

TTIC 31230, Fundamentals of Deep Learning

David McAllester, Winter 2022

Variational Auto-Encoders (VAEs)

Image Compression and Image Generation

Suppose that we want to model a population distribution on y , for example the distribution of “natural images”.

Shannon’s source coding theorem implies that there exists a coding function with an inverse decoding function such that decoding a random string samples an image from the population distribution on images.

While we cannot optimally compress images, it can be useful to represent a population distribution on y in terms of a latent (unobserved) variable $z(y)$ loosely analogous to a compressed form.

The Encoder, Decoder and the Prior

Consider a probabilistic encoder algorithm $P_{\text{enc}}(z|y)$ — perhaps a stochastic image compression algorithm.

The encoder $P_{\text{enc}}(z|y)$ defines a joint probability distribution on pairs (y, z) .

We will also introduce a decoder model (decompressor) $P_{\text{dec}}(y|z)$ and a prior probability model $P_{\text{pri}}(z)$ which are to be trained using a cross-entropy loss to $P(y|z)$ and $P(z)$ as defined by the encoder.

The ELBO

$$H(y, z) = H(y) + H(z|y) = H(z) + H(y|z)$$

$$H(y) = H(z) + H(y|z) - H(z|y)$$

$$\leq CE(P(z), P_{\text{pri}}(z)) + CE(P(y|z), P_{\text{dec}}(y|z)) - H_{\text{enc}}(z|y)$$

$$= E_{y \sim \text{Pop}, z \sim P_{\text{enc}}(z|y)} - \ln \frac{P_{\text{pri}}(z) P_{\text{dec}}(y|z)}{P_{\text{enc}}(z|y)}$$

The last line is the (negative) ELBO. ELBO stands for “Evidence Lower Bound” but this terminology is obscure and unhelpful.

The ELBO Loss

We now interpret the ELBO as a loss on a given value y .

$$\begin{aligned} H(y) &\leq E_{y \sim \text{Pop}, z \sim P_{\text{enc}}(z|y)} - \ln \frac{P_{\text{pri}}(z)P_{\text{dec}}(y|z)}{P_{\text{enc}}(z|y)} \\ &= E_{y \sim \text{Pop}} \mathcal{L}_E(y) \end{aligned}$$

$$\mathcal{L}_E(y) = E_{z \sim P_{\text{enc}}(z|y)} - \ln \frac{P_{\text{pri}}(z)P_{\text{dec}}(y|z)}{P_{\text{enc}}(z|y)}$$

A Third Fundamental Equation

$$\text{pri}^*, \text{dec}^* = \underset{\text{pri}, \text{dec}}{\operatorname{argmin}} \quad E_y \mathcal{L}_E(y)$$

$$E_y \mathcal{L}_E(y)$$

$$= E_{y \sim \text{Pop}, z \sim P_{\text{enc}}(z|y)} \ln \frac{P_{\text{enc}}(z|y)}{P_{\text{pri}}(z)P_{\text{dec}}(y|z)}$$

$$= E_y [-\ln \text{Pop}(y)] + KL(\text{Pop}(y)P_{\text{enc}}(z|y), P_{\text{pri}}(z)P_{\text{dec}}(y|z))$$

$$= H(y) + KL(\text{Pop}(y)P_{\text{enc}}(z|y), P_{\text{pri}}(z)P_{\text{dec}}(y|z))$$

The ELBO Loss

$$E_y \mathcal{L}_E(y) = H(y) + KL(\text{Pop}(y)P_{\text{enc}}(z|y), P_{\text{pri}}(z)P_{\text{dec}}(y|z))$$

Minimization occurs when the prior and the docoder satisfy

$$P_{\text{pri}}(z)P_{\text{dec}}(y|z) = \text{Pop}(y)P_{\text{enc}}(z|y)$$

The ELBO Loss

$$E_y \mathcal{L}_E(y) = H(y) + KL(\text{Pop}(y)P_{\text{enc}}(z|y), P_{\text{pri}}(z)P_{\text{dec}}(y|z))$$

Minimizing gives

$$P_{\text{pri}}(z)P_{\text{dec}}(y|z) = \text{Pop}(y)P_{\text{enc}}(z|y) = P(y, z)$$

and hence

$$\mathcal{L}_E(y) = E_z \ln \frac{P(z|y)}{P(z, y)} = E_z \ln \frac{P(z, y)/P(y)}{P(z, y)} = -\ln \text{Pop}(y)$$

After optimization one can interpret $\mathcal{L}_E(y)$ as $-\ln \text{Pop}(y)$.

Optimizing the Encoder

Although in principle the encoder need not be trained, it is sometimes jointly optimized with the prior and the decoder.

$$\text{pri}^*, \text{dec}^*, \text{enc}^* = \underset{\text{pri}, \text{dec}, \text{enc}}{\operatorname{argmin}} E_{y, z \sim P_{\text{enc}}(z|y)} \left[-\ln \frac{P_{\text{pri}}(z) P_{\text{dec}}(y|z)}{P_{\text{enc}}(z|y)} \right]$$

This is necessary if want to interpret z as some kind of “understanding” of the distribution on y that facilitates representing the prior and the decoder.

Optimizing the Encoder

$$\text{pri}^*, \text{dec}^*, \text{enc}^* = \underset{\text{pri}, \text{dec}, \text{enc}}{\operatorname{argmin}} E_{y, z \sim P_{\text{enc}}(z|y)} \left[-\ln \frac{P_{\text{pri}}(z) P_{\text{dec}}(y|z)}{P_{\text{enc}}(z|y)} \right]$$

Gradient descent on the encoder parameters must take into account the fact that we are sampling from the encoder.

To handle this we sample noise ϵ from a fixed noise distribution and replace z with a deterministic function $z_{\text{enc}}(y, \epsilon)$

$$\text{enc}^*, \text{pri}^*, \text{dec}^* = \underset{\text{enc}, \text{pri}, \text{dec}}{\operatorname{argmin}} E_{y, \epsilon, z = z_{\text{enc}}(y, \epsilon)} \left[-\ln \frac{P_{\text{pri}}(z) P_{\text{dec}}(y|z)}{P_{\text{enc}}(z|y)} \right]$$

The Re-Parameterization Trick

$$\text{enc}^*, \text{pri}^*, \text{dec}^* = \underset{\text{enc}, \text{pri}, \text{dec}}{\text{argmin}} \quad E_{y, \epsilon, z=z_{\text{enc}}(y, \epsilon)} \left[-\ln \frac{P_{\text{pri}}(z) P_{\text{dec}}(y|z)}{P_{\text{enc}}(z|y)} \right]$$

To get gradients we must have that $z_{\text{enc}}(y, \epsilon)$ is a smooth function of the encoder parameters and all probabilities must be a smooth function of z .

Most commonly $\epsilon \in \mathbb{R}^d$ with $\epsilon \sim \mathcal{N}(0, I)$ and

$$z_{\text{enc}}^i(y, \epsilon) = \hat{z}_{\text{enc}}^i(y) + \sigma^i \epsilon^i.$$

Optimizing the encoder is tricky for discrete z . Discrete z is handled effectively in EM algorithms and in VQ-VAEs.

EM is Alternating Optimization of the ELBO Loss

Expectation Maximimization (EM) applies in the (highly special) case where the exact posterior $P_{\text{pri},\text{dec}}(z|y)$ is samplable and computable. EM alternates exact optimization of enc and the pair (pri, dec) in:

$$\text{VAE:} \quad \text{pri}^*, \text{dec}^* = \underset{\text{pri}, \text{dec}}{\operatorname{argmin}} \min_{\text{enc}} E_{y, z \sim P_{\text{enc}}(z|y)} - \ln \frac{P_{\text{pri}, \text{dec}}(z, y)}{P_{\text{enc}}(z|y)}$$

$$\text{EM:} \quad \text{pri}^{t+1}, \text{dec}^{t+1} = \underset{\text{pri}, \text{dec}}{\operatorname{argmin}} E_{y, z \sim P_{\text{pri}^t, \text{dec}^t}(z|y)} - \ln P_{\text{pri}, \text{dec}}(z, y)$$

Inference
(E Step)

$$P_{\text{enc}}(z|y) = P_{\text{pri}^t, \text{dec}^t}(z|y)$$

Update
(M Step)

Hold $P_{\text{enc}}(z|y)$ fixed

END