

Variational Inference

Rafael Rafailov

October 9, 2021

Definition

Consider two distributions p and q over a set X . The KL-divergence $D_{KL}[p||q]$ is defined as

$$\begin{aligned} D_{KL}[p||q] &= \int_X p(x) \log \frac{q(x)}{p(x)} dx \\ &= E_x \left[\log \frac{p(x)}{q(x)} \right] \end{aligned}$$

Properties of the KL divergence

- 1 The KL divergence is not symmetric $D_{KL}[p||q] \neq D_{KL}[q||p]$

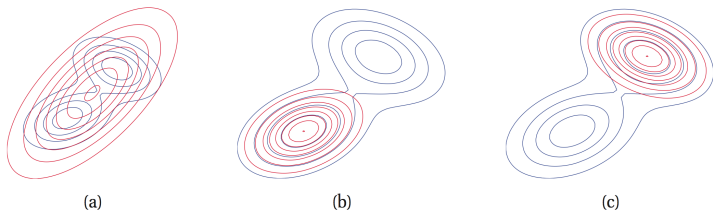


Figure: Fitting a unimodal approximating distribution q (red) to a multimodal p (blue). Using $KL(p||q)$ leads to **mode-covering** (a). However, using $KL(q||p)$ forces q to be **mode-seeking** (b, c)

¹Image credit: CS 236

- 1 The KL-divergence between two distributions (when it is defined) is always non-negative.

Proof: We have

$$\begin{aligned} D_{KL}[q||p] &= \int p(x) \log \frac{q(x)}{p(x)} dx \\ &= \int p(x) \log \frac{q(x)}{p(x)} dx \\ &= \int \log \frac{q(x)}{p(x)} p(x) dx = \int \log 1 dx = 0 \end{aligned}$$

where the above follows from Jensen's inequality.

- 2 From the previous point $D_{KL}[q||p] = 0$ if and only if $p = q$ (modulo some measure-theoretic considerations).

Generative models over data



$$\mathbf{x}_i \sim P_{\text{data}}$$
$$i = 1, 2, \dots, n$$

Figure: We want to fit the probability distribution of the data

¹Image credit: CS 236

Generative models over data



$$x_i \sim P_{\text{data}} \\ i = 1, 2, \dots, n$$

Figure: We want to fit the probability distribution of the data

$$\arg \min_p D_{KL}(p_{\text{data}} \| p) = \max_p \frac{1}{|D|} \sum_{x_i \in D} \log p(x_i) \quad (0.1)$$

¹Image credit: CS 236

Generative models over data



$$x_i \sim P_{\text{data}} \\ i = 1, 2, \dots, n$$

Figure: We want to fit the probability distribution of the data

$$\arg \min_p D_{KL}(p_{\text{data}} \| p) = \max_p \frac{1}{J D_j} \prod_{x_i \in D} \log p(x_i) \quad (0.1)$$

We just need to train a maximum likelihood model!

¹Image credit: CS 236

Evaluating likelihoods over high dimensions is hard!

Latent variable models

¹Image credit: Jeremy Jordan's blog post.

Latent variable models

¹Image credit: Jeremy Jordan's blog post.

- 1 Data is governed by a simple latent distribution $p(Z)$.

¹Image credit: Jeremy Jordan's blog post.

- 1 Data is governed by a simple latent distribution $p(Z)$.
- 2 The observed data X is generated by a conditional distribution $p(X|Z)$.

¹Image credit: Jeremy Jordan's blog post.

Data likelihood under latent variable models

General solution: Introduce inference distribution $q(z|x)$.

General solution: Introduce inference distribution
 $q(z|x)$.

Intuition:

General solution: Introduce inference distribution $q(z|x)$.

Intuition:

- 1 Guess the likely z given x_i and use those to compute likelihood.

General solution: Introduce inference distribution $q(z|x)$.

Intuition:

Guess the likely z given x_i and use those to compute likelihood.

Evaluate uncertainty through a distribution over z - $q(z|x)$.

General solution: Introduce inference distribution $q(z|x)$.

Intuition:

Guess the likely z given x_i and use those to compute likelihood.

Evaluate uncertainty through a distribution over z - $q(z|x)$.

Approach is similar to the EM algorithm.

$\log_p (x) =$

$$\log p(x) = \log \int_Z p(x|z)p(z)dz =$$

$$\log p(x) = \log \int_Z p(x|z)p(z)dz =$$

$$\log \int_Z \frac{q(z|x)}{\int_Z q(z|x)} p(x|z)p(z)dz =$$

$$\int_Z \frac{q(z|x)}{\int_Z q(z|x)} p(x|z)p(z)dz =$$

$$\log p(x) = \int_Z p(x|z)p(z)dz =$$

$$\int_Z \frac{q(z|x)}{\underbrace{q(z|x)}_{=1}} p(x|z)p(z)dz = \int_Z q(z|x) \frac{p(x; z)}{q(z|x)} dz =$$

$$\begin{aligned} \log p(x) &= \log \int_Z p(x|z)p(z)dz = \\ \log \int_Z \frac{q(z|x)}{\int_Z q(z|x)} p(x|z)p(z)dz &= \log \int_Z q(z|x) \frac{p(x; z)}{q(z|x)} dz = \\ \log E_Z \frac{p(x; z)}{q(z|x)} & \end{aligned}$$

$$\begin{aligned}
 \log p(x) &= \int_Z p(x|z)p(z) dz = \\
 \log \int_Z \frac{q(z|x)}{\int_Z q(z|x)} p(x|z)p(z) dz &= \log \int_Z q(z|x) \frac{p(x; z)}{q(z|x)} dz = \\
 \log E_Z \left[\frac{p(x; z)}{q(z|x)} \right] & \underset{\text{Jensen}}{\leq} E_Z \left[\log \frac{p(x; z)}{q(z|x)} \right] =
 \end{aligned}$$

$$\begin{aligned}
\log p(x) &= \int_Z p(x|z)p(z) dz = \\
\log \int_Z \underbrace{q(z|x)}_{=1} p(x|z)p(z) dz &= \int_Z q(z|x) \frac{p(x; z)}{q(z|x)} dz = \\
\log E_Z \left[\frac{p(x; z)}{q(z|x)} \right] &\stackrel{\text{Jensen}}{\leq} E_Z \left[\log \frac{p(x; z)}{q(z|x)} \right] = \\
E_Z \left[\log p(x|z) \right] + E_Z \left[\log \frac{p(z)}{q(z|x)} \right] &=: L(x) \\
&= D_{KL}(q(z|x) || p(z))
\end{aligned}$$

Why exactly do we sample from $q(z|x)$?

.

Why exactly do we sample from $q(z|x)$?

Proposition

For any distribution $q(z)$ we have:

$$\log p(x) - D_{\text{KL}}(q(z) \parallel p(z|x)) = \mathbb{E}_{z \sim q(z)}[\log p(x|z)] - D_{\text{KL}}(q(z) \parallel p(z))$$

Why exactly do we sample from $q(z|x)$?

Proposition

For any distribution $q(z)$ we have:

$$\log p(x) - D_{\text{KL}}(q(z) \parallel p(z|x)) = \mathbb{E}_{z \sim q(z)}[\log p(x|z)] - D_{\text{KL}}(q(z) \parallel p(z))$$

The Evidence Lower Bound (ELBO) is a lower bound on the data log-likelihood under any sampling distribution $q(z)$.

Why exactly do we sample from $q(z|x)$?

Proposition

For any distribution $q(z)$ we have:

$$\log p(x) - D_{\text{KL}}(q(z) \parallel p(z|x)) = \mathbb{E}_{z \sim q(z)}[\log p(x|z)] - D_{\text{KL}}(q(z) \parallel p(z))$$

The Evidence Lower Bound (ELBO) is a lower bound on the data log-likelihood under any sampling distribution $q(z)$.

From the properties of KL-divergences, equality is achieved only when $q(z) = p(z|x)$.

Why exactly do we sample from $q(z|x)$?

Proposition

For any distribution $q(z)$ we have:

$$\log p(x) - D_{\text{KL}}(q(z) \parallel p(z|x)) = \mathbb{E}_{z \sim q(z)}[\log p(x|z)] - D_{\text{KL}}(q(z) \parallel p(z))$$

The Evidence Lower Bound (ELBO) is a lower bound on the data log-likelihood under any sampling distribution $q(z)$.

From the properties of KL-divergences, equality is achieved only when $q(z) = p(z|x)$.

To minimize the ELBO gap we choose $q(z) = p(z|x)$.

The ELBO consists of two terms: reconstruction and a KL regularization:

$$L ; = \underbrace{E_z \text{ } q(z|x) [\log p(x|z)]}_{\text{reconstruction}} \quad \underbrace{D_{\text{KL}}(q(z|x) || p(z))}_{\text{KL regularization}}$$

The ELBO consists of two terms: reconstruction and a KL regularization:

$$L = \underbrace{E_z \left[\log p(x|z) \right]}_{\text{reconstruction}} - \underbrace{D_{\text{KL}}(q(z|x) \parallel p(z))}_{\text{KL regularization}}$$

Set $q(z|x)$ to be multivariate normal distribution parameterized by a neural network, i.e. $q(z|x) = N(z; \mu(x); \Sigma(x))$, where $\Sigma(x) = \text{diag}(\sigma_1^2(x); \dots; \sigma_K^2(x))$ is a diagonal matrix.

The ELBO consists of two terms: reconstruction and a KL regularization:

$$L; = \underbrace{E_z q(z|x) [\log p(x|z)]}_\text{reconstruction} - \underbrace{D_{KL}(q(z|x)||p(z))}_\text{KL regularization}$$

Set $q(z|x)$ to be multivariate normal distribution parameterized by a neural network, i.e. $q(z|x) = N(z; \mu(x); \Sigma(x))$, where $\Sigma(x) = \text{diag}(\sigma_1^2(x); \dots; \sigma_k^2(x))$ is a diagonal matrix.

Proposition

Let $q = N(\mu_1; \Sigma_1)$ and $p = N(\mu_2; \Sigma_2)$, then

$$D_{KL}(q||p) = \frac{1}{2} \text{tr} \left(\Sigma_2^{-1} \Sigma_1 + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) \right) + \log \frac{\det \Sigma_2}{\det \Sigma_1}$$

The ELBO consists of two terms: reconstruction and a KL regularization:

$$L; = \underbrace{E_z q(z|x) [\log p(x|z)]}_\text{reconstruction} - \underbrace{D_{KL}(q(z|x)||p(z))}_\text{KL regularization}$$

Set $q(z|x)$ to be multivariate normal distribution parameterized by a neural network, i.e. $q(z|x) = N(z; \mu(x); \Sigma(x))$, where $\Sigma(x) = \text{diag}(\sigma_1^2(x); \dots; \sigma_k^2(x))$ is a diagonal matrix.

Proposition

Let $q = N(\mu_1; \Sigma_1)$ and $p = N(\mu_2; \Sigma_2)$, then

$$D_{KL}(q||p) = \frac{1}{2} \text{tr} \left(\Sigma_2^{-1} \Sigma_1 + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) \right) + k + \log \frac{\det \Sigma_2}{\det \Sigma_1}$$

This is a differentiable function!

The ELBO consists of two terms: reconstruction and a KL regularization:

$$L = \underbrace{E_z \left[\log p(x|z) \right]}_{\text{reconstruction}} - \underbrace{D_{\text{KL}}(q(z|x) \parallel p(z))}_{\text{KL regularization}}$$

The ELBO consists of two terms: reconstruction and a KL regularization:

$$L = \underbrace{E_z \left[\log p(x|z) \right]}_{\text{reconstruction}} - \underbrace{D_{\text{KL}}(q(z|x) \parallel p(z))}_{\text{KL regularization}}$$

Optimizing the ELBO

The ELBO consists of two terms: reconstruction and a KL regularization:

$$L ; = \underbrace{E_z q(z|x) [\log p(x|z)]}_{\text{reconstruction}} + \underbrace{D_{KL}(q(z|x) || p(z))}_{\text{KL regularization}}$$

Need to compute

$$r ; E_z q(z|x) [\log p(x|z)]$$

Optimizing the ELBO

The ELBO consists of two terms: reconstruction and a KL regularization:

$$L ; = \underbrace{E_z q(z|x) [\log p(x|z)]}_{\text{reconstruction}} - \underbrace{D_{KL}(q(z|x) || p(z))}_{\text{KL regularization}}$$

Need to compute

$$r ; E_z q(z|x) [\log p(x|z)]$$

Both the expectation and the likelihood are functions of model parameters!

Optimizing the ELBO

The ELBO consists of two terms: reconstruction and a KL regularization:

$$L ; = \underbrace{E_z q(z|x) [\log p(x|z)]}_{\text{reconstruction}} \quad \underbrace{D_{KL}(q(z|x) || p(z))}_{\text{KL regularization}}$$

Need to compute

$$r ; E_z q(z|x) [\log p(x|z)] \quad r ; \hat{E}_z q(z|x) [\log p(x|z)]$$

Optimizing the ELBO

The ELBO consists of two terms: reconstruction and a KL regularization:

$$L ; = \underbrace{E_z q(z|x) [\log p(x|z)]}_{\text{reconstruction}} - \underbrace{D_{KL}(q(z|x) || p(z))}_{\text{KL regularization}}$$

Need to compute

$$r ; E_z q(z|x) [\log p(x|z)] \quad r ; \hat{E}_z q(z|x) [\log p(x|z)]$$

Reparameterization Trick:

$$z = \mu(x) + \sigma(x); \text{ where } N(0; I)$$

Optimizing the ELBO

The ELBO consists of two terms: reconstruction and a KL regularization:

$$L ; = \underbrace{E_z q(z|x) [\log p(x|z)]}_{\text{reconstruction}} - \underbrace{D_{KL}(q(z|x) || p(z))}_{\text{KL regularization}}$$

Need to compute

$$r ; E_z q(z|x) [\log p(x|z)] \quad r ; \hat{E}_z q(z|x) [\log p(x|z)]$$

Reparameterization Trick:

$$z = \mu(x) + \sigma(x); \text{ where } N(0; I)$$

$$r ; E_z q(z|x) [\log p(x|z)] \quad r ; \frac{1}{M} \sum_{j=1}^M \log p(x| \underbrace{\mu(x) + \sigma(x)}_{z^{(j)}})$$

Optimizing the ELBO

The ELBO consists of two terms: reconstruction and a KL regularization:

$$L ; = \underbrace{E_z q(z|x) [\log p(x|z)]}_{\text{reconstruction}} - \underbrace{D_{KL}(q(z|x) || p(z))}_{\text{KL regularization}}$$

Need to compute

$$r ; E_z q(z|x) [\log p(x|z)] \quad r ; \hat{E}_z q(z|x) [\log p(x|z)]$$

Reparameterization Trick:

$$z = \mu(x) + \sigma(x) \epsilon; \text{ where } \epsilon \sim N(0; I)$$

$$r ; E_z q(z|x) [\log p(x|z)] \quad r ; \frac{1}{M} \sum_{j=1}^M \log p(x| \underbrace{\mu(x) + \sigma(x) \epsilon^{(j)}}_{z^{(j)} \sim N(\mu(x); \sigma(x))})$$

In practice usually $M = 1$.

Variational Auto Encoder (VAE)

Variational Auto Encoder (VAE)

Variational Auto Encoder (VAE)

$$\max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p(x_i | z_i(\theta)) + \mathbb{E}_{z \sim N(0, I)} \mathbb{E}_{x \sim p(x|z)} \mathbb{D}_{KL}[N(x; \mu(x)) || N(0; I)]$$

What does the VAE actually do

$$\max_{\theta} = \underbrace{\mathbb{E}_z \underbrace{q(z|x) [\log p(x|z)]}_{\text{reconstruction}}}_{\text{reconstruction}} \quad \underbrace{D_{KL}(q(z|x) \parallel p(z))}_{\text{KL regularization}}$$

¹Image credit: Jeremy Jordan's blog post.

What does the VAE actually do

$$\max_{\theta} = \underbrace{\mathbb{E}_z \underbrace{q(z|x) [\log p(x|z)]}_{\text{reconstruction}}}_{\text{reconstruction}} \quad \underbrace{D_{KL}(q(z|x) \parallel p(z))}_{\text{KL regularization}}$$

¹Image credit: Jeremy Jordan's blog post.

What does the VAE actually do

$$\max_{\theta} = \underbrace{\mathbb{E}_z \left[q(z|x) \log p(x|z) \right]}_{\text{reconstruction}} - \underbrace{\mathbb{D}_{KL}(q(z|x) || p(z))}_{\text{KL regularization}}$$

The VAE objective arranges data on a compact manifold (we can sample from) in a continuous smooth way.

¹Image credit: Jeremy Jordan's blog post.

Example: MNIST VAE

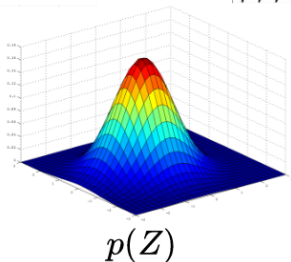
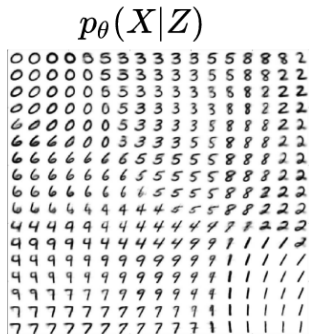
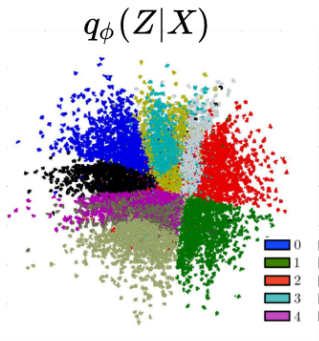


Figure: Samples from CVAE trained on SVHN

Figure: Samples from CVAE trained on SVHN

$$\max_{\theta} \mathbb{E}_{z \sim q(z|x; y)} [\log p(x|z; y)] - D_{KL}(q(z|x; y) \parallel p(z|y))$$

Figure: Samples from CVAE trained on SVHN

$$\max_z \mathbb{E}_{z \sim q(z|x;y)} [\log p(x|z;y)] - D_{KL}(q(z|x;y) \parallel p(z|y))$$

We have (optional) additional conditional-specific prior $p(z|y)$.

Questions?