

2 most hot topics on NIPS 2016?

2 most hot topics on NIPS 2016

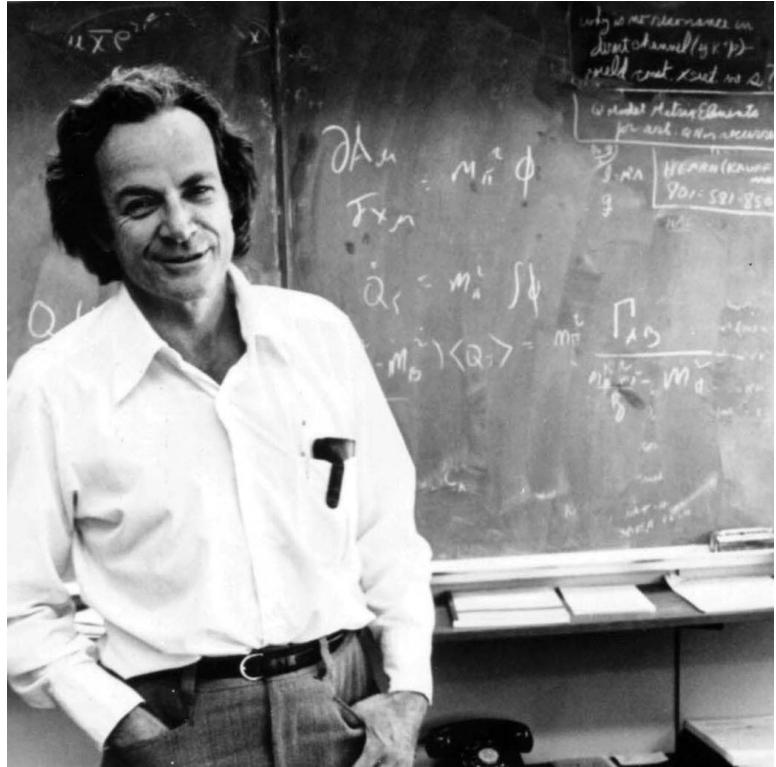
- Reinforcement Learning
- Generative models

Generative Models. Why?

Generative Models. Why?

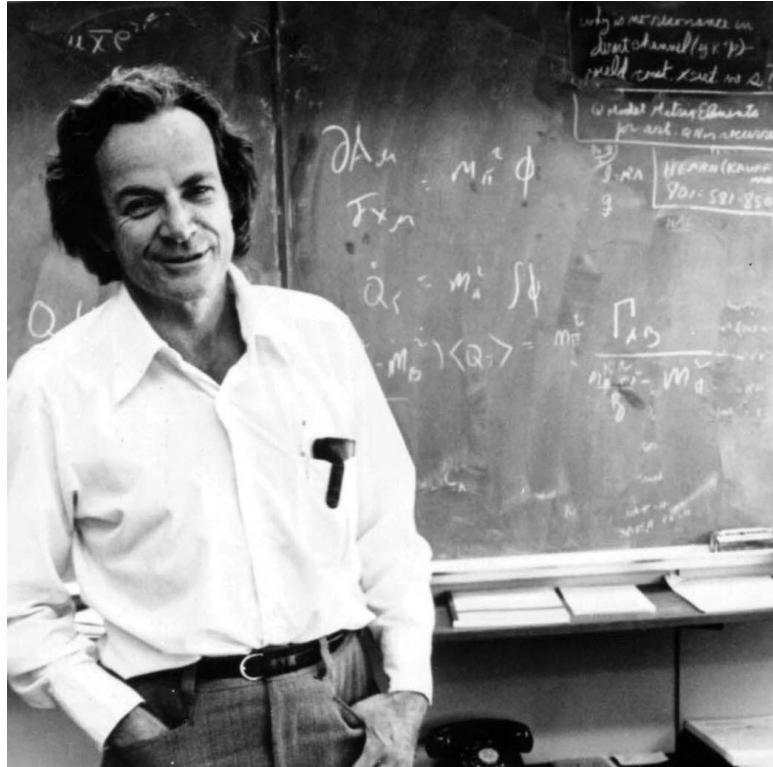
*What I cannot create,
I do not understand.*

Generative Models. Why?



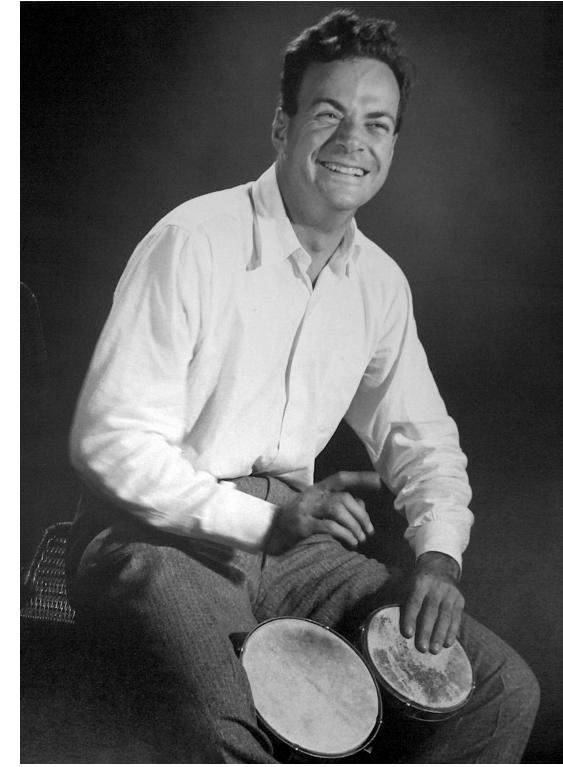
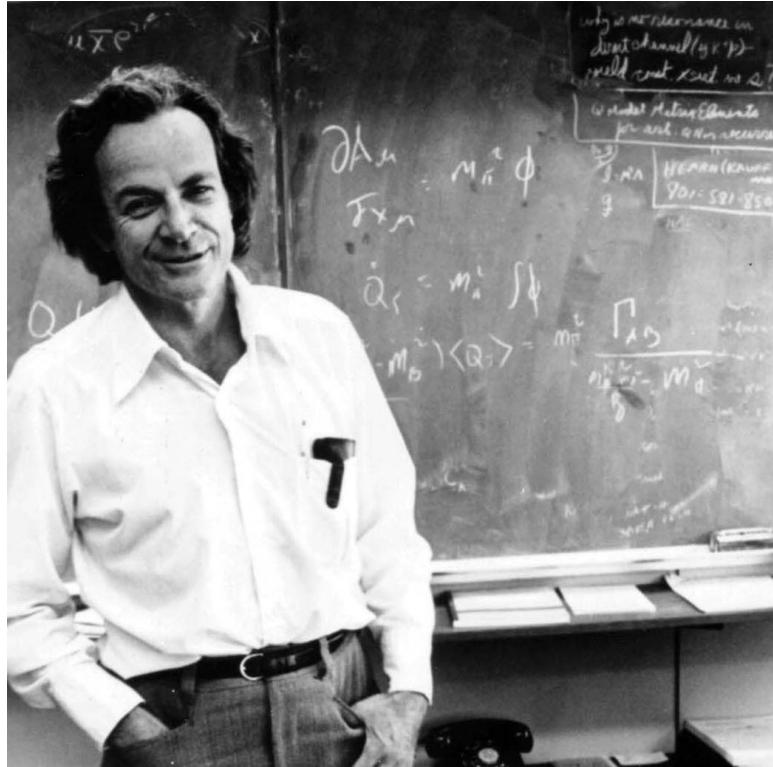
*What I cannot create,
I do not understand.
Richard Feynman*

Generative Models. Why?



*What I cannot create,
I do not understand.
Richard Feynman*

Generative Models. Why?



AI Cake

Reinforcement Learning



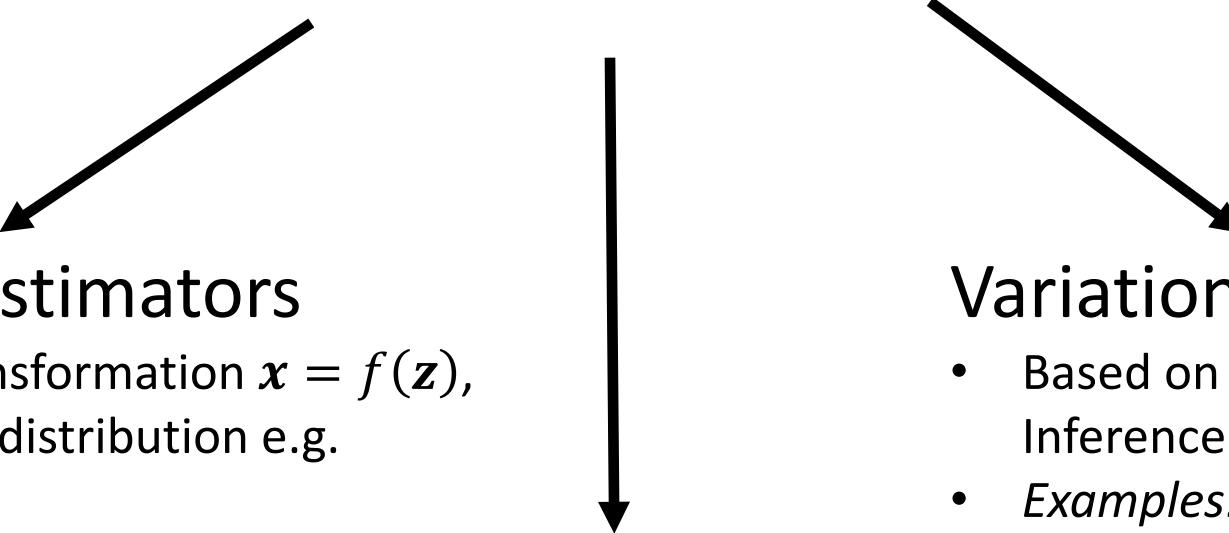
Supervised Learning

Unsupervised Learning

Why?

- Most of the data is unsupervised
- Humans learn in an unsupervised way

Generative Models of Images



Density Estimators

- Train Inversible Transformation $x = f(\mathbf{z})$, where \mathbf{z} have fixed distribution e.g. $\mathcal{N}(0,1)$
- $p(x) = p(\mathbf{z}) \det \left| \frac{\partial x}{\partial \mathbf{z}} \right|^{-1}$
- Examples: NICE, Real NVP...
– Must have $\dim(\mathbf{z}) = \dim(\mathbf{x})$

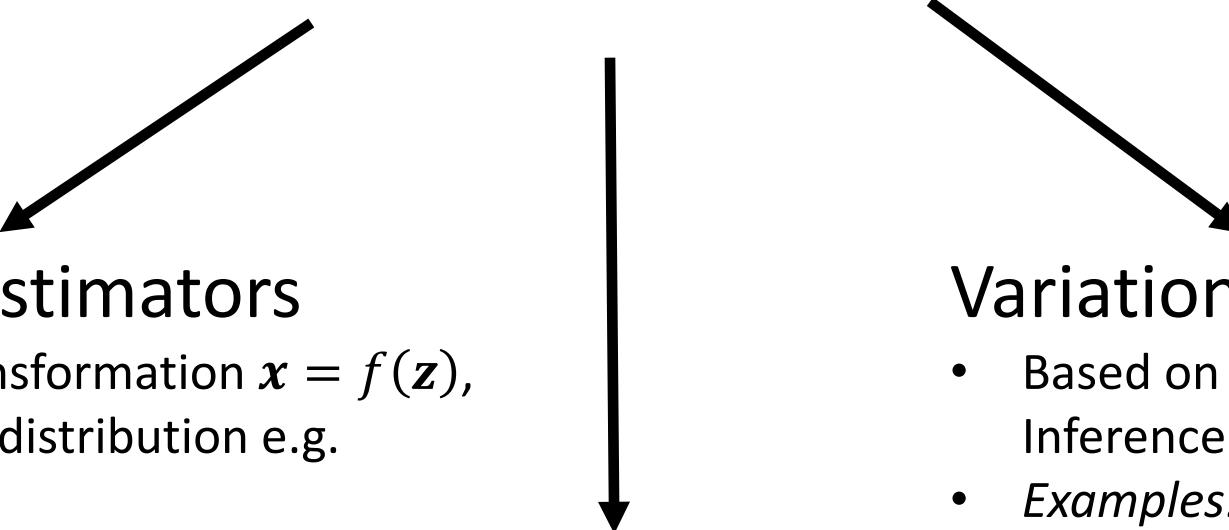
GANs

- Based on adversarial training
- Examples: GAN, InfoGAN, StackGAN ...
– Tricky to train
– Tricky to evaluate

Variational Autoencoders

- Based on approximate Bayesian Inference
- Examples: VAE, DRAW...
– Currently cannot achieve good samples quality (blurry samples)

Generative Models of Images



Density Estimators

- Train Inversible Transformation $x = f(\mathbf{z})$, where \mathbf{z} have fixed distribution e.g. $\mathcal{N}(0,1)$
- $p(x) = p(\mathbf{z}) \det \left| \frac{\partial x}{\partial \mathbf{z}} \right|^{-1}$
- Examples: NICE, Real NVP...
– Must have $\dim(\mathbf{z}) = \dim(\mathbf{x})$

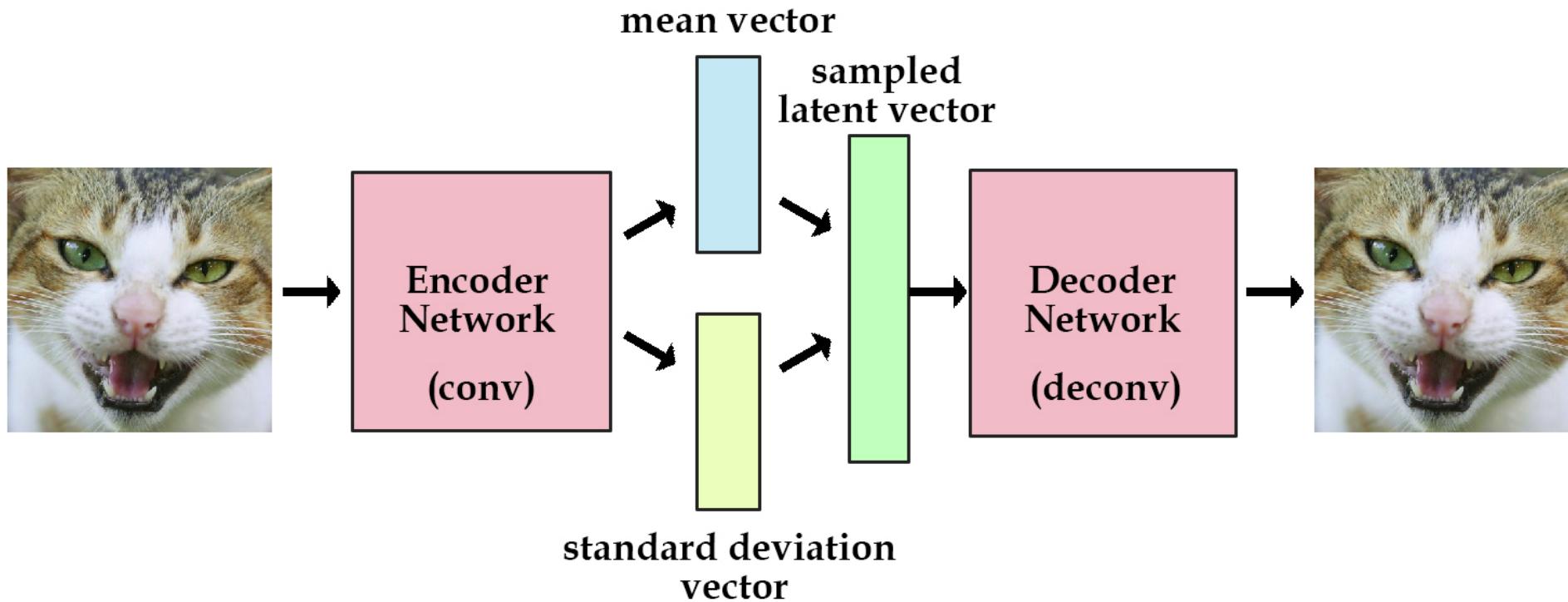
GANs

- Based on adversarial training
- Examples: GAN, InfoGAN, StackGAN ...
 - Tricky to train
 - Tricky to evaluate

Variational Autoencoders

- Based on approximate Bayesian Inference
- Examples: VAE, DRAW...
 - Currently cannot achieve good samples quality (blurry samples)

VAE – Variational AutoEncoder



VAE = Variational Autoencoder

Why variational?

VAE = Variational Autoencoder

Why variational?

Variational inference!

Bayesian ML in a nutshell

Recall Bayes formula

$$P(a \mid b) =$$

Recall Bayes formula

$$P(a \mid b) = \frac{P(b \mid a) \cdot P(a)}{P(b)}$$

Recall Bayes formula

$$P(a \mid b) = \frac{P(b \mid a) \cdot P(a)}{P(b)} = \frac{P(b \mid a) \cdot P(a)}{\int_a P(b \mid a) \cdot P(a) da}$$

Recall Bayes formula

$$P(a \mid b) = \frac{P(b \mid a) \cdot P(a)}{\int_a P(b \mid a) \cdot P(a) da}$$

Recall Bayes formula

$$P(\mathbf{z} \mid \mathbf{x}) = \frac{P(\mathbf{x} \mid \mathbf{z}) \cdot P(\mathbf{z})}{\int_{\mathbf{z}} P(\mathbf{x} \mid \mathbf{z}) \cdot P(\mathbf{z}) d\mathbf{z}}$$

Holds when $\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{R}^k$

Recall Bayes formula

$$P(\mathbf{z} \mid \mathbf{x}) = \frac{P(\mathbf{x} \mid \mathbf{z}) \cdot P(\mathbf{z})}{\int_{\mathbf{z}} P(\mathbf{x} \mid \mathbf{z}) \cdot P(\mathbf{z}) d\mathbf{z}}$$

Holds when $\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{R}^k$

- \mathbf{x} – observed data
- \mathbf{z} – unobserved latent variables

Recall Bayes formula

$$P(z | x) = \frac{P(x | z) \cdot P(z)}{\int_z P(x | z) \cdot P(z) dz}$$

- $P(z | x)$ – posterior
- $P(x | z)$ – likelihood
- $P(z)$ – prior

How to be Bayesian in 3(4) easy steps:

1. Make up your likelihood $P(\mathbf{x} \mid \mathbf{z})$
2. Make up your prior on $P(\mathbf{z})$ (regularization)
3. Infer $P(\mathbf{z} \mid \mathbf{x})$

How to be Bayesian in 3(4) easy steps:

1. Make up your likelihood $P(\mathbf{x} \mid \mathbf{z})$
2. Make up your prior on $P(\mathbf{z})$ (regularization)
3. Infer $P(\mathbf{z} \mid \mathbf{x})$
4. (optional) Find best model $P(\mathbf{x} \mid \mathbf{z}), P(\mathbf{z})$
maximizing $\int_{\mathbf{z}} P(\mathbf{x} \mid \mathbf{z}) \cdot P(\mathbf{z}) d\mathbf{z}$

How to be R easy: 6 steps:

1. Make a $P(z)$
2. Make a $P(z)$ (generalization)
3. Infer $P(z)$
4. Implement $P(z)$

Recall Bayes formula

$$P(z | x) = \frac{P(x | z) \cdot P(z)}{\int_z P(x | z) \cdot P(z) dz}$$

- What about denominator $\int_z P(x | z) \cdot P(z) dz$?

Recall Bayes formula

$$P(z | x) = \frac{P(x | z) \cdot P(z)}{\int_z P(x | z) \cdot P(z) dz}$$

- What about denominator $\int_z P(x | z) \cdot P(z) dz$?
- We **don't** need it if we only need Maximum a-posteriori (MAP)

Recall Bayes formula

$$P(\mathbf{z} \mid \mathbf{x}) = \frac{P(\mathbf{x} \mid \mathbf{z}) \cdot P(\mathbf{z})}{\int_{\mathbf{z}} P(\mathbf{x} \mid \mathbf{z}) \cdot P(\mathbf{z}) d\mathbf{z}}$$

- What about denominator $\int_{\mathbf{z}} P(\mathbf{x} \mid \mathbf{z}) \cdot P(\mathbf{z}) d\mathbf{z}$?
- We **don't** need it if we only need Maximum a-posteriori (MAP)
- However, it's essential if we want to
 1. Sample from $P(\mathbf{z} \mid \mathbf{x})$
 2. Compute mean/variance/almost anything about $P(\mathbf{z} \mid \mathbf{x})$

How to compute denominator?

$$\int_{\mathbf{z}} P(\mathbf{x} \mid \mathbf{z}) \cdot P(\mathbf{z}) d\mathbf{z} \quad - ???$$

How to compute denominator?

$$\int_{\mathbf{z}} P(\mathbf{x} \mid \mathbf{z}) \cdot P(\mathbf{z}) d\mathbf{z} \quad - ???$$

- Analytically
 - + Exact
 - Works only in some cases

How to compute denominator?

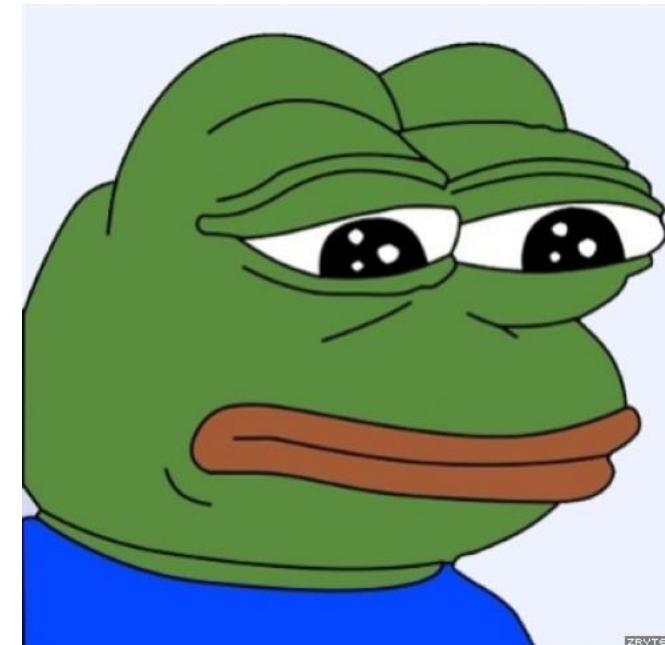
$$\int_{\mathbf{z}} P(\mathbf{x} \mid \mathbf{z}) \cdot P(\mathbf{z}) d\mathbf{z} \quad - ???$$

- Analytically
 - + Exact
 - Works only in some cases
- Numerically
 - + Always possible
 - Takes **AGES** if $k > 10$

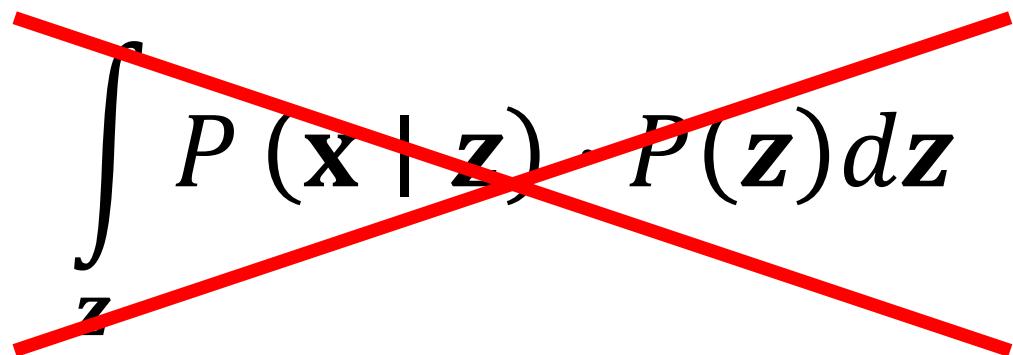
How to compute denominator?

$$\int_z P(\mathbf{x} | \mathbf{z}) \cdot P(\mathbf{z}) d\mathbf{z} \quad - ???$$

- Analytically
 - + Exact
 - Works only in some cases
- Numerically
 - + Always possible
 - Takes **AGES** if $k > 10$



We can't compute posterior...

$$\int_z^1 P(\mathbf{x} | \mathbf{z}) \cdot P(\mathbf{z}) d\mathbf{z}$$
A large red 'X' is drawn across the entire equation, indicating that it is incorrect or cannot be computed.

We can't compute posterior... Approximate it!

$$\int_z^1 P(\mathbf{x} | \mathbf{z}) \cdot P(\mathbf{z}) d\mathbf{z}$$

$$q(\mathbf{z}) \approx p(\mathbf{z} | \mathbf{x})$$

- $q(\mathbf{z})$ – easy to work with
- $p(\mathbf{z} | \mathbf{x})$ – true posterior

Two ideas how to approximate $P(\mathbf{z} \mid \mathbf{x})$

- MCMC – Markov Chain Monte Carlo
 - *Idea:* Sample a lot of stuff that looks like $P(\mathbf{z} \mid \mathbf{x})$ somehow
 - *Examples:* Metropolis-Hastings, Gibbs sampling
- VI – Variational Inference
 - *Idea:* Find the best approximation in some parametric family $q(\mathbf{z}) = q_\phi(\mathbf{z})$
 - *Examples:* Mean-Field, SGVB

Two ideas how to approximate $P(\mathbf{z} \mid \mathbf{x})$

- MCMC – Markov Chain Monte Carlo
 - *Idea:* Sample a lot of stuff that looks like $P(\mathbf{z} \mid \mathbf{x})$ somehow
 - *Examples:* Metropolis-Hastings, Gibbs sampling
- VI – Variational Inference
 - *Idea:* Find the best approximation in some parametric family $q(\mathbf{z}) = q_\phi(\mathbf{z})$
 - *Examples:* Mean-Field, SGVB

Two ideas how to approximate $P(\mathbf{z} \mid \mathbf{x})$

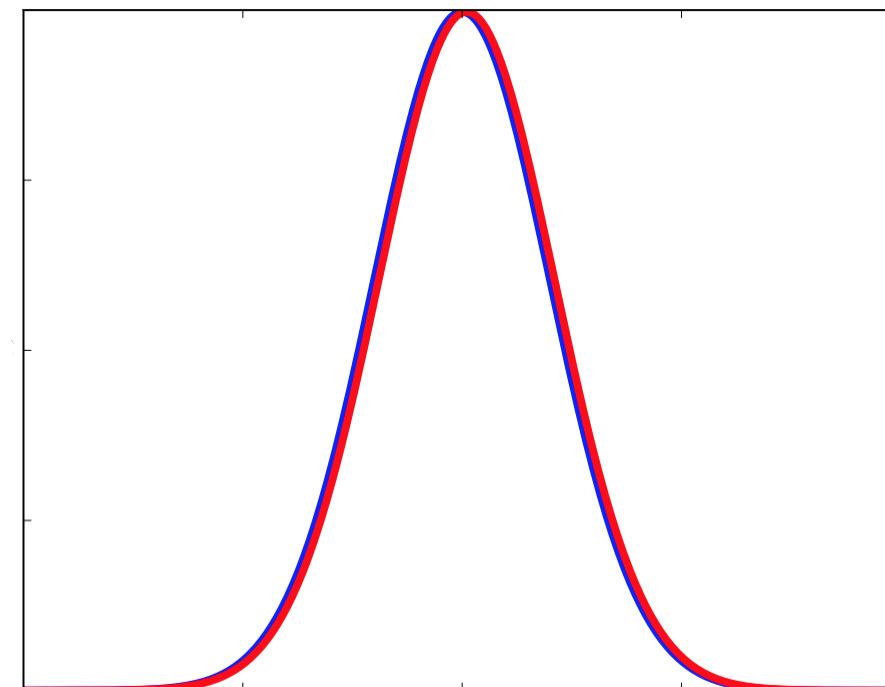
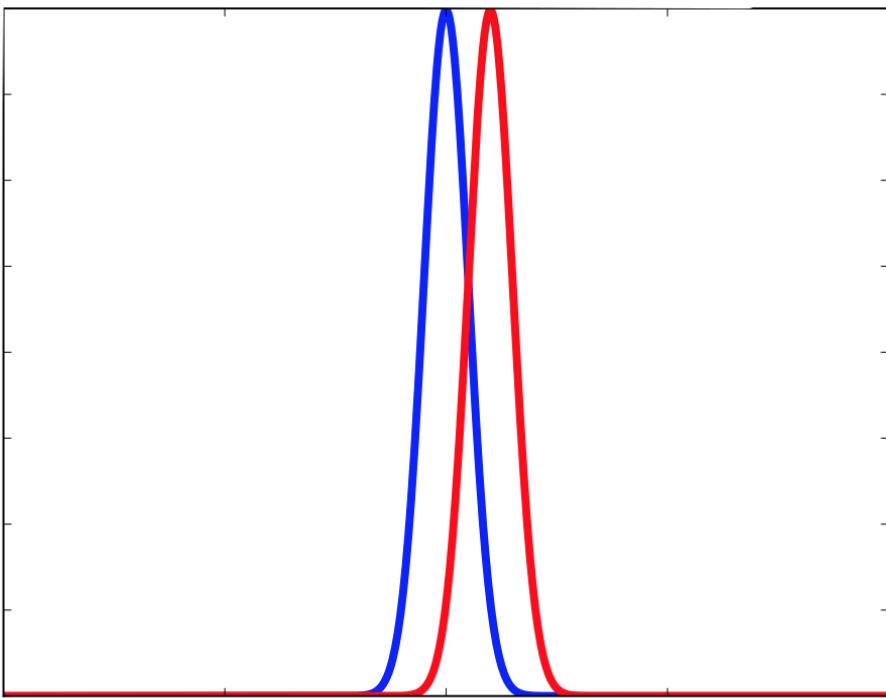
- MCMC – Markov Chain Monte Carlo
 - *Idea:* Sample a lot of stuff that looks like $P(\mathbf{z} \mid \mathbf{x})$ somehow
 - *Examples:* Metropolis-Hastings, Gibbs sampling
- VI – Variational Inference
 - *Idea:* Find the best approximation in some parametric family $q(\mathbf{z}) = q_\phi(\mathbf{z})$
 - *Examples:* Mean-Field, SGVB

Two ideas how to approximate $P(\mathbf{z} \mid \mathbf{x})$

- MCMC – Markov Chain Monte Carlo
 - *Idea:* Sample a lot of stuff that looks like $P(\mathbf{z} \mid \mathbf{x})$ somehow
 - *Examples:* Metropolis-Hastings, Gibbs sampling
- VI – Variational Inference
 - *Idea:* Find the best approximation in some parametric family $q(\mathbf{z}) = q_\phi(\mathbf{z})$
 - *Examples:* Mean-Field, SGVB

How to evaluate quality of approximation?

$$q(\mathbf{z}) \approx p(\mathbf{z} | \mathbf{x})$$



How to evaluate quality of approximation?

Kullback-Leibler divergence:

$$KL(p_1 \parallel p_2) = \mathbb{E}_{z \sim p_1} \log \frac{p_1(\mathbf{z})}{p_2(\mathbf{z})} = \int_{\mathbf{z}} p_1(\mathbf{z}) \log \frac{p_1(\mathbf{z})}{p_2(\mathbf{z})} d\mathbf{z}$$

- $KL(p_1 \parallel p_2) \geq 0$
- $KL(p_1 \parallel p_2) = 0 \Leftrightarrow p_1 = p_2$
- $KL(p_1 \parallel p_2) \neq KL(p_2 \parallel p_1)$

How to evaluate quality of approximation?

Kullback-Leibler divergence:

$$KL(q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})) = \mathbb{E}_q \log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x})} = \int_{\mathbf{z}} q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x})} d\mathbf{z}$$

ELBO: Evidence Lower Bound

$\log p(x)$

ELBO: Evidence Lower Bound

$\log p(x) = \mathbb{E}_{z \sim q} \log p(x) // \text{we integrate over } z \text{ here. } \int_z q(\mathbf{z}) d\mathbf{z} = 1$

ELBO: Evidence Lower Bound

$$\begin{aligned}\log p(\mathbf{x}) &= \mathbb{E}_{\mathbf{z} \sim q} \log p(\mathbf{x}) = \\ &= \mathbb{E}_{\mathbf{z} \sim q} \log \left[\frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \frac{q(\mathbf{z})}{q(\mathbf{z})} \right] \quad // p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})}\end{aligned}$$

ELBO: Evidence Lower Bound

$$\begin{aligned}\log p(\mathbf{x}) &= \mathbb{E}_{\mathbf{z} \sim q} \log p(\mathbf{x}) = \\&= \mathbb{E}_{\mathbf{z} \sim q} \log \left[\frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \frac{q(\mathbf{z})}{q(\mathbf{z})} \right] = \\&= \mathbb{E}_{\mathbf{z} \sim q} \left[\log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x})} + \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right]\end{aligned}$$

ELBO: Evidence Lower Bound

$$\begin{aligned}\log p(\mathbf{x}) &= \mathbb{E}_{\mathbf{z} \sim q} \log p(\mathbf{x}) = \\&= \mathbb{E}_{\mathbf{z} \sim q} \log \left[\frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \frac{q(\mathbf{z})}{q(\mathbf{z})} \right] = \\&= \mathbb{E}_{\mathbf{z} \sim q} \left[\log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x})} + \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right] = \\&= KL(q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})) + \mathbb{E}_{\mathbf{z} \sim q} [\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z})]\end{aligned}$$

ELBO: Evidence Lower Bound

$$\begin{aligned}\log p(\mathbf{x}) &= \mathbb{E}_{\mathbf{z} \sim q} \log p(\mathbf{x}) = \\&= \mathbb{E}_{\mathbf{z} \sim q} \log \left[\frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \frac{q(\mathbf{z})}{q(\mathbf{z})} \right] = \\&= \mathbb{E}_{\mathbf{z} \sim q} \left[\log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x})} + \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right] = \\&= KL(q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})) + \mathbb{E}_{\mathbf{z} \sim q} [\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z})] \\&\geq 0 && \mathcal{L}(q)\end{aligned}$$

ELBO: Evidence Lower Bound

$$\log p(\mathbf{x}) = KL(q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})) + \mathcal{L}(q)$$

$$\mathcal{L}(q) \rightarrow \max_q \quad \Leftrightarrow \quad KL(q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})) \rightarrow \min_q$$

ELBO: Evidence Lower Bound

$$\log p(\mathbf{x}) = KL(q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})) + \mathcal{L}(q)$$

$$\mathcal{L}(q) \rightarrow \max_q \iff KL(q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})) \rightarrow \min_q$$

Let's maximize ELBO instead of minimizing KL!

- Now we don't need to know how $p(\mathbf{z}|\mathbf{x})$ looks like

ELBO: Evidence Lower Bound

$$\log p(\mathbf{x}) = KL(q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})) + \mathcal{L}(q)$$

$$\mathcal{L}(q) \rightarrow \max_q \iff KL(q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})) \rightarrow \min_q$$

Let's maximize ELBO instead of minimizing KL!

- Now we don't need to know how $p(\mathbf{z}|\mathbf{x})$ looks like
- Note 1: we can maximize \mathcal{L} wrt $q(\mathbf{z})$ and $p(\mathbf{x}|\mathbf{z})$ finding best model and making inference simultaneously!

ELBO: Evidence Lower Bound

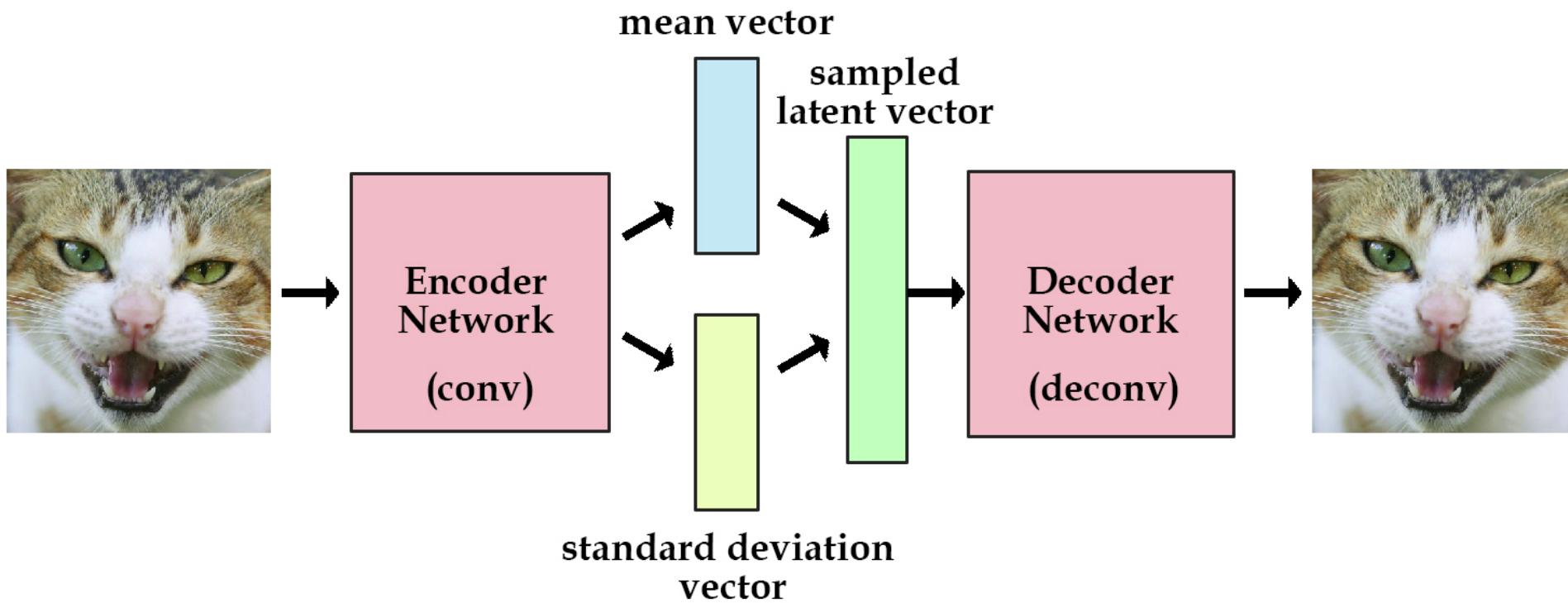
$$\log p(\mathbf{x}) = KL(q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})) + \mathcal{L}(q)$$

$$\mathcal{L}(q) \rightarrow \max_q \iff KL(q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x})) \rightarrow \min_q$$

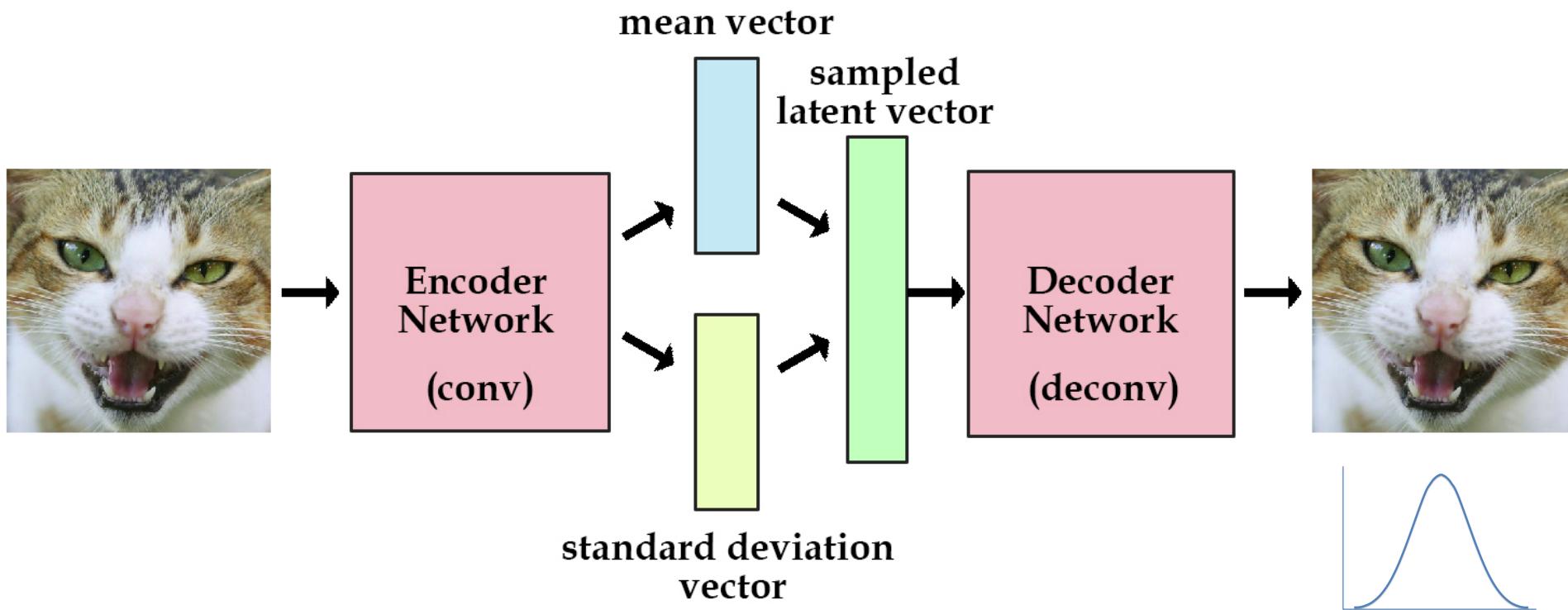
Let's maximize ELBO instead of minimizing KL!

- Now we don't need to know how $p(\mathbf{z}|\mathbf{x})$ looks like
- Note 1: we can maximize \mathcal{L} wrt $q(\mathbf{z})$ and $p(\mathbf{x}|\mathbf{z})$ finding best model and making inference simultaneously!
- Note 2: $q(\mathbf{z})$ can be parameterized as a function of \mathbf{x} : $q_\phi(\mathbf{z}|\mathbf{x})$

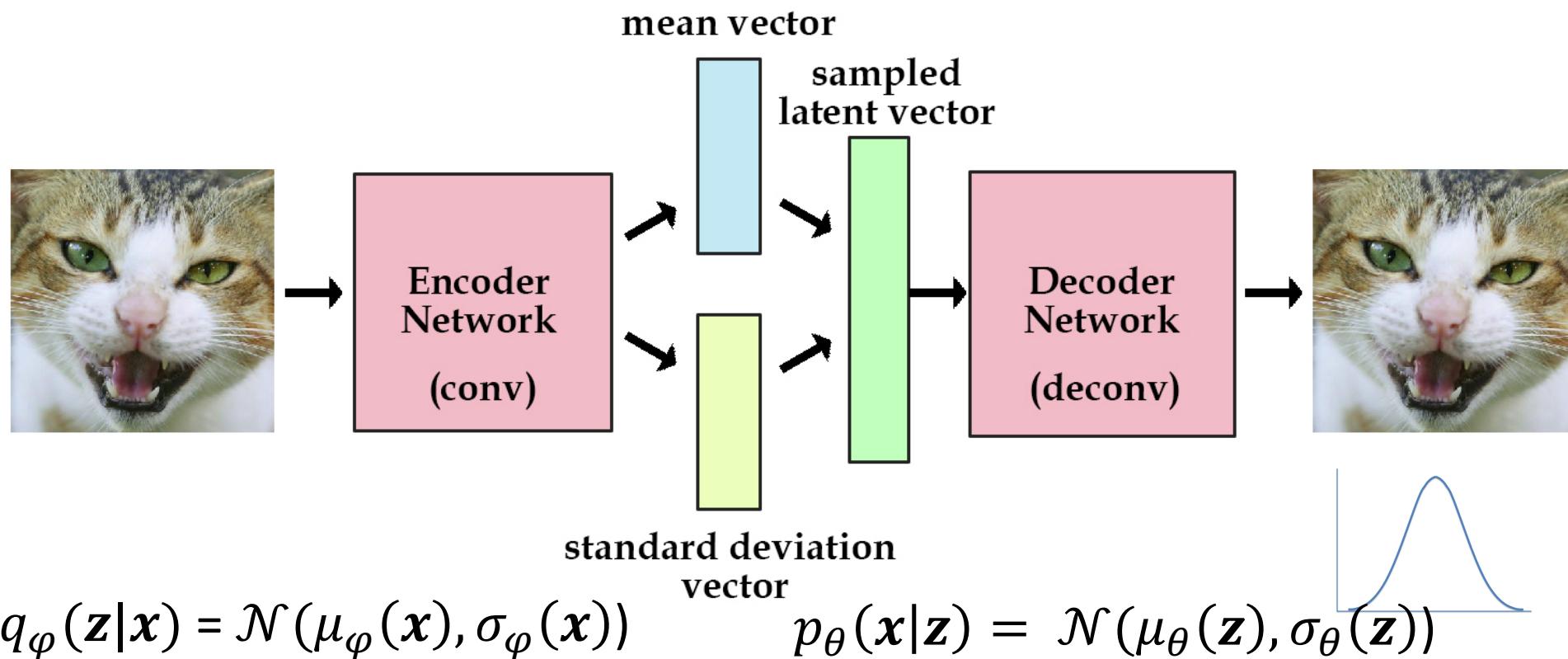
Variational Autoencoder



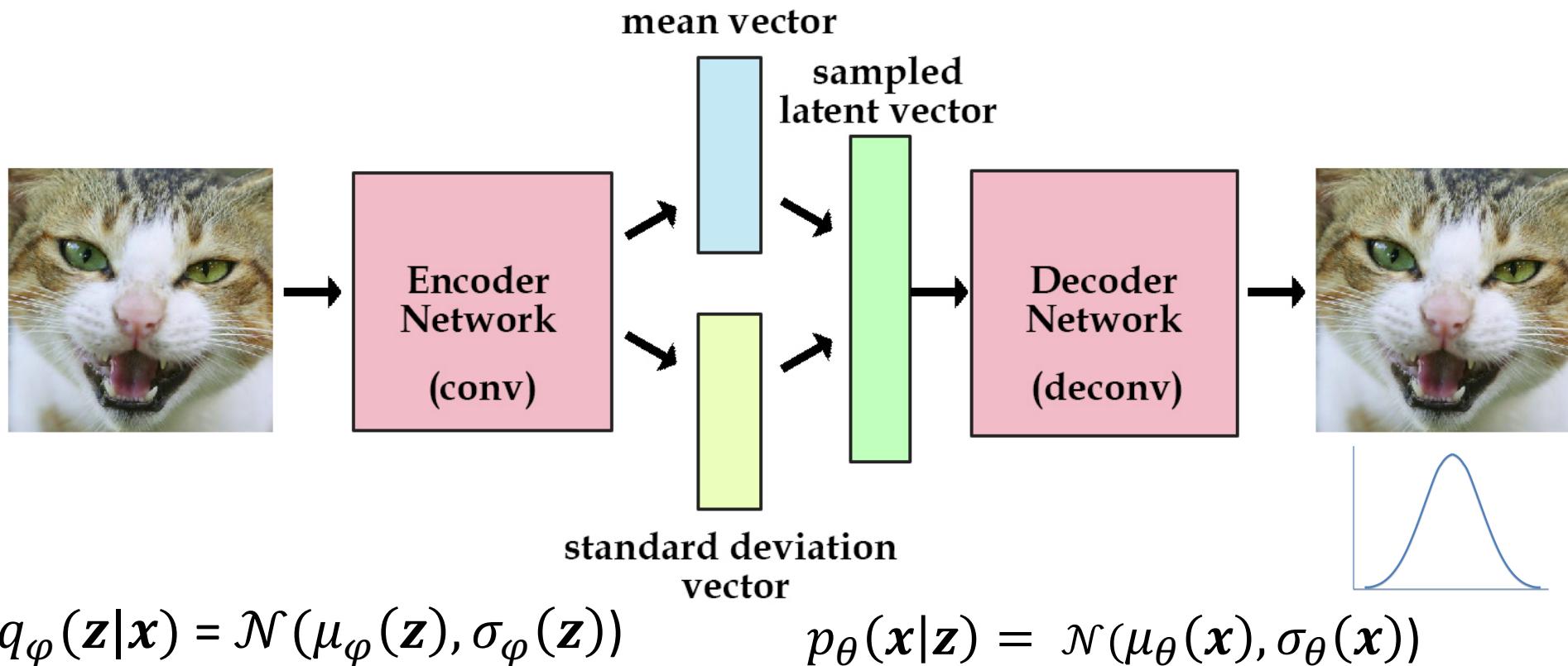
Variational Autoencoder



Variational Autoencoder



$$\mathcal{L}(q) = \mathbb{E}_{z \sim q} [\log p_{\theta}(x, z) - \log q_{\varphi}(z)] \rightarrow \max_{\varphi, \theta}$$



$$\mathcal{L}(q) = \mathbb{E}_{z \sim q} \left[\log p_{\theta}(x, z) - \log q_{\varphi}(z|x) \right] \rightarrow \max_{\varphi, \theta}$$

$$\mathcal{L}(q) = \mathbb{E}_{z \sim q} \left[\log p_{\theta}(x|z) + \log p_{\theta}(z) - \log q_{\varphi}(z|x) \right]$$

$$q_{\varphi}(z|x)=\mathcal{N}(\mu_{\varphi}(z),\sigma_{\varphi}(z))$$

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

$$\mathcal{L}(q) = \mathbb{E}_{z \sim q} [\log p_\theta(x, z) - \log q_\varphi(z|x)] \rightarrow \max_{\varphi, \theta}$$

$$\begin{aligned}\mathcal{L}(q) &= \mathbb{E}_{z \sim q} [\log p_\theta(x|z) + \log p_\theta(z) - \log q_\varphi(z|x)] = \\ &= \mathbb{E}_{z \sim q} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q} [\log q_\varphi(z|x) - \log p_\theta(z)]\end{aligned}$$

$$q_\varphi(z|x) = \mathcal{N}(\mu_\varphi(z), \sigma_\varphi(z))$$

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

$$\mathcal{L}(q) = \mathbb{E}_{z \sim q} [\log p_\theta(x, z) - \log q_\varphi(z|x)] \rightarrow \max_{\varphi, \theta}$$

$$\begin{aligned}\mathcal{L}(q) &= \mathbb{E}_{z \sim q} [\log p_\theta(x|z) + \log p_\theta(z) - \log q_\varphi(z|x)] = \\ &= \mathbb{E}_{z \sim q} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q} [\log q_\varphi(z|x) - \log p_\theta(z)] = \\ &= \mathbb{E}_{z \sim q} [\log p_\theta(x|z)] - KL(q_\varphi(z|x) \parallel p_\theta(z))\end{aligned}$$

$$q_\varphi(z|x) = \mathcal{N}(\mu_\varphi(z), \sigma_\varphi(z))$$

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

$$\mathcal{L}(q) = \mathbb{E}_{z \sim q} [\log p_\theta(x, z) - \log q_\varphi(z|x)] \rightarrow \max_{\varphi, \theta}$$

$$\begin{aligned}\mathcal{L}(q) &= \mathbb{E}_{z \sim q} [\log p_\theta(x|z) + \log p_\theta(z) - \log q_\varphi(z|x)] = \\ &= \mathbb{E}_{z \sim q} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q} [\log q_\varphi(z|x) - \log p_\theta(z)] = \\ &= \mathbb{E}_{z \sim q} [\log p_\theta(x|z)] - KL(q_\varphi(z|x) \parallel p_\theta(z)) \rightarrow \max_{\varphi, \theta}\end{aligned}$$

$$q_\varphi(z|x) = \mathcal{N}(\mu_\varphi(z), \sigma_\varphi(z))$$

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

$$\mathcal{L}(q) = \mathbb{E}_{z \sim q} [\log p_\theta(x, z) - \log q_\varphi(z|x)] \rightarrow \max_{\varphi, \theta}$$

$$\begin{aligned}\mathcal{L}(q) &= \mathbb{E}_{z \sim q} [\log p_\theta(x|z) + \log p_\theta(z) - \log q_\varphi(z|x)] = \\ &= \mathbb{E}_{z \sim q} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q} [\log q_\varphi(z|x) - \log p_\theta(z)] = \\ &= \mathbb{E}_{z \sim q} [\log p_\theta(x|z)] - KL(q_\varphi(z|x) \parallel p_\theta(z)) \rightarrow \max_{\varphi, \theta}\end{aligned}$$

Reconstruction error *Regularization term*

$$q_\varphi(z|x) = \mathcal{N}(\mu_\varphi(z), \sigma_\varphi(z)) \quad p_\theta(x|z) = \mathcal{N}(\mu_\theta(x), \sigma_\theta(x))$$

$$\mathcal{L}(\varphi, \theta) = \mathbb{E}_{z \sim q} [\log p_\theta(x, z) - \log q_\varphi(z)] \rightarrow \max_{\varphi, \theta}$$

Training algorithm:

1. Take batch of data x . Sample batch of gaussian noise ε .
2. Pass x through encoder to get $\mu_\varphi(x), \sigma_\varphi(x)$.
3. $z = \sigma_\varphi(x)\varepsilon + \mu_\varphi(x)$
4. Pass z through decoder to get $\mu_\theta(z), \sigma_\theta(z)$
5. Compute and backpropagate estimate of loss $\mathcal{L}(\varphi, \theta)$:
 $\log p_\theta(x, z) - \log q_\varphi(z)$
6. Make gradient updates for decoder and encoder with your favourite SGD-like algorithm

Thank you for your attention!