TTIC 31230, Fundamentals of Deep Learning David McAllester, Autumn 2024

Variational Auto-Encoders (VAEs)

Widely Used VAEs

Diffusion Models: A diffusion model is a hierarchical Gaussian VAE.

Vector Quantized VAEs: A 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 modelled with a transformer. This is done in GPT-40.

Auto-Regressive Language Models: An auto-regressive language model, such as a transformer, is mathematically equivalent to a hierarchical VAE.

VAEs

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

- An encoder distribution $P_{\text{enc}}(z|y)$.
- A decoder distribution $P_{\text{dec}}(y|z)$
- A "prior" distribution $P_{\rm pri}(z)$

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

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

Two Joint Distributions

A VAE defines two joint distributions on y and z, namely $P_{\text{gen}}(y, z)$ and $P_{\text{enc}}(y, z)$ defined by

 $P_{\text{gen}}(y, z) = P_{\text{pri}}(z)P_{\text{dec}}(y|z)$

 $P_{\text{enc}}(y, z) = \operatorname{Pop}(y)P_{\text{enc}}(z|y)$

Training the Generator

Fix the encoder arbitrarily and train P_{gen} by cross entropy.

gen^{*} = argmin
$$E_{(y,z)\sim P_{\text{enc}}(y,z)} \left[-\ln P_{\text{gen}}(y,z) \right]$$

Under universality we have $P_{\text{gen}^*} = P_{\text{enc}}$ and hence the generator distribution on y defined by gen^{*} matches the population distribution.

In a diffusion model the encoder just adds noise. The encoder is not trained.

The ELBO Loss

$$\operatorname{Pop}(y) = \frac{\operatorname{Pop}(y)P_{\operatorname{enc}}(z|y)}{P_{\operatorname{enc}}(z|y)} = \frac{P_{\operatorname{enc}}(y,z)}{P_{\operatorname{enc}}(z|y)}$$
$$E_{(y,z)\sim P_{\operatorname{enc}}}\left[-\ln\operatorname{Pop}(y)\right] = E_{(y,z)\sim P_{\operatorname{enc}}}\left[-\ln\frac{P_{\operatorname{enc}}(y,z)}{P_{\operatorname{enc}}(z|y)}\right]$$
$$H(y) \leq E_{(y,z)\sim P_{\operatorname{enc}}}\left[-\ln\frac{P_{\operatorname{gen}}(y,z)}{P_{\operatorname{enc}}(z|y)}\right]$$

The right hand side of the last line is called the **ELBO Loss**.

The Variational Bayes Interpretation

The generator is interpreted as a Bayesian model where y is evidence for z.

$$\ln P_{\text{gen}}(y) = \ln \frac{P_{\text{gen}}(y)P_{\text{gen}}(z|y)}{P_{\text{gen}}(z|y)}$$

$$= E_{z \sim P_{\text{enc}}(z|y)} \left[\ln \frac{P_{\text{gen}}(y,z)}{P_{\text{gen}}(z|y)} \right]$$

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

Hence the name **evidence lower bound** or **ELBO**.

Fundamental Equations of Deep Learning

• Cross Entropy Loss: $\Phi^* =$

$$\underset{\Phi}{\operatorname{argmin}} E_{(x,y)\sim \operatorname{Pop}} \left[-\ln P_{\Phi}(y|x) \right]$$

• GAN: $gen^* =$

 $\underset{\text{gen disc}}{\operatorname{argmax}} \min_{\text{disc}} E_{i \sim \{-1,1\}, y \sim P_i} \left[-\ln P_{\text{disc}}(i|y) \right]$

• VAE:
$$pri^*$$
, dec^* , $enc^* =$

argmin
$$E_{(y,z)\sim P_{\text{enc}}} \left[-\ln \frac{P_{\text{gen}}(y,z)}{P_{\text{enc}}(z|y)} \right]$$

Training the Encoder

In diffusion models the encoder simply adds noise and is not trained.

In a VQ-VAE the encoder is trained. A naive approach to training the encoder is to optimize the ELBO loss.

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

Training the Encoder

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

Unfortunately this optimization involves optimizing a sampling distribution (the encoder). As with GANs, optimizing a sampling distribution (such as a GAN generator) is subject to mode collapse — the loss function is very forgiving of a failure of the sampling distribution to cover the desired space of values.

Gaussian VAEs

Gaussian VAEs, which are the original version from 2014, train the encoder.

Gaussian VAEs solve the problem of optimizing a sampling distribution by sampling instead from fixed Gaussian noise. This further allows expressing the ELBO loss as a "closed form" L_2 loss which avoids the need to even sample the noise.

However, non-hierarchical Gaussian VAEs (with a single Gaussian latent variable) produce poor results in practice. A diffusion model is a hierarchical Gaussian VAEs which does not train the encoder. Hierarchical Gaussian VAEs which train the encoder can also produce good results.

A Non-Hierarchical Gaussian VAE

$$P_{ ext{pri}}(z) = \mathcal{N}(0, I)$$

 $P_{ ext{enc}}(z|y) = \mathcal{N}(\hat{z}(y), I)$
 $P_{ ext{dec}}(y|z) = \mathcal{N}(\hat{y}(z), I)$

In general we can use arbitrary Gaussians but this example makes the math simple.

Gaussian VAEs

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

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

$$= E_{y \sim \text{Pop}} \left[\frac{1}{2} || \hat{z}_{\text{enc}}(y) ||^2 + E_{\epsilon} \left[\frac{1}{2} || y - \hat{y}_{\text{dec}}(\hat{z}_{\text{enc}}(y) + \epsilon)) ||^2 \right] \right]$$

Training the Encoder

In a VQ-VAE the encoder is traned jointly with the decoder $P_{\text{dec}}(y|z)$ but is trained independently of $P_{\text{pri}}(z)$. $P_{\text{pri}}(z)$ is trained later using a transformer model. The encoder of a VQ-VAE is closely related to K-means clustering. In a VQ-VAE the encoder converts vectors to tokens so that a transformer can be applied.

This minimal training of the encoder again exploits the fact that under universality Pop(y) can be modelled fully for any encoder.

A different approach to training the encoder, an ME-VAE, is discussed below.

Hierarchical VAEs

Hierarchical Gaussian VAEs which train the encoder are explored both theoretically and empirically by Vahdat and Kautz.

NVAE: A Deep Hierarchical Variational Autoencoder, Arash Vahdat, Jan Kautz (NVIDIA, January 2021)

But diffusion models and autoregressive models are also instances of hierarchical VAEs.

Hierachical VAEs



 $[Sally talked to John] \stackrel{\rightarrow}{\leftarrow} [Sally talked to] \stackrel{\rightarrow}{\leftarrow} [Sally talked] \stackrel{\rightarrow}{\leftarrow} [Sally] \stackrel{\rightarrow}{\leftarrow} []$

$$y \stackrel{\rightarrow}{\leftarrow} z_1 \stackrel{\rightarrow}{\leftarrow} \cdots \stackrel{\rightarrow}{\leftarrow} z_N$$

Hierarchical VAEs

$$y \stackrel{\rightarrow}{\leftarrow} z_1 \stackrel{\rightarrow}{\leftarrow} \cdots \stackrel{\rightarrow}{\leftarrow} z_N$$

Encoder: Pop(y), $P_{enc}(z_1|y)$, $P_{enc}(z_2|z_1)$, ..., $P(z_N|z_{N-1})$.

Generator:
$$P_{\text{pri}}(z_N), P_{\text{dec}}(z_{N-1}|z_N), \dots, P_{\text{dec}}(z_1|z_2), P_{\text{dec}}(y|z_1)$$

The encoder and the decoder define distributions $P_{\text{enc}}(y, z_1, \ldots, z_N)$ and $P_{\text{gen}}(y, z_1, \ldots, z_N)$ respectively.

Hierarchical ELBO Loss

$$H(y) = E_{(y,z_1,...,z_n) \sim P_{\text{enc}}} \left[-\ln \frac{P_{\text{gen}}(y,z_1,...,z_n)}{P_{\text{enc}}(z_1,...,z_N|y)} \right]$$

EM-VAEs

The use of minimal encoder training may reflect the mode collapse problem of training a sampling distribution, such as a GAN generator or a VAE encoder.

The situation might be different if a better method were available for training the encoder. Here I will propose a method for training the encoder that avoids the mode collapse problem.

EM-VAEs

We start with the following "optimum encoder" inequality.

$$E_{y \sim \text{Pop}, z \sim P_{\text{gen}}(z|y)} \left[-\ln \frac{P_{\text{gen}}(y, z)}{P_{\text{gen}}(z|y)} \right] \le E_{(y, z) \sim P_{\text{enc}}} \left[-\ln \frac{P_{\text{gen}}(y, z)}{P_{\text{enc}}(z|y)} \right]$$

This implies $P_{\text{enc}}^*(z|y) = P_{\text{gen}}(z|y)$ and universality gives

$$\operatorname{enc}^* = \operatorname{argmin}_{\operatorname{enc}} E_{(y,z)\sim P_{\operatorname{gen}}} - \ln P_{\operatorname{enc}}(z|y)$$

EM-VAE

E: enc⁺ = argmin
$$E_{(y,z)\sim P_{\text{gen}}} - \ln P_{\text{enc}}(z|y)$$

M: gen^{*} = argmin $E_{(y,z)\sim P_{\text{enc}}(y,z)} \left[-\ln P_{\text{gen}}(y,z) \right]$

*

The classical EM algorithm is the case where we alternate optimizing the encoder (the E step) and the generator (the M step) and where the E step yields $P_{\text{enc}}(z|y) = P_{\text{gen}}(z|y)$ exactly and where the M step cannot fully fit the population.

Here we can use SGD on these two objectives independent of details of the models.

Derivation of the Encoder Optimum

$$E_{(y,z)\sim P_{enc}} \left[-\ln \frac{P_{gen}(y,z)}{P_{enc}(z|y)} \right]$$

$$= E_{(y,z)\sim P_{enc}} \left[-\ln \frac{P_{gen}(y,z)}{P_{gen}(z|y)} \right] + E_{y\sim Pop} KL(P_{enc}(z|y), P_{gen}(z|y))$$

$$\geq E_{(y,z)\sim P_{enc}} \left[-\ln \frac{P_{gen}(y,z)}{P_{gen}(z|y)} \right]$$

$$= E_{(y,z)\sim P_{enc}} \left[-\ln P_{gen}(y) \right]$$

$$= E_{y\sim Pop,z\sim P_{gen}(z|y)} \left[-\ln P_{gen}(y) \right]$$

$$= E_{y \sim \text{Pop}, z \sim P_{\text{gen}}(z|y)} \left[-\ln \frac{P_{\text{gen}}(y, z)}{P_{\text{gen}}(z|y)} \right]$$

END