

The variational autoencoder

Clayton W. Seitz

April 17, 2022

Outline

References

The logic of generative modeling

Say we have a set of variables $\mathbf{x} = (x_1, x_2, \dots, x_n)$ which might have some statistical dependence

The variable \mathbf{x} might be an amino acid sequence, gene expression data, microscopy image, etc.

- ▶ Often we are handed a batch of empirical samples $\{\mathbf{x}_i\}_{i=1}^N$
- ▶ We want to know the generating distribution $p(\mathbf{x})$

In supervised **generative learning**, we try to explicitly learn the joint distribution $p(\mathbf{x}) = \prod_{i=1}^{N-1} p(x_i | x_{i+1:N}) p(x_N)$, which is generally more difficult than discriminative learning.

Perks of generative modeling

- ▶ Fitting complete multivariate distributions $p(\mathbf{x})$ goes beyond correlation-based or clustering approaches
- ▶ Correlations cannot discover partial correlation in the context of other neighbors
- ▶ Fitting $p(\mathbf{x})$ permits sampling based inference

Why generative modeling is difficult

When describing a distribution over multiple variables, we may not know the proper normalization Z . That is,

$$p(\mathbf{x}) = \frac{1}{Z} \tilde{p}(\mathbf{x})$$

This **very important** situation arises in several contexts:

1. In **Bayesian inference** where $p(x_1|x_2) = p(x_2|x_1)p(x_1)/p(x_2)$ is intractable due to $Z = p(x_2) = \int p(x_2|x_1)p(x_1)dx_1$. This integral can be very difficult or impossible to compute.
2. In models from statistical physics, e.g. the Ising model, we only know $\tilde{p}(\mathbf{x}) = e^{-H(\mathbf{x})}$ where $H(\mathbf{x})$ is the Hamiltonian

Bayesian inference

The variable \mathbf{x} has a latent representation or code \mathbf{z} . We often say that \mathbf{z} is the *causal source* of \mathbf{x} . Ultimately, we would like to know the distribution $P_\phi(\mathbf{x})$

$$P_\phi(\mathbf{x}) = \frac{P_\phi(\mathbf{x}|\mathbf{z})P_\phi(\mathbf{z})}{Q_\psi(\mathbf{z}|\mathbf{x})}$$

in order to find the model parameters that maximize the likelihood of the observed data:

$$\Phi^* = \underset{\Phi}{\operatorname{argmin}} -\log P_\Phi(\mathbf{x})$$

but we generally do not know $P_\psi(\mathbf{z}|\mathbf{x})$ due to the intractable integral $Z = \int P_\phi(\mathbf{x}|\mathbf{z})P_\phi(\mathbf{z})d\mathbf{z}$ (see slide 5)

Computing the evidence

We can rewrite the evidence as

$$\begin{aligned}P_{\phi}(\mathbf{x}) &= \int P_{\phi}(\mathbf{z})P_{\phi}(\mathbf{x}|\mathbf{z})d\mathbf{z} \\&= \int P_{\phi}(\mathbf{z})P_{\phi}(\mathbf{x}|\mathbf{z})\frac{P_{\phi}(\mathbf{z}|\mathbf{x})}{P_{\phi}(\mathbf{z}|\mathbf{x})}d\mathbf{z} \\&= \mathbb{E}_{\mathbf{z}\sim P_{\phi}(\mathbf{z}|\mathbf{x})}\frac{P_{\phi}(\mathbf{z})P_{\phi}(\mathbf{x}|\mathbf{z})}{P_{\phi}(\mathbf{z}|\mathbf{x})}\end{aligned}$$

where $P_{\phi}(\mathbf{z}|\mathbf{x})$ is our model "encoder"

The evidence lower bound (ELBO)

$$\begin{aligned}\log P_{\phi}(\mathbf{x}) &= \log \int_z P(x, z) dx \\&= \log \int_z P(x, z) \frac{Q(z|x)}{Q(z|x)} dz \\&= \log \mathbb{E}_{\mathbf{z} \sim P_{\phi}(\mathbf{z}|\mathbf{x})} \frac{P(x|z)P(z)}{Q(z|x)} \\&\geq \mathbb{E}_{\mathbf{z} \sim P_{\phi}(\mathbf{z}|\mathbf{x})} \log \frac{Q(x|z)}{P(z)} + \log P(x|z) \\-\log P_{\phi}(\mathbf{x}) &\leq \mathbb{E}_{\mathbf{z} \sim P_{\phi}(\mathbf{z}|\mathbf{x})} \log \frac{Q(x|z)}{P(z)} - \log P(x|z)\end{aligned}$$

The ELBO objective

$$\begin{aligned}\Phi^* &= \mathcal{L}(\Phi) \\ &= \operatorname{argmin}_{\Phi} \mathbb{E}_{\mathbf{x} \sim P_{\text{OP}}, \mathbf{z} \sim P_{\Phi}(\mathbf{z}|\mathbf{x})} \log \frac{Q_{\Psi}(\mathbf{z}|\mathbf{x})}{P(\mathbf{z})} - \log P(\mathbf{x}|\mathbf{z})\end{aligned}$$

The ELBO can be rewritten in terms of a KL-divergence and population entropy. We often think of Φ as “position” and the loss \mathcal{L} as an “energy”

References I