

TTIC 31230, Fundamentals of Deep Learning

David McAllester, Autumn 2024

Variational Auto-Encoders (VAEs)

Fundamental Equations of Deep Learning

- Cross Entropy Loss: $\Phi^* = \operatorname{argmin}_{\Phi} E_{(x,y) \sim P_{\text{op}}} [-\ln P_{\Phi}(y|x)]$.
- GAN: $\text{gen}^* = \operatorname{argmax}_{\text{gen}} \min_{\text{disc}} E_{i \sim \{-1,1\}, y \sim P_i} [-\ln P_{\text{disc}}(i|y)]$.
- VAE (including diffusion models)
 $\text{pri}^*, \text{dec}^*, \text{enc}^*$
$$= \operatorname{argmin}_{\text{pri}, \text{dec}, \text{enc}} E_{y \sim P_{\text{op}}, z \sim P_{\text{enc}}(z|y)} \left[-\ln \frac{P_{\text{pri}}(z) P_{\text{dec}}(y|z)}{P_{\text{enc}}(z|y)} \right]$$

VAEs

A variational autoencoder (VAE) is defined by three parts:

- An encoder distribution $P_{\text{enc}}(z|y)$.
- A “prior” distribution $P_{\text{pri}}(z)$
- A generator distribution $P_{\text{dec}}(y|z)$

VAE generation uses $P_{\text{pri}}(z)$ and $P_{\text{dec}}(y|z)$ (like a GAN).

VAE training uses an encoder $P_{\text{enc}}(z|y)$.

Fixing an Arbitray Encoder

Fix the encoder arbitrarily and train P_{pri} and P_{dec} by cross entropy loss.

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

Universality gives

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

Sampling from $P_{\text{pri}^*}(z) P_{\text{dec}^*}(y|z)$ now samples y from the population.

Training the Encoder — the ELBO

In practice the choice of encoder matters.

$$P(y) = \frac{P_{\text{pop}}(y)P_{\text{enc}}(z|y)}{P_{\text{enc}}(z|y)} = \frac{P_{\text{enc}}(z)P_{\text{enc}}(y|z)}{P_{\text{enc}}(z|y)}$$

$$H(y) \leq E_{y \sim P_{\text{pop}}, z \sim P_{\text{enc}}(z|y)} \left[-\ln \frac{P_{\text{pri}}(z)P_{\text{dec}}(y|z)}{P_{\text{enc}}(z|y)} \right]$$

The inequality follows from the fact that cross-entropy (using the models P_{pri} and P_{dec}) upper bounds entropy.

This upper bound on $H(y)$ is called **the ELBO** (Acronym described later).

Difficulties in Training the Encoder

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

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

Training a sampling distribution typically suffers from mode collapse (as in GANs).

Often the encoder collapses to fixing $z = 0$. $P_{\text{dec}}(y|z)$ can always just ignore z . We are then back to standard cross-entropy loss. This is called posterior collapse.

Types of VAEs

In a **Gaussian VAE** we have $P_{\text{pri}}(z)$ and $P_{\text{enc}}(z|y)$ are both Gaussian distributions on R^d . A diffusion model involves a Gaussian VAE at each incremental step of diffusion.

A Vector Quantized VAE (VQ-VAE) defines $P_{\text{enc}}(z|y)$ in terms of vector quantization analogous to K -means clustering. VQ-VAEs provide a translation from continuous data, such as images, to token data that can be modeled with a transformer. This is done in Chat-GPT 4o and other multi-modal language models.

We will consider each these approaches.

Gaussian VAEs

We sample noise ϵ from a Gaussian distribution on R^d .

$$\begin{aligned}\text{enc}^* &= \underset{\text{enc}}{\operatorname{argmin}} \quad E_{y, \epsilon \sim \mathcal{N}(0, \sigma I)} \left[-\ln \frac{P_{\text{pri}}(z) P_{\text{dec}}(y|z)}{P_{\text{enc}}(z|y)} \right] \quad z = \hat{z}(y) + \epsilon \\ &= \underset{\text{enc}}{\operatorname{argmin}} \quad KL(P_{\text{enc}}(z|y), P_{\text{pri}}(z)) + E_{z \sim P_{\text{enc}}(z|y)} [-\ln P_{\text{dec}}(y|z)] \\ &= \underset{\text{enc}}{\operatorname{argmin}} \quad \frac{||\hat{z}_{\text{enc}}(y) - \hat{z}_{\text{pri}}||^2}{2\sigma^2} + E_{\epsilon} \frac{||y - \hat{y}_{\text{dec}}(\hat{z}_{\text{enc}}(y, \epsilon))||^2}{2\sigma^2}\end{aligned}$$

A closed-form expression for the KL term avoids sampling noise.

VAEs Evolved from Variational Bayesian Inference

$P_{\text{pri}}(z)$ is interpreted as the Bayesian prior on “hypothesis” z .

$P_{\text{dec}}(y|z)$ is interpreted as the propability of the “evidence” y given hypothesis z .

We consider the Bayesian distribution defined by $P_{\text{pri}}(z)$ and $P_{\text{dec}}(y|z)$ and we want to compute $P_{\text{pri,dec}}(z|y)$ under this Bayesian distribution. I will write this $P_{\text{Bayes}}(z|y)$

$P_{\text{enc}}(z|y)$ is interpreted as an approximation for $P_{\text{pri,dec}}(z|y)$.

Bayesian Interpretation

$$\begin{aligned}\ln P_{\text{Bayes}}(y) &= \ln \frac{P_{\text{Bayes}}(y)(z)P_{\text{Bayes}}(z|y)}{P_{\text{Bayes}}(z|y)} \\ &= E_{z \sim P_{\text{enc}}(z|y)} \left[\ln \frac{P_{\text{pri}}(z)P_{\text{dec}}(y|z)}{P_{\text{Bayes}}(z|y)} \right] \\ &\geq E_{z \sim P_{\text{enc}}(z|y)} \left[\ln \frac{P_{\text{pri}}(z)P_{\text{dec}}(y|z)}{P_{\text{enc}}(z|y)} \right]\end{aligned}$$

Here we have replaced a cross-entropy by an entropy.

Bayesian Interpretation

$$\ln P_{\text{Bayes}}(y) \geq E_{z \sim P_{\text{enc}}(z|y)} \left[\ln \frac{P_{\text{pri}}(z) P_{\text{dec}}(y|z)}{P_{\text{enc}}(z|y)} \right]$$

Here y is the evidence about z under the Bayesian model.

This is the **evidence lower bound** or **ELBO**.

Expectation Maximization (EM)

Expectation Maximization (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:

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

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

Inference
(E Step)

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

slidePosterior Collapse

Update

(M Step)

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

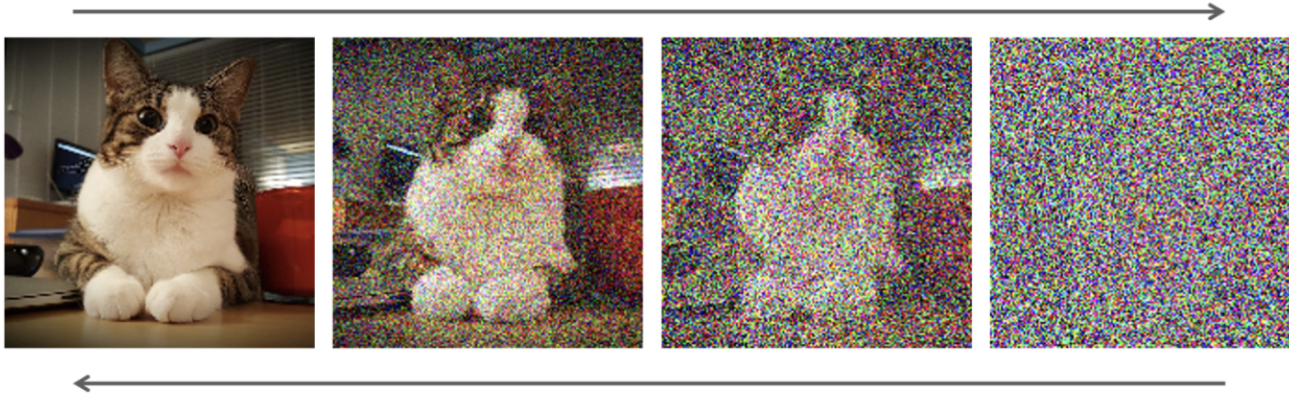
$$P_{\text{pri}}(z)P_{\text{dec}}(y|z) = \text{Pop}(y)P_{\text{enc}}(z|y)$$

Any joint distribution on (y, z) with the desired marginal on y optimizes the bound.

This allows the prior and the encoder (the posterior) to both degenerate to having no mutual information with y .

This often happens in language modeling.

Hierarchical VAEs



[Sally talked to John] $\xleftrightarrow{\quad}$ [Sally talked to] $\xleftrightarrow{\quad}$ [Sally talked] $\xleftrightarrow{\quad}$ [Sally] $\xleftrightarrow{\quad}$ []

$$y \xleftrightarrow{\quad} z_1 \xleftrightarrow{\quad} \dots \xleftrightarrow{\quad} z_N$$

Hierarchical VAEs

$$y \overset{\rightarrow}{\leftarrow} z_1 \overset{\rightarrow}{\leftarrow} \dots \overset{\rightarrow}{\leftarrow} z_N$$

Encoder: $\text{Pop}(y)$, $P_{\text{enc}}(z_1|y)$, and $P_{\text{enc}}(z_{\ell+1}|z_\ell)$.

Generator: $P_{\text{pri}}(z_N)$, $P_{\text{dec}}(z_{\ell-1}|z_\ell)$, $P_{\text{dec}}(y|z_1)$.

The encoder and the decoder define distributions $P_{\text{enc}}(y, \dots, z_N)$ and $P_{\text{dec}}(y, \dots, z_N)$ respectively.

Hierarchical VAEs

$$y \overset{\rightarrow}{\leftarrow} z_1 \overset{\rightarrow}{\leftarrow} \dots \overset{\rightarrow}{\leftarrow} z_N$$

- autoregressive models
- diffusion models

Hierarchical (or Diffusion) ELBO

$$\begin{aligned}
H(y) &= E_{\text{enc}} \left[-\ln \frac{P_{\text{enc}}(y) P_{\text{enc}}(z_1, \dots, z_N | y)}{P_{\text{enc}}(z_1, \dots, z_N | y)} \right] \\
&= E_{\text{enc}} \left[-\ln \frac{P_{\text{enc}}(y|z_1) P_{\text{enc}}(z_1|z_2) \cdots P_{\text{enc}}(z_{N-1}|z_N) P_{\text{enc}}(z_N)}{P_{\text{enc}}(z_1|z_2, y) \cdots P_{\text{enc}}(z_{N-1}|z_N, y) P_{\text{enc}}(z_N|y)} \right] \\
&\leq E_{\text{enc}} \left[-\ln \frac{P_{\text{dec}}(y|z_1) P_{\text{dec}}(z_1|z_2) \cdots P_{\text{dec}}(z_{N-1}|z_N) P_{\text{dec}}(z_N)}{P_{\text{enc}}(z_1|z_2, y) \cdots P_{\text{enc}}(z_{N-1}|z_N, y) P_{\text{enc}}(z_N|y)} \right] \\
&= \begin{cases} E_{\text{enc}} [-\ln P_{\text{dec}}(y|z_1)] \\ + \sum_{i=2}^N E_{\text{enc}} KL(P_{\text{enc}}(z_{i-1}|z_i, y), P_{\text{dec}}(z_{i-1}|z_i)) \\ + E_{\text{enc}} KL(P_{\text{enc}}(Z_N|y), p_{\text{dec}}(Z_N)) \end{cases}
\end{aligned}$$

END