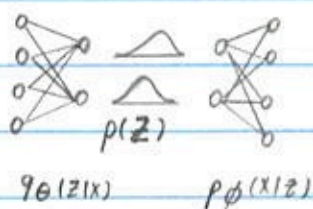# Evidence Lower Bound (ELBO) for VAE

**Goal:** A VAE is an autoencoder so we need our network to reconstruct the data. However, we also need our network to learn to map the latent space to a chosen prior. Thus, we need an objective function that does both of these things.

**Let:**



$$q_\theta(z|x) \qquad p_\phi(x|z)$$
$$p(z)$$

**Recall:** Baye's Rule,

$$p(z|x) = p(z) \cdot \frac{p(x|z)}{\int p(x|z) p(z) dz}$$

- $p(z|x)$ — posterior
- $p(z)$ — prior
- $\int p(x|z)p(z)dz$ — margin likelihood
- $p(x|z)$ — likelihood

**ELBO for VAE:** 1. Let $p(z|x_i)$ be a chosen prior distribution to express latent space.

2. To make the encoder learn to map to our chosen prior, we can minimize, $D_{KL}(q_\theta(z|x_i) \| p(z|x_i))$

3. We shall proceed to show that $D_{KL}(q_\theta(z|x_i) \| p(z|x_i))$ can be used to construct an objective function that maximizes the log-likelihood of our data (the goal of any probabilistic model) by optimizing a lower-bound that takes into account the VAE's objectives of mapping to a chosen prior for the latent and reconstruction of data.

**4.** $$D_{KL}(q_\theta(z|x_i) \| p(z|x_i)) = -\int q_\theta(z|x_i) \log\left(\frac{p(z|x_i)}{q_\theta(z|x_i)}\right) dz \geq 0$$

$$= -\int q_\theta(z|x_i) \log\left(\frac{p_\phi(x_i|z)\, p(z)}{q_\theta(z|x_i)\, p(x_i)}\right) dz \geq 0 \qquad // \text{apply bayes} \atop \text{to } p(z|x_i)$$

$$= -\int q_\theta(z|x_i)\left[\log\left(\frac{p_\phi(x_i|z)\, p(z)}{q_\theta(z|x_i)}\right) - \log(p(x_i))\right] dz \geq 0 \qquad // \log\left(\frac{a}{b}\right) = \log(a) - \log(b)$$

$$= \left(-\int q_\theta(z|x_i)\log\left(\frac{p_\phi(x_i|z)\, p(z)}{q_\theta(z|x_i)}\right) dz\right) + \left(\log(p(x_i)) \underbrace{\int q_\theta(z|x_i)\, dz}_{\text{integrates to } \mathbf{1}}\right) \geq 0$$

$$= \log(p(x_i)) - \int q_\theta(z|x_i) \log\left(\frac{p_\phi(x_i|z)\, p(z)}{q_\theta(z|x_i)}\right) dz \geq 0$$

$$\log(p(x_i)) \geq \int q_\theta(z|x_i) \log\left(\frac{p_\phi(x_i|z)\, p(z)}{q_\theta(z|x_i)}\right) dz$$

$$\geq \int q_\theta(z|x_i) \log\left(\frac{p(z)}{q_\theta(z|x_i)}\right) dz + \int q_\theta(z|x_i) \log(p_\phi(x_i|z)) dz$$

$$\geq \underbrace{-D_{KL}(q_\theta(z|x_i) \| p(z)) + \mathop{\mathbb{E}}_{\sim q_\theta(z|x_i)}[\log(p_\phi(x_i|z))]}_{\text{"ELBO"}}$$

**5.** $$\log(p(x)) \geq \underbrace{-D_{KL}(q_\theta(z|x_i) \| p(z))}_{\nearrow} + \underbrace{\mathop{\mathbb{E}}_{\sim q_\theta(z|x_i)}[\log(p_\phi(x_i|z))]}_{\nwarrow}$$

• This term causes the encoder to map to a chooser prior, $p(z)$.

• This term causes the decoder to map the latent back to our data's space.

• Also, by maximizing ELBO, we maximize the likelihood of our data being observed by our network. Thus, our network learns in respect to our data. □

**Using ELBO for a VAE in practice:** we have shown that maximizing ELBO fits a VAE yet, in practice, we must choose a prior, $p(z)$, for the latent to map the data too. We then need to manipulate ELBO to learn the parameters of the chosen prior.

**Closed Form VAE Loss, Gaussian Latent:** we shall show ELBO for when the prior, $p(z)$, is gaussian.

1. Let $Z \sim N(\mu, \sigma^2)$, then $p(z) = \frac{1}{\sqrt{2\pi\sigma_p^2}} \exp\left(\frac{-(x-\mu_p)^2}{2\sigma_p^2}\right)$

   and $q_\theta(z|x_i) = \frac{1}{\sqrt{2\pi\sigma_q^2}} \exp\left(\frac{-(x-\mu_q)}{2\sigma_q^2}\right)^2$

2. Plug in $p(z)$ and $q_\theta(z|x_i)$ into the $KL$ term of ELBO,

$$-D_{KL}(q_\theta(z|x_i) \| p(z)) = \int q_\theta(z|x_i) \log\left(\frac{q_\theta(z|x_i)}{p(z)}\right) dz$$

$$= \int \frac{1}{\sqrt{2\pi\sigma_q^2}} \exp\left(\frac{-(x-\mu_q)^2}{2\sigma_q^2}\right) \log\left(\frac{\frac{1}{\sqrt{2\pi\sigma_p^2}}\exp\left(\frac{-(x-\mu_p)^2}{2\sigma_p}\right)}{\frac{1}{\sqrt{2\pi\sigma_p^2}}\exp\left(\frac{-(x-\mu_q)}{2\sigma_q^2}\right)}\right) dz$$

3. We will work the current form of $D_{KL}$ to be in terms of $E(\cdot)$ and $V(\cdot)$

$$-D_{KL}(q_\theta(z|x_i) \| p(z)) = \int \frac{1}{\sqrt{2\pi\sigma_q^2}} \exp\left(\frac{-(x-\mu_q)^2}{2\sigma_q^2}\right) \cdot \left[ \left(-\tfrac{1}{2}\right)\log(2\pi) - \log(\sigma_p) - \frac{(x-\mu_p)^2}{2\sigma_p^2} + \left(\tfrac{1}{2}\right)\log(2\pi) + \log(\sigma_q) + \frac{(x-\mu_q)^2}{2\sigma_q^2} \right] dz$$

$$= \frac{1}{\sqrt{2\pi\sigma_q^2}} \int \exp\left(\frac{-(x-\mu_q)^2}{2\sigma_q^2}\right)\left[-\log(\sigma_p) - \frac{(x-\mu_p)^2}{2\sigma_p^2} + \log(\sigma_q) + (\tfrac{1}{2})\log(2\pi) \right.$$
$$\left. + \log(\sigma_q) + \frac{(x-\mu_q)^2}{2\sigma_q^2}\right]dz$$

$$= \frac{1}{\sqrt{2\pi\sigma_q^2}} \int \exp\left(\frac{-(x-\mu_q)^2}{2\sigma_q^2}\right)\left(\log\left(\frac{\sigma_q}{\sigma_p}\right) - \frac{(x-\mu_p)^2}{2\sigma_p^2} + \frac{(x-\mu_q)^2}{2\sigma_q^2}\right)dz$$

$$\Rightarrow -D_{KL}(q_\theta(z|x;) \| p(z)) = E_q\left[\log\left(\frac{\sigma_q}{\sigma_p}\right) - \frac{(x-\mu_p)^2}{2\sigma_p^2} + \frac{(x-\mu_q)^2}{2\sigma_q^2}\right]$$

$$= \log\left(\frac{\sigma_q}{\sigma_p}\right) + E_q\left[\frac{-(x-\mu_p)^2}{2\sigma_p^2} + \frac{(x-\mu_q)^2}{2\sigma_q^2}\right]$$

$$= \log\left(\frac{\sigma_q}{\sigma_p}\right) - \frac{1}{2\sigma_p^2}E_q[(x-\mu_p)^2] + \frac{1}{2\sigma_q^2}E_q\{(x-\mu_q)^2\}$$

4. Recall $V(\cdot) = \sigma^2 = E[(x-\mu)^2]$, thus

$$-D_{KL}(q_\theta(z|x;) \| p(z)) = \log\left(\frac{\sigma_q}{\sigma_p}\right) - \frac{1}{2\sigma_p^2}E_q[(x-\mu_p)^2] + \frac{\sigma_q^2}{2\sigma_q^2}$$

$$= \log\left(\frac{\sigma_q}{\sigma_p}\right) - \frac{1}{2\sigma_p^2}E_q[(x-\mu_p)^2] + \frac{1}{2}$$

$$= \log\left(\frac{\sigma_q}{\sigma_p}\right) - \frac{1}{2\sigma_p^2}E_q[\underbrace{(x-\mu_q)}_{a} + \underbrace{\mu_q - \mu_p)}_{b}^2] + \frac{1}{2}$$

recall, $(a+b)^2 = a^2 + 2ab + b^2$

$$= \log\left(\frac{\sigma_q}{\sigma_p}\right) - \frac{1}{2\sigma_p^2}E_q[(x-\mu_q)^2 + 2(x-\mu_q)(\mu_q-\mu_p) + (\mu_q-\mu_p)^2] + \frac{1}{2}$$

$$= \log\left(\frac{\sigma_q}{\sigma_p}\right) - \frac{\sigma_q^2 + (\mu_q-\mu_p)^2}{2\sigma_p^2} + \frac{1}{2}$$

5. Assume prior, $p(z)$, is $N(0,1)$. Thus, $\sigma_p = 1$ and $\mu_p = 0$.

$$\therefore \quad -D_{KL}(q_\theta(z|x_i) \| p(z)) = \log(\sigma_q) - \frac{\sigma_q^2 + \mu_q^2}{2} + \frac{1}{2}$$

$$= (\frac{1}{2})[1 + \log(\sigma_p^2) - \sigma_q^2 - \mu_q^2]$$

6. Thus, ELBO with a standard normal prior is

$$\log(p(x_i)) \geq \frac{1}{2}[1 + \log(\sigma_j^2) - \sigma_j^2 - \mu_j^2] + \underset{\sim q_\theta(z|x_i)}{\mathbb{E}}[\log(p_\phi(x_i|z))]$$

For a training batch, ELBO with $N(0,1)$ prior is

$$G = \sum_{j=1}^{J}(\frac{1}{2})[\log(\sigma_j^2) - \sigma_j^2 - \mu_j^2 + 1] + \frac{1}{L}\sum_{\ell}\underset{\sim q_\theta(z|x_i)}{\mathbb{E}}[\log(p(x_i|z^{(i,\ell)}))]$$

where $J$ is the dimension of the latent vector $z$, and $L$ is number of samples drawn according to reparameterization trick.

Re-parameterization Trick : In ELBO, $\underset{\sim q_\theta(z|x_i)}{\mathbb{E}}[\log(p_\phi(x_i|z))]$ is not differentiable. Thus, SGD does not work. By using Stochastic Gradient Variational Bayes (SGVB), we can approximate $\underset{\sim q_\theta(z|x_i)}{\mathbb{E}}[\log(p_\phi(x_i|z))]$.

To do so, we will reparameterize the random variable $\tilde{z} \sim q_\phi(z|x)$ using a differentiable function $g_\phi(\varepsilon, x)$ where $\varepsilon$ is a noise variable

$$\tilde{z} = g_\phi(\varepsilon, x) \quad \text{with} \quad \varepsilon \sim p(\varepsilon)$$

We can form Monte Carlo estimates of expectations of some function $f(z)$ w.r.t. $q_\phi(z|x)$ as follows:

$$\mathbb{E}_{q_\phi(z|x_i)}[f(z)] = \mathbb{E}_{p(\epsilon)}[f(g_\phi(\epsilon, x_i))] \simeq \frac{1}{L}\sum_{\ell=1}^{L} f(g_\phi(\epsilon^{(\ell)}, x^{(i)}))$$

$$\text{where } \epsilon^{(\ell)} \sim p(\epsilon)$$

Applying to lower bound $\mathbb{E}_{\sim q_\theta(z|x_i)}[\log(p(x_i|z^{(i,\ell)}))]$

$$\hat{\mathcal{L}}^A(\theta, \phi; x^{(i)}) \simeq \mathcal{L}(\theta, \phi; x^{(i)})$$

$$\hat{\mathcal{L}}^A(\theta, \phi; x^{(i)}) = \frac{1}{L}\sum_{\ell=1}^{L} \log(p_\theta(x^{(i)}, z^{(i,\ell)})) - \log(q_\phi(z^{(i,\ell)}|x^{(i)}))$$

$$\text{where } z^{(i,\ell)} = g_\phi(\epsilon^{(i,\ell)}, x^{(i)}) \text{ and } \epsilon^{(\ell)} \sim p(\epsilon)$$

Proof:
1. Given the deterministic mapping $z = g_\phi(\epsilon, x)$ we know that
$$q_\phi(z|x) \prod_i dz_i = p(\epsilon) \prod_i d\epsilon_i$$

2. $\therefore \int q_\phi(z|x) f(z) dz = \int p(\epsilon) f(z) d\epsilon = \int p(\epsilon) f(g_\phi(\epsilon, x)) d\epsilon$

3. Thus, a differentiable estimator can be made
$$\int q_\phi(z|x) f(z) dz \simeq \frac{1}{L}\sum_{\ell=1}^{L} f(g_\phi(x, \epsilon^{(\ell)})) \text{ where}$$
$$\epsilon^{(\ell)} \sim p(\epsilon)$$

## VAE ELBO Re-param Gauss:

$$\log(p(x_j) \geq (\tfrac{1}{2})\sum_j [1 + \log(\sigma_j^2) - \sigma_j^2 - \mu_j^2] + \mathbb{E}_{\sim q_\theta(z|x_j)}[\log(p_\phi(x_j|z_j))]$$

$$\text{where } z_j = \mu_j + \sigma_j \odot \epsilon \quad \text{and} \quad \epsilon \sim N(0,1)$$