

Variational Autoencoders (VAEs)

Sargur N. Srihari

srihari@cedar.buffalo.edu

Topics in VAE

1. Generative Model
2. Black-box variational inference
3. The reparameterization trick
4. Choosing q and p
5. Autoencoding Variational Bayes (AEVB)
6. Variational autoencoder (VAE)
7. VAE: The neural network perspective
<https://cedar.buffalo.edu/~srihari/CSE676/21.2-VAE-NeuralNets.pdf>
8. Examples

Generative Models

- *A generative model is*
 - Able to generate samples like those in training data
 - With MNIST: samples would be synthetic images of digits
- Any of these samples can be decoded
 - Into a reasonable image of a handwritten digit

VAE is a Generative Model

- VAE is a directed generative model
- It has both observed and latent variables
 - Gives us a latent space to sample from
- It has emerged as a popular unsupervised learning method

VAE and Deep Learning

- A VAE models a data distribution using a collection of independent latent variables
 - A deep learning technique for learning latent representations
 - A hybrid of deep learning and variational inference

Uses of VAEs

- They have been used to
 1. Draw images
 2. State-of-the-art results in semi-supervised learning
 3. Interpolate between sentences

A simple generative model

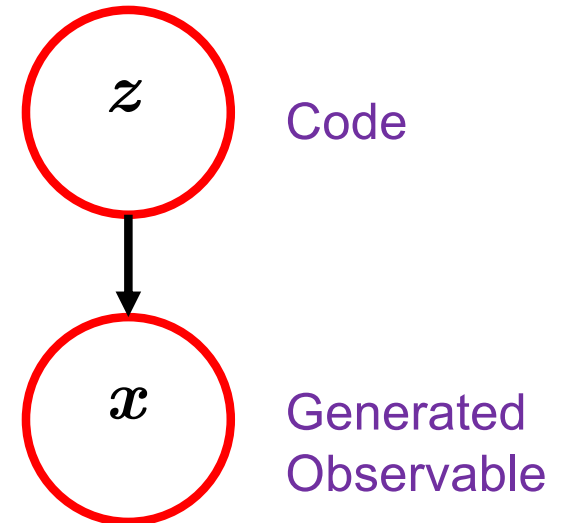
- Consider a directed model:

$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})$$

- \mathbf{x} represents observed data
- \mathbf{z} is a collection of latent variables
 - $\mathbf{x} \in \mathcal{X}$, where \mathcal{X} continuous or discrete
 - latent $\mathbf{z} \in \mathbb{R}^k$

- Example

- \mathbf{x} can be an image (e.g., human face)
- \mathbf{z} : latent factors not seen during training that explain features of face
 - E.g., a coordinate of \mathbf{z} encodes happy or sad, another male or female



VAE applied to face images



Variational autoencoder $p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$
applied to face images (modeled by \mathbf{x})

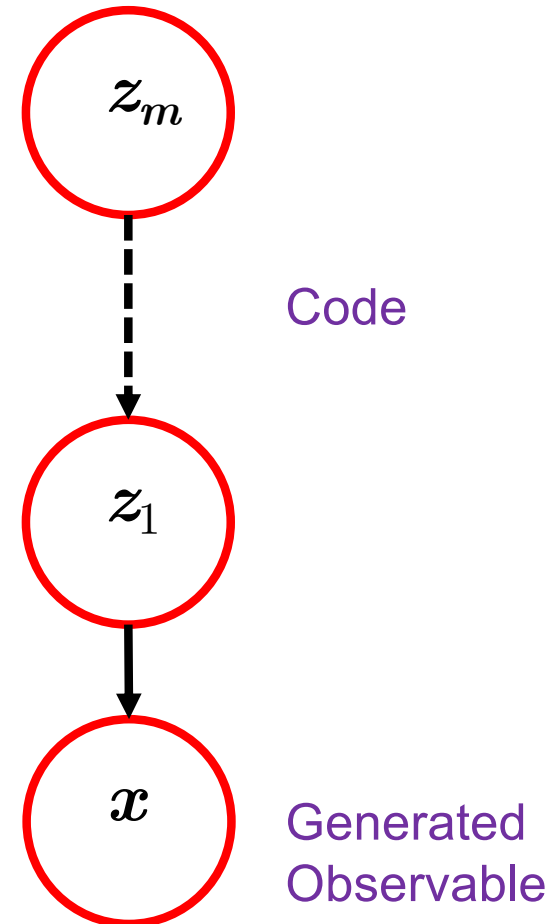
The learned latent space \mathbf{z} can be used
to interpolate between facial expressions

Deep Generative Model

- A model with many layers

$$p(x|z_1)p(z_1|z_2)p(z_2|z_3)\cdots p(z_{m-1}|z_m)p(z_m)$$

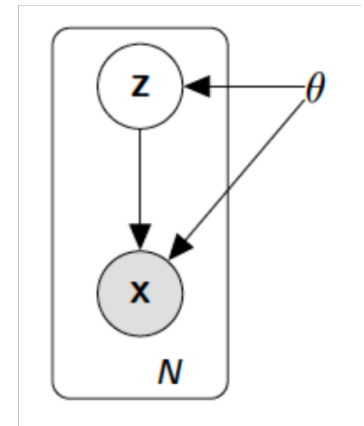
- This is a hierarchy of latent representations



Our Objective

- The setup

- \mathbf{x} : observed variables
- \mathbf{z} : latent variables (continuous)
- θ : model parameters
- A specified joint pdf $p_{\theta}(\mathbf{x}, \mathbf{z})$
- Posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$ is intractable



- Objectives

- Infer an (approximate) posterior over latent \mathbf{z} :
 - Given \mathbf{x} , what are its latent factors? Determine $p(\mathbf{z}|\mathbf{x})$
- Learn the parameters θ (obtain a MAP estimate)
OR an appropriate posterior over the parameters
 - We are given a dataset $D = \{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n\}$

Ex: Discrete Child with Continuous Parent

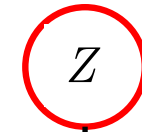
• Thermostat

$$P(x^1) = \begin{cases} 0.9 & z \leq 65 \\ 0.05 & \text{otherwise} \end{cases}$$

$$P(x^0) = 1 - P(x^1)$$



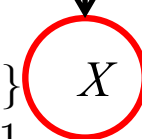
Continuous



Binary

$$\text{Val}(X) = \{x^0, x^1\}$$

$$P(x^0) + P(x^1) = 1$$



– Z is the temperature in Fahrenheit and X is the thermostat turning the heater on

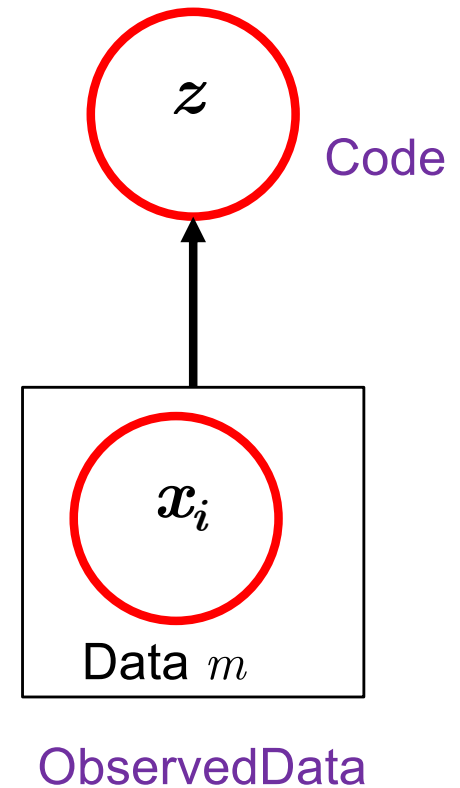
- Threshold model has abrupt change in probability with Z , i.e., from 0.9 to 0.05

– Can use the logistic model to fix this

$$P(x^1 | z) = \text{sigmoid}(w_0 + w_1 z)$$

Additional assumptions made

- Intractability:
 - computing the posterior probability $p_{\theta}(z|x)$ is intractable
- Big data:
 - the dataset D is too large to fit in memory; we can only work with small, subsampled batches of D



Trying out standard approaches

1. *EM algorithm* to learn latent-variable models
 2. *Mean Field Inference* to perform approximate inference
 3. *Sampling-based* approximate inference
- Let us see what the drawbacks of each are next

EM to learn latent variables

- E step requires computing the approximate posterior $p_{\theta}(z|x)$, which we have assumed to be intractable.
- In the M step, we learn the θ by looking at the entire dataset
 - Note, however, that there exists a generalization called online EM, which performs the M -step over mini-batches, which is going to be too large to hold in memory.

Mean Field Inference

- To perform approximate inference, $p_{\theta}(z|x)$ consider mean field
- One step of mean field requires us to compute an expectation
 - whose time complexity scales exponentially with the size of the Markov blanket of the target variable

Size of Markov Blanket

If we assume at least one component of x depends on each component of z , then this introduces a V-structure into the graph of our model (the x , which are observed, are explaining away the differences among the z)

Thus, we know that all the z variables will depend on each other and the Markov blanket of some z_i will contain all the other z -variables.

This will make mean-field intractable

An equivalent (and simpler) explanation is that p will contain a factor $p(x_i|z_1, \dots, z_k)$, in which all the z variables are tied

Sampling-based methods

- Compare VAE against sampling-based algorithms
 - Sampling-based methods don't scale well to large datasets.
- Techniques such as Metropolis-Hastings require a hand-crafted proposal distribution, which might be difficult to choose
 - <https://arxiv.org/abs/1312.6114>

Auto Encoding Variational Bayes

- AEVB algorithm can efficiently solve our inference and learning tasks
 - VAE will be one instantiation of this algorithm
 - AEVB is based on ideas from variational inference
- We first look at the variational lower bound

1. In Gibbs, $-KL(q \parallel \tilde{p}) = E_q [\log \tilde{p}(X) - \log q(X)]$ is lower bound for $\log Z(\theta)$

- Note: Z is the denominator in computing $p(X)$:

$$p(\mathbf{x}; \theta) = \frac{\tilde{p}(\mathbf{x}; \theta)}{Z(\theta)} = \frac{1}{Z(\theta)} \prod_{k \in C} \phi_k(x_k; \theta)$$

2. With observed and latent variables $p(X, Z)$:

$$L(p, q) = E_q [\log p(X, Z) - \log q(Z)] \text{ is lower bound for } \log p(X)$$

- Note: $p(X)$ is the denominator in computing $p(Z|X)$:
 $p(Z|X) = p(X, Z) / p(X)$

Inferring $p(\mathbf{x})$: Lower Bound for $Z(\theta)$

- Suppose distribution p over $\mathbf{x} = [x_1, \dots, x_n]$ is

$$p(\mathbf{x}; \theta) = \frac{\tilde{p}(\mathbf{x}; \theta)}{Z(\theta)} = \frac{1}{Z(\theta)} \prod_{k \in C} \phi_k(x_k; \theta)$$

$$Z(\theta) = \sum_{\mathbf{x}} \prod_{k \in C} \phi_k(x_k; \theta)$$

- It is intractable because of $Z(\theta)$

- Consider $J(q) = KL(q \parallel \tilde{p}) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{\tilde{p}(\mathbf{x})}$ which is tractable

– Its negation is the lower bound for $Z(\theta)$ as follows

$$\begin{aligned} J(q) &= KL(q \parallel \tilde{p}) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{\tilde{p}(\mathbf{x})} \\ &= \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} - \log Z(\theta) \\ &= KL(q \parallel p) - \log Z(\theta) \end{aligned}$$

$$\Rightarrow \log Z(\theta) = KL(q \parallel p) - J(q) \geq -J(q)$$

Note:
 $\tilde{p}(x) = p(x)Z(\theta)$

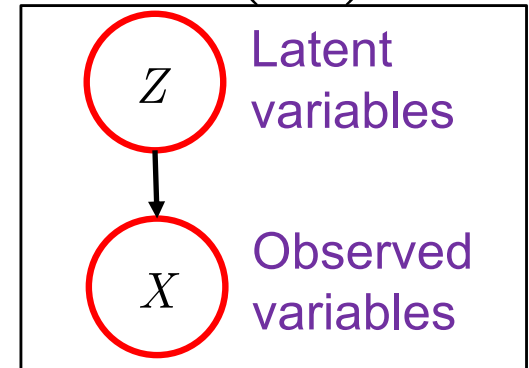
- Tractable - $J(q) = -KL(q \parallel \tilde{p}) = E_q[\log \tilde{p}(X) - \log q(X)]$ is a lower bound for $\log Z(\theta)$

Inferring $p(Z|X)$: lower bound for $p(X)$

- Given $p(X, Z) = p(Z)p(X|Z)$
- We need $p(Z|X) = p(X, Z) / p(X)$

- But $p(Z|X) = \frac{p(X|Z)p(Z)}{p(X)} = \frac{p(X|Z)p(Z)}{\int_Z p(X, Z)}$ is intractable

- Approximate $p(Z|X)$ by $q(Z)$ with minimum $KL[q(Z)||p(Z|X)]$



$$\begin{aligned}
 KL[q(Z)||p(Z|X)] &= \int_Z q(Z) \log \frac{q(Z)}{p(Z|X)} \\
 &= - \int_Z q(Z) \log \frac{p(Z|X)}{q(Z)} \\
 &= - \left(\int_Z q(Z) \log \frac{p(X, Z)}{q(Z)} - \int_Z q(Z) \log p(X) \right) \\
 &= - \int_Z q(Z) \log \frac{p(X, Z)}{q(Z)} + \log p(X) \int_Z q(Z) \\
 &= -L + \log p(X)
 \end{aligned}$$

where

$$L(p, q) = E_q \left[\log p(X, Z) - \log q(Z) \right]$$

Known as evidence lower bound: ELBO

Rearranging $\log p(X) = KL[q(Z)||p(Z|X)] + \mathcal{L}$

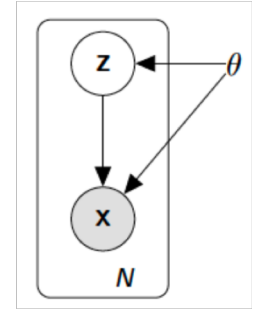
– Thus $\log p(X) \geq \mathcal{L}$

- Or $L(p, q) = E_q \left[\log p(X, Z) - \log q(Z) \right]$ is a lower bound for $\log p(X)$
- Maximizing ELBO is equivalent to minimizing KL
- This is the form useful for VAE

AEVB Variational Inference

- Distribution with latent variables and parameters:

$$p_{\theta}(\mathbf{x}, \mathbf{z}) = \prod_{i=1}^N p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{(i)}) p_{\theta}(\mathbf{z}^{(i)})$$

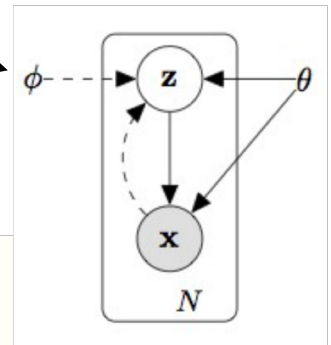


- We use conditional $q_{\phi}(z|x)$

- Effectively we choose a different q for every x
- Marginal likelihood is

$$\log p_{\theta}(x) = KL(q_{\phi}(z|x) || p(z|x)) + \mathcal{L}(p_{\theta}, q_{\phi})$$

- where $\mathcal{L}(p_{\theta}, q_{\phi}) = \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x, z) - \log q_{\phi}(z|x)]$
- i.e., the variational lower bound of marginal likelihood
- We maximize this bound wrt ϕ and θ (jointly)



- Doubly stochastic optimization

Short-coming of mean field

- How exactly do we optimize over $q(\mathbf{z}|\mathbf{x})$?
- Consider using mean field
 - Where $q(\mathbf{x}) = q_1(x_1)q_2(x_2)\cdots q_n(x_n)$ that approximates p in terms of $KL(q\|p)$
 - After each step of coordinate descent, we increase the variational lower bound, tightening it around $\log Z(\theta)$
 - But this is inaccurate for our purpose
 - The assumption that q is fully factored is too strong
 - And the coordinate descent optimization algorithm is too simplistic

Black-box variational inference

- A general purpose approach for optimizing q
 - for large classes of q more complex than mean field
- Maximize ELBO using gradient descent over ϕ
 - Instead of coordinate descent in mean field
- Assume q_ϕ is differentiable in its parameters ϕ
 - Instead of only doing inference, simultaneously perform learning via gradient descent on both ϕ, θ (jointly)
 - Optimization over ϕ will keep ELBO tight around $\log p(\mathbf{x})$
 - Optimization over θ will keep pushing the lower bound (and hence $\log p(\mathbf{x})$) up
 - This is somewhat similar to how the EM algorithm optimizes a lower bound on the marginal likelihood.

Score function gradient estimator

- We need to compute the gradient of the ELBO

$$\nabla_{\theta, \phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(x, z) - \log q_{\phi}(z)]$$

- No closed form for expectation wrt q
- Can obtain Monte Carlo estimates of gradient by sampling from q
 - This is easy for the gradient of p :
 - swap gradient and expectation and estimate the following expression via Monte Carlo: $\mathbb{E}_{q_{\phi}(z)} [\nabla_{\theta} \log p_{\theta}(x, z)]$
 - However, gradient wrt q is more difficult
 - Can use score function estimator

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(x, z) - \log q_{\phi}(z)] = \mathbb{E}_{q_{\phi}(z)} [(\log p_{\theta}(x, z) - \log q_{\phi}(z)) \nabla_{\phi} \log q_{\phi}(z)]$$

Disadvantage of Score Function

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(x, z) - \log q_{\phi}(z)] = \mathbb{E}_{q_{\phi}(z)} [(\log p_{\theta}(x, z) - \log q_{\phi}(z)) \nabla_{\phi} \log q_{\phi}(z)]$$

- It has a high variance
- What does this mean?
 - Suppose we use Monte Carlo to estimate an expectation whose mean is 1
 - If samples are 0.9, 1.1, 0.96, 1.05, ... All close to 1, then after a few samples, we get a good estimate true expectation.
 - If on the other hand we sample zero 99 times out of 100, and you sample 100 once, the expectation is still correct
 - Variance = $(1/100) (99 + 99^2)$
 - We will need a very large number of samples to figure out that the true expectation is actually one
 - We refer to the second case as being high variance

Key contribution of VAE

- This is the kind of problem we often run into with the score function estimator. In fact, its variance is so high, that we cannot use it to learn many models
- VAE proposes a much better estimator
 - This is done in two steps:
 - Reformulate ELBO so that parts of it can be computed in closed form (without Monte Carlo), and
 - Then use an alternative gradient estimator, based on the so-called reparametrization trick

The SGVB estimator

- The reformulation of the ELBO is as follows:

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)] - KL(q_\phi(z|x) || p(z))$$

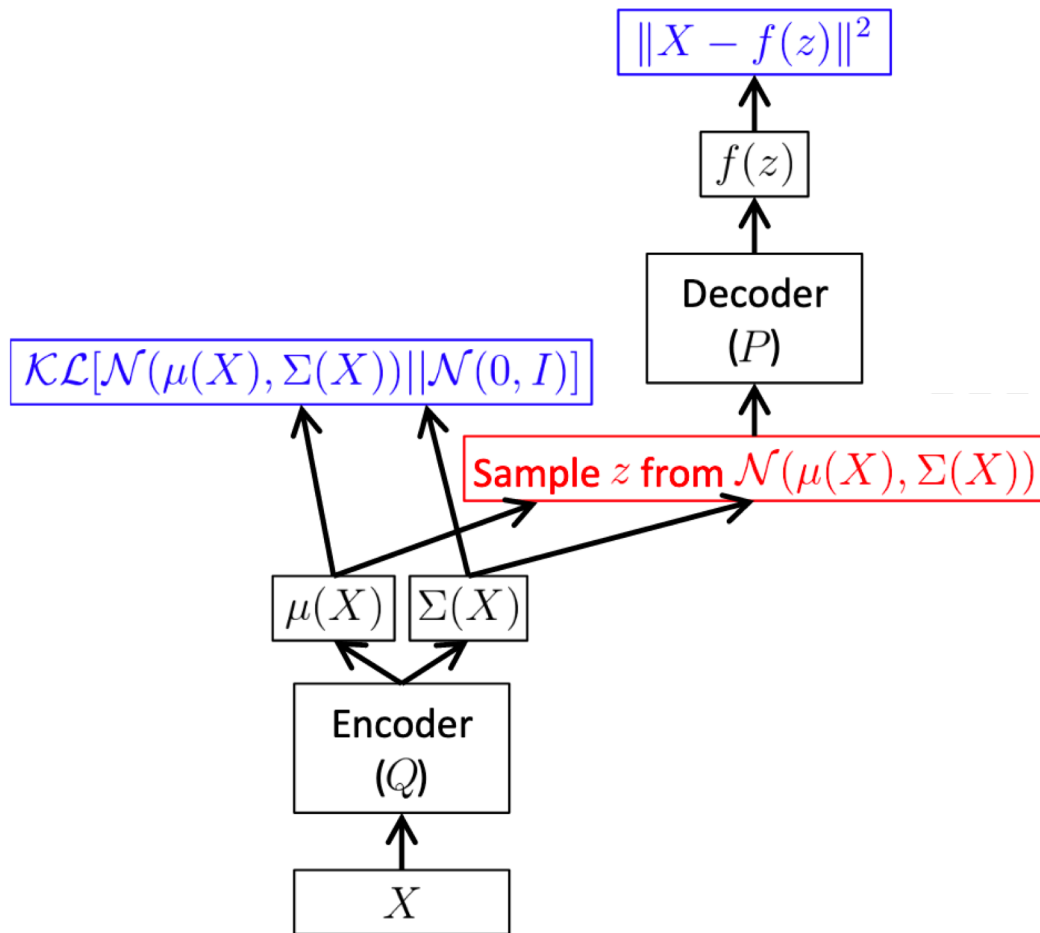
- x is an observed data point
- The rhs consists of two terms
 - both involve taking a sample $z \sim q(z|x)$, which we can interpret as a code describing x
 - We are also going to call q the *encoder*

Two terms of ELBO

$$\log p(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)] - KL(q_\phi(z|x) || p(z))$$

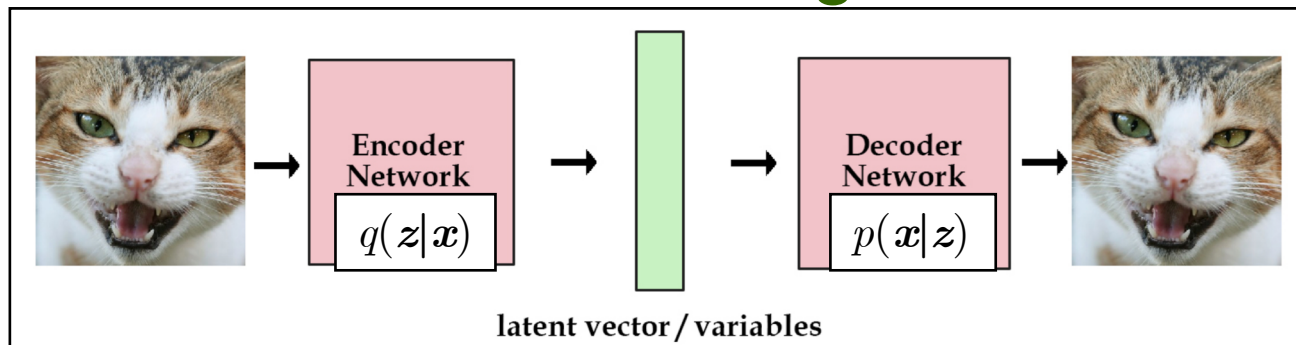
1. $\log p(x|z)$ is the log-likelihood of the observed x given the code z that we have sampled
 - Maximized when $p(x|z)$ assigns high probability to x
 - It is trying to reconstruct x given the code z ; for that reason $p(x|z)$ is the *decoder* network and the term is the *reconstruction error*
2. Divergence between $q(z|x)$ and prior $p(z)$, which we will fix to be a unit *Normal*
 - It encourages codes z to look Gaussian
 - We call it the *regularization* term
 - It prevents $q(z|x)$ from simply encoding an identity mapping, and instead forces it to learn some more interesting representation (e.g. facial features in first example)

VAE Training (without reparameterization)



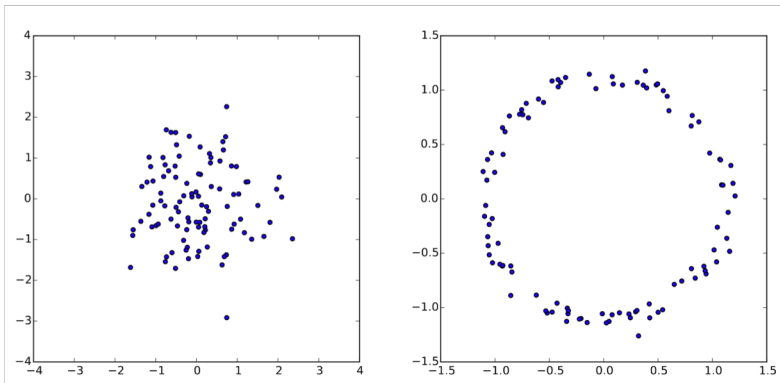
Optimization Objective: Autoencoding

- Fit a $q(z|x)$ to map x into a useful latent space z from which we can reconstruct x via $p(x|z)$
 - Reminiscent of *auto-encoder* neural networks
 - An auto-encoder is a pair of neural networks f, g that are composed as $x' = f(g(x))$
 - They are trained to minimize reconstruction error $\|x' - x\|$
 - In practice, $g(x)$ learns to embed x in a latent space that often has an intuitive interpretation
 - This is where the AEVB algorithm takes its name



The reparameterization trick

- As seen earlier, optimizing our objective requires a good estimate of the gradient
- The main technical contribution of VAE is a low-variance gradient estimator based on the *reparameterization trick*
 - Given a random variable z with one distribution we can create another random variable $\mathbf{X}=g(z)$ with a completely different distribution



$$g(z) = z/10 + z/||z||$$

Two-step reparameterization

- Under mild conditions, we may express $q_\phi(z|x)$ as the following two-step generative process

1. We sample a noise variable ϵ

- from a simple distribution $p(\epsilon)$ such as $N(0,1)$

$$\epsilon \sim p(\epsilon)$$

2. We apply a deterministic transformation $g_\phi(\epsilon, x)$ to map the random noise into a more complex distribution

$$z = g_\phi(\epsilon, x)$$

- For many interesting classes of q_ϕ
 - it is possible to choose a $g_\phi(\epsilon, x)$, such that $z = g_\phi(\epsilon, x)$ will be distributed according to $q_\phi(z|x)$

Transformation to get a Gaussian

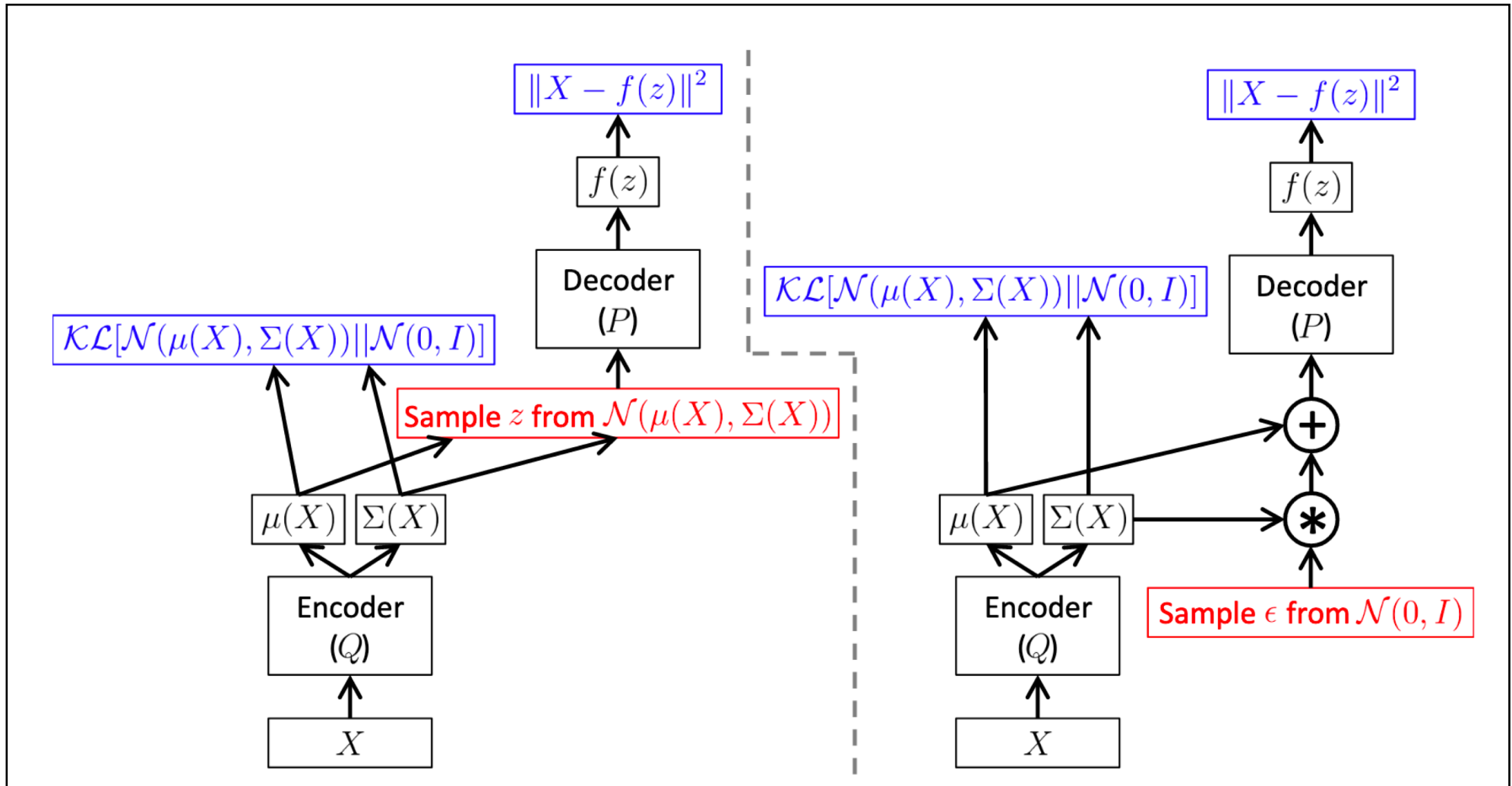
- Gaussian variables provide the simplest example of the reparametrization trick
- Instead of writing $z \sim q_{\mu,\sigma}(z) = N(\mu, \sigma)$, we write

$$z = g_{\mu,\sigma}(\epsilon) = \mu + \epsilon \cdot \sigma$$

where $\epsilon \sim N(0,1)$

- It is easy to check that the two ways of expressing the random variable z lead to the same distribution.

VAE Training (without/with reparameterization)



The feedforward behavior of these networks is identical, but backpropagation can be applied only to the right network

Advantage is reduced variance

- Biggest advantage of this approach is that we many now write the gradient of an expectation with respect to $q(\mathbf{z})$ (for any f) as

$$\nabla_{\phi} \mathbb{E}_{z \sim q(z|x)} [f(x, z)] = \nabla_{\phi} \mathbb{E}_{\epsilon \sim p(\epsilon)} [f(x, g(\epsilon, x))] = \mathbb{E}_{\epsilon \sim p(\epsilon)} [\nabla_{\phi} f(x, g(\epsilon, x))]$$

- Gradient is inside the expectation
 - We may use Monte Carlo to estimate right-hand term
- Approach has much lower variance than the score function estimator
 - Enables learning more complex models

Computing the gradient

- The KL divergence between two Gaussians can be computed in closed form as

$$KL[N(\mu(X), \Sigma(X)) || N(0, I)] = 1/2(\text{tr}(\Sigma(X)) + (\mu(X))^T \mu(X) - k - \log \det(\Sigma(X)))$$

- We want to optimize, over dataset D

$$E_{X \sim D}[\log P(X) - KL[Q(z|X) || P(z|X)]] = E_{X \sim D}[E_{z \sim Q}[\log P(X|z)] - KL[Q(z|X) || P(z)]]$$

- The equation we actually take the gradient of is:

$$E_{X \sim D}[E_{\epsilon \sim N(0, I)}[\log P(X|z = \mu(X) + \Sigma^{1/2}(X) * \epsilon)] - KL[Q(z|X) || P(z)]]$$

- Given a fixed X and ϵ , this function is deterministic and continuous in the parameters of P and Q , meaning backpropagation can compute a gradient that will work for SGD

Choice of p and q

- We have not specified exact form of p or q
 - How to parameterize these distributions?
- The best $q(\mathbf{z}|\mathbf{x})$ should be able to approximate the true posterior $p(\mathbf{z}|\mathbf{x})$
- Similarly, $p(\mathbf{x})$ should be flexible enough to represent the richness of the data.

Parameterizing by neural networks

- For these reasons, we are going to parameterize q and p by *neural networks*
 - Which are extremely expressive
 - Efficiently optimized over large datasets
- Choice is a bridge between classical machine learning (approximate Bayesian inference) and modern deep learning

Form of neural network

- Assuming that $q(\mathbf{z}|\mathbf{x})$ and $p(\mathbf{x}|\mathbf{z})$ are Gaussian

$$q(\mathbf{z}|\mathbf{x}) = N(\mathbf{z}; \boldsymbol{\mu}(\mathbf{x}), \boldsymbol{\sigma}(\mathbf{x}) \odot \mathbf{I})$$

where $\boldsymbol{\mu}(\mathbf{x})$, $\boldsymbol{\sigma}(\mathbf{x})$ are deterministic vector-valued functions of \mathbf{x} parametrized by an arbitrary complex neural network

\odot is element-wise multiplication (Hadamard multiplication)

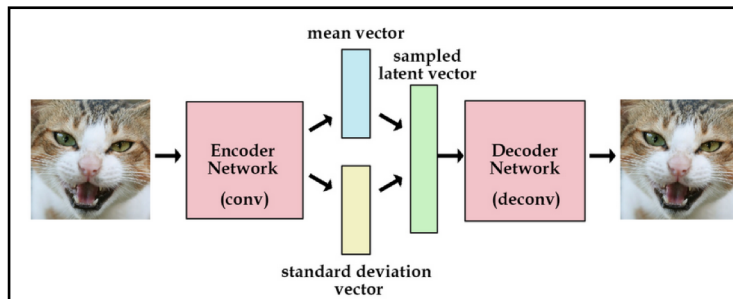
- More generally, the same technique can be applied to any exponential family distribution by parameterizing the sufficient statistics by a function of \mathbf{x}

VAE Training

- Backprop proceeds through parameters of the latent distribution
 - Called reparameterization trick

$$N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \boldsymbol{\mu} + \boldsymbol{\Sigma} N(0, \mathbf{I})$$

where the covariance matrix $\boldsymbol{\Sigma}$ is diagonal



- Due to randomness involved, training is called Stochastic Gradient Variational Bayes (SGVB)

AEVB Algorithm

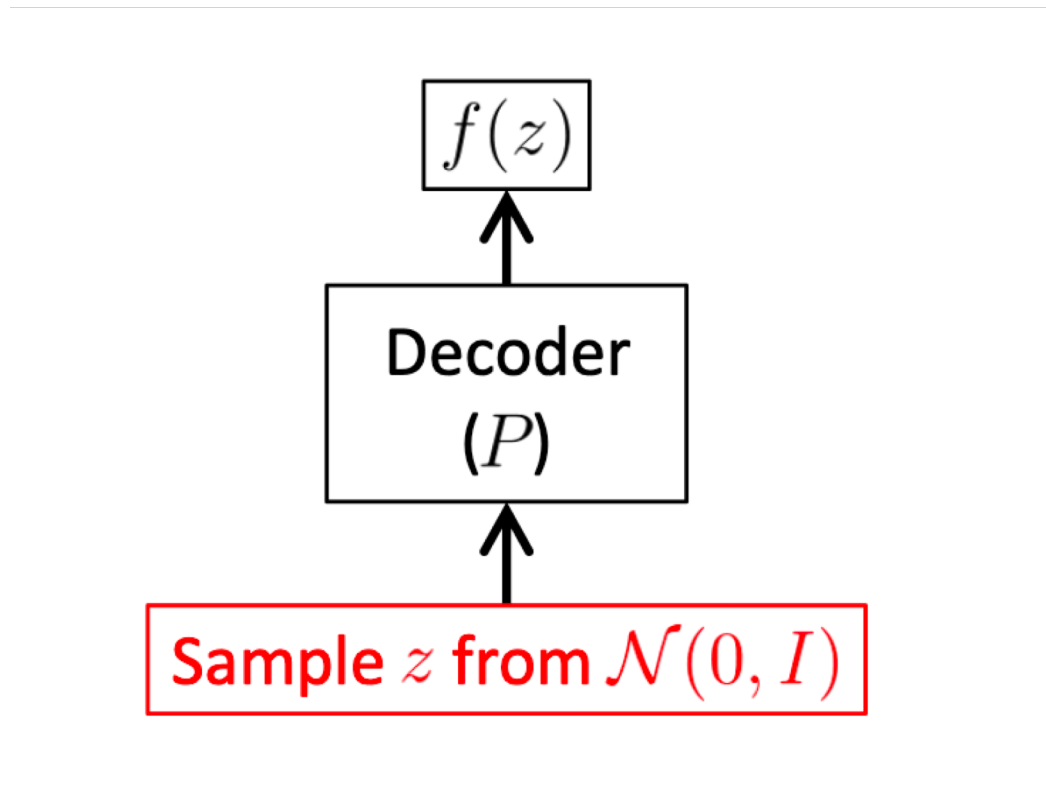
- The AEVB algorithm is the combination of
 1. Auto-encoding ELBO reformulation
 2. Black-box variational inference approach, and
 3. Reparameterization-based low-variance gradient estimator
 - It optimizes the auto-encoding ELBO using black-box variational inference with the reparametrized gradient estimator
- AEVB is applicable to any deep generative model p_{θ} with latent variables that is differentiable in θ

VAE Algorithm

- A VAE uses the AEVB algorithm to learn a specific model p using a particular encoder q
 - The model p is parametrized as
$$p(\mathbf{x}|\mathbf{z})=N(\mathbf{x}; \boldsymbol{\mu}(\mathbf{z}), \boldsymbol{\sigma}(\mathbf{z})\odot\mathbf{I})p(\mathbf{z})=N(\mathbf{z};0,\mathbf{I})$$
where $\boldsymbol{\mu}(\mathbf{z}),\boldsymbol{\sigma}(\mathbf{z})$ are parameterized by a neural network (typically, two dense hidden layers of 500 units each)
 - The model q is similarly parametrized as
$$q(\mathbf{z}|\mathbf{x})=N(\mathbf{z}; \boldsymbol{\mu}(\mathbf{x}),\boldsymbol{\sigma}(\mathbf{x})\odot\mathbf{I})$$
- This p and q simplifies auto-encoding ELBO
 - We can use a closed form expression to compute the regularization term, and we only use Monte-Carlo estimates for the reconstruction term

Testing time VAE

- A remarkably simple setup
 - The values of μ and σ are part of the decoder

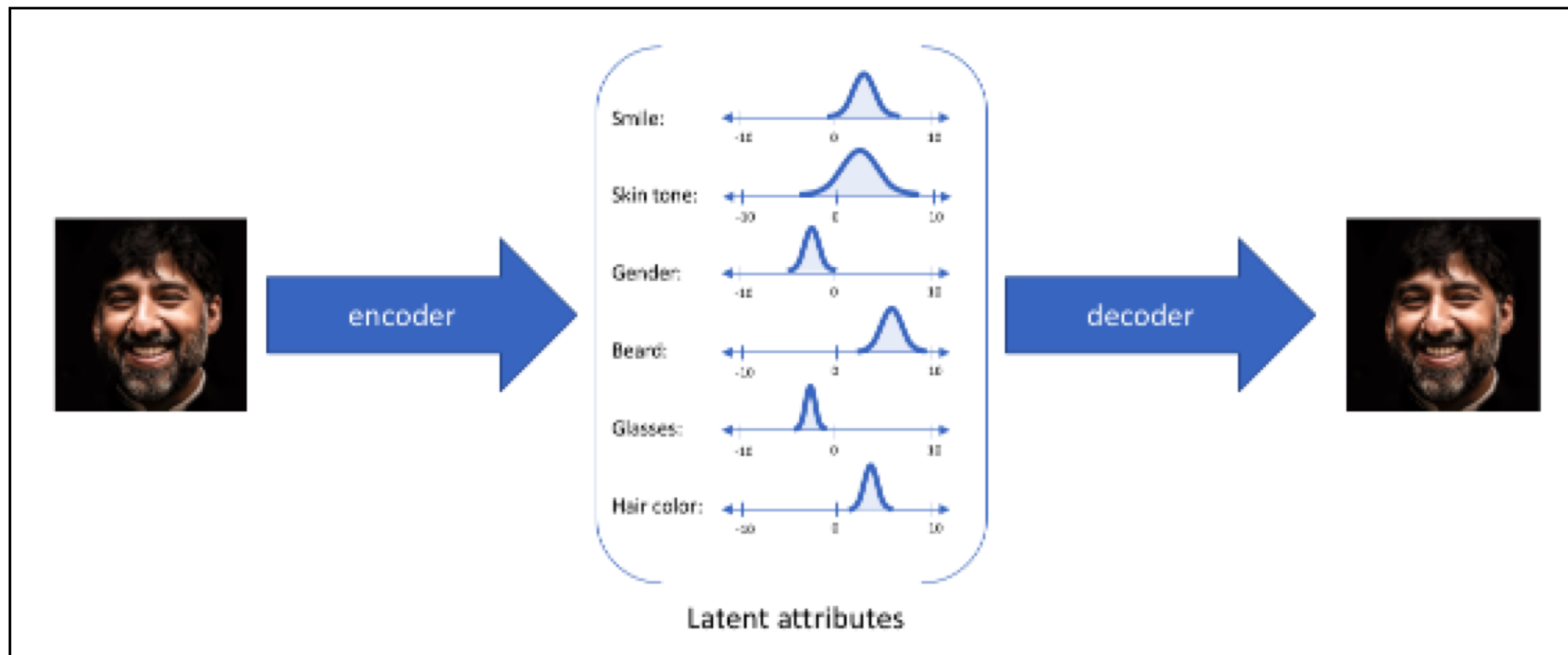


VAE as a PGM

- We may interpret the VAE as a directed latent-variable PGM
- We may also view it as a particular objective for training an auto-encoder neural network
 - unlike previous approaches, this objective derives reconstruction and regularization terms from a more principled, Bayesian perspective

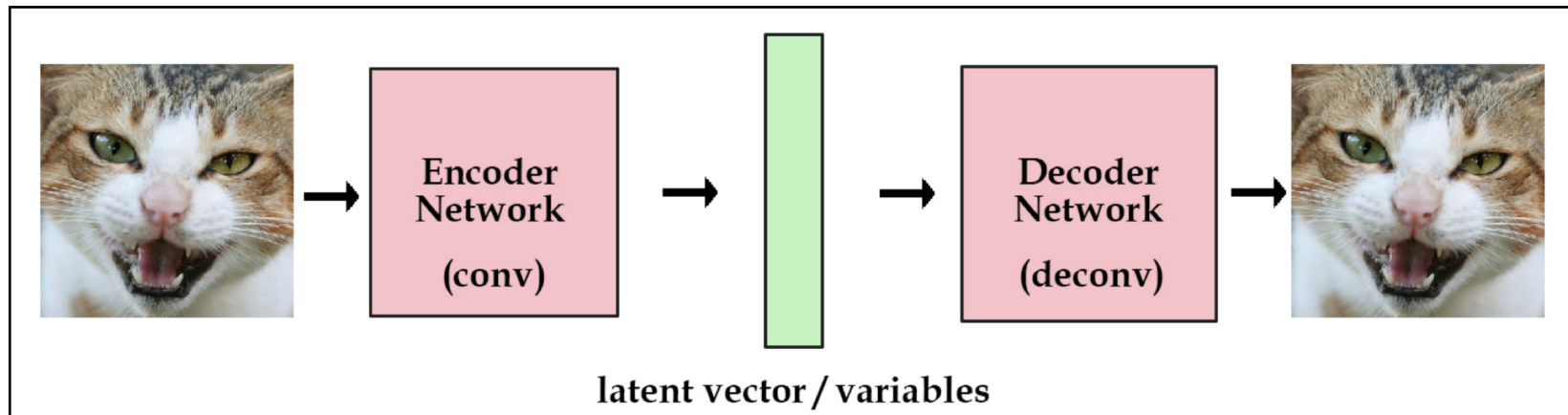
Autoencoding using a VAE

- Move from samples x_i to latent space z_i and reconstruct \hat{x}_i

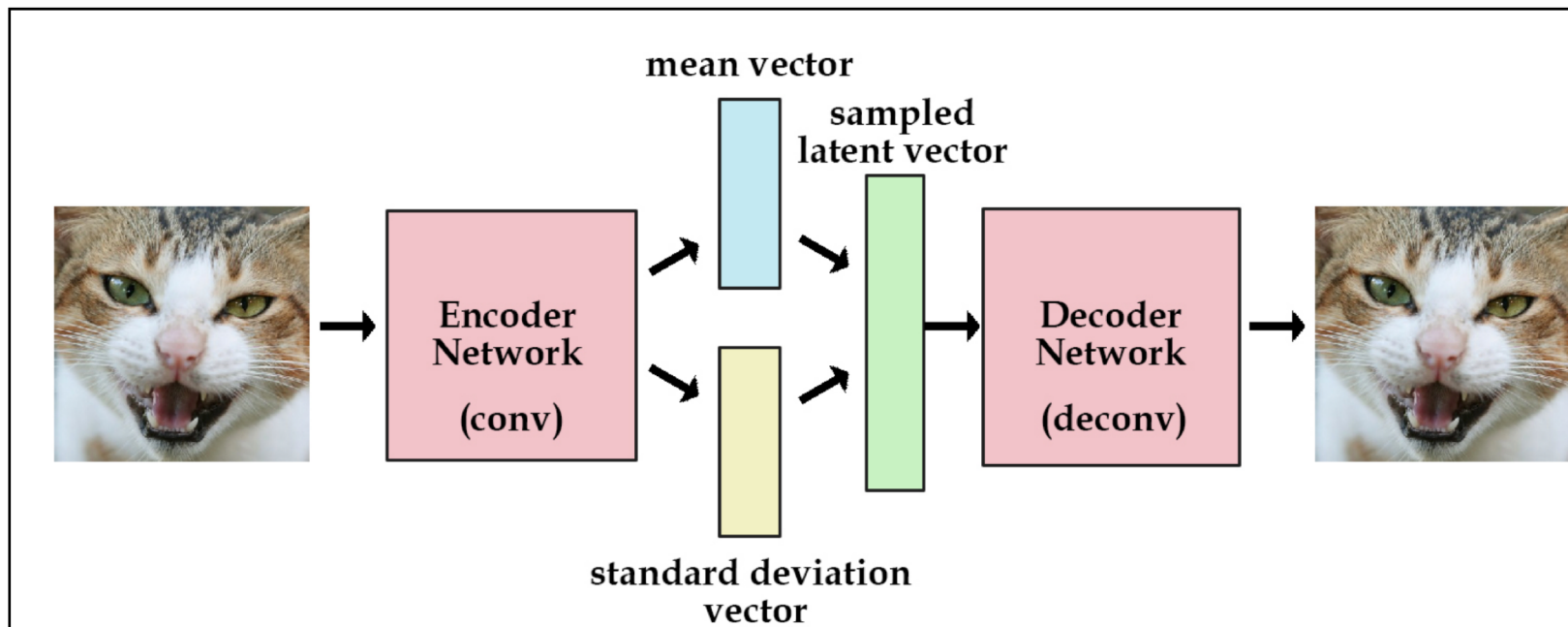


Source: *Auto-Encoding Variational Bayes*, Diederik P Kingma, Max Welling

Autoencoder



Variational Autoencoder



Parameters of the VAE model

- **Generative model:** $p(x, z) = p(x|z)p(z)$
- $p(x|z)$ is parameterized by a deep neural network (decoder)
- Components of z are independent Bernoulli or Gaussian
- Learned approx inference trained using gradient descent
 - $q(z|x) = N(\mu, \sigma^2 I)$ whose parameters are given by another deep network (encoder)
- Thus we have $z \sim \text{Enc}(x) = q(z|x)$ and $y \sim \text{Dec}(z) = p(x|z)$

Key insight of VAE

- They can be trained by maximizing variational lower bound $\mathcal{L}(q)$ associated with data point \mathbf{x}

$$\begin{aligned}\mathcal{L}(q) &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log p_{\text{model}}(\mathbf{z}, \mathbf{x}) + \mathcal{H}(q(\mathbf{z} | \mathbf{x})) \\ &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log p_{\text{model}}(\mathbf{x} | \mathbf{z}) - D_{\text{KL}}(q(\mathbf{z} | \mathbf{x}) || p_{\text{model}}(\mathbf{z})) \\ &\leq \log p_{\text{model}}(\mathbf{x}).\end{aligned}$$

- where $\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log p_{\text{model}}(\mathbf{z}, \mathbf{x})$ is the joint log-likelihood of the visible and hidden variables under the approximate posterior over the latent variables
 - ♦ $\mathcal{H}(q(\mathbf{z}|\mathbf{x}))$ is the entropy of the approximate posterior
- When q is chosen to be Gaussian with noise added to a predicted mean, maximizing this entropy term encourages increasing σ

Summary of VAE

- VAE is a directed model that uses
 - Learned approximate inference
 - Trained purely with gradient-based methods
- VAE generates a sample from the model,
 - First draw sample z from code distribution $p_{\text{model}}(z)$.
 - Sample is then run through a differentiable generator network $g(z)$
 - \mathbf{x} is sampled from distribution $p_{\text{model}}(\mathbf{x}; g(z)) = p_{\text{model}}(\mathbf{x} | g(z))$
 - However during training the approximate inference network (or encoder) $q(z | \mathbf{x})$ is used to obtain z and $p_{\text{model}}(\mathbf{x} | z)$ is viewed as a decoder network

References

- Autoencoding Variational Bayes (Kingma and Welling):
 - <https://arxiv.org/pdf/1312.6114.pdf>
- Variational Autoencoder (S. Ermon):
 - <https://ermongroup.github.io/cs228-notes/extras/vae/>
- Blackbox Variational Inference (Ranganath, et.al.):
 - <http://www.cs.columbia.edu/~blei/papers/RanganathGerriShBlei2014.pdf>
- Stochastic Backpropagation (Rezende, et. al.):
 - <https://arxiv.org/pdf/1401.4082.pdf>
- Tutorial on VAE (C. Doersch):
<https://arxiv.org/pdf/1606.05908.pdf>