

Variational Autoencoders

Kingma et al. (2014)

JE Starling, 10-2017

Fall 2017

Review of Variational Inference

Setting:

x is our data.

z is our latent variable.

Joint density is $p(x, z) = p(x|z) \cdot p(z)$

Goal: Approximate the posterior $p(z|x) = \frac{p(x|z)p(z)}{p(x)}$,
where marginal $p(x) = \int p(x|z)p(z)dz$ is intractable.

Strategy: We pose a family of approximations, Q , and choose a member of that family, $q(z) \in Q$, to minimize $KL [q(z)||p(z|x)]$.

Review of Variational Inference (2)

We want to find the best approximation:

$$q^*(z) = \arg \min_{q(z) \in Q} KL [q(z) || p(z|x)]$$

This objective is intractable because it involves $p(x)$:

$$\begin{aligned} KL [q(z) || p(z|x)] &= E_{q(z)} \left[\log \left(\frac{q(z)}{p(z|x)} \right) \right] \\ &= E_{q(z)} [\log (q(z))] - E_{q(z)} [\log (p(x, z))] + \log (p(x)) \end{aligned}$$

Review of Variational Inference (3)

We can maximize another quantity which is equivalent to minimizing the KL divergence. (The $\log(p(x))$ term is constant wrt $q(z)$.)

$$KL[q(z)||p(z|x)] = \underbrace{E_{q(z)}[\log(q(z))] - E_{q(z)}[\log(p(x, z))]}_{-ELBO(q)} + \log(p(x))$$

We can write $ELBO(q)$ as

$$\begin{aligned} ELBO(q) &= E_{q(z)}[\log(q(z))] - E_{q(z)}[\log(p(x, z))] \\ &= E_{q(z)}[\log(p(z))] + E_{q(z)}[\log(p(x|z))] - E_{q(z)}[\log(q(z))] \\ &= E_{q(z)}[\log(p(x|z))] - KL[q(z)||p(z)] \end{aligned}$$

The $ELBO(q)$ is a lower bound on $\log(p(x))$.

Variational Auto Encoders: Notation Note

Going forward, we will write $q(z)$ as $q(z|x)$.

We will also add a subscript to indicate that $q(z|x)$ is parameterized by variational parameters labeled ϕ .

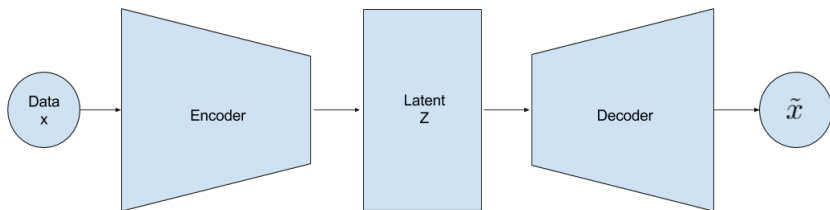
We write: $q_{\phi}(z|x)$

Purpose of VAEs

What do VAEs do?

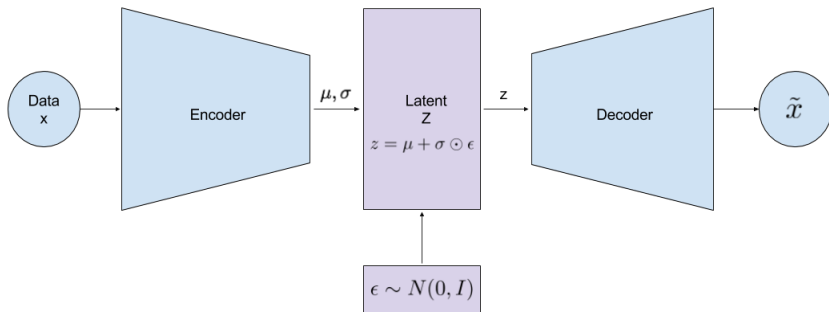
- Fit a generative model to a large dataset.
- Model the underlying distribution of the data and sample new data from the distribution.
- Image re-generation, and generating new images like existing ones.

Basic VAE Setup



In neural nets language, a VAE consists of an encoder, a decoder, and a loss function.

The Latent Variables



- Latent z is a continuous, lower-dimensional; dimension is a user input.
- Let the prior $p(z)$ be isotropic multivariate Gaussian: $p(z) \sim N(z; \mathbf{0}, \mathbf{I})$.

The Latent Variables: Reparameterization Trick

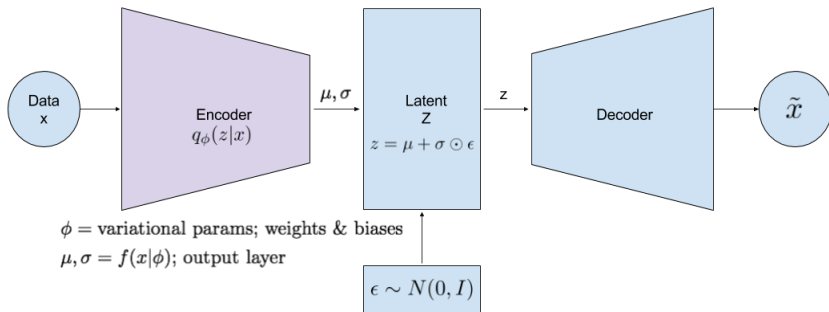
Why use $z = \mu + \sigma \odot \epsilon, \epsilon \sim N(0, I)$ instead of sampling $z \sim N(\mu, \sigma I)$?

We must remove the randomness from the neural networks in order to use back-propagation with stochastic gradient descent to train.

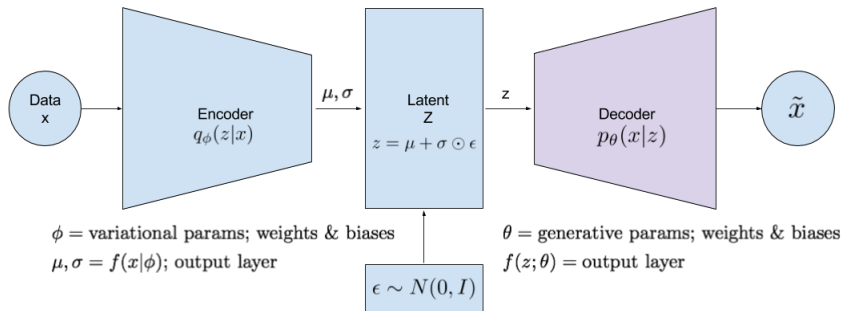
Interpretation of σ :

- Think of z as encoding information in a lower-dimensional space. How to prevent encoding an infinite amount of info?
- What if we have two x 's which are very different, mapped to two z 's which are close?
- σ acts as noise, which prevents reconstructing two different \tilde{x} values from z 's which are less than σ apart. Regulates how much info you can encode.

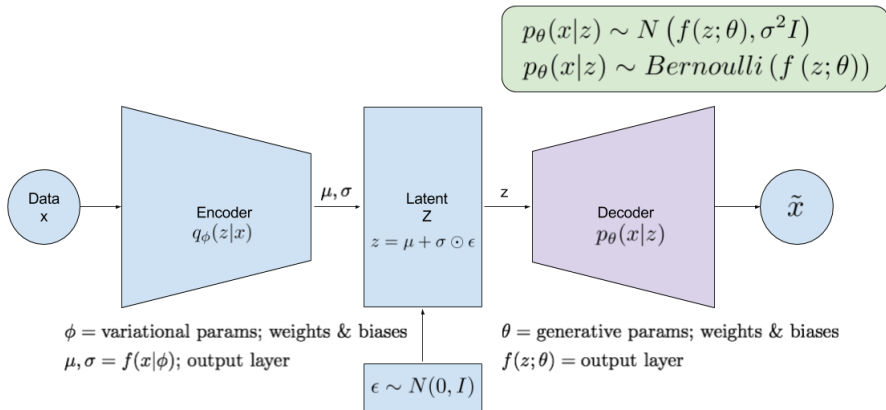
The Encoder



The Decoder



The Decoder



Training the Model

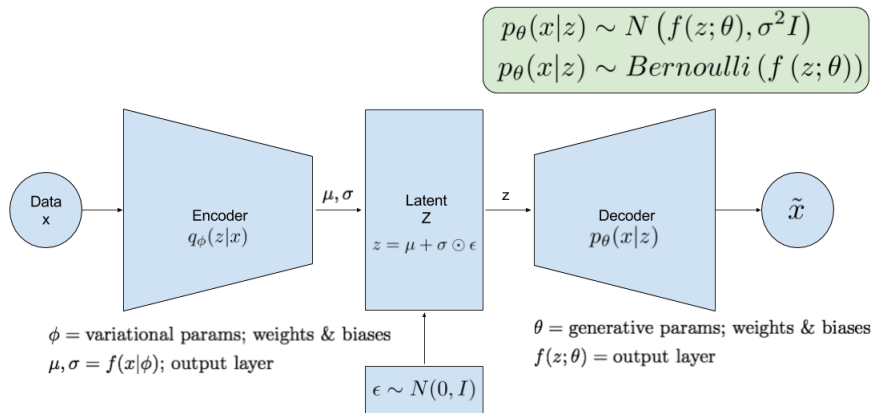
Model parameters (ϕ, θ) are trained using stochastic gradient descent with back-propagation.

- (ϕ, θ) are learned jointly.
- With step size ρ , encoder and decoder parameters are updated as

$$\begin{aligned}\phi_{\text{new}} &= \phi_{\text{old}} - \rho \frac{\partial L_i(\phi, \theta, x_i)}{\partial \phi} \\ \theta_{\text{new}} &= \theta_{\text{old}} - \rho \frac{\partial L_i(\phi, \theta, x_i)}{\partial \theta}\end{aligned}$$

where $L_i(\phi, \theta, x_i)$ is a loss function.

The Loss Function



$$L(\phi, \theta, x) = -ELBO_q(\phi, \theta, x) = -E_{q_\phi(z|x)} [\log p_\theta(x|z)] + KL[q_\phi(z|x) || p(z)]$$

Tasks we can accomplish after training

- (1) Inference on parameters (the Gaussian means, or Bernoulli probabilities).
- (2) Efficient approximation of posterior of latent z given observed x for a set of parameters. (Coding tasks, data representation)
- (3) Efficient approximate marginal inference of x :
 - Reconstruct data sets we have seen.
 - Generate new data sets that are like one we have already seen.

Sampling of papers

- Variational Autoencoder for Deep Learning of Images, Labels and Captions
- Alternative priors for Deep Generative Models
- Least Squares VAE with Regularization
- Stein VAE
- Ladder VAEs

Example on github

In the example on github:

- Uses the MNIST data set.
- Uses a gaussian encoder and a bernoulli decoder.
- Both neural nets have two hidden layers.
- Uses drop-outs (`tf.nn.dropout`) to avoid overfitting.