

# Useful Facts About Latent-Variable Generative Models

Karl Stratos

In this note,  $p_X$  always refers to the marginal density of a latent-variable generative model  $p_Z$  and  $p_{X|Z}$  where  $Z \in \mathcal{Z}$ .<sup>1</sup> So “learning  $p_X$ ” means learning  $p_Z$  and  $p_{X|Z}$ .

## 1 Three Types of Estimation

We are interested in estimating the population density  $\mathbf{pop}_X$  with the marginal  $p_X$  by minimizing some divergence  $D(\mathbf{pop}_X, p_X)$ . Three common types of the estimation problem are

$$\begin{aligned} \min_{p_X} D(\mathbf{pop}_X, p_X) & \quad \text{(direct marginalization)} \\ \min_{p_X} \min_q U(\mathbf{pop}_X, p_X, q) & \quad \text{(variational optimization)} \\ \min_{p_X} \max_q L(\mathbf{pop}_X, p_X, q) & \quad \text{(variational adversarial optimization)} \end{aligned}$$

**Direct marginalization.** Sometimes it is possible to directly calculate the marginal. It is trivial if  $\mathcal{Z}$  is small. If  $\mathcal{Z}$  is a set of discrete structures such as sequences or trees, it is possible with conditional independence assumptions [3, 12]. In this case it is natural to minimize the KL divergence. Using  $D_{\text{KL}}(\mathbf{pop}_X || p_X) = H(\mathbf{pop}_X, p_X) - H(\mathbf{pop}_X)$

$$\min_{p_X} D_{\text{KL}}(\mathbf{pop}_X, p_X) \equiv \min_{p_X} H(\mathbf{pop}_X, p_X) = \max_{p_X} \mathbf{E}_{x \sim \mathbf{pop}_X} [\log p_X(x)]$$

**Variational optimization.** Generally direct marginalization is intractable. In this case we can consider estimating an equivalent objective by minimizing an upper bound  $U(\mathbf{pop}_X, p_X, q)$  with a variational model  $q$  such that (1) it is easy to compute, and (2) it is tight for an optimal  $q$  (henceforth tightable). For instance, EM minimizes the KL divergence by minimizing the minimization of a tightable upper bound  $-\text{ELBO}(\mathbf{pop}_X, p_X, q_{Z|X})$  on  $H(\mathbf{pop}_X, p_X)$  where  $q_{Z|X}$  estimates the intractable posterior  $p_{Z|X}$  [6, 13]:

$$\min_{p_X} D_{\text{KL}}(\mathbf{pop}_X, p_X) \equiv \min_{p_X} \left( \min_{q_{Z|X}} -\text{ELBO}(\mathbf{pop}_X, p_X, q_{Z|X}) \right)$$

The consistency of  $p_X$  follows from the fact that the bound is tightable (assuming universality). The fact that the objective remains a minimization is convenient in practice.

**Variational adversarial optimization.** There are cases we maximize a lower bound  $L(\mathbf{pop}_X, p_X, q)$  where  $q$  enjoys similar properties for estimating other divergence measures. For instance, GAN minimizes the Jensen-Shannon divergence by minimizing the maximization of a tightable upper bound  $\text{NCE}(\mathbf{pop}_X, p_X, q_{A|X})$  on  $2\text{JSD}(\mathbf{pop}_X || p_X) - \log 4$  where  $q_{A|X}$  is tasked with discriminating between  $\mathbf{pop}_X$  and  $p_X$  [8]:

$$\min_{p_X} \text{JSD}(\mathbf{pop}_X || p_X) \equiv \min_{p_X} \left( \max_{q_{A|X}} \text{NCE}(\mathbf{pop}_X, p_X, q_{A|X}) \right)$$

The divergence measure has been generalized to  $f$ -divergences [16] and the Wasserstein metric [1]. The consistency of  $p_X$  again follows from the fact that the bound is tightable (assuming universality). However, the objective is adversarial and more difficult to optimize in practice.

---

<sup>1</sup> We write “density” or “distribution” interchangeably to denote a probability function and  $\sum$  to denote marginalization whether the considered variable is discrete or continuous.

## 2 Forms of ELBO

The evidence lower bound (ELBO) emerges through an effort to replace the posterior  $p_{Z|X}$  with a variational model  $q_{Z|X}$  in the expected log-likelihood:<sup>2</sup>

$$\mathbf{E}_{x \sim \mathbf{pop}_X} [\log p_X(x)] = \mathbf{E}_{\substack{x \sim \mathbf{pop}_X \\ z \sim q_{Z|X}(\cdot|x)}} \left[ \log \frac{p_{XZ}(x, z)}{p_{Z|X}(z|x)} \frac{q_{Z|X}(z|x)}{q_{Z|X}(z|x)} \right] = \underbrace{\mathbf{E}_{\substack{x \sim \mathbf{pop}_X \\ z \sim q_{Z|X}(\cdot|x)}} \left[ \log \frac{p_{XZ}(x, z)}{q_{Z|X}(z|x)} \right]}_{\text{ELBO}(\mathbf{pop}_X, p_X, q_{Z|X})} + \underbrace{D_{\text{KL}}(q_{Z|X} \| p_{Z|X})}_{\geq 0}$$

One can write ELBO in various forms by manipulating terms:

$$\text{ELBO}(\mathbf{pop}_X, p_X, q_{Z|X}) = \mathbf{E}_{x \sim \mathbf{pop}_X} [\log p_X(x)] - D_{\text{KL}}(q_{Z|X} \| p_{Z|X}) \quad (\text{EM})$$

$$= \mathbf{E}_{\substack{x \sim \mathbf{pop}_X \\ z \sim q_{Z|X}(\cdot|x)}} [\log p_{XZ}(x, z)] + H(q_{Z|X})$$

$$= \mathbf{E}_{\substack{x \sim \mathbf{pop}_X \\ z \sim q_{Z|X}(\cdot|x)}} [\log p_{X|Z}(x|z)] - D_{\text{KL}}(q_{Z|X} \| p_Z) \quad (\text{VAE})$$

$$= \mathbf{E}_{\substack{x \sim \mathbf{pop}_X \\ z \sim q_{Z|X}(\cdot|x)}} [\log p_{X|Z}(x|z)] - D_{\text{KL}}(q_{Z|X} \| q_Z) - D_{\text{KL}}(q_Z \| p_Z) \quad (\text{DVAE})$$

$$= \mathbf{E}_{\substack{x \sim \mathbf{pop}_X \\ z \sim q_{Z|X}(\cdot|x)}} [\log p_{X|Z}(x|z)] + H(q_{Z|X}) - H(q_Z, p_Z) \quad (\text{MAXENT})$$

The first two forms yield the traditional alternating optimization steps in the EM algorithm. The VAE form is the standard VAE objective in which the reconstruction term is estimated by sampling and the KL term is estimated in closed form. The DVAE form is the VAE form further decomposed; one can easily check that  $D_{\text{KL}}(q_{Z|X} \| p_Z) = D_{\text{KL}}(q_{Z|X} \| q_Z) + D_{\text{KL}}(q_Z \| p_Z)$  where  $q_Z(z) = \mathbf{E}_{x \sim \mathbf{pop}_X} [q_{Z|X}(z|x)]$ . MAXENT is easily derived from DVAE. DVAE can be used for the following additional interpretations of VAE.

**Rate-distortion autoencoders.** First note that  $D_{\text{KL}}(q_{Z|X} \| q_Z) = I(X, Z; \mathbf{pop}_{XZ}^q)$  is the mutual information between  $X$  and  $Z$  under the joint density  $\mathbf{pop}_{XZ}^q(x, z) = \mathbf{pop}_X(x)q_{Z|X}(z|x)$ .<sup>3</sup> We can view DVAE as a nested maximization over  $p_Z$ ,  $p_{X|Z}$ , and  $q_{Z|X}$ . The optimal prior is always  $p_Z = q_Z$  which eliminates the last KL term. The resulting optimization problem is equivalent to

$$\min_{p_{X|Z}, q_{Z|X}} I(X, Z; \mathbf{pop}_{XZ}^q) + H(\mathbf{pop}_{X|Z}^q, p_{X|Z})$$

This shows the standard rate-distortion tradeoff where we want to limit the channel capacity of the encoder  $q_{Z|X}$  for light-weight communication while limiting distortion  $H(\mathbf{pop}_{X|Z}^q, p_{X|Z}) \leq H_{\text{max}}$ .

**Disentanglement.** Assume  $Z = (Z_1 \dots Z_m)$  and the model prior  $p_Z = \prod_i p_{Z_i}$  is component-wise independent (usually the case). Let  $q_{Z_i}$  denote the marginal of  $q_Z$  for the  $i$ -th variable; note that  $q_Z$  may still be a complicated joint density. It is easy to verify that

$$D_{\text{KL}}(q_Z \| p_Z) = D_{\text{KL}}\left(q_Z \left\| \prod_{i=1}^m q_{Z_i}\right.\right) + \sum_{i=1}^m D_{\text{KL}}(q_{Z_i} \| p_{Z_i})$$

where the first term is also known as total correlation (a multivariate generalization of mutual information). Thus minimizing the KL term in VAE involves minimizing dependencies between latent components under  $q_Z$ . Implicit disentanglement observed in VAE (e.g.,  $Z_{\text{gender}}$  vs  $Z_{\text{mustache}}$ ) is attributed to this term. Several weighting schemes have been proposed to control the level of disentanglement [9, 5].

<sup>2</sup>One related lower bound on the log-likelihood is

$$\mathbf{E}_{\substack{x \sim \mathbf{pop}_X \\ z \sim p_Z}} [\log p_{X|Z}(x|z)] \leq \mathbf{E}_{x \sim \mathbf{pop}_X} \left[ \log \mathbf{E}_{z \sim p_Z} [p_{X|Z}(x|z)] \right] = \mathbf{E}_{x \sim \mathbf{pop}_X} [\log p_X(x)]$$

which follows from Jensen's inequality. This objective is useful if we want to sample  $z$  from the model's own prior during learning; it can be optimized with REINFORCE in that case (Section 4) [21].

<sup>3</sup>An aside: the posterior collapse in VAE  $D_{\text{KL}}(q_{Z|X} \| p_Z) = 0$  implies  $I(X, Z; \mathbf{pop}_{XZ}^q) = 0$  since

$$I(X, Z; \mathbf{pop}_{XZ}^q) = D_{\text{KL}}(q_{Z|X} \| q_Z) = D_{\text{KL}}(q_{Z|X} \| p_Z) - D_{\text{KL}}(q_Z \| p_Z) = -D_{\text{KL}}(q_Z \| p_Z)$$

and both mutual information and KL divergence are nonnegative [7].

## 2.1 Direct Estimation of the Log-Likelihood

ELBO is just a lower bound on the log-likelihood (LL) which is what we really care about. The gap between ELBO and LL is the KL divergence between  $q_{Z|X}$  and  $p_{Z|X}$ , which means when they are not equal a Monte Carlo estimate of ELBO (assuming fixed data  $x$  for simplicity)  $(1/K) \log(p_{XZ}(x, z^{(k)})/q_{Z|X}(z^{(k)}|x))$  where  $z^{(1)} \dots z^{(K)} \sim q_{Z|X}(\cdot|x)$  may still be far smaller than  $\log p_X(x)$  even if  $K \rightarrow \infty$ . However, we can simply compute a Monte Carlo estimate of  $p_X(x)$  directly, in particular using importance sampling with  $q_{Z|X}$  as the proposal distribution. That is,

$$\log p_X(x) = \log \mathbf{E}_{z \sim q_{Z|X}(\cdot|x)} \left[ \frac{p_{XZ}(x, z)}{q_{Z|X}(z|x)} \right] \approx \log \left( \frac{1}{K} \sum_{k=1}^K \frac{p_{XZ}(x, z^{(k)})}{q_{Z|X}(z^{(k)}|x)} \right)$$

Clearly as  $K \rightarrow \infty$  the estimate converges to LL assuming bounded  $\mathbf{E}_{z \sim q_{Z|X}(\cdot|x)} [p_{XZ}(x, z)/q_{Z|X}(z|x)]$ . The corresponding population-level objective defined for each value of  $K$  is actually a lower bound on LL by Jensen’s inequality:

$$\mathcal{L}_K = \mathbf{E}_{z^{(1)} \dots z^{(K)} \sim q_{Z|X}(\cdot|x)} \left[ \log \left( \frac{1}{K} \sum_{k=1}^K \frac{p_{XZ}(x, z^{(k)})}{q_{Z|X}(z^{(k)}|x)} \right) \right] \leq \log p_X(x)$$

Note that  $\mathcal{L}_1$  coincides with ELBO: thus  $\mathcal{L}_K$  can be viewed as a multi-sample lower bound that, unlike ELBO, converges to LL as  $K \rightarrow \infty$ .  $\mathcal{L}_K$  can be used as an alternative training objective for learning  $p_{XZ}$  and  $q_{Z|X}$  (importance weighted autoencoders (IWAs) [4]) or as a way to estimate true LL value (instead of ELBO) after learning a VAE for evaluation purposes.

## 3 Prior

The choice of prior can be important from an optimization perspective as well as a modeling perspective. Recall that the VAE objective is, written explicitly as a function of  $p_Z$  and  $p_{X|Z}$ ,

$$\text{ELBO}(\mathbf{pop}_X, p_Z, p_{X|Z}, q_{Z|X}) = \mathbf{E}_{\substack{x \sim \mathbf{pop}_X \\ z \sim q_{Z|X}(\cdot|x)}} [\log p_{X|Z}(x|z)] - D_{\text{KL}}(q_{Z|X} \| p_Z)$$

$p_Z$  affects the decoder  $q_{Z|X}$  through the regularization term, which in turn affects the reconstruction term. A richer prior such as multimodal instead of unimodal can help achieve better objective value.

**Mixture prior.** One common approach to enriching the prior is to make it a mixture distribution by introducing an additional latent variable  $C \in \{1 \dots K\}$  and define  $p_Z(z) = \sum_{c=1}^K p_{Z|C}(z|c)p_C(c)$ . We assume  $X \perp\!\!\!\perp C|Z$  so that the model defines the joint density  $p_{XZC}(x, z, c) = p_{X|Z}(x|z)p_{Z|C}(z|c)p_C(c)$ . If we further define the variational posterior  $q_{ZC|X}(z, c|x) = q_{Z|X}(z|x)q_{C|X}(c|x)$  with the assumption  $Z \perp\!\!\!\perp C|X$ , the first term in ELBO does not change and only the regularization term changes to

$$\begin{aligned} D_{\text{KL}}(q_{ZC|X} \| p_{ZC}) &= D_{\text{KL}}(q_{C|X} \| p_C) + D_{\text{KL}}(q_{Z|X} \| p_{Z|C}) \\ &= \mathbf{E}_{\substack{x \sim \mathbf{pop}_X \\ c \sim q_{C|X}(\cdot|x)}} \left[ \log \frac{q_{C|X}(c|x)}{p_C(c)} + D_{\text{KL}}(q_{Z|X}(\cdot|x) \| p_{Z|C}(\cdot|c)) \right] \end{aligned}$$

which is easy to calculate assuming small  $K$  and a closed-form solution for the KL term over  $Z$  as usual. It can be viewed as more fine-grained regularization in which we make  $q_{Z|X} \approx p_{Z|C}$  where  $C$  is with respect to  $q_{C|X} \approx p_C$ .

**Mixture prior from the variational posterior.** Choosing  $p_Z$  to be a mixture distribution can be justified in terms of the optimal prior. Recall from the DVAE form of ELBO that the optimal prior is given by  $\mathbf{E}_{x \sim \mathbf{pop}_X} [q_{Z|X}(z|x)]$  for any fixed  $q_{Z|X}$ . Thus the optimal prior for an empirical estimate of ELBO based on iid samples  $x_1 \dots x_N \sim \mathbf{pop}_X$  is actually the mixture distribution  $(1/N) \sum_{i=1}^N q_{Z|X}(z|x_i)$ . The mixture prior  $p_Z(z) = \sum_{c=1}^K p_{Z|C}(z|c)p_C(c)$  can be seen as approximating this optimal prior which becomes exact when  $p_C$  is uniform over  $K = N$  components and  $p_{Z|C}(z|c) = q_{Z|X}(z|x_c)$ . We can consider a more direct approximation by explicitly using  $q_{Z|X}$  to define  $p_Z$ , for instance  $p_Z(z) = (1/K) \sum_{c=1}^K q_{Z|X}(z|\tilde{x}_c)$  where  $\tilde{x}_1 \dots \tilde{x}_K$  are either random samples or learnable parameters that represent “pseudo-inputs”. This parameter sharing between the prior and the decoder is shown to be potentially useful [17].

**Hierarchical prior.** There may be cases in which there is a natural structure to the latent variable  $Z$ . For instance, if  $X$  represents a sentence, we may think of a topic  $Z_1$ , think of facts about the topic  $Z_2$ , and then generate  $X$ . In this case we can model the joint density as  $p_{XZ_1Z_2}(x, z_1, z_2) = p_{X|Z_1Z_2}(x|z_1, z_2)p_{Z_2|Z_1}(z_2|z_1)p_{Z_1}(z_1)$  and define the variational posterior as  $q_{Z_1Z_2|X}(z_1, z_2|x) = q_{Z_2|XZ_1}(z_2|x, z_1)q_{Z_1|X}(z_1|x)$ . This is almost the same as having a mixture prior except that we do not make conditional independence assumptions. ELBO can still be optimized with a suitable parameterization, for instance  $Z_1$  and  $Z_2$  are isotropic Gaussians.

**Compartmentalized prior.** There may also be cases in which we want to compartmentalize  $Z = (Z_1, Z_2)$  with  $Z_1 \perp\!\!\!\perp Z_2$ . For instance, Kingma et al. (2014) define  $Z_1 \in \mathbb{R}^d$  as the style and  $Z_2 \in \{1 \dots L\}$  as the label of an MNIST digit image  $X$  and consider the model  $p_{XZ_1Z_2}(x, z_1, z_2) = p_{X|Z_1Z_2}(x|z_1, z_2)p_{Z_2}(z_2)p_{Z_1}(z_1)$  where the decoder  $p_{X|Z_1Z_2}$  is a continuous-discrete hybrid [14]. The variational posterior  $q_{Z_1Z_2|X}(z_1, z_2|x) = q_{Z_2|X}(z_2|x)q_{Z_1|XZ_2}(z_1|x, z_2)$  provides a label classifier  $q_{Z_2|X}$ . The model can be trained in a semi-supervised manner by jointly optimizing the ELBO of  $\log p_X(x)$  on unlabeled images and  $\log p_{XZ_2}(x, z_2)$  on labeled images.

**Conditional prior.** Assume a joint density  $\mathbf{pop}_{XY}$  over  $X$  and  $Y$  as the population density. Consider the model  $p_{XYZ}(x, y, z) = p_{X|Z}(x|z)p_{Z|Y}(z|y)\mathbf{pop}_Y(y)$  representing the generative story  $Y \rightarrow Z \rightarrow X$  in which  $X$  is conditionally independent of  $Y$  given the bottleneck variable  $Z$ . We minimize  $D_{\text{KL}}(\mathbf{pop}_{XY} || p_{XY}(x, y))$  over the model parameters. Equivalently we maximize the conditional log likelihood  $\mathbf{E}_{(x,y) \sim \mathbf{pop}_{XY}}[\log p_{X|Y}(x|y)]$ . We do this by introducing a variational posterior  $q_{Z|XY}$  and maximizing the ELBO lower bound

$$\max_{p_{Z|Y}, p_{X|Z}, q_{Z|XY}} \mathbf{E}_{\substack{(x,y) \sim \mathbf{pop}_{XY} \\ z \sim q_{Z|XY}(\cdot|x,y)}} [\log p_{X|Z}(x|z)] - D_{\text{KL}}(q_{Z|XY} || p_{Z|Y})$$

## 4 Types of Non-Differentiability

**Deterministic operation.** Consider the step function  $\text{STEP} : \mathbb{R} \rightarrow \{0, 1\}$  which outputs 1 iff the input value is nonnegative. Since it is non-differentiable at 0 and has derivative 0 almost everywhere, it is meaningless to talk about a gradient. One heuristic to obtain a meaningful gradient signal is to linearize  $\text{STEP}(a) \approx a$  (which preserves the sign) in the backward pass so that

$$\frac{\partial J(\text{STEP}(f(\theta)))}{\partial \theta} = \frac{\partial J(\text{STEP}(f(\theta)))}{\partial \text{STEP}(f(\theta))} \frac{\partial \text{STEP}(f(\theta))}{\partial f(\theta)} \frac{\partial f(\theta)}{\partial \theta} \approx \frac{\partial J(\text{STEP}(f(\theta)))}{\partial \text{STEP}(f(\theta))} \frac{\partial f(\theta)}{\partial \theta}$$

We can also consider a “multidimensional step function”. Define  $\text{SNAP} : \mathbb{R}^K \rightarrow \{e_1 \dots e_K\}$  by  $\text{SNAP}(u) = e_{k^*}$  where  $k^* = \arg \max_{k=1}^K u_k$  and  $e_1 \dots e_K \in \{0, 1\}^K$  are standard basis elements. It is non-differentiable along  $(K - 1)$ -dimensional manifolds and has gradient  $0^K$  almost everywhere, but we can linearize  $\text{SNAP}(u) \approx u$  (which preserves the argmax) in the backward pass so that

$$\frac{\partial J(\text{SNAP}(f(\theta)))}{\partial \theta} \approx \frac{\partial J(\text{SNAP}(f(\theta)))}{\partial \text{SNAP}(f(\theta))} \frac{\partial f(\theta)}{\partial \theta} \tag{1}$$

**Stochastic operation.** Consider a stochastic function which outputs a certain value with a certain probability. It is only meaningful to talk about the differentiability of such a function with respect to its expectation. Let  $p_Z^\theta$  denote a differentiable function of  $\theta$  that defines a density over some variable  $Z$ . Given an objective function  $J(z)$ <sup>4</sup>, we can consider unbiased gradient estimators such as

$$\begin{aligned} \frac{\partial}{\partial \theta} \mathbf{E}_{z \sim p_Z^\theta} [J(z)] &= \sum_{z \in \mathcal{Z}} J(z) \frac{\partial}{\partial \theta} p_Z^\theta(z) && \text{(direct marginalization)} \\ &= \mathbf{E}_{z \sim p_Z^\theta} \left[ J(z) \frac{\partial}{\partial \theta} \log p_Z^\theta(z) \right] && \text{(score function estimator)} \\ &= \mathbf{E}_{\epsilon \sim p_\epsilon} \left[ \frac{\partial}{\partial \theta} J(\pi^\theta(\epsilon)) \right] && \text{(reparameterization trick)} \end{aligned}$$

Direct marginalization is an option if computationally possible (e.g.,  $\mathcal{Z}$  is a small discrete set) [14]. The other two gradient estimates are based on sampling. The score function estimator (aka. REINFORCE [19]) is high-variance

<sup>4</sup>In general  $J(\theta, z)$  can depend on  $\theta$  through other connections.

and normally requires additional techniques to reduce variance (control variates). The reparameterization trick is an option if  $z \sim p_Z^\theta$  is distributed as  $z = \pi^\theta(\epsilon)$  where  $\pi^\theta(\epsilon)$  is a differentiable function of  $\theta$  and  $\epsilon \sim p_\epsilon$  is some random variable that does not depend on  $\theta$ . An isotropic Gaussian density is a classical example:  $z \sim \mathcal{N}(\mu_\theta, \text{diag}(\sigma_\theta^2) I_d)$  is distributed as  $z = \mu_\theta + \sigma_\theta \odot \epsilon$  where  $\epsilon \sim \mathcal{N}(0, I_d)$  [13].

## 4.1 Backpropagation Through Discrete Sampling

When  $Z$  is discrete and neither direct marginalization nor the score function estimator is a good option (due to computational costs or high variance), we can consider biased gradient estimators.

**Bernoulli variable.** Let  $Z \in \{0, 1\}$  with  $p_Z^\theta(1) = f(\theta)$ . Combining the reparameterization trick  $z = \text{STEP}(f(\theta) - \epsilon)$  where  $\epsilon \sim \mathcal{U}(0, 1)$  and the linear approximation  $\text{STEP}(a) \approx a$  in the backward pass, we derive the straight-through gradient estimator [10, 2]:

$$\frac{\partial}{\partial \theta} \mathbf{E}_{z \sim p_Z^\theta} [J(z)] = \mathbf{E}_{\epsilon \sim p_\epsilon} \left[ \frac{\partial}{\partial \theta} J(\text{STEP}(f(\theta) - \epsilon)) \right] \approx \mathbf{E}_{\epsilon \sim p_\epsilon} \left[ \frac{\partial J(\text{STEP}(f(\theta) - \epsilon))}{\partial \text{STEP}(f(\theta) - \epsilon)} \frac{\partial f(\theta)}{\partial \theta} \right] = \mathbf{E}_{z \sim p_Z^\theta} \left[ \frac{\partial J(z)}{\partial z} \frac{\partial f(\theta)}{\partial \theta} \right]$$

Of course in this case direct marginalization is trivial, but this is applicable for  $Z \in \{0, 1\}^d$  where  $p_Z^\theta$  is a product distribution. In that case direct marginalization has complexity  $O(2^d)$  whereas sampling has complexity  $O(d)$ .

**Categorical variable: Gumbel-Softmax.** Let  $Z \in \{1 \dots K\}$  where  $K > 2$ . Without loss of generality, let  $Z \in \{e_1 \dots e_K\}$  be represented as a  $K$ -dimensional standard basis element with  $p_Z^\theta = f(\theta) \in \Delta^{K-1}$ . It can be shown that  $z \sim p_Z^\theta$  (which is a vertex in  $\Delta^{K-1}$ ) is distributed as  $z = \text{softmax}((\log f(\theta) + \epsilon)/\tau)$  where  $\epsilon \sim \text{Gumbel}^K(0, 1)$  and  $\tau > 0$  as  $\tau$  goes to zero (Appendix A), yielding the Gumbel-Softmax (GS) estimator [11, 15]

$$\frac{\partial}{\partial \theta} \mathbf{E}_{z \sim p_Z^\theta} [J(z)] = \lim_{\tau \rightarrow 0} \mathbf{E}_{\epsilon \sim \text{Gumbel}^K(0, 1)} \left[ \frac{\partial}{\partial \theta} J \left( \text{softmax} \left( \frac{\log f(\theta) + \epsilon}{\tau} \right) \right) \right]$$

In practice  $\tau$  is fixed (e.g., 0.9) or annealed, so GS is biased. GS involves a  $K$ -dimensional sample and does not seem to offer any computational advantage over direct marginalization over  $K$  values. But consider:

1. Suppose the objective  $J(z, z')$  involves nested sampling  $z \sim \text{Cat}(f(\theta))$  and  $z' \sim \mathcal{N}(\mu(z), \text{diag}(\sigma^2(z)) I_d)$ . Direct marginalization requires drawing  $K$  conditional samples of  $Z'$  and thus  $O(Kd)$  time, whereas GS gives

$$\lim_{\tau \rightarrow 0} \mathbf{E}_{\epsilon \sim \text{Gumbel}^K(0, 1)} \left[ \mathbf{E}_{\epsilon' \sim \mathcal{N}(0, I_d)} \left[ \frac{\partial}{\partial \theta} J(z_\theta^{\epsilon, \tau}, \mu(z_\theta^{\epsilon, \tau}) + \sigma(z_\theta^{\epsilon, \tau}) \odot \epsilon') \right] \right]$$

where  $z_\theta^{\epsilon, \tau} = \text{softmax}((\log f(\theta) + \epsilon)/\tau) \in \Delta^{K-1}$ . This estimator requires one conditional sample of  $Z'$  and thus  $O(d)$  time.

2. Suppose  $Z \in \{1 \dots K\}^d$  and  $p_Z^\theta$  is a product distribution over  $d$  dimensions. Then direct marginalization has complexity  $O(K^d)$  whereas sampling has complexity  $O(Kd)$ .

We can consider a straight-through version of GS for cases in which we need a discrete sample by snapping in the forward pass and softmaxing in the backward pass:

$$\frac{\partial}{\partial \theta} \mathbf{E}_{\epsilon \sim \text{Gumbel}^K(0, 1)} \left[ J \left( \text{SNAP} \left( \frac{\log f(\theta) + \epsilon}{\tau} \right) \right) \right] \approx \mathbf{E}_{\epsilon \sim \text{Gumbel}^K(0, 1)} \left[ \frac{\partial}{\partial \theta} J \left( \text{softmax} \left( \frac{\log f(\theta) + \epsilon}{\tau} \right) \right) \right]$$

**Categorical variable: vector quantization** Let  $f(\theta) \in \mathbb{R}^d$  and assume  $C \in \mathbb{R}^{K \times d}$ . We discretize  $f(\theta) \in \mathbb{R}^d$  into  $\{1 \dots K\}$  by treating the rows of  $C$  as centroids in  $k$ -means. Estimating the gradient in this case using the straight-through estimator for snapping (1) is called vector quantization (VQ):

$$\frac{\partial J(\text{SNAP}(u(\theta, C)))}{\partial \theta} \approx \frac{\partial J(\text{SNAP}(u(\theta, C)))}{\partial \text{SNAP}(u(\theta, C))} \frac{\partial f(\theta)}{\partial \theta}$$

where  $u_i(\theta, C) = -\|f(\theta) - C_i\|$ . When we want to update  $C$  as well we can add additional objectives such as minimizing  $\|f(\theta) - C_{i^*}\|^2$  where  $i^* = \arg \max_{i=1}^K u_i(\theta, C)$  [18].

## References

- [1] Arjovsky, M., Chintala, S., and Bottou, L. (2017). Wasserstein generative adversarial networks. In D. Precup and Y. W. Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 214–223, International Convention Centre, Sydney, Australia. PMLR.
- [2] Bengio, Y., Léonard, N., and Courville, A. (2013). Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*.
- [3] Berg-Kirkpatrick, T., Bouchard-Côté, A., DeNero, J., and Klein, D. (2010). Painless unsupervised learning with features. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 582–590. Association for Computational Linguistics.
- [4] Burda, Y., Grosse, R., and Salakhutdinov, R. (2015). Importance weighted autoencoders. *arXiv preprint arXiv:1509.00519*.
- [5] Chen, T. Q., Li, X., Grosse, R. B., and Duvenaud, D. K. (2018). Isolating sources of disentanglement in variational autoencoders. In *Advances in Neural Information Processing Systems*, pages 2610–2620.
- [6] Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, **39**(1), 1–22.
- [7] Dieng, A. B., Kim, Y., Rush, A. M., and Blei, D. M. (2019). Avoiding latent variable collapse with generative skip models. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2397–2405. PMLR.
- [8] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680.
- [9] Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., and Lerchner, A. (2017). beta-vae: Learning basic visual concepts with a constrained variational framework. *Iclr*, **2**(5), 6.
- [10] Hinton, G. (2012). Neural networks for machine learning. coursera,[video lectures].
- [11] Jang, E., Gu, S., and Poole, B. (2016). Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*.
- [12] Kim, Y., Dyer, C., and Rush, A. M. (2019). Compound probabilistic context-free grammars for grammar induction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2369–2385.
- [13] Kingma, D. P. and Welling, M. (2014). Auto-encoding variational bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*.
- [14] Kingma, D. P., Mohamed, S., Rezende, D. J., and Welling, M. (2014). Semi-supervised learning with deep generative models. In *Advances in neural information processing systems*, pages 3581–3589.
- [15] Maddison, C. J., Mnih, A., and Teh, Y. W. (2016). The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*.
- [16] Nowozin, S., Cseke, B., and Tomioka, R. (2016). f-gan: Training generative neural samplers using variational divergence minimization. In *Advances in neural information processing systems*, pages 271–279.
- [17] Tomczak, J. and Welling, M. (2018). Vae with a vampprior. In *International Conference on Artificial Intelligence and Statistics*, pages 1214–1223.
- [18] van den Oord, A., Vinyals, O., *et al.* (2017). Neural discrete representation learning. In *Advances in Neural Information Processing Systems*, pages 6306–6315.
- [19] Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, **8**(3-4), 229–256.
- [20] Wolpert, R. L. (2014). Extremes.

[21] Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R., and Bengio, Y. (2015). Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057.

## A The Gumbel Distribution

The Gumbel distribution with location  $\mu \in \mathbb{R}$  and scale  $\beta > 0$  has the CDF over  $\mathbb{R}$ :

$$\overline{\text{Gumbel}}(\mu, \beta)(x) = \int_{-\infty}^x \text{Gumbel}(\mu, \beta)(t) dt = \exp\left(-\exp\left(-\frac{x-\mu}{\beta}\right)\right) \stackrel{\mu=0}{\stackrel{\beta=1}{\equiv}} \exp(-\exp(-x))$$

where  $\mu = 0$  and  $\beta = 1$  yields the standard form. Differentiating the CDF gives the PDF:<sup>5</sup>

$$\text{Gumbel}(\mu, \beta)(x) = \frac{1}{\beta} \exp\left(-\frac{x-\mu}{\beta} - \exp\left(-\frac{x-\mu}{\beta}\right)\right) \stackrel{\mu=0}{\stackrel{\beta=1}{\equiv}} \exp(-x - \exp(-x))$$

which is admittedly complicated with nested exponentials.<sup>6</sup> Gumbel is closed under linear transformation: if  $Z \sim \text{Gumbel}(0, 1)$ , then  $\mu + \beta Z \sim \text{Gumbel}(\mu, \beta)$ .<sup>7</sup> If  $z \sim \text{Gumbel}(\mu, \beta)$  and  $z' \sim \text{Gumbel}(\mu', \beta)$ , then  $z - z' \sim \text{Logistic}(\mu - \mu', \beta)$ . Gumbel is the limiting distribution over the maximum of (properly normalized)  $N$  iid variables as  $N \rightarrow \infty$  [20]. Let  $X_1 \dots X_N \geq 0$  denote iid variables distributed as  $\text{Exp}_\lambda$ .<sup>8</sup> Define  $Z_N^{\max} := \lambda \max_{i=1}^N X_i - \log N$ . The CDF of  $Z_N^{\max}$  is

$$\begin{aligned} \Pr(Z_N^{\max} \leq x) &= \Pr(\lambda X_i - \log N \leq x)^N && \text{(iid assumption on } X_1 \dots X_N) \\ &= \Pr\left(X_i \leq \frac{1}{\lambda}(x + \log N)\right)^N \\ &= (1 - \exp(-x - \log N))^N && \text{(CDF of } \text{Exp}_\lambda) \\ &= \left(1 - \frac{1}{N} \exp(-x)\right)^N \\ &\stackrel{N \rightarrow \infty}{\equiv} \exp(-\exp(-x)) = \overline{\text{Gumbel}}(0, 1)(x) && \text{(Lemma A.2)} \end{aligned}$$

Hence  $Z_N^{\max} \sim \text{Gumbel}(0, 1)$  with median  $\bar{z} = -\log \log 2$  (i.e.,  $\Pr(Z_N^{\max} \leq \bar{z}) = 1/2$ ) as  $N \rightarrow \infty$ , which implies that there is a 50/50 chance that  $\max_{i=1}^N X_i \leq \frac{1}{\lambda}(\log N - \log \log 2)$ . More generally, the maximum under a distribution with an exponentially thin tail (Gaussian, Gamma, etc., with appropriate normalization) behaves like Gumbel with the median growing like  $O((\log N)^p)$  for some  $p$ . Think of it as a “max version” of the central limit theorem which states that  $Z_N^{\text{mean}} := \sqrt{N/\sigma^2}(\bar{X}_N - \mu) \sim \mathcal{N}(0, 1)$  as  $N \rightarrow \infty$  where the variance of  $\bar{X}_N = \sum_{i=1}^N X_i/N$  shrinks like  $O(1/N)$ .

### A.1 The Gumbel-Max Trick

**Theorem A.1.** Let  $u \in \mathbb{R}^K$  where  $K \geq 1$ . Pick any  $\beta > 0$  and define the categorical random variable  $X \in \{1 \dots K\}$  with distribution  $\text{Cat}(\text{softmax}(u/\beta))$ , that is

$$\Pr(X = k) = \frac{\exp(u_k/\beta)}{\sum_{l=1}^K \exp(u_l/\beta)} \quad \forall k = 1 \dots K$$

Now define

$$\epsilon_1 \dots \epsilon_K \stackrel{\text{iid}}{\sim} \text{Gumbel}(0, \beta) \quad Y = \arg \max_{k=1}^K u_k + \epsilon_k$$

(the argmax is unique with probability 1 since  $\epsilon_1 \dots \epsilon_K$  are drawn iid from a continuous distribution). Then  $\Pr(X = k) = \Pr(Y = k)$  for all  $k = 1 \dots K$ .

<sup>5</sup>If  $F(x) = \int_c^x f(t) dt$  is a CDF of a PDF  $f$  with support on  $x \geq c$ , and  $G$  is any antiderivative of  $f$  (i.e.,  $G'(x) = f(x)$ ), the fundamental theorem of calculus says  $F(x) = G(x) - G(c)$ . Therefore,  $F'(x) = \frac{\partial}{\partial x}(G(x) + G(c)) = G'(x) = f(x)$ .

<sup>6</sup>It can be verified that the mean is  $\mu + \gamma\beta$  where  $\gamma = 0.5772 \dots$  is the Euler-Mascheroni constant, the variance is  $\frac{\pi^2}{6}\beta^2$ .

<sup>7</sup>The CDF of  $\mu + \beta Z$  is  $\Pr(\mu + \beta Z \leq x) = \Pr(Z \leq \frac{x-\mu}{\beta}) = \exp\left(-\exp\left(-\frac{x-\mu}{\beta}\right)\right)$ . But this is exactly the CDF of  $\overline{\text{Gumbel}}(\mu, \beta)$ .

<sup>8</sup>Recall  $\text{Exp}_\lambda(x) = \mathbb{1}[x \geq 0] \lambda \exp(-\lambda x)$  is the continuous version of the geometric distribution that represents how long we have to wait until an event with rate  $\lambda$  happens. It has the CDF:  $1 - \exp(-\lambda x)$ .

The proof is trivial for  $K = 1$  and simple for  $K = 2$  with  $\beta = 1$ .<sup>9</sup> The proof for the general case is less simple.

*Proof.* WLOG let  $k = 1$ .

$$\begin{aligned}
\Pr(Y = 1) &= \Pr_{\epsilon_1 \dots \epsilon_K \sim \text{Gumbel}(0, \beta)} (u_1 + \epsilon_1 \geq u_k + \epsilon_k \ \forall k > 1) \\
&= \mathbf{E}_{\epsilon_1 \sim \text{Gumbel}(0, \beta)} \left[ \prod_{k > 1} \Pr_{\epsilon_k \sim \text{Gumbel}(0, \beta)} (u_1 + \epsilon_1 \geq u_k + \epsilon_k) \right] \\
&= \mathbf{E}_{\epsilon_1 \sim \text{Gumbel}(0, \beta)} \left[ \prod_{k > 1} \Pr_{\epsilon_k \sim \text{Gumbel}(0, \beta)} (\epsilon_k \leq u_1 - u_k + \epsilon_1) \right] \\
&= \mathbf{E}_{\epsilon_1 \sim \text{Gumbel}(0, \beta)} \left[ \prod_{k > 1} \exp \left( - \exp \left( \frac{-u_1 + u_k - \epsilon_1}{\beta} \right) \right) \right] \\
&= \int_{-\infty}^{\infty} \frac{1}{\beta} \exp \left( - \frac{\epsilon_1}{\beta} - \exp \left( \frac{-\epsilon_1}{\beta} \right) \right) \prod_{k > 1} \exp \left( - \exp \left( \frac{-u_1 + u_k - \epsilon_1}{\beta} - \frac{\epsilon_1}{\beta} \right) \right) d\epsilon_1
\end{aligned}$$

Let  $x = \epsilon_1/\beta$  which yields the infinitesimal relation  $d\epsilon_1 = \beta dx$ . Plugging it in (the integrated region is still  $\mathbb{R}$ ), we have

$$\begin{aligned}
&\int_{-\infty}^{\infty} \exp(-x - \exp(-x)) \prod_{k > 1} \exp \left( - \exp \left( \frac{-u_1 + u_k - x}{\beta} - x \right) \right) dx \\
&= \int_{-\infty}^{\infty} \exp \left( -x - \exp(-x) - \sum_{k > 1} \exp \left( \frac{-u_1 + u_k}{\beta} \right) \exp(-x) \right) dx \\
&= \int_{-\infty}^{\infty} \exp \left( -x - \sum_{k=1}^K \exp \left( \frac{-u_1 + u_k}{\beta} \right) \exp(-x) \right) dx
\end{aligned}$$

By <https://www.integral-calculator.com>, for  $C > 0$ ,

$$\int_{-\infty}^{\infty} \exp(-x - C \exp(-x)) dx = \frac{\exp(-C \exp(-x))}{C} \Big|_{-\infty}^{\infty} = \frac{1}{C}$$

Thus the above integral evaluates to

$$\frac{1}{\sum_{k=1}^K \exp((-u_1 + u_k)/\beta)} = \frac{\exp(u_1/\beta)}{\sum_{k=1}^K \exp(u_k/\beta)} = \Pr(X = 1)$$

□

**Gumbel-Softmax.** We observe that for any  $\epsilon \in \mathbb{R}^K$  the *distribution*

$$\delta_\tau := \text{softmax} \left( \frac{u + \epsilon}{\tau} \right) \in [0, 1]^K$$

converges to the one-hot vector representation of  $k^* = \arg \max_{k=1}^K u_k + \epsilon_k$  as  $\tau \rightarrow 0^+$ . Thus if  $\epsilon_1 \dots \epsilon_K \sim \text{Gumbel}(0, 1)$ , then  $\delta_\tau$  is distributed as the one-hot vector representation of  $X \sim \text{Cat}(\text{softmax}(u))$  as  $\tau \rightarrow 0^+$  by Theorem A.1.

<sup>9</sup> WLOG let  $k = 1$ . Using the fact that the CDF of Logistic(0, 1) is the sigmoid function, we have

$$\Pr(Y = 1) = \Pr_{\epsilon_1, \epsilon_2 \sim \text{Gumbel}(0, 1)} (u_1 + \epsilon_1 \geq u_2 + \epsilon_2) = \Pr_{z \sim \text{Logistic}(0, 1)} (z \leq u_1 - u_2) = \sigma(u_1 - u_2)$$

By the usual relation between the sigmoid and the softmax when  $K = 2$ , we have

$$\sigma(u_1 - u_2) = \frac{1}{1 + \exp(u_2 - u_1)} = \frac{\exp(u_1)}{\exp(u_1) + \exp(u_2)} = \Pr(X = 1)$$

**Lemma A.2.** For any  $c > 0$ ,

$$\lim_{N \rightarrow \infty} \left(1 - \frac{1}{N}c\right)^N = \lim_{\epsilon \rightarrow 0^-} (1 + \epsilon c)^{-\frac{1}{\epsilon}} = e^{-c}$$

*Proof.* The first equality follows by the change of variable  $N = -\frac{1}{\epsilon}$  and taking the one-sided limit  $\epsilon \rightarrow 0^-$ . Taking negative log on both sides of the second equality gives

$$-\log \lim_{\epsilon \rightarrow 0^-} (1 + \epsilon c)^{-\frac{1}{\epsilon}} = \lim_{\epsilon \rightarrow 0^-} \frac{\log(1 + \epsilon c)}{\epsilon} = c \quad (2)$$

where the first equality holds since log is continuous at  $(1 + \epsilon c)^{-\frac{1}{\epsilon}} > 0$  for all  $\epsilon_{\min} < \epsilon < 0$  for some  $\epsilon_{\min}$ . Define  $f(x) = \log(1 + xc)$  and note that

$$f'(x) := \lim_{\epsilon \rightarrow 0} \frac{f(x + \epsilon) - f(x)}{\epsilon} = \lim_{\epsilon \rightarrow 0} \frac{\log(1 + (x + \epsilon)c) - \log(1 + xc)}{\epsilon} \quad \Rightarrow \quad f'(0) = \lim_{\epsilon \rightarrow 0} \frac{\log(1 + \epsilon c)}{\epsilon}$$

$f(x)$  is uniformly continuous at  $x = 0$  so the one-sided limit is the same as the two-sided limit. Hence the claim (2) is equivalent to  $f'(0) = c$ . This follows since  $f'(x) = c/(1 + xc)$  by the chain rule on log.  $\square$