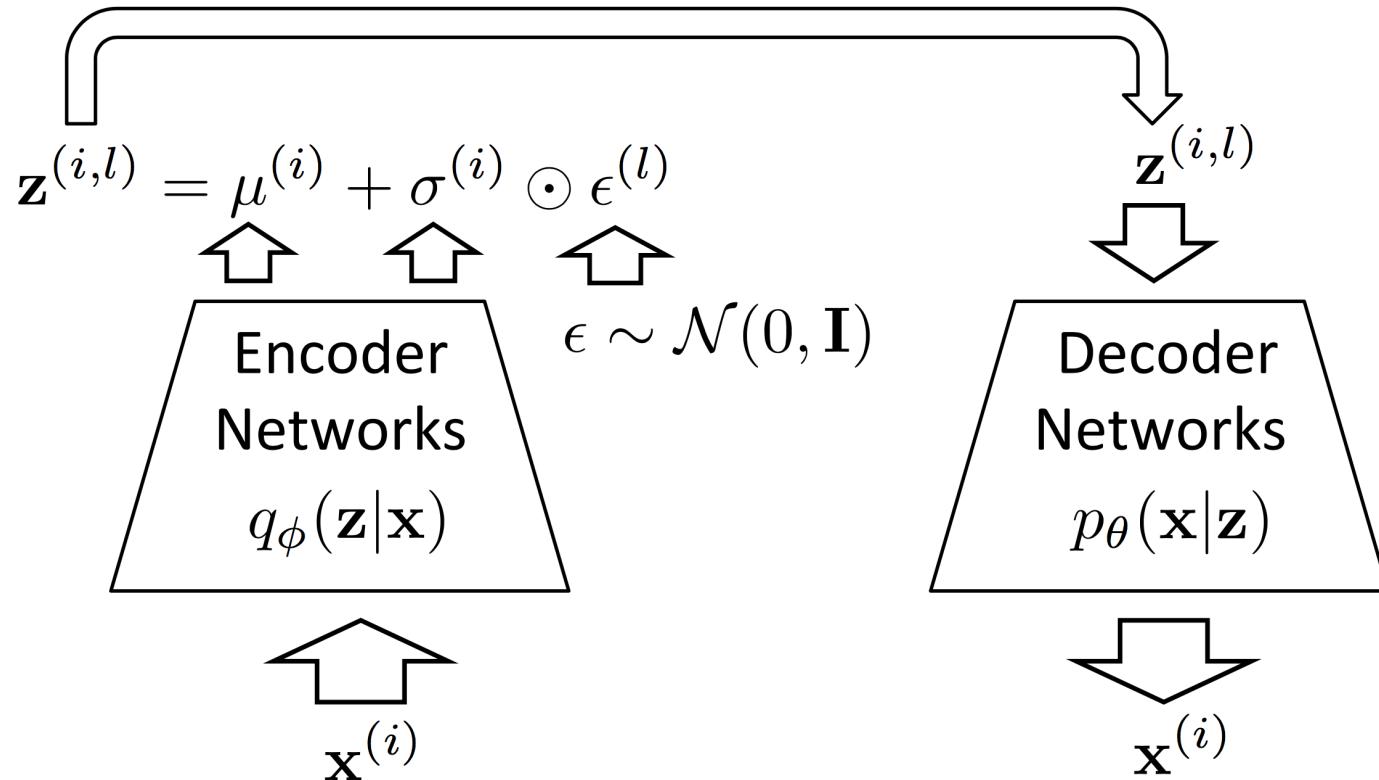


ELBO:

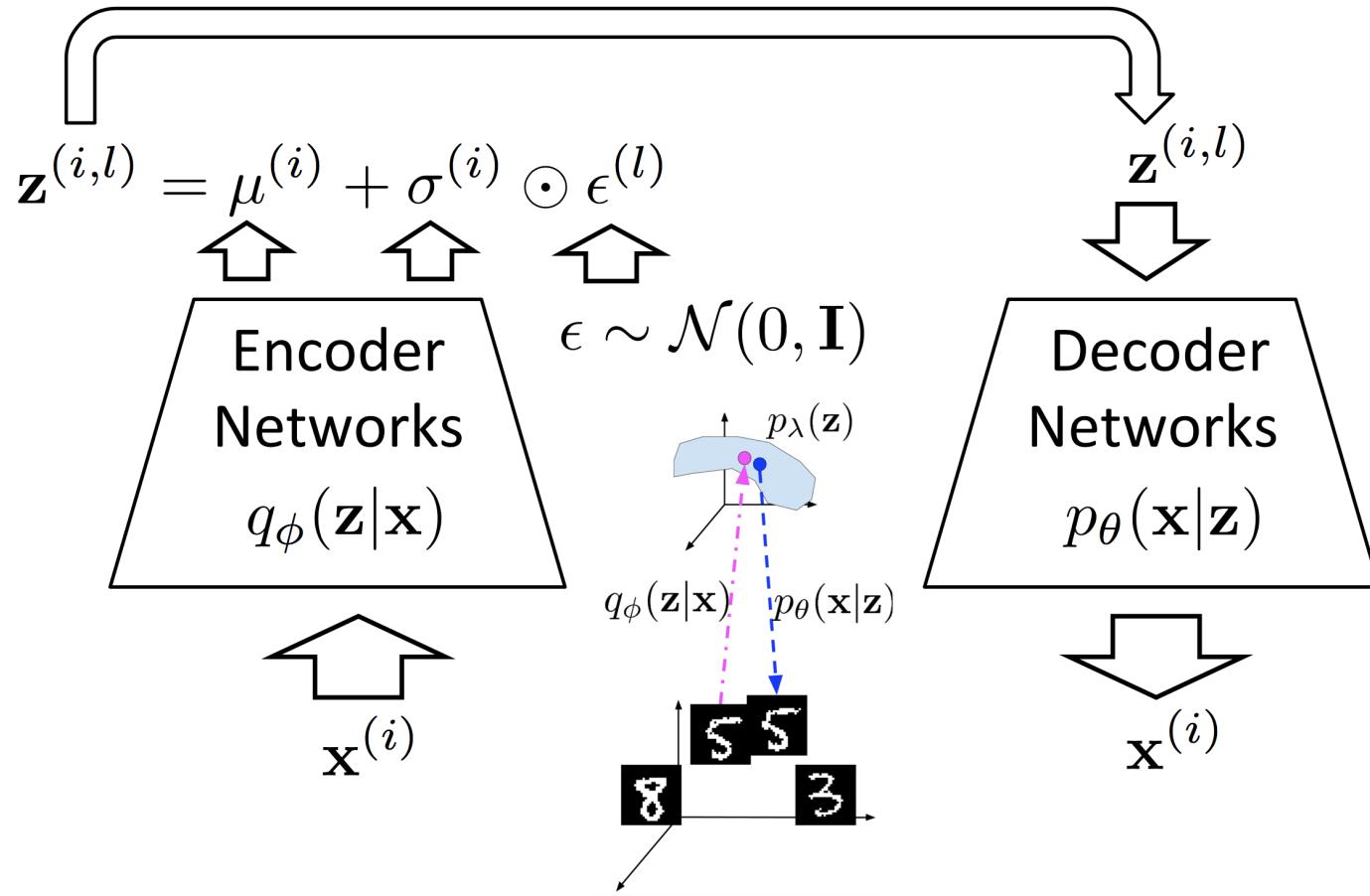
$$\begin{aligned} \sum_{i=1}^n \log p_{\theta}(x_i) &\geq \\ &\geq \sum_{i=1}^n L(\theta, \phi, x_i) \rightarrow \max_{\theta, \phi} \end{aligned}$$



Variational AutoEncoders: Reparameterization Trick



Variational AutoEncoders: Reparameterization Trick



Empirical variational lower-bound:

Empirical variational lower-bound:

$$p(z) = \mathcal{N}(z|0, \mathbf{I})$$



Empirical variational lower-bound:

$$p(z) = \mathcal{N}(z|0, \mathbf{I})$$

$$z_{i,l} \sim q_\phi(z|x_i)$$

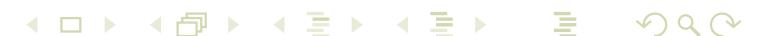


Empirical variational lower-bound:

$$p(z) = \mathcal{N}(z|0, \mathbf{I})$$

$$z_{i,l} \sim q_\phi(z|x_i)$$

$$\begin{aligned}\hat{L}(\theta, \phi, x_i) &= \frac{1}{2} \sum_{j=1}^d [1 + \log \sigma_{j,\phi}^2(x_i) - \mu_{j,\phi}^2(x_i) - \sigma_{j,\phi}^2(x_i)] \\ &+ \frac{1}{L} \sum_{l=1}^L \log p_\theta(x_i|z_{i,l}) \rightarrow \max_{\theta, \phi}\end{aligned}$$



Empirical variational lower-bound:

$$p(z) = \mathcal{N}(z|0, \mathbf{I})$$

$$z_{i,l} \sim q_\phi(z|x_i)$$

$$\begin{aligned}\hat{L}(\theta, \phi, x_i) &= \frac{1}{2} \sum_{j=1}^d [1 + \log \sigma_{j,\phi}^2(x_i) - \mu_{j,\phi}^2(x_i) - \sigma_{j,\phi}^2(x_i)] \\ &\quad + \frac{1}{L} \sum_{l=1}^L \log p_\theta(x_i|z_{i,l}) \rightarrow \max_{\theta, \phi}\end{aligned}$$

Regularization!



Empirical variational lower-bound:

$$p(z) = \mathcal{N}(z|0, \mathbf{I})$$

$$z_{i,l} \sim q_\phi(z|x_i)$$

$$\begin{aligned}\hat{L}(\theta, \phi, x_i) &= \frac{1}{2} \sum_{j=1}^d [1 + \log \sigma_{j,\phi}^2(x_i) - \mu_{j,\phi}^2(x_i) - \sigma_{j,\phi}^2(x_i)] \\ &\quad + \frac{1}{L} \sum_{l=1}^L \log p_\theta(x_i|z_{i,l}) \rightarrow \max_{\theta, \phi}\end{aligned}$$

↑
Reconstruction error!
↓
Regularization!

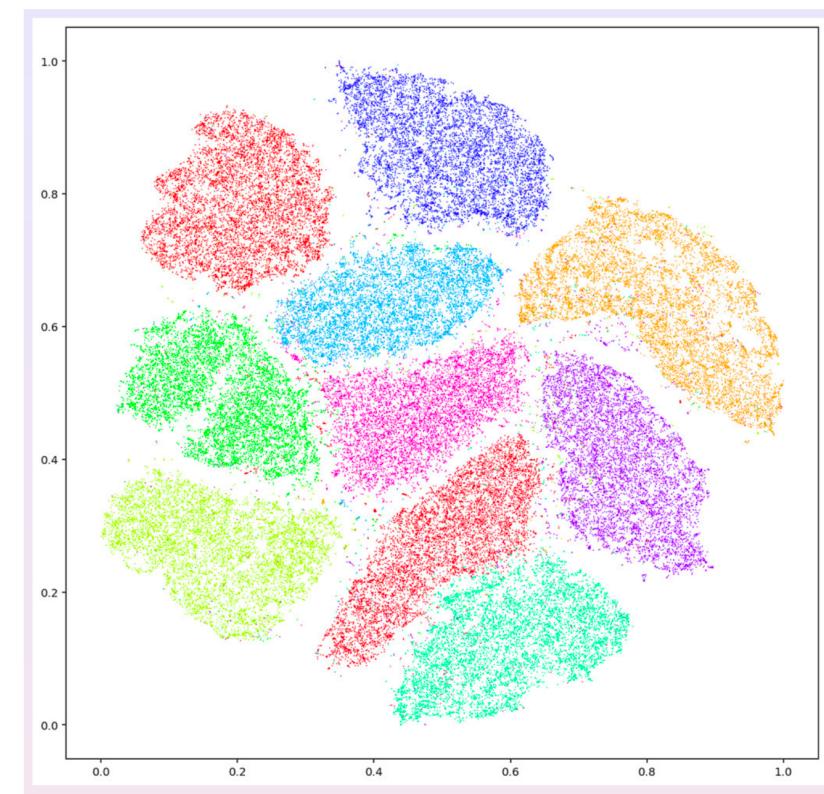


Example: MNIST

Use decoder network



Use decoder network



Each digit corresponds to a distribution on the manifold

Use decoder network. Now sample z from prior!

Use decoder network. Now sample \mathbf{z} from prior!

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

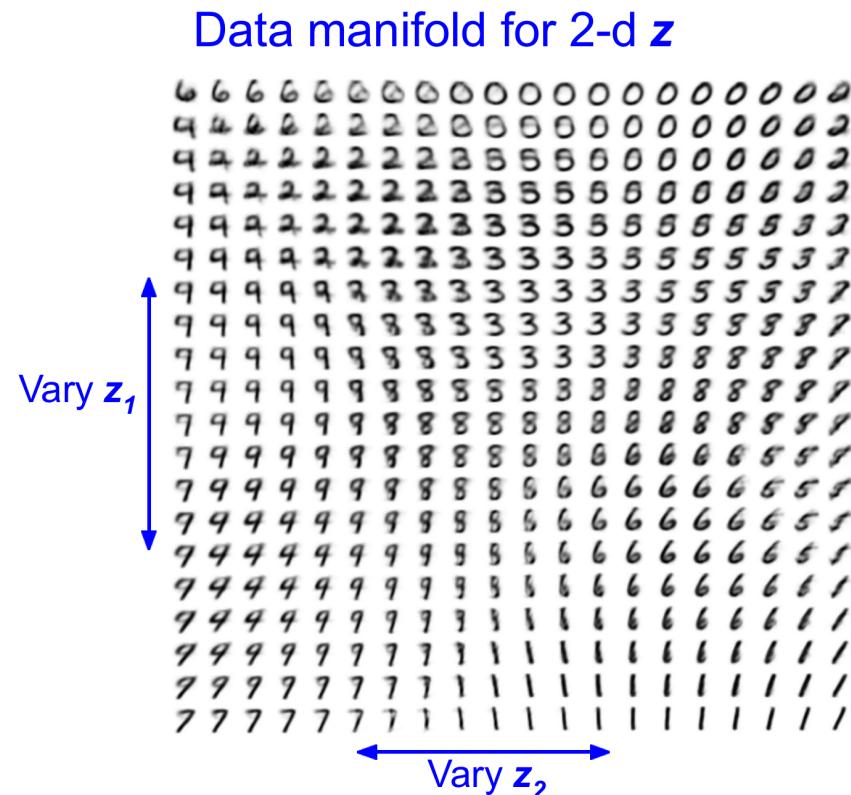
Use decoder network. Now sample \mathbf{z} from prior!

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

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

Example: MNIST

Use decoder network. Now sample \mathbf{z} from prior!

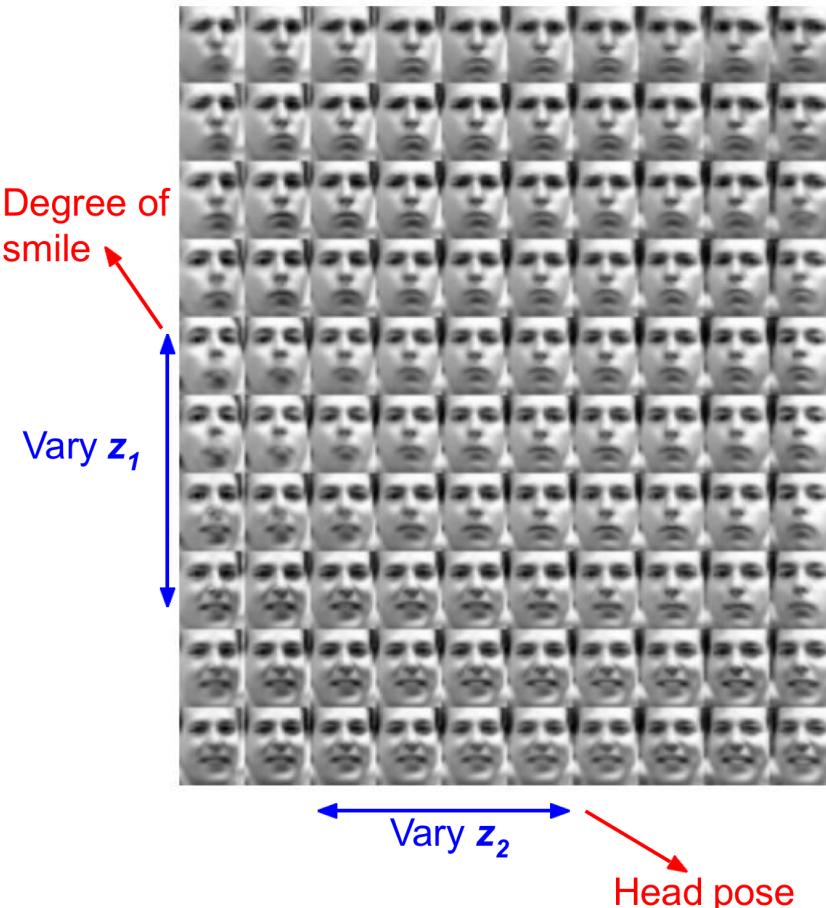
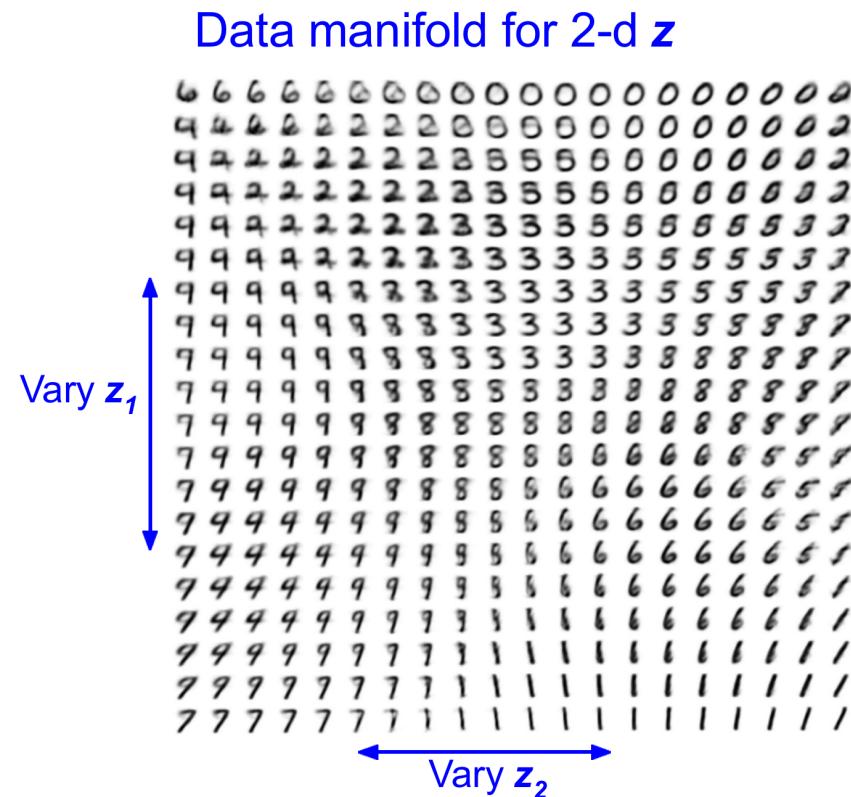


Diagonal prior on \mathbf{z}
=> independent latent
variables

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

Example: Faces

Use decoder network. Now sample \mathbf{z} from prior!



Example: Faces

Use decoder network. Now sample z from prior!



Labeled Faces in the Wild

