

Deep Generative Models

Changyou Chen

Department of Computer Science and Engineering
University at Buffalo, SUNY
changyou@buffalo.edu

April 4, 2019

Variational Inference

$$\mathcal{F}(y, q) = \mathbb{E}_{q(z)}[\log p(y|z)] - KL[q(z)\|p(z)]$$

Approx. Posterior Reconstruction Penalty

Some comments on q

- **Integration is now optimization:** optimize for $q(z)$ directly:
 z typically depends on data y , but write $q(z)$ for notation simplicity.
Easy convergence assessment.
- **Variational parameters:** parameters of $q(z)$:
e.g., if a Gaussian, the parameters are mean and variance.
Optimization allows us to tighten the bound and get as close as possible to the true marginal likelihood.

Free-form and Fixed-form Solutions

Free-form

- Solves for the exact distribution $q(z)$ by setting the functional derivative to zero via **calculus of variations**:

$$\frac{\delta \mathcal{F}(y, q)}{\delta q(z)} = 0, \quad s.t. \quad \int q(z)dz = 1$$
$$\Rightarrow q(z) \propto p(z)p(y|z, \theta)$$

- The optimal solution is the true posterior distribution, but solving for the normalization is our original problem.

Fixed-form

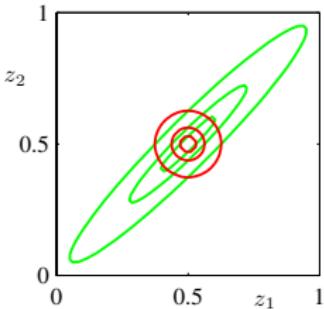
- Specify an explicit form of $q(z)$, e.g., a normal distribution.
- Parameter in q is called the variational parameter.

Mean-field Variational Inference

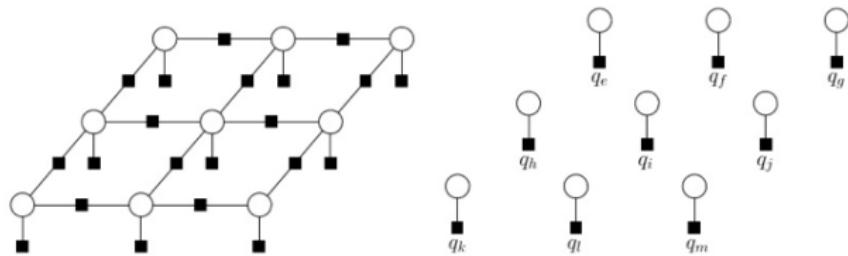
Mean-field methods assume that the distribution is factorised.

$$q(z) = \prod_i q_i(z_i)$$

Restricted class of approximations: every dimension (or subset of dimensions) of the posterior is independent.



$$q(z) = \prod_i \mathcal{N}(z_i | \mu_i, \sigma_i^2)$$



Example: Mean-field for Latent Gaussian Models

Generative model:

$$\mathbf{z} \sim \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}), \quad \mathbf{y} \sim p(\mathbf{y} | f_{\theta}(\mathbf{z}))$$

Variational distribution:

$$q(\mathbf{z}) = \prod_i \mathcal{N}(z_i; \mu_i, \sigma_i^2)$$

Example: Mean-field for Latent Gaussian Models

$$\begin{aligned}\mathcal{F}(\mathbf{y}, q) &= \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{y} | \mathbf{z})] - KL [q(\mathbf{z}) \| p(\mathbf{z})] \\&= \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{y} | \mathbf{z})] - \sum_i KL [q(z_i) \| p(z_i)] \\&= \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{y} | \mathbf{z})] - \sum_i KL \left[\mathcal{N}(z_i; \mu_i, \sigma_i^2) \| \mathcal{N}(z_i; 0, 1) \right] \\&= \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{y} | \mathbf{z})] - \frac{1}{2} \sum_i \left(\sigma_i^2 + \mu_i^2 - 1 - \log \sigma_i^2 \right)\end{aligned}$$

KL between two Gaussians

Let $p(x) = \mathcal{N}(x; \mu_1, \sigma_1^2)$, $q(x) = \mathcal{N}(x; \mu_2, \sigma_2^2)$, then

$$KL(p(x) \| q(x)) = \log \frac{\sigma_2^2}{\sigma_1^2} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

Example: Mean-field for Latent Gaussian Models

$$\begin{aligned}\mathcal{F}(\mathbf{y}, q) &= \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{y} | \mathbf{z})] - KL [q(\mathbf{z}) \| p(\mathbf{z})] \\&= \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{y} | \mathbf{z})] - \sum_i KL [q(z_i) \| p(z_i)] \\&= \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{y} | \mathbf{z})] - \sum_i KL \left[\mathcal{N}(z_i; \mu_i, \sigma_i^2) \| \mathcal{N}(z_i; 0, 1) \right] \\&= \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{y} | \mathbf{z})] - \frac{1}{2} \sum_i \left(\sigma_i^2 + \mu_i^2 - 1 - \log \sigma_i^2 \right)\end{aligned}$$

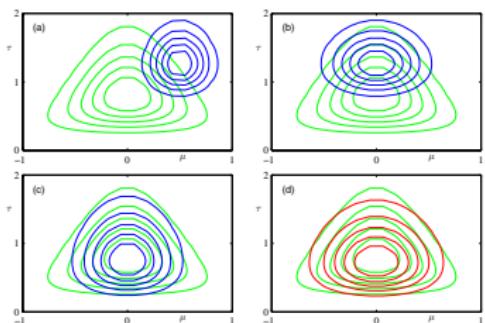
Optimize parameters of $q(z)$ by gradient descent.

Optimization for the Variational Bound

$$\max_{q,\theta} \mathcal{F}(y, q) = \mathbb{E}_{q(z)}[\log p(y|z, \theta)] - KL[q(z)\|p(z)]$$

Approx. Posterior Reconstruction Penalty

- *Variational EM*
- *Stochastic Variational Inference*
- *Doubly Stochastic Variational Inference*
- *Amortised Inference*

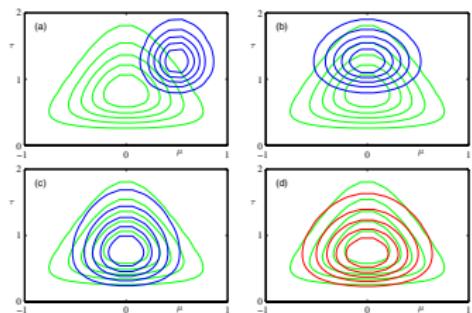


Optimization for the Variational Bound

$$\max_{q, \theta} \mathcal{F}(y, q) = \mathbb{E}_{q(z)}[\log p(y|z, \theta)] - KL[q(z)\|p(z)]$$

Approx. Posterior Reconstruction Penalty

- *Variational EM*
- *Stochastic Variational Inference*
- *Doubly Stochastic Variational Inference*
- *Amortised Inference*



Dealt with big data and complicated $q(z)$

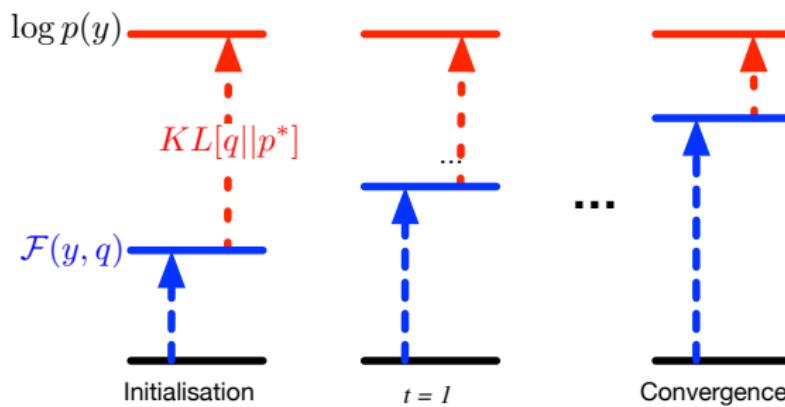
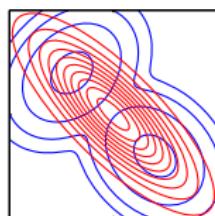
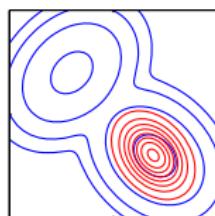
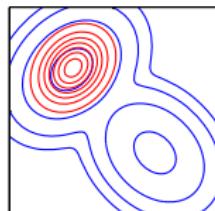
Variational Expectation Maximization

Alternating optimization for the variational parameters and then model parameters.

Repeat:

E-step $\phi \propto \nabla_\phi \mathcal{F}(y, q)$ *Var. params*

M-step $\theta \propto \nabla_\theta \mathcal{F}(y, q)$ *Model params*



Variational Expectation Maximization

Alternating optimization for the variational parameters and then model parameters.

Repeat:

E-step

(Inference)

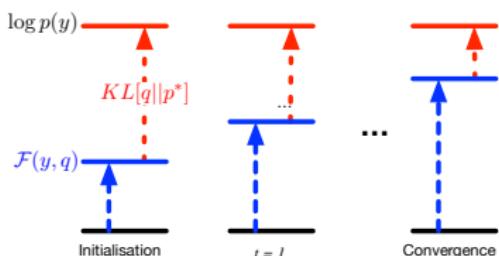
For $i = 1, \dots, N$

$$\phi_n \propto \nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(y_n | z_n)] - \nabla_{\phi} KL[q(z_n) \| p(z_n)]$$

M-step

(Parameter Learning)

$$\theta \propto \frac{1}{N} \sum_n \mathbb{E}_{q_{\phi}(z)} [\nabla_{\theta} \log p_{\theta}(y_n | z_n)]$$

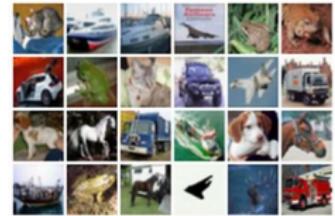
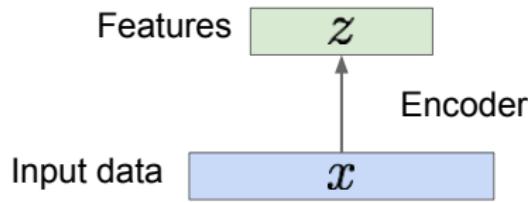


Autoencoder¹

¹ Partially adapted from http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture13.pdf

Background: Autoencoder

Learning a lower-dimensional feature representation from unlabeled training data



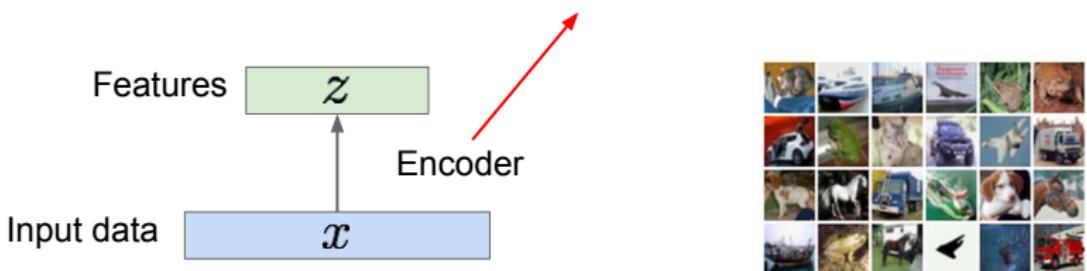
Background: Autoencoder

Learning a lower-dimensional feature representation from unlabeled training data

Originally: Linear +
nonlinearity (sigmoid)

Later: Deep, fully-connected

Later: ReLU CNN

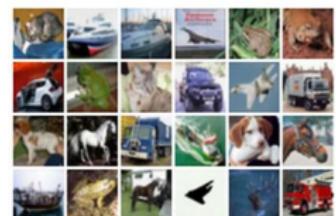
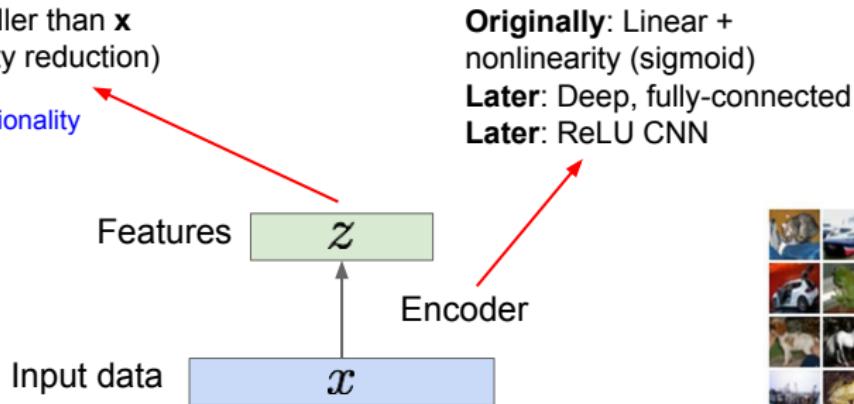


Background: Autoencoder

Learning a lower-dimensional feature representation from unlabeled training data

z usually smaller than x
(dimensionality reduction)

Q: Why dimensionality reduction?



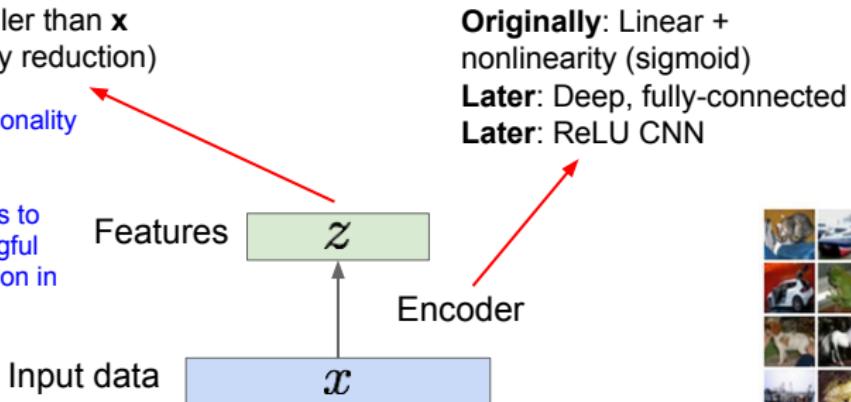
Background: Autoencoder

Learning a lower-dimensional feature representation from unlabeled training data

z usually smaller than x
(dimensionality reduction)

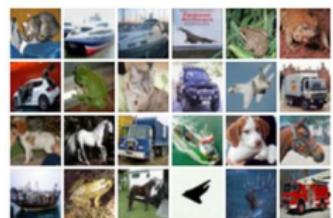
Q: Why dimensionality reduction?

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



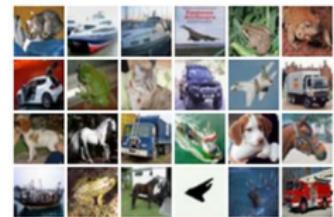
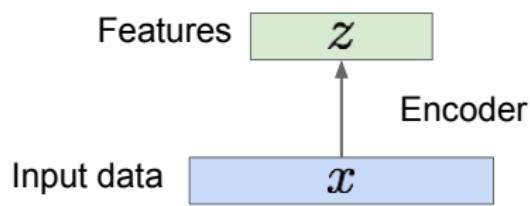
Originally: Linear +
nonlinearity (sigmoid)

Later: Deep, fully-connected
Later: ReLU CNN



Background: Autoencoder

How to learning this feature representation?

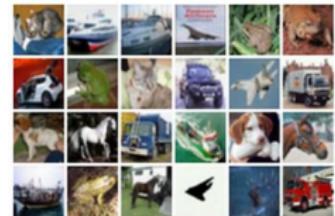
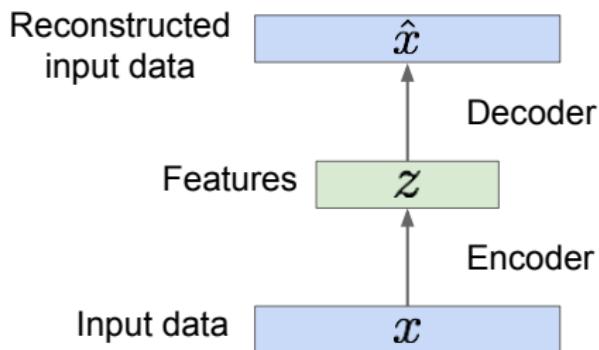


Background: Autoencoder

How to learn this feature representation?

We want that feature can be used to reconstruct itself:

- “Autoencoding” – encoding itself.

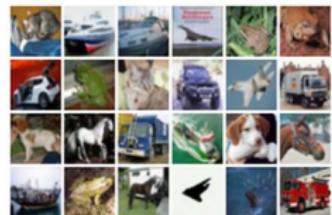
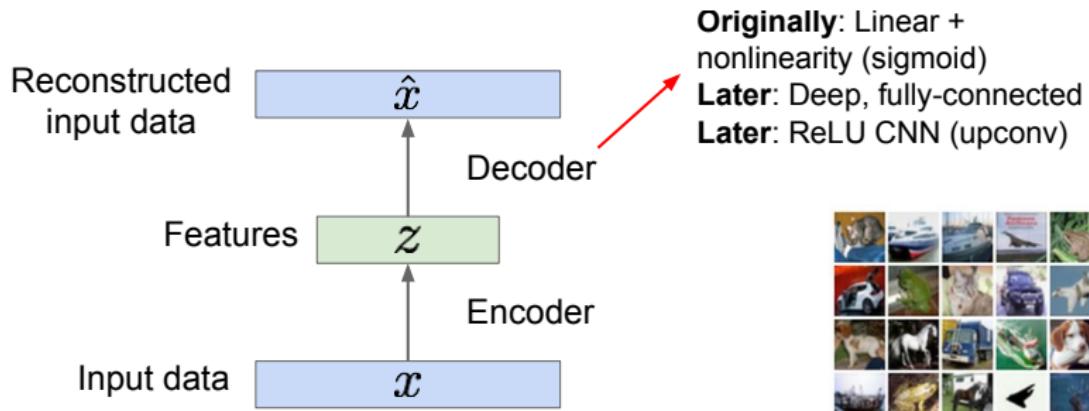


Background: Autoencoder

How to learn this feature representation?

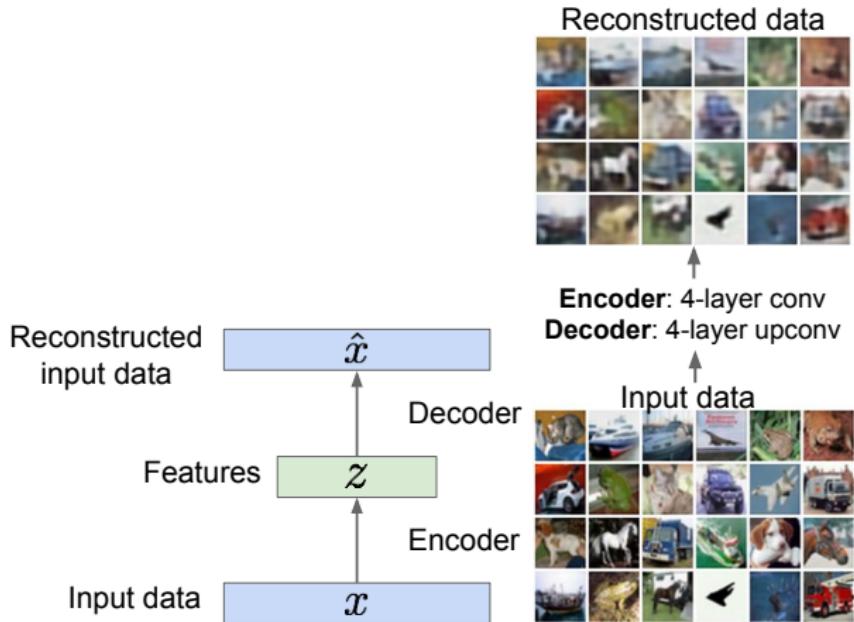
We want that feature can be used to reconstruct itself:

- “Autoencoding” – encoding itself.



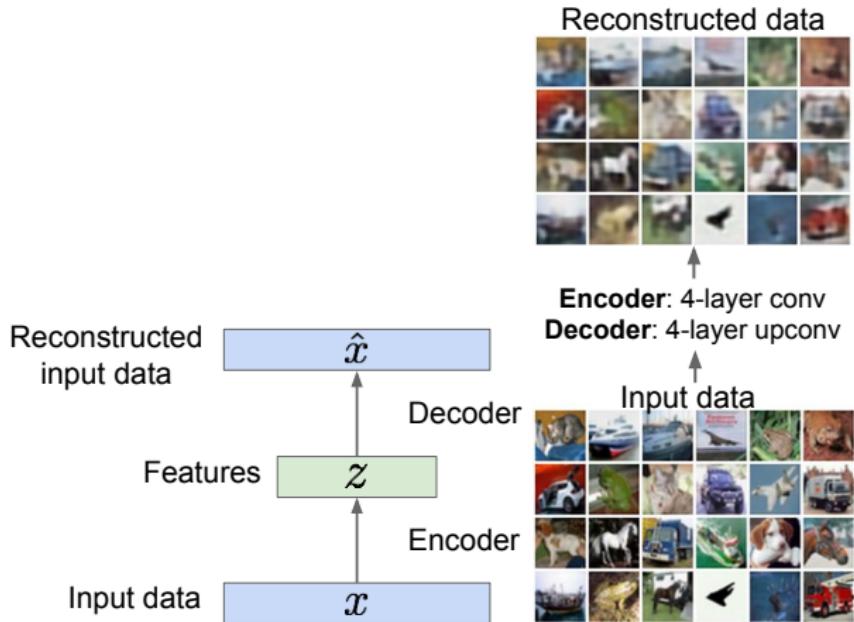
Background: Autoencoder

- We want that feature can be used to reconstruct itself
- Use L^2 loss:
 $\|x - \hat{x}\|^2$
- Doesn't use labels



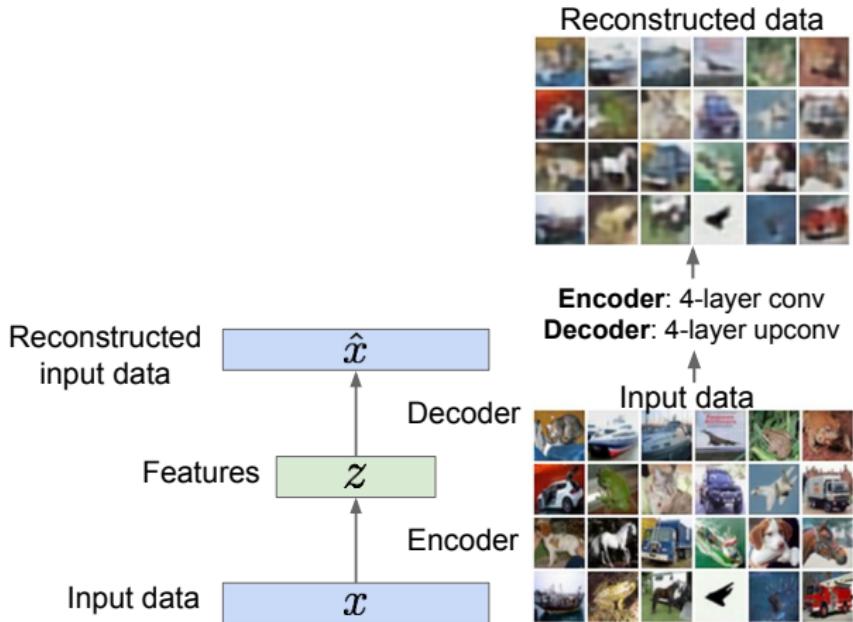
Background: Autoencoder

- We want that feature can be used to reconstruct itself
- Use L^2 loss:
 $\|\mathbf{x} - \hat{\mathbf{x}}\|^2$
- Doesn't use labels



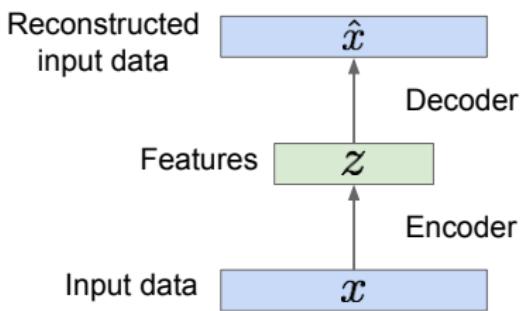
Background: Autoencoder

- We want that feature can be used to reconstruct itself
- Use L^2 loss:
 $\|\mathbf{x} - \hat{\mathbf{x}}\|^2$
- Doesn't use labels



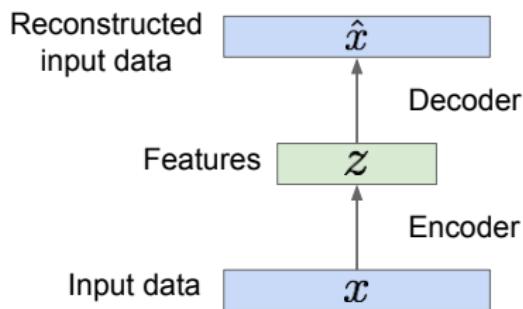
Background: Autoencoder

- What can an autoencoder do?
 - Reconstruct data.
 - Pretrain for supervised learning.



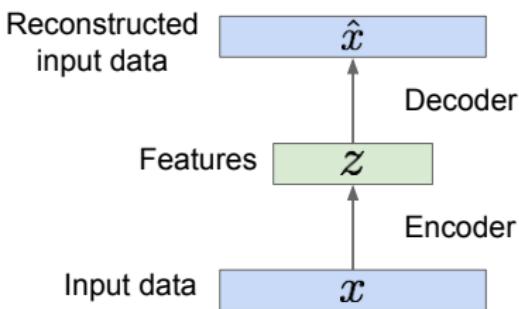
Background: Autoencoder

- What can an autoencoder do?
 - Reconstruct data.
 - Pretrain for supervised learning.



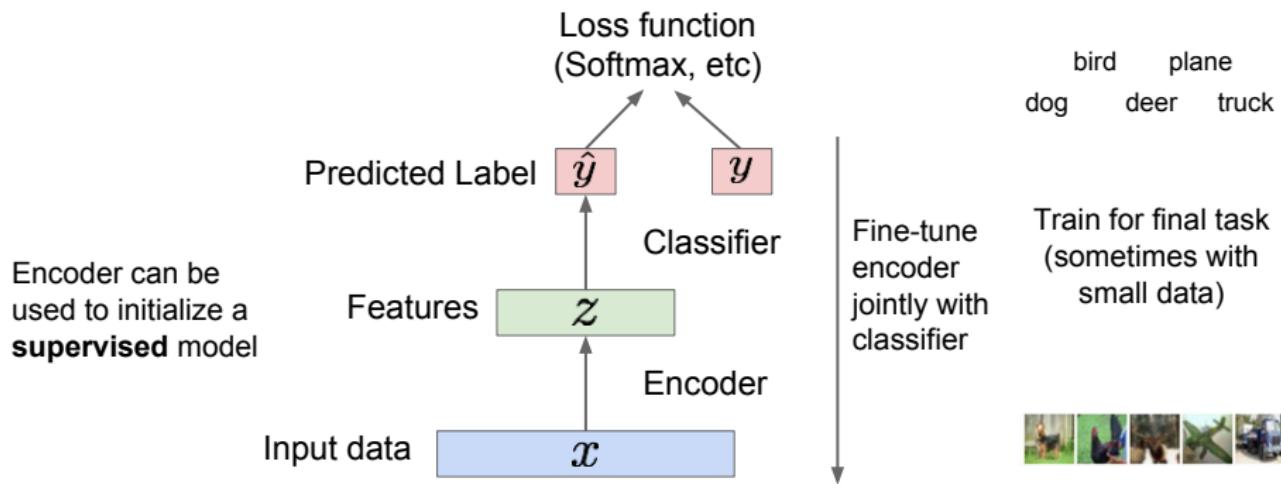
Background: Autoencoder

- What can an autoencoder do?
 - Reconstruct data.
 - Pretrain for supervised learning.



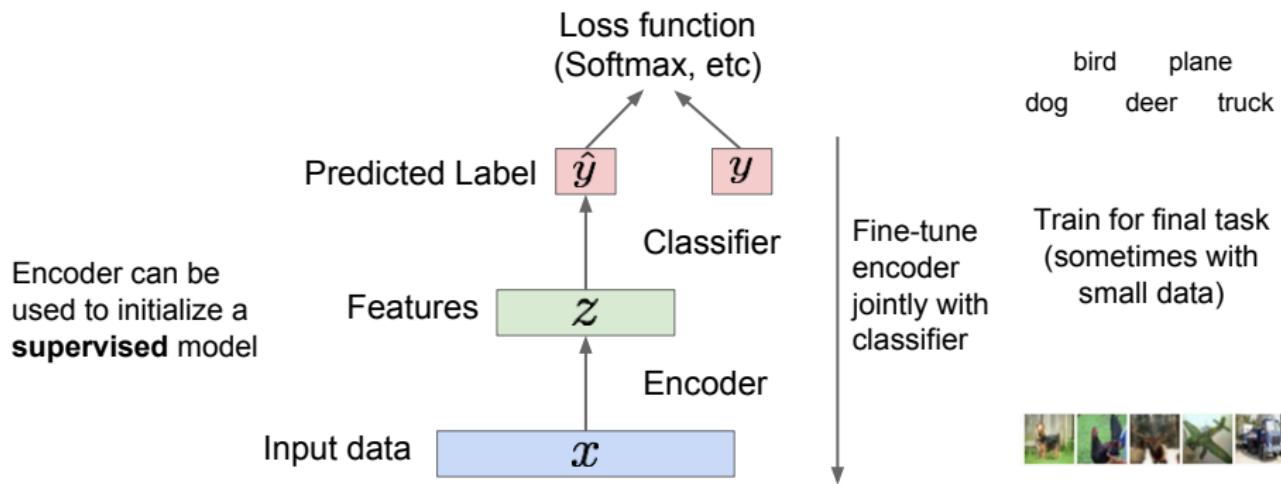
Background: Autoencoder

After training, throw away decoder.



Background: Autoencoder

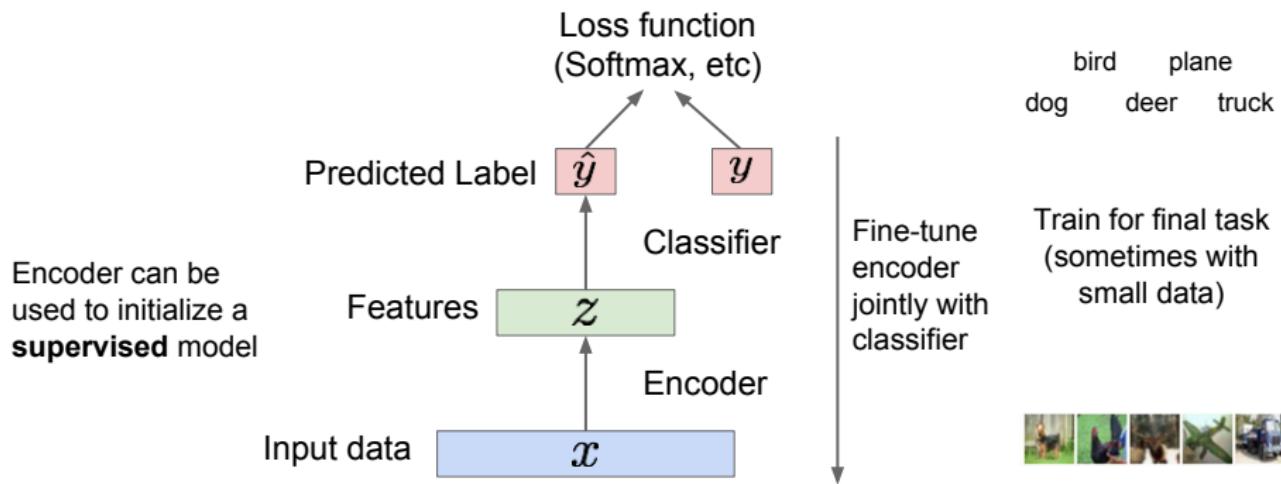
After training, throw away decoder.



- Autoencoders can reconstruct data, and can learn features to initialize a supervised model.
- However, can we generate new images from it?
 - NO

Background: Autoencoder

After training, throw away decoder.



- Autoencoders can reconstruct data, and can learn features to initialize a supervised model.
- However, can we generate new images from it?
 - NO

Variational Autoencoder

Variational Autoencoder \Rightarrow Probabilistic Autoencoder

- ① Define a generative probability model for the decoder.
- ② Assume training data $\{\mathbf{x}^{(i)}\}_{i=1}^N$ is generated from underlying unobserved (latent) representation \mathbf{z} .

Variational Autoencoder \Rightarrow Probabilistic Autoencoder

- 1 Define a generative probability model for the decoder.
- 2 Assume training data $\{\mathbf{x}^{(i)}\}_{i=1}^N$ is generated from underlying unobserved (latent) representation \mathbf{z} .

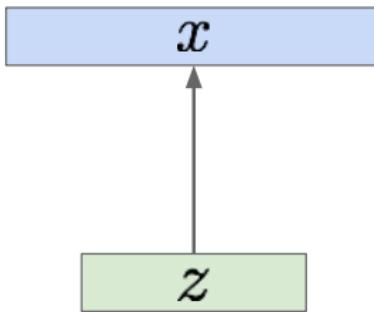
- Sample latent \mathbf{z} from true prior

$$p_{\theta^*}(\mathbf{z})$$

- Sample data \mathbf{x} from true

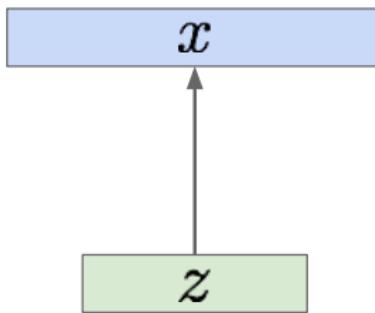
conditional

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



Variational Autoencoder \Rightarrow Probabilistic Autoencoder

- 1 Define a generative probability model for the deconder.
- 2 Assume training data $\{\mathbf{x}^{(i)}\}_{i=1}^N$ is generated from underlying unobserved (latent) representation \mathbf{z} .
- Sample latent \mathbf{z} from true prior $p_{\theta^*}(\mathbf{z})$
- Sample data \mathbf{x} from true conditional $p_{\theta^*}(\mathbf{x} \mid \mathbf{z})$



Intuition

\mathbf{x} is an image, \mathbf{z} is latent factors used to generate \mathbf{x} , e.g., attributes, orientation, etc

Variational Autoencoder \Rightarrow Probabilistic Autoencoder

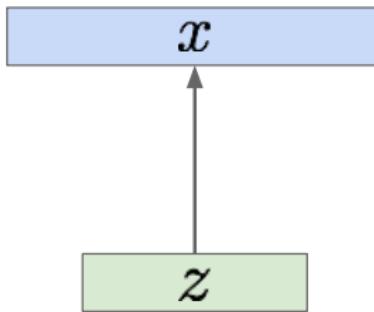
- ① Define a generative probability model for the decoder.
- ② Assume training data $\{\mathbf{x}^{(i)}\}_{i=1}^N$ is generated from underlying unobserved (latent) representation \mathbf{z} .

- Sample latent \mathbf{z} from true prior

$$p_{\theta^*}(\mathbf{z})$$

- Sample data \mathbf{x} from true conditional

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

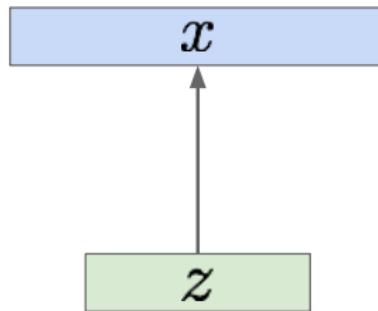


We want to estimate the true parameters θ^* of this generative model.

Variational Autoencoder \Rightarrow Probabilistic Autoencoder

- ① Define a generative probability model for the deconder.
- ② Assume training data $\{\mathbf{x}^{(i)}\}_{i=1}^N$ is generated from underlying unobserved (latent) representation \mathbf{z} .

- Sample latent \mathbf{z} from true prior $p_{\theta^*}(\mathbf{z})$



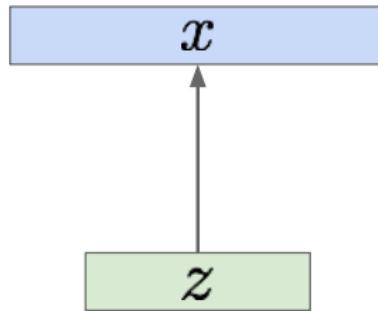
- Sample data \mathbf{x} from true conditional $p_{\theta^*}(\mathbf{x} | \mathbf{z})$

How should we represent this model?

Variational Autoencoder \Rightarrow Probabilistic Autoencoder

- 1 Define a generative probability model for the decoder.
- 2 Assume training data $\{\mathbf{x}^{(i)}\}_{i=1}^N$ is generated from underlying unobserved (latent) representation \mathbf{z} .

- Sample latent \mathbf{z} from true prior $p_{\theta^*}(\mathbf{z})$



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

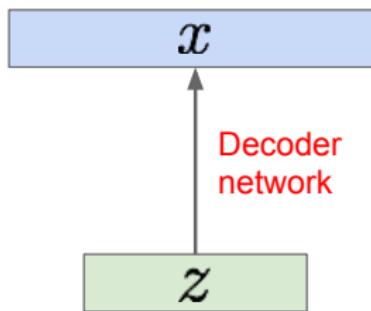
How should we represent this model?

- Choose prior $p(z)$ to be simple, e.g. Gaussian.

Variational Autoencoder \Rightarrow Probabilistic Autoencoder

- 1 Define a generative probability model for the deconder.
- 2 Assume training data $\{\mathbf{x}^{(i)}\}_{i=1}^N$ is generated from underlying unobserved (latent) representation \mathbf{z} .

- Sample latent \mathbf{z} from true prior $p_{\theta^*}(\mathbf{z})$
- Sample data \mathbf{x} from true conditional $p_{\theta^*}(\mathbf{x} | \mathbf{z})$

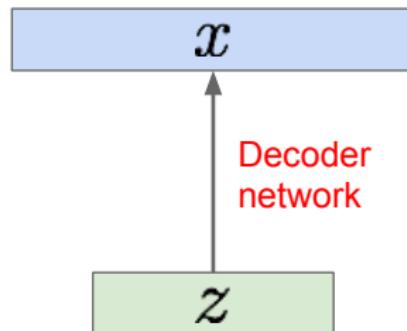


How should we represent this model?

- Choose prior $p(\mathbf{z})$ to be simple, e.g. Gaussian.
- Choose $p(\mathbf{x} | \mathbf{z})$ to be complex (generates image) \rightarrow represented with a neural network.

Variational Autoencoder \Rightarrow Probabilistic Autoencoder

- 1 Define a generative probability model for the deconder.
- 2 Assume training data $\{\mathbf{x}^{(i)}\}_{i=1}^N$ is generated from underlying unobserved (latent) representation \mathbf{z} .



- Sample latent \mathbf{z} from true prior $p_{\theta^*}(\mathbf{z})$
- Sample data \mathbf{x} from true conditional $p_{\theta^*}(\mathbf{x} \mid \mathbf{z})$

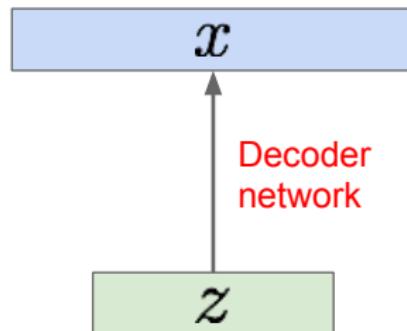
How to train the model?

- Maximum likelihood?

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x} \mid \mathbf{z}) p_{\theta}(\mathbf{z}) d\mathbf{z}$$

Variational Autoencoder \Rightarrow Probabilistic Autoencoder

- 1 Define a generative probability model for the deconder.
- 2 Assume training data $\{\mathbf{x}^{(i)}\}_{i=1}^N$ is generated from underlying unobserved (latent) representation \mathbf{z} .



- Sample latent \mathbf{z} from true prior $p_{\theta^*}(\mathbf{z})$
- Sample data \mathbf{x} from true conditional $p_{\theta^*}(\mathbf{x} \mid \mathbf{z})$

How to train the model?

- Maximum likelihood?

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x} \mid \mathbf{z}) p_{\theta}(\mathbf{z}) d\mathbf{z}$$

Intractable \Rightarrow Variational Inference!

Kingma & Welling, ICLR 2014

Variational Autoencoder: Intractability

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

Variational Autoencoder: Intractability

- Data likelihood: $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x} | \mathbf{z}) \underbrace{p_{\theta}(\mathbf{z})}_{\text{simple Gaussian prior}} d\mathbf{z}$

Variational Autoencoder: Intractability

- Data likelihood: $p_{\theta}(\mathbf{x}) = \int \underbrace{p_{\theta}(\mathbf{x} | \mathbf{z})}_{\text{decoder neural network}} p_{\theta}(\mathbf{z}) d\mathbf{z}$
- typically model as:

$$p_{\theta}(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x}; \mu_{\theta}(\mathbf{z}), \Sigma_{\theta}(\mathbf{z})) ,$$

where μ_{θ} and Σ_{θ} are two neural networks parameterized by θ .

Variational Autoencoder: Intractability

- Data likelihood: $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x} | \mathbf{z}) p_{\theta}(\mathbf{z}) d\mathbf{z}$
intractable \mathbf{x}

Variational Autoencoder: Intractability

- Data likelihood: $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x} | \mathbf{z}) p_{\theta}(\mathbf{z}) d\mathbf{z}$
intractible \mathbf{x}
- Posterior also intractable: $p_{\theta}(\mathbf{z} | \mathbf{x}) = \underbrace{p_{\theta}(\mathbf{x} | \mathbf{z})}_{\checkmark} \underbrace{p_{\theta}(\mathbf{z})}_{\checkmark} / \underbrace{p_{\theta}(\mathbf{x})}_{\times}$

Variational Autoencoder: Intractability

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

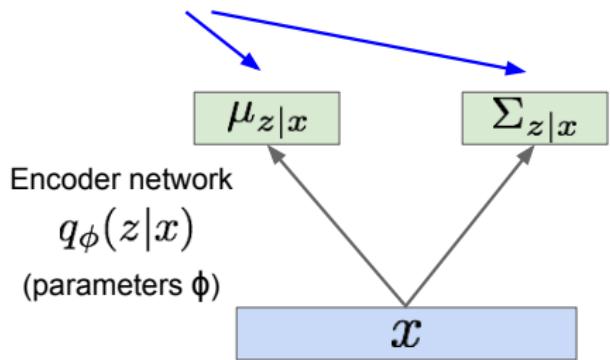
- Posterior also intractable: $p_{\theta}(\mathbf{z} | \mathbf{x}) = \underbrace{p_{\theta}(\mathbf{x} | \mathbf{z})}_{\checkmark} \underbrace{p_{\theta}(\mathbf{z})}_{\checkmark} / \underbrace{p_{\theta}(\mathbf{x})}_{\times}$
- Solution:
 - In addition to define the decoder network $p_{\theta}(\mathbf{x} | \mathbf{z})$, define an additional *encoder network* $q_{\phi}(\mathbf{z} | \mathbf{x})$ that approximates the posterior $p_{\theta}(\mathbf{z} | \mathbf{x}) \Rightarrow$ **variational inference with variational distribution $q_{\phi}(\mathbf{z} | \mathbf{x})$!**
 - We will see this allows us to derive a lower bound on the data likelihood, which can be optimized tractably.

Variational Autoencoder

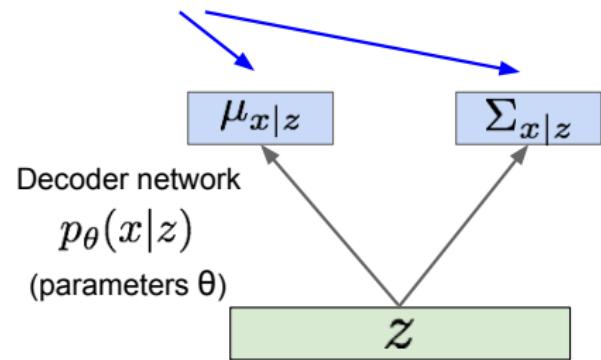
Encoder and decoder networks are probabilistic

- Encoder also called inference/recognition network.
- Decoder also called generation network.

Mean and (diagonal) covariance of $z | x$



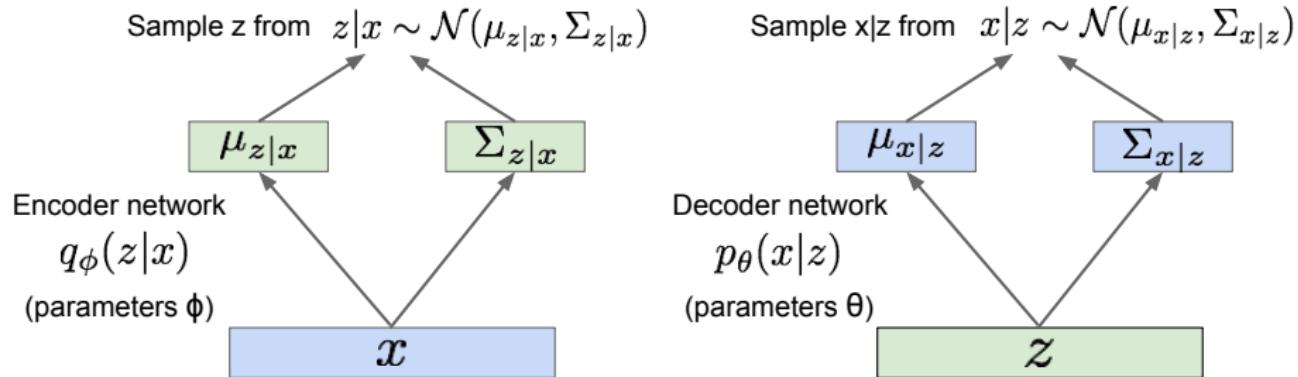
Mean and (diagonal) covariance of $x | z$



Variational Autoencoder

Encoder and decoder networks are probabilistic

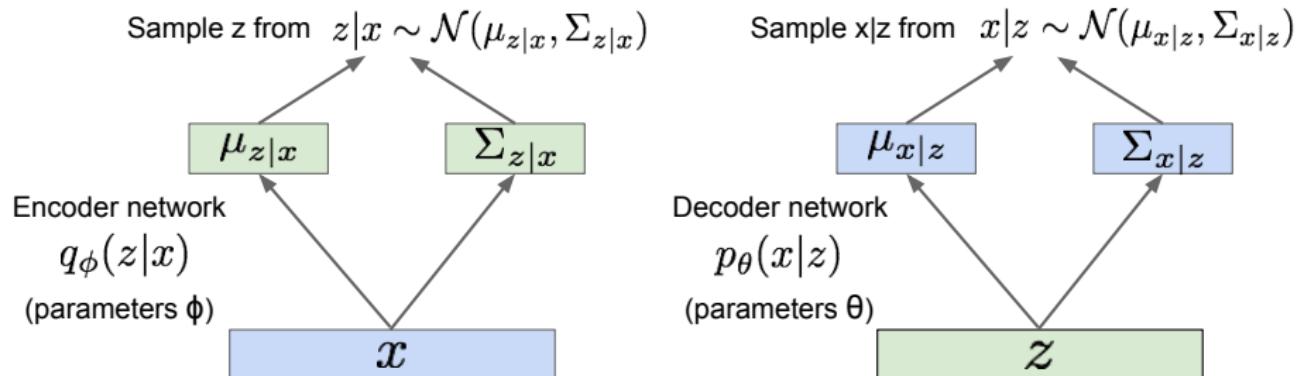
- Encoder also called inference/recognition network.
- Decoder also called generation network.



Variational Autoencoder

Encoder and decoder networks are probabilistic

- Encoder also called inference/recognition network.
- Decoder also called generation network.



Want to match $q_\phi(z|x)$ and $p(z|x) \propto p(z)p(x|z)$.

Variational Autoencoder: Log Data Likelihood

Derive from a different but equivalent way as in variational inference:

- Don't need to apply Jensen's inequality, or equivalently, prove Jensen's inequality.

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

Variational Autoencoder: Log Data Likelihood

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



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

Variational Autoencoder: Log Data Likelihood

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

Variational Autoencoder: Log Data Likelihood

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

Variational Autoencoder: Log Data Likelihood

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

Variational Autoencoder: Log Data Likelihood

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

Variational Autoencoder: Log Data Likelihood

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

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

Variational Autoencoder: Log Data Likelihood

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

↑
Decoder network gives $p_\theta(x|z)$, can
compute estimate of this term through
sampling. (Sampling differentiable
through reparam. trick, see paper.)

↑
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 Autoencoder: Log 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 Autoencoder: Log Data Likelihood

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

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

Variational lower bound ("ELBO")

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

Training: Maximize lower bound

Variational Autoencoder: Log Data Likelihood

Reconstruct the input data

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

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

Variational lower bound ("ELBO")

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

Training: Maximize lower bound

Variational Autoencoder: Log Data Likelihood

Putting it all together: maximizing the likelihood lower bound

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

Variational Autoencoder: Log Data Likelihood

Putting it all together: maximizing the likelihood lower bound

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

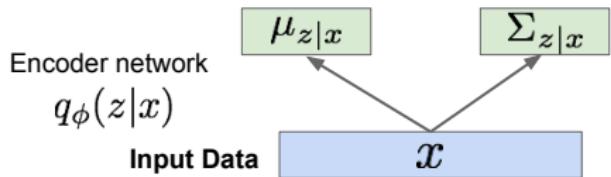
Let's look at computing the bound (forward pass) for a given minibatch of input data



Variational Autoencoder: Log Data Likelihood

Putting it all together: maximizing the likelihood lower bound

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



Variational Autoencoder: Log Data Likelihood

Putting it all together: maximizing the likelihood lower bound

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

Make approximate posterior distribution close to prior

Encoder network

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

Input Data

$$\mu_{z|x}$$

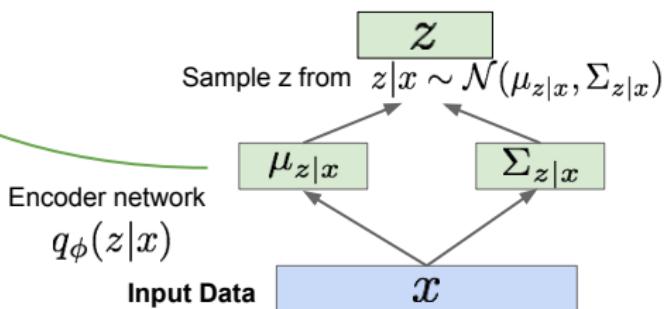
$$\Sigma_{z|x}$$

Variational Autoencoder: Log Data Likelihood

Putting it all together: maximizing the likelihood lower bound

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

Make approximate posterior distribution close to prior

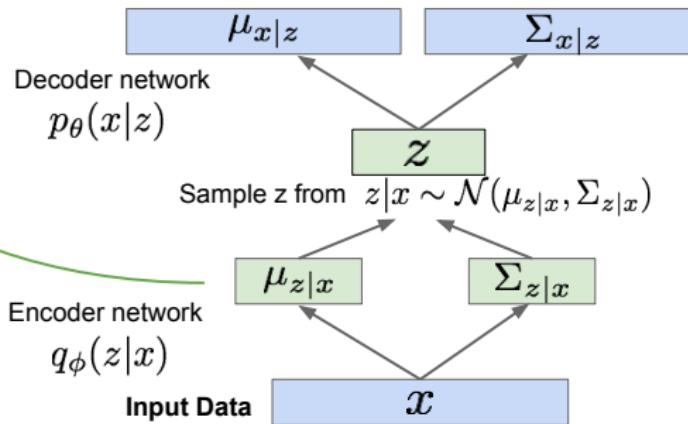


Variational Autoencoder: Log Data Likelihood

Putting it all together: maximizing the likelihood lower bound

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

Make approximate posterior distribution close to prior



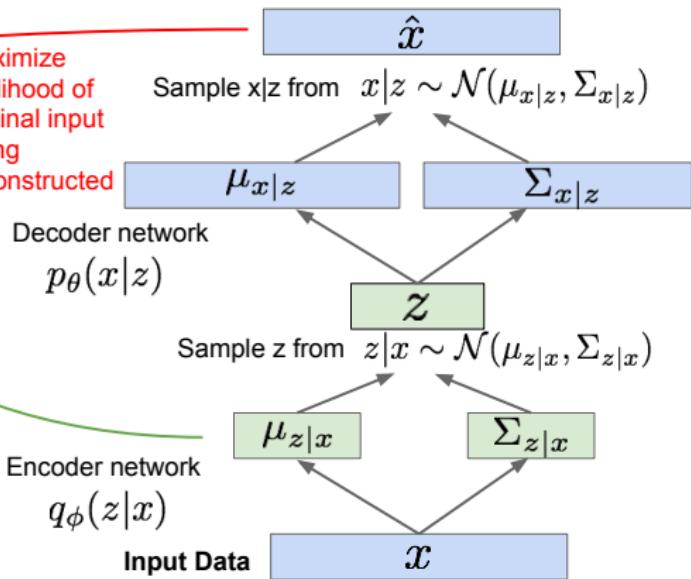
Variational Autoencoder: Log Data Likelihood

Putting it all together: maximizing the likelihood lower bound

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

Make approximate posterior distribution close to prior

Maximize likelihood of original input being reconstructed



Variational Autoencoder: Log Data Likelihood

Putting it all together: maximizing the likelihood lower bound

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

Make approximate posterior distribution close to prior

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

Maximize likelihood of original input being reconstructed

