

TTIC 31230, Fundamentals of Deep Learning

David McAllester, Winter 2022

Variational Auto-Encoders (VAEs)

Rate-Distortion Autoencoders

Consider image compression where we compress an image y into a compressed file z .

We will assume a stochastic compression algorithm which we will call the “encoder” $P_{\text{enc}}(z|y)$.

The number of bits needed for the compressed file is given by $H(z)$. $H(z)$ is the “rate” (bits per image) for transmitting compressed images.

The number of unknown additional bits needed to exactly recover y is $H(y|z)$. $H(y|z)$ is a measure of the “distortion” of y when y is decoded without the missing bits.

Rate-Distortion Autoencoders

In practice we model $H(z)$ with a “prior model” $P_{\text{pri}}(z)$ and model $H(y|z)$ with a “decoder model” $P_{\text{dec}}(y|z)$.

So the rate-distortion auto-encoder has three parts $P_{\text{enc}}(z|y)$, $P_{\text{pri}}(z)$, and $P_{\text{dec}}(y|z)$.

The **variational autoencoder (VAE)** with latent variable z is mathematically the same as a rate-distortion autoencoder with compressed form z .

An “Encoder First” Treatment of VAEs

Fix an arbitrary encoder model $P_{\text{enc}}(z|y)$.

For $y \sim \text{Pop}$ and $z \sim P_{\text{enc}}(z|y)$ train models pri and dec.

$$\text{Prior Model: } \text{pri}^* = \underset{\text{pri}}{\text{argmin}} \quad E_{y,z} \quad - \ln P_{\text{pri}}(z)$$

$$\text{Decoder Model: } \text{dec}^* = \underset{\text{dec}}{\text{argmin}} \quad E_{y,z} \quad - \ln P_{\text{dec}}(y|z)$$

For any $P_{\text{enc}}(z|y)$ the universality assumption for pri^* and dec^* gives

$$\text{Pop}(y) = \sum_z P_{\text{pri}^*}(z) P_{\text{dec}^*}(y|z)$$

Upper Bounding $H(y)$

Cross-entropy upper bounds $H(y)$ and equals $H(y)$ assuming universality.

The ELBO plays the role of cross-entropy for latent variable models.

The negative ELBO upper bounds $H(y)$ and equals $H(y)$ assuming universality.

Deriving the ELBO

Sample $y \sim \text{Pop}$ and $z \sim P_{\text{enc}}(z|y)$.

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

Solving for $H(y)$ gives

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

An Encoder First Treatment of VAEs

Replace the first two entropies by cross entropies.

$$H(y) = H(z) + H(y|z) - H_{\text{enc}}(z|y)$$

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

$$H_{\text{pri}}(z) = E_{y,z} [-\ln P_{\text{pri}}(z)]$$

$$H_{\text{dec}}(y|z) = E_{y,z} [-P_{\text{dec}}(y|z)]$$

VAE

$$\begin{aligned} H(y) &= H(z) + H(y|z) - H_{\text{enc}}(z|y) \\ &\leq H_{\text{pri}}(z) + H_{\text{dec}}(y|z) - H_{\text{enc}}(z|y) \\ &= E_{y,z} \left[-\ln \frac{P_{\text{pri}}(z)P_{\text{dec}}(y|z)}{P_{\text{enc}}(z|y)} \right] \end{aligned}$$

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

The Re-Parameterization Trick

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

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 VAE

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:

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

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

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