

Control Variate

Richard Xu

December 4, 2020

1 Control variate motivations

variance reduction works by modifying “function of a random variable” so that its expectation remains same, but variance reduces

in this article, our primary interest is in estimating a derivative:

$$\nabla_{\theta} \mathbb{E}_{p(b|\theta)}[f(b)] \quad (1)$$

1.1 control variate C

We illustrate this through Reinforcement Learning, letting $b \rightarrow \tau$, $p \rightarrow \pi$ and $f(b) \rightarrow R(\tau)$:

1.1.1 Reinforcement Learning

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)}[R(\tau)] \\ \implies \nabla_{\theta} J(\theta) &= \nabla_{\theta} \int_{\tau} \pi_{\theta}(\tau) R(\tau) = \int_{\tau} \nabla_{\theta} \pi_{\theta}(\tau) R(\tau) \\ &= \int_{\tau} \nabla_{\theta} [\log(\pi_{\theta}(\tau))] \pi_{\theta}(\tau) R(\tau) \\ &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log(\pi_{\theta}(\tau))] R(\tau) \end{aligned}$$

now adding a control variate C :

$$\begin{aligned} &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log(\pi_{\theta}(\tau)) (R(\tau) - C)] \\ &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log(\pi_{\theta}(\tau)) R(\tau)] - \underbrace{\mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log(\pi_{\theta}(\tau)) C]}_{=0} \end{aligned} \quad (2)$$

so what does the RHS of Eq. (2) = 0?

$$\begin{aligned} &\mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log(\pi_{\theta}(\tau)) C] \\ &= \int_{\tau} \nabla_{\theta} \log(\pi_{\theta}(\tau)) C \times \pi_{\theta}(\tau) \\ &= C \int_{\tau} \underbrace{\nabla_{\theta} [\log(\pi_{\theta}(\tau)) \pi_{\theta}(\tau)]}_{\nabla_{\theta} \pi_{\theta}(\tau)} \\ &= C \nabla_{\theta} \int_{\tau} \pi_{\theta}(\tau) = 0 \end{aligned} \quad (3)$$

This means $\times C$ won't matter. This means that a constant C , i.e., **independent** of both θ and b will generate an unbiased control variate.

1.1.2 when C is a function of θ

however, when C is independent of τ , but dependent of θ :

$$\begin{aligned}\nabla_{\theta} \int_{\tau} C(\theta) \pi(\tau|\theta) d\tau &= \nabla_{\theta} C(\theta) \int_{\tau} \pi(\tau|\theta) d\tau \\ &= C'(\theta)\end{aligned}\tag{4}$$

2 More sophisticated Control Variate

We can add a more generic, biased control variate C . Since we know that when C is a function of b , then $\mathbb{E}_{p(b|\theta)} [C \nabla_{\theta} \log p(b|\theta)]$ is unbiased. Therefore, the biased C must be a function of b . Since it's biased, we need to delete it:

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{p(b|\theta)} [f(b)] &= \nabla_{\theta} \mathbb{E}_{p(b|\theta)} [f(b) - C + C] \\ &= \nabla_{\theta} \mathbb{E}_{p(b|\theta)} [f(b) - C] + \nabla_{\theta} \mathbb{E}_{p(b|\theta)} [C] \\ &= \mathbb{E}_{p(b|\theta)} \left[(f(b) - C) \nabla_{\theta} \log p(b|\theta) \right] + \nabla_{\theta} \mathbb{E}_{p(b|\theta)} [C] \\ &= \mathbb{E}_{p(b|\theta)} \left[f(b) \nabla_{\theta} \log p(b|\theta) \right] - \mathbb{E}_{p(b|\theta)} \left[C \nabla_{\theta} \log p(b|\theta) \right] + \underbrace{\nabla_{\theta} \mathbb{E}_{p(b|\theta)} [C]}_{\mathbb{E}_{p(b|\theta)} [C \nabla_{\theta} \log p(b|\theta)]} \\ &= \mathbb{E}_{p(b|\theta)} \left[f(b) \nabla_{\theta} \log p(b|\theta) \right] - \mathbb{E}_{p(b|\theta)} \left[C \nabla_{\theta} \log p(b|\theta) \right] + \mathbb{E}_{p(b|\theta)} [C \nabla_{\theta} \log p(b|\theta)]\end{aligned}\tag{5}$$

substitute $C \equiv C_{\phi}(z)$:

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{p(b|\theta)} [f(b)] &= \mathbb{E}_{p(b|\theta)} \left[f(b) \nabla_{\theta} \log p(b|\theta) \right] - \mathbb{E}_{p(b|\theta)} \left[C_{\phi}(z) \nabla_{\theta} \log p(b|\theta) \right] + \nabla_{\theta} \mathbb{E}_{p(b|\theta)} [C_{\phi}(z)] \\ &= \underbrace{\mathbb{E}_{p(b|\theta)} \left[f(b) \nabla_{\theta} \log p(b|\theta) \right]}_{\nabla_{\theta} \hat{f}_{\text{Reinforce}}} - \underbrace{\mathbb{E}_{p(b|\theta)} \left[C_{\phi}(z) \nabla_{\theta} \log p(b|\theta) \right]}_{\nabla_{\theta} \hat{C}_{\text{Reinforce without re-parameterization}}} + \underbrace{\nabla_{\theta} \mathbb{E}_{p(z)} [C_{\phi}(z)]}_{\nabla_{\theta} \hat{C}_{\text{Reinforce with re-parameterization}}}\end{aligned}\tag{6}$$

3 More sophisticated control variate: more the correlation, the better!

(Paisley ICML12) introduce a control variate $g(x)$ approximates $f(x)$ well when closed-form $\mathbb{E}[g(\theta)]$ isn't possible, a low variance approximate is used:

$$\hat{f}(x) = f(x) - h \left(\underbrace{g(x) - \mathbb{E}[g(x)]}_{\mathbb{E}[g(x) - \mathbb{E}[g(x)]] = 0} \right)\tag{7}$$

knowing:

$$\text{Var} \left(\sum_{i=1}^n a_i X_i \right) = \sum_{i=1}^n a_i^2 \text{Var}(X_i) + 2 \sum_{i=1}^n \sum_{j: j>i}^n a_i a_j \text{Cov}(X_i, X_j) \quad (8)$$

try it on $\text{Var}(\hat{f})$:

$$\begin{aligned} \text{Var}(\hat{f}) &= \text{Var}(f) - 2h \text{Cov}(f, g) + h^2 \text{Var}(g) \\ \implies \nabla_h \text{Var}(\hat{f}) &= -2 \text{Cov}(f, g) + 2h \text{Var}(g) = 0 \\ \implies h^* &= \frac{\text{Cov}(f, g)}{\text{Var}(g)} \end{aligned} \quad (9)$$

substitute h^* :

$$\begin{aligned} \text{Var}(\hat{f}) &= \text{Var}(f) - 2 \frac{\text{Cov}(f, g)}{\text{Var}(g)} \text{Cov}(f, g) + \frac{\text{Cov}(f, g)^2}{\text{Var}(g)^2} \text{Var}(g) \\ &= \text{Var}(f) - 2 \frac{\text{Cov}(f, g)^2}{\text{Var}(g)} + \frac{\text{Cov}(f, g)^2}{\text{Var}(g)} \\ &= \text{Var}(f) - \frac{\text{Cov}(f, g)^2}{\text{Var}(g)} \\ \implies \frac{\text{Var}(\hat{f})}{\text{Var}(f)} &= 1 - \frac{\text{Cov}(f, g)^2}{\text{Var}(g)^2} = 1 - \text{Corr}(f, g) \end{aligned} \quad (10)$$

meaning, higher the correlation between f and g , the less the variance of \hat{f}

4 RELAX algorithm:

the control variate we have added is:

$$C \equiv \mathbb{E}_{p(\tilde{z}|b, \theta)} [c_\phi(\tilde{z})] \quad (11)$$

note that C is a function of b , through $p(\tilde{z}|b, \theta)$:

$$\begin{aligned} \nabla_\theta \mathbb{E}_{p(b|\theta)} [f(b)] &= \mathbb{E}_{p(b|\theta)} \left[f(b) \nabla_\theta \log p(b|\theta) \right] - \mathbb{E}_{p(b|\theta)} \left[\mathbb{E}_{p(\tilde{z}|b, \theta)} [c_\phi(\tilde{z})] \nabla_\theta \log p(b|\theta) \right] \\ &\quad - \mathbb{E}_{p(b|\theta)} \left[\nabla_\theta \mathbb{E}_{p(\tilde{z}|b, \theta)} [c_\phi(\tilde{z})] \right] + \nabla_\theta \mathbb{E}_{p(z|\theta)} [c_\phi(z)] \end{aligned} \quad (12)$$

we can equally write $\mathbb{E}_{p(\tilde{z}|b, \theta)} [c_\phi(\tilde{z})]$ as $\mathbb{E}_{p(z|b, \theta)} [c_\phi(z)]$, however, \tilde{z} emphasis that it is sampled from posterior $p(\tilde{z}|b)$ condition on b , whereas z is from prior $p(z)$

the posterior term $p(\tilde{z}|b, \theta)$ highlights that \tilde{z} and b are correlated! The same does **not** seen in prior $p(z)$

note that RELAX estimator does not mean it correspond to a particular expression for C in Eq.(6)

4.1 Why the above is unbiased

imagine we can prove:

$$\mathbb{E}_{p(b|\theta)} \left[\mathbb{E}_{p(\tilde{z}|b,\theta)} [c_\phi(\tilde{z})] \nabla_\theta \log p(b|\theta) \right] = -\mathbb{E}_{p(b|\theta)} \left[\nabla_\theta \mathbb{E}_{p(\tilde{z}|b,\theta)} [c_\phi(\tilde{z})] \right] + \nabla_\theta \mathbb{E}_{p(z|\theta)} [c_\phi(z)] \quad (13)$$

or,

$$\nabla_\theta \mathbb{E}_{p(z|\theta)} [c_\phi(z)] = \mathbb{E}_{p(b|\theta)} \left[\nabla_\theta \mathbb{E}_{p(\tilde{z}|b,\theta)} [c_\phi(\tilde{z})] \right] + \mathbb{E}_{p(b|\theta)} \left[\mathbb{E}_{p(\tilde{z}|b,\theta)} [c_\phi(\tilde{z})] \nabla_\theta \log p(b|\theta) \right] \quad (14)$$

then Eq.(25) is unbiased estimator and we are all done.

To show this, we first need under **Deterministic function**: it has a very special property when $p(b|z)$ is deterministic:

$$\begin{aligned} p(z) &\equiv p(z|b)p(b) \quad \text{or} \\ p(z|\theta) &\equiv p(z|b, \theta)p(b|\theta) \end{aligned} \quad (15)$$

with that

$$\begin{aligned} \nabla_\theta \mathbb{E}_{p(z|\theta)} [c_\phi(z)] &= \mathbb{E}_{p(z|\theta)} [\nabla_\theta \log(p(z|\theta)) c_\phi(z)] \\ &= \mathbb{E}_{p(b|\theta)} \left[\mathbb{E}_{p(\tilde{z}|b,\theta)} [c_\phi(\tilde{z}) \nabla_\theta \log(p(z|\theta))] \right] \quad \text{using Eq.(15)} \\ &= \mathbb{E}_{p(b|\theta)} \left[\mathbb{E}_{p(\tilde{z}|b,\theta)} [c_\phi(\tilde{z}) \nabla_\theta \log(p(z|\theta) + p(z|b, \theta))] \right] \quad \text{using Eq.(15) again!} \\ &= \mathbb{E}_{p(b|\theta)} \left[\mathbb{E}_{p(\tilde{z}|b,\theta)} [c_\phi(\tilde{z}) \nabla_\theta \log p(z|\theta)] \right] + \mathbb{E}_{p(b|\theta)} \left[\mathbb{E}_{p(\tilde{z}|b,\theta)} [c_\phi(\tilde{z}) \nabla_\theta \log p(z|b, \theta)] \right] \\ &= \mathbb{E}_{p(b|\theta)} \left[\nabla_\theta \mathbb{E}_{p(\tilde{z}|b,\theta)} [c_\phi(\tilde{z})] \right] + \mathbb{E}_{p(b|\theta)} \left[\mathbb{E}_{p(\tilde{z}|b,\theta)} [c_\phi(\tilde{z})] \nabla_\theta \log p(b|\theta) \right] \end{aligned} \quad (16)$$

also, for $\mathbb{E}_{p(b|\theta)} \left[f(b) \nabla_\theta \log p(b|\theta) \right]$:

$$\begin{aligned}
\mathbb{E}_{p(z)} \left[f(H(z)) \nabla_{\theta} \log p(z) \right] &= \mathbb{E}_{p(b)} \left[\mathbb{E}_{p(z|b)} \left[f(H(z)) \nabla_{\theta} \log p(z) \right] \right] \\
&= \mathbb{E}_{p(b)} \left[\mathbb{E}_{p(z|b)} \left[f(b) \nabla_{\theta} \log p(z) \right] \right] \\
&= \mathbb{E}_{p(b)} \left[f(b) \mathbb{E}_{p(z|b)} \left[\nabla_{\theta} \log p(z) \right] \right] \\
&= \mathbb{E}_{p(b)} \left[f(b) \mathbb{E}_{p(z|b)} \left[\nabla_{\theta} \log(p(z|b)) + \nabla_{\theta} \log(p(b)) \right] \right] \\
&= \mathbb{E}_{p(b)} \left[f(b) \underbrace{\mathbb{E}_{p(z|b)} \left[\nabla_{\theta} \log(p(z|b)) \right]}_{=0} \right] + \mathbb{E}_{p(b)} \left[f(b) \nabla_{\theta} \log(p(b)) \right] \\
&= \mathbb{E}_{p(b)} \left[f(b) \underbrace{\int_z \nabla_{\theta} p(z|b)}_{=0} \right] + \nabla_{\theta} \mathbb{E}_{p(b)} \left[f(b) \right] \\
&= \nabla_{\theta} \mathbb{E}_{p(b)} \left[f(b) \right] \\
&= \mathbb{E}_{p(b|\theta)} \left[f(b) \nabla_{\theta} \log p(b|\theta) \right]
\end{aligned}$$

(17)

Add Baseline trick to the problem

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [R(\tau)] \\ \implies \nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) R(\tau)] - \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [C \nabla_{\theta} \log \pi_{\theta}(\tau)] \end{aligned} \quad (18)$$

replace $\pi_{\theta}(\tau)$ with $p(b|\theta)$ and $r(\tau)$ with $f(b)$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{b \sim p(b|\theta)} [\nabla_{\theta} \log p(b|\theta) f(b)] - \mathbb{E}_{b \sim p(b|\theta)} [C \nabla_{\theta} \log p(b|\theta)] \quad (19)$$

what should a good expression of C be?

Variance Reduction through control variate

apply this sophisticated control variate:

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{p(b,z)} [f(b)] &= \nabla_{\theta} \left(\mathbb{E}_{p(b,z)} [f(b) - g(z)] + \mathbb{E}_{p(b,z)} [g(z)] \right) \\ &= \mathbb{E}_{p(b,z)} \left[(f(b) - g(z)) \nabla_{\theta} \log p(b) \right] + \nabla_{\theta} \mathbb{E}_{p(b,z)} [g(z)] \\ \nabla_{\theta} \mathbb{E}_{p(b)} [f(b)] &= \mathbb{E}_{p(b)} \left[f(b) \nabla_{\theta} \log p(b) \right] - \mathbb{E}_{p(b,z)} \left[g(z) \nabla_{\theta} \log p(b) \right] + \nabla_{\theta} \mathbb{E}_{p(z)} [g(z)] \end{aligned} \quad (20)$$

a good choice of $g(z)$ is important, remember $\text{corr}(f, g)$ needs to be high

$$g(z) \equiv \mathbb{E}_{p(z|b)} [f(\sigma_{\lambda}(z))] \quad (21)$$

after some simplification:

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{p(b)} [f(b)] &= \mathbb{E}_{p(b)} \left[(f(b) \nabla_{\theta} \log p(b)) \right] - \mathbb{E}_{p(b)} \left[\mathbb{E}_{p(z|b)} [f(\sigma_{\lambda}(z))] \nabla_{\theta} \log p(b) \right] + \nabla_{\theta} \underbrace{\mathbb{E}_{p(z|b)} [f(\sigma_{\lambda}(z))]}_{p(z)} \end{aligned} \quad (22)$$

change of variables
this is:

$$p(z) = \left| \frac{d\epsilon}{dz} \right| p(\epsilon) \implies |p(z) dz| = |p(\epsilon) d\epsilon| \quad (23)$$

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{p(z;\theta)} [f(z)] &= \nabla_{\theta} \int p(z; \theta) f(z) dz \\ &= \nabla_{\theta} \int p(\epsilon) f(z) d\epsilon = \nabla_{\theta} \int p(\epsilon) f(g(\epsilon, \theta)) d\epsilon \\ &= \nabla_{\theta} \mathbb{E}_{p(\epsilon)} [f(g(\epsilon, \theta))] = \mathbb{E}_{p(\epsilon)} [\nabla_{\theta} f(g(\epsilon, \theta))] \end{aligned} \quad (24)$$

5 REBAR algorithm

5.1 Put REBAR in action

we want to estimate $\nabla_{\theta} \mathbb{E}_{p(b|\theta)} [f(b)]$, the estimator used is:

$$\begin{aligned} & \nabla_{\theta} \mathbb{E}_{p(b|\theta)} [f(b)] \\ &= \mathbb{E}_{p(b|\theta)} \left[f(b) \nabla_{\theta} \log p(b|\theta) \right] - \mathbb{E}_{p(b|\theta)} \left[\mathbb{E}_{p(\tilde{z}|b, \theta)} [c_{\phi}(\tilde{z})] \nabla_{\theta} \log p(b|\theta) \right] \\ & \quad - \mathbb{E}_{p(b|\theta)} \left[\nabla_{\theta} \mathbb{E}_{p(\tilde{z}|b, \theta)} [c_{\phi}(\tilde{z})] \right] + \nabla_{\theta} \mathbb{E}_{p(z|\theta)} [c_{\phi}(z)] \end{aligned} \quad (25)$$

remember you may substitute any function of $c_{\phi}(\tilde{z})$. In REBAR, one substitute:

$$c_{\phi}(\tilde{z}) \equiv \sigma_{\lambda}(\tilde{z}) \quad (26)$$

5.1.1 What is $\sigma_{\lambda}(\tilde{z}(\theta, u))$

In **un-relaxed** version, you can obtain b from $z(\theta, u)$ by:

$$b = H(z(\theta, u)) \quad (27)$$

In REBAR, you can obtain a **relaxed** version:

$$\tilde{b} = \sigma_{\lambda}(z(\theta, u)) = \frac{1}{1 + \exp(-z/\lambda)} \quad (28)$$

$$\begin{aligned} & \nabla_{\theta} \mathbb{E}_{p(b)} [f(b)] \\ &= \mathbb{E}_{p(b)} \left[(f(b) \nabla_{\theta} \log p(b)) \right] - \mathbb{E}_{p(b)} \left[\mathbb{E}_{p(z|b)} [f(\sigma_{\lambda}(z))] \nabla_{\theta} \log p(b) \right] \\ & \quad - \mathbb{E}_{p(b)} \nabla_{\theta} \mathbb{E}_{p(z|b)} [f(\sigma_{\lambda}(z))] + \mathbb{E}_{p(b)} [\mathbb{E}_{p(z|b)} [f(\sigma_{\lambda}(z)) \nabla_{\theta} \log(p(b))]] \end{aligned} \quad (29)$$

5.2 REBAR in action

5.2.1 Under a Bernoulli setting

illustrate the case where we need to estimate $\nabla_{\theta} \mathbb{E}_{p(b)} [f(b)]$ where $b \sim \text{Bernoulli}(\theta)$

from detailed derivation in re-parameterization section, we see re-parameterization of Logistic random variable can be done by:

$$\begin{aligned} & u \sim \mathcal{U}(0, 1) \\ & z(\theta, u) = \log \left(\frac{\theta}{1 - \theta} \right) + \log \left(\frac{u}{1 - u} \right) \end{aligned} \quad (30)$$

so REBAR version is:

$$\begin{aligned}\tilde{b} &= \sigma_\lambda(z(\theta, u)) = \frac{1}{1 + \exp(-z/\lambda)} \\ &= \frac{1}{1 + \exp\left[-\left(\log\left(\frac{\theta}{1-\theta}\right) - \log\left(\frac{u}{1-u}\right)\right)/\lambda\right]}\end{aligned}\quad (31)$$

after **re-parameterization**:

$$\begin{aligned}\nabla_\theta \mathbb{E}_{p(b, z)}[f(b)] \\ &= \mathbb{E}_{u \sim \mathcal{U}}[f(H((z))\nabla_\theta \log(p(z)))] - \mathbb{E}_{u \sim \mathcal{U}}\left[\mathbb{E}_{v \sim \mathcal{U}}[f(\sigma_\lambda(\tilde{z}))]\nabla_\theta \log p(b)\right] \\ &\quad - \underbrace{\mathbb{E}_{p(b)}\left[\mathbb{E}_{v \sim \mathcal{U}}[\nabla_\theta f(\sigma_\lambda(\tilde{z}))]\right]}_{\text{using } \tilde{z}} + \underbrace{\nabla_\theta \mathbb{E}_{p(b)}\left[\mathbb{E}_{v \sim \mathcal{U}}[f(\sigma_\lambda(z))]\right]}_{\text{using } z}\end{aligned}\quad (32)$$

5.2.2 REBAR algorithm

the algorithm is like, at each iteration, the derivative is calculated as:

$$\begin{aligned}\nabla_\theta \mathbb{E}_{p(b)}[f(b)] \\ = \mathbb{E}_{p(u, v)}\left[\left(f(b) - f(\sigma_\lambda(\tilde{z}))\right)\nabla_\theta \log p(b|\theta) - (\nabla_\theta f(\sigma_\lambda(\tilde{z})) - \nabla_\theta f(\sigma_{\tilde{z}}))\right]\end{aligned}\quad (33)$$

basically, $u \rightarrow z \rightarrow b$, then having $(b, v) \rightarrow \tilde{z}$

6 re-parameterization of $p(\mathbf{z}|\mathbf{b})$ in general case

we look at even more generically case of k , instead of binary case
 this means knowing $\mathbf{z} = [z_1, \dots, z_n]$, a vector of Gumbel random variables,

1. one-hot vector $\mathbf{b} = [0, 0, \dots, \underbrace{1}_{k^{\text{th}}}, 0, \dots]$

2. probability vector \mathbf{p}

3. CDF sampling variables \mathbf{v}

what is the linking function \tilde{g} which outputs $\tilde{\mathbf{z}}$?

$$\tilde{\mathbf{z}} = \tilde{g}(\mathbf{v}, \mathbf{b}, \mathbf{p}) \quad (34)$$

obviously the z_k correspond to 1 in one-hot vector should be the largest!

6.1 Gumbel distribution properties

1. Gumbel's first wonderful property derived from its CDF:

$$\mathcal{GA}(g; \phi) = \exp(-\exp(-g + \phi))$$

then:

$$\begin{aligned} \Rightarrow \mathcal{GA}(g; \phi) \mathcal{GA}(g; \gamma) &= \mathcal{GA}\left(g; \log(\exp(\phi) + \exp(\gamma))\right) \\ &\quad \text{product of two Gumbel CDFs is a Gumbel CDF!} \\ \Rightarrow \prod_{i=1}^n \mathcal{GA}(g; \phi_i) &= \mathcal{GA}\left(g; \log\left(\underbrace{\sum_{i=1}^n \exp(\phi_i)}_{\mathcal{Z}}\right)\right) = \mathcal{GA}(g; \log(\mathcal{Z})) \end{aligned} \quad (35)$$

in words, product of multiple Gumbel random variables are also Gumbel, with parameters are the **logSumExp** of individual parameters

2. Gumbel's second wonderful property:

$$\underbrace{\mathcal{ga}(g; \phi)}_{\text{pdf}} = \exp(-g + \phi) \underbrace{\mathcal{GA}(g; \phi)}_{\text{cdf}} \quad (36)$$

6.2 Formulation for $p(\mathbf{z}|\mathbf{b})$

We know from Gumbel-Max trick that:

Given a set of Gumbel RVs:

$$\mathbf{z} = \{z_1 \sim \text{ga}(\log(p_1)), \dots, z_n \sim \text{ga}(\log(p_n))\} \quad (37)$$

probability that k^{th} element is their maximum is equal probability of sampling:

$$k \sim (p_1, \dots, p_n), \text{i.e., } p_k \quad (38)$$

condition on the prior:

$$p(\mathbf{z}) = \prod_{i=1}^n \text{ga}(z_i; \log(p_i)) \quad (39)$$

the conditional density that the k^{th} of one-hot vector \mathbf{b} is one, i.e., $b_k = 1$, $b_{i \neq k} = 0 \forall i$ is:

$$b_i = H_i(\mathbf{z}) \begin{cases} 1 & \text{if } z_i \geq z_j \forall j \neq i \\ 0 & \text{otherwise} \end{cases} \quad (40)$$

where H is the “argmax” operator

$b_k = 1$ means one-hot \mathbf{b} vector having k^{th} element one, the conditional:

$$p(b_k = 1|\mathbf{z}) = \mathbf{1}(z_k \geq z_i) \quad (41)$$

computing posterior using usual:

$$\begin{aligned} p(\mathbf{z}|b_k = 1) &= \frac{p(\mathbf{z}, b_k = 1)}{p(b_k = 1)} \\ &= \frac{p(b_k = 1|\mathbf{z})p(\mathbf{z})}{p(b_k = 1)} \end{aligned} \quad (42)$$

note that this is different to computing softmax, where we need to integrate out z_k

$$\begin{aligned}
p(b_k = 1|\mathbf{z})p(\mathbf{z}) &= \underbrace{\prod_{i=1}^n \text{ga}(z_i; \log(p_i))}_{p(\mathbf{z})} \underbrace{\mathbf{1}(z_k \geq z_i)}_{p(b_k=1|\mathbf{z})} \\
&= \frac{\mathcal{GA}(z_k; \log(1-p_k))}{\mathcal{GA}(z_k; \log(1-p_k))} \prod_{i=1}^n \text{ga}(z_i; \log(p_i)) \mathbf{1}(z_k \geq z_i) \\
&= \frac{\text{ga}(z_k; \log(p_k)) \mathcal{GA}(z_k; \log(1-p_k))}{\mathcal{GA}(z_k; \log(1-p_k))} \prod_{i \neq k}^n \text{ga}(z_i; \log(p_i)) \mathbf{1}(z_k \geq z_i) \\
&= \text{ga}(z_k; \log(p_k)) \mathcal{GA}(z_k; \log(1-p_k)) \prod_{i \neq k}^n \frac{\text{ga}(z_i; \log(p_i)) \mathbf{1}(z_k \geq z_i)}{\mathcal{GA}(z_k; \log(1-p_k))}
\end{aligned} \tag{43}$$

using **Property 1**: Eq.(35):

$$\begin{aligned}
\mathcal{GA}(z_k; \log(1-p_k)) &= \mathcal{GA}\left(z_k; \log\left(\sum_{i \neq k} p_i\right)\right) \\
&= \mathcal{GA}\left(z_k; \log \underbrace{\sum_{i \neq k} \exp(\underbrace{\log(p_i)}_{\phi_i})}_{\log \text{SumExp}}\right) \\
&= \prod_{i \neq k}^n \mathcal{GA}(z_k; \log(p_i))
\end{aligned} \tag{44}$$

$$\begin{aligned}
&= \underbrace{\text{ga}(z_k; \log(p_k))}_{\text{Property 2: Eq.(36) } \text{ga}(g; \phi) = \exp(-g + \phi) \mathcal{GA}(g; \phi)} \mathcal{GA}(z_k; \log(1 - p_k)) \prod_{i \neq k}^n \frac{\text{ga}(z_i; \log(p_i)) \mathbf{1}(z_k \geq z_i)}{\mathcal{GA}(z_k; \log(p_i))} \\
&= \underbrace{\exp(-z_k + \log(p_k)) \mathcal{GA}(z_k; \log(p_k))}_{\text{Property 2: Eq.(36) } \text{ga}(g; \phi) = \exp(-g + \phi) \mathcal{GA}(g; \phi)} \mathcal{GA}(z_k; \log(1 - p_k)) \prod_{i \neq k}^n \frac{p(z_i; \log(p_i)) \mathbf{1}(z_k \geq z_i)}{\mathcal{GA}(z_k; \log(p_i))} \\
&= \exp(-z_k + \log(p_k)) \underbrace{\mathcal{GA}(z_k; \log(p_k)) \mathcal{GA}(z_k; \log(1 - p_k))}_{\text{Property 1: Eq.(35)}} \prod_{i \neq k}^n \frac{p(z_i; \log(p_i)) \mathbf{1}(z_k \geq z_i)}{\mathcal{GA}(z_k; \log(p_i))} \\
&= \underbrace{p_k \exp(-z_k)}_{\text{Property 1: Eq.(35)}} \underbrace{\mathcal{GA}(z_k; \log(1))}_{\text{Property 1: Eq.(35)}} \prod_{i \neq k}^n \frac{p(z_i; \log(p_i)) \mathbf{1}(z_k \geq z_i)}{\mathcal{GA}(z_k; \log(p_i))} \\
&= p_k \underbrace{\exp(-z_k) \mathcal{GA}(z_k; 0)}_{\text{Property 1: Eq.(35)}} \prod_{i \neq k}^n \frac{\text{ga}(z_i; \log(p_i)) \mathbf{1}(z_k \geq z_i)}{\mathcal{GA}(z_k; \log(p_i))} \quad \log(1) = 0 \\
&= p_k \underbrace{\underbrace{\text{ga}(z_k; 0)}_{\text{Eq.(36) } \text{ga}(g; \phi) = \exp(-g + \phi) \mathcal{GA}(g; \phi)}}_{\text{Eq.(36) } \text{ga}(g; \phi) = \exp(-g + \phi) \mathcal{GA}(g; \phi)} \prod_{i \neq k}^n \frac{\text{ga}(z_i; \log(p_i)) \mathbf{1}(z_k \geq z_i)}{\mathcal{GA}(z_k; \log(p_i))} \\
&= \underbrace{p_k}_{p(b_k=1)} \underbrace{p_0(z_k) \prod_{i \neq k}^n \frac{\text{ga}(z_i; \log(p_i)) \mathbf{1}(z_k \geq z_i)}{\mathcal{GA}(z_k; \log(p_i))}}_{p(\mathbf{z}|b_k=1)} \\
&= p(\mathbf{z}|p_k = 1) \Pr(b_k = 1)
\end{aligned} \tag{45}$$

Gumbel-max trick tells us that $\Pr(b_k = 1) = p_k$, so we have the conditional density required for the linking function:

6.2.1 sample $p(\mathbf{z}|p_k = 1)$ via re-parameterization

from uniform random variables \mathbf{v} , the link function to obtain \mathbf{z} is:

$$\mathbf{z} = \tilde{g}(\mathbf{v}, \mathbf{b}, \mathbf{p}) \tag{46}$$

we can obtain the link function by looking at the conditional:

$$p(\mathbf{z}|b_k = 1) = \text{ga}(z_k; 0) \prod_{i \neq k}^n \frac{\text{ga}(z_i; \log(p_i)) \mathbf{1}(z_k \geq z_i)}{\mathcal{GA}(z_k; \log(p_i))} \tag{47}$$

link function is different for z_k and $z_{k \neq i}$:

$$z_i = \begin{cases} \mathcal{GA}^{-1}(v_i, 0) & \text{if } i = k \\ \mathcal{GA}_{\text{Truncated}}^{-1}(v_i, \log(p_i), z_k) & \text{otherwise} \end{cases} \tag{48}$$

$$\mathcal{GA}^{-1}(z_k; v_k, 0) = -\log(-\log(v_k)) \tag{49}$$

and we know:

$$\mathcal{GA}_{\text{truncate}}^{-1}(u, \phi, T) = \phi - \log(\exp(\phi - T) - \log(u)) \quad (50)$$

let

$$\begin{aligned} \phi &= \log(p_i) \\ T &= z_k = -\log(-\log(v_k)) \end{aligned} \quad (51)$$

then:

$$\begin{aligned} &\mathcal{GA}_{\text{truncate}}^{-1}(z_i; v_i, \log(p_i), z_k) \\ &= \log(p_i) - \log(\exp(\log(p_i) + \log(-\log(v_k))) - \log(u)) \\ &= \log(p_i) - \log(\exp(\log(\frac{p_i}{\log(v_k)})) - \log(u)) \\ &= -\log\left(\frac{1}{p_i}\right) - \log\left(\frac{p_i}{\log(v_k)} - \log(u)\right) \\ &= -\log\left(\left(\frac{1}{p_i}\right)\frac{p_i}{\log(v_k)} - \left(\frac{1}{p_i}\right)\log(u)\right) \\ &= -\log\left(-\log(v_k) - \frac{\log(u)}{p_i}\right) \end{aligned} \quad (52)$$

6.2.2 relationship with how Gumbel-max trick is caculated

key point to remember, the above require the explicit values of $z_{i \neq k}$, therefore, it uses $p(\log(p_i))$ whereas Gumbel-max trick is calculated without the explicit value of $z_{i \neq k}$, therefore it uses $\Pr(\log(p_i))$

6.2.3 Formulation for $p(z|b)$ under recursive truncation τ

another problem is to add in further truncation τ (this is not used in this setting), in here, no requirement to obtain explicit value of $z_{i \neq k}$, so we use $\mathcal{GA}(\cdot)$

and we perform the same trick, i.e., write down joint density:

$$\begin{aligned}
& p(\text{max element} = \textcolor{red}{b}, \text{max value} = \textcolor{blue}{z} | \tau) \\
&= p\left(G(b) = z, G(b) \geq \max_{i \neq b} G(i) | \tau\right) = \textcolor{red}{p}(z) \textcolor{red}{p}(b|z) \\
&= \underbrace{\frac{\textcolor{red}{ga}(z; \phi_k) \mathbf{1}(z \leq \tau)}{\textcolor{green}{GA}(\tau; \phi_k)}}_{\Pr(G(i^* \equiv b) = z | \tau)} \times \underbrace{\prod_{i \neq b} \frac{\mathcal{GA}(z; \phi_i)}{\textcolor{blue}{GA}(\tau; \phi_i)}}_{\Pr(G(b) \geq \max_{i \neq b} G(i) | G(b) = z, \tau)} \\
&= \exp(-z + \phi_k) \mathbf{1}(z \leq \tau) \times \left(\frac{\textcolor{red}{GA}(z; \phi_k)}{\textcolor{green}{GA}(\tau; \phi_k)} \times \prod_{i \neq b} \frac{\mathcal{GA}(z; \phi_i)}{\textcolor{blue}{GA}(\tau; \phi_i)} \right) \text{ using } p_\phi(z) = \exp(-z + \phi) \mathcal{GA}(z; \phi) \\
&= \frac{\exp(\phi_k)}{\textcolor{blue}{Z}} \exp(-z + \textcolor{blue}{\log Z}) \mathbf{1}(z \leq \tau) \frac{\mathcal{GA}(z; \log(\mathcal{Z}))}{\textcolor{blue}{GA}(\tau; \log(\mathcal{Z}))} \quad \text{using } \prod_{i=1}^n \mathcal{GA}(z; \phi_i) = \mathcal{GA}(z; \log(\mathcal{Z})) \\
&= \frac{\exp(\phi_k)}{\mathcal{Z}} \overbrace{\frac{\exp(-z + \log \mathcal{Z}) \mathcal{GA}(z; \log(\mathcal{Z}))}{\mathcal{GA}(\tau; \log(\mathcal{Z}))} \mathbf{1}(z \leq \tau)}^{p_\phi(z) = \exp(-z + \phi) \mathcal{GA}(z; \phi)} \\
&= \underbrace{\frac{\exp(\phi_k)}{\mathcal{Z}}}_{\textcolor{red}{Pr}(b)} \underbrace{\frac{\textcolor{red}{ga}(z; \log(\mathcal{Z}))}{\textcolor{blue}{GA}(\tau; \log(\mathcal{Z}))} \mathbf{1}(z \leq \tau)}_{\textcolor{red}{p}(z|b, \tau)} \\
&= \Pr(b) p(z|b, \tau)
\end{aligned} \tag{53}$$

6.3 Re-parameterization in Bernoulli variable

just like the Gumbel-Max trick

$$\Pr(k : k \sim (\text{softmax}(\Phi_1, \dots, \Phi_k))) = \Pr(k = \textcolor{red}{\max}(\Phi_1 + z_1, \dots, \Phi_k + z_k)) \tag{54}$$

where $z_i \sim \text{ga}(0, 1)$

in a similar way, instead of sampling $b \sim \text{Bernoulli}(\theta)$, we can use **re-parameterization** of $u \sim \mathcal{U}(0, 1)$

$$\underbrace{\Pr(b = 1 | \theta)}_{=\theta} = \Pr\left(b = \textcolor{red}{H}\left(\underbrace{\log\left(\frac{\theta}{1-\theta}\right) + \log\left(\frac{u}{1-u}\right)}_{z(\theta, u)}\right)\right) \tag{55}$$

where $u \sim U(0, 1)$ $H(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$

in fact, the above is a specialized Gumbel-Max trick of binary cases, where H is the “max” operator

6.3.1 Re-parameterization of binary case via Logistic Distribution

Logistic Distribution has the following properties:

$$\text{PDF : } p(\mu, s) = \frac{\exp^{-\frac{x-\mu}{s}}}{s \left(1 + \exp^{-\frac{x-\mu}{s}}\right)^2} \quad (56)$$

$$\text{CDF : } \Pr(\mu, s) = \frac{1}{1 + \exp^{-\frac{x-\mu}{s}}} \quad (57)$$

$$\begin{aligned} \text{CDF}^{-1} : \quad u &= \frac{1}{1 + \exp^{-\frac{x-\mu}{s}}} \\ \implies \frac{1}{u} &= 1 + \exp^{-\frac{x-\mu}{s}} \implies \frac{1}{u} - 1 = \exp^{-\frac{x-\mu}{s}} \\ \implies \log\left(\frac{1}{u} - 1\right) &= -\frac{x-\mu}{s} \\ \implies x &= -\log\left(\frac{1}{u} - 1\right)s + \mu = -\log\left(\frac{1-u}{u}\right)s + \mu \\ x &= (\log u - \log(1-u))s + \mu \end{aligned} \quad (58)$$

6.3.2 Obtain linking function $z = g(\theta, u)$ via Logistic Distribution

property: if $z_1 \sim \text{ga}(\mu_1, \beta)$ and $z_2 \sim \text{ga}(\mu_2, \beta)$, then:

$$z = z_1 - z_2 \sim \text{Logistic}(\mu_1 - \mu_2, \beta) \quad (59)$$

difference of two Gumbel RVs is a Logistic RV!

in binary case, i.e., $K = 2$, max of K Gumbels, becomes “max of two Gumbel random variables” with locations $\mu_1 = \log \alpha_1$ and $\mu_2 = \log \alpha_2$ respectively, and also let $\beta = 1$

$$U \sim \mathcal{U}(0, 1), \text{ then } z = z_1 - z_2 = \log U - \log(1 - U) + \log \alpha_1 - \log \alpha_2 \quad (60)$$

$$b = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (61)$$

it is obvious that:

$$H(z) = H\left(\log U - \log(1 - U) + \log\left(\frac{\alpha_1}{\alpha_2}\right)\right) \quad (62)$$

where H is the unit step function

$$\begin{aligned}
\Pr([b \equiv H(z)] = 1) &= P(z_1 \geq z_2) \\
&= P(z_1 - z_2 \geq 0) \\
&= P\left(\underbrace{\log U - \log(1 - U) + \log \frac{\alpha_1}{\alpha_2}}_z \geq 0\right) \\
&= P(z \geq 0)
\end{aligned} \tag{63}$$

in case $\alpha_1 = \theta$ and $\alpha_2 = 1 - \theta$:

$$z = \log \frac{\theta}{1 - \theta} + \log \frac{U}{1 - U} \tag{64}$$

6.3.3 verify re-parameterization of u :

in here, we work backwards to verify:

$$\begin{aligned}
H(z) &= 1 \quad \text{if } z > 0 \\
\implies b = H(z(\theta, u)) &= 1 \quad \text{if } z(\theta, u) > 0
\end{aligned} \tag{65}$$

so to find $\Pr(b = 1)$ we just need to find $p(z(\theta, u) \geq 0)$, it turns out that $z(\theta, u)$ is monotonically increasing, so, let $u_0 \in (0, 1)$

$$\begin{aligned}
\Pr(b = 1) &\equiv \Pr(z(\theta, u) \geq 0) = \Pr\left(u \geq u_0 \mid \underbrace{z(\theta, u_0) = 0}_{u_0 \text{ is some cut-off}}\right) \\
&= 1 - u_0
\end{aligned} \tag{66}$$

need to find cut-off u_0 , such that $z(\theta, u_0) = 0$:

$$\begin{aligned}
z(\theta, u_0) = 0 &\implies \log\left(\frac{\theta}{1 - \theta}\right) + \log\left(\frac{u_0}{1 - u_0}\right) = 0 \\
&\implies \log\left(\frac{u_0}{1 - u_0}\right) = -\log\left(\frac{\theta}{1 - \theta}\right) \implies \log\left(\frac{u_0}{1 - u_0}\right) = \log\left(\frac{1 - \theta}{\theta}\right) \\
&\implies \frac{u_0}{1 - u_0} = \frac{1 - \theta}{\theta} \\
&\implies u_0 = \frac{1 - \theta}{\theta} \times (1 - u_0) \implies u_0 \times \left(1 + \frac{1 - \theta}{\theta}\right) = \frac{1 - \theta}{\theta} \\
&\implies u_0 \times \frac{1}{\theta} = \frac{1 - \theta}{\theta} \implies u_0 = 1 - \theta
\end{aligned} \tag{67}$$

to summarize:

$$\begin{aligned}
\Pr(b = 1|\theta) &= \Pr(z(\theta, u) > 0|\theta) = \Pr(u > u_0) = 1 - u_0 = \theta \\
&\text{where } u \sim U(0, 1)
\end{aligned} \tag{68}$$

which shows it's a Bernoulli distribution

6.3.4 $p(z|b)$ under binary case

it may be tempting to solve it in the same way as you do for multi-class case, i.e., looking at the joint density and hope to deduce $p(z|b)$ of the something like:

$$\begin{aligned}
 p(b_k = 1|z)p(z) &= \underbrace{\text{Logistic}(z; \log(\theta) - \log(1 - \theta))}_{p(z)} \underbrace{\mathbf{1}(z \geq 0)}_{p(b=1|z)} \\
 &= \frac{\exp(-(z - \mu))}{(1 + \exp(-(z - \mu)))^2} \mathbf{1}(z \geq 0) \\
 &= \frac{\exp(-(z - \log \theta + \log(1 - \theta)))}{(1 + \exp(-(z - \log \theta + \log(1 - \theta))))^2} \mathbf{1}(z \geq 0)
 \end{aligned} \tag{69}$$

it is difficult to see where this leads, but there is an easier way to solve it, looking at [?]:

$$\begin{aligned}
 b = 1 &\implies z = \log \frac{\theta}{1 - \theta} + \log \frac{u}{1 - u} \geq 0 \\
 &\implies \log \frac{u}{1 - u} \geq -\log \frac{\theta}{1 - \theta} \\
 &\quad \geq \log \frac{1 - \theta}{\theta} \\
 &\implies \frac{u}{1 - u} \geq \frac{1 - \theta}{\theta} \\
 &\quad \implies u \theta \geq (1 - u)(1 - \theta) \quad \text{valid as } \theta \geq 0 \quad u \geq 0 \\
 &\quad \geq 1 - u - \theta + u\theta \\
 &\quad \implies u \geq 1 - \theta
 \end{aligned} \tag{70}$$

in words, it meant the unit line is divided by the point $1 - \theta$:

$$b = \begin{cases} 1 & u \in (1 - \theta, 1) & \text{this interval has length } \theta \\ 0 & u \in (0, 1 - \theta) & \text{this interval has length } 1 - \theta \end{cases} \tag{71}$$

so instead of sampling $u \sim \mathcal{U}(0, 1)$ in the prior case, for sampling $p(z|b)$ we transform u conditionally by:

$$b = \begin{cases} 1 & v' = (1 - \theta) + \theta u \\ 0 & v' = (1 - \theta)u \end{cases} \tag{72}$$

then compute the inverse of Logistic distribution using v' :

$$z|b = \log \frac{\theta}{1 - \theta} + \log \frac{v'}{1 - v'} \tag{73}$$

7 Generalised Re-parameterization

7.1 a quick revision on change of variables in probability

7.1.1 continuous PDF

for re-parameterization knowing:

$$\begin{aligned} q_\epsilon(\epsilon) &= \left| \det \frac{\partial z}{\partial \epsilon} \right| q(z) \\ q_z(z) &= \left| \det \frac{\partial \epsilon}{\partial z} \right| q_\epsilon(\epsilon) \end{aligned} \quad (74)$$

because let $z = g(\epsilon) \implies \epsilon = g^{-1}(z)$:

$$\begin{aligned} \text{Pr}_Z(z) &= \Pr_Z(Z \leq z) \\ &= \Pr_Z(g(\mathcal{E}) \leq z) \\ &= \Pr_{\mathcal{E}}(\mathcal{E} \leq g^{-1}(z)) \\ &= \text{Pr}_{\mathcal{E}}(\mathcal{E} \leq \epsilon) \end{aligned} \quad (75)$$

$$\begin{aligned} p_Z(z) &= \frac{d \text{Pr}_Z(z)}{dz} = \frac{d \text{Pr}_{\mathcal{E}}(\epsilon)}{dz} \\ &= \frac{d \text{Pr}_{\mathcal{E}}(\epsilon)}{d\epsilon} \frac{d\epsilon}{dz} \\ &= p_\epsilon(\epsilon) \frac{d\epsilon}{dz} \\ &= p_\epsilon(\epsilon) \left| \frac{d\epsilon}{dz} \right| \quad \text{because of probability} \\ &\equiv p_\epsilon(\epsilon) \left| \det \frac{\partial \epsilon}{\partial z} \right| \quad \text{for multi-dimension} \end{aligned}$$

when computing relationship between $d\epsilon$ and dz :

$$\left| \det \frac{\partial \epsilon}{\partial z} \right| dz = d\epsilon \quad (76)$$

when computing expectation:

$$\begin{aligned} \mathbb{E}_{q(z)}[f(z)] &= \int_z f(z) q(z) dz \\ &= \int_z f(z) \left(\left| \det \frac{\partial \epsilon}{\partial z} \right| q_\epsilon(\epsilon) \right) dz \\ &= \int_\epsilon f(g(\epsilon)) q_\epsilon(\epsilon) \left(\left| \det \frac{\partial \epsilon}{\partial z} \right| dz \right) \\ &= \int_\epsilon f(g(\epsilon)) q_\epsilon(\epsilon) d\epsilon \end{aligned} \quad (77)$$

looking at the relationship:

$$\int_z f(z)q(z)dz = \int_\epsilon f(g(\epsilon))q_\epsilon(\epsilon)d\epsilon \quad (78)$$

one may find that the Jacobian does not appear in the above equation. Therefore, as long as one may **establish** relationship that:

$$q(\epsilon) = \left| \det \frac{\partial z}{\partial \epsilon} \right| q(z)$$

then, we can use the Eq.(77) for that. Note that without $q(z) = \left| \det \frac{\partial \epsilon}{\partial z} \right| q(\epsilon)$ being established then equation Eq.(77) can not be true, as:

$$\begin{aligned} \mathbb{E}_{q(z)}[f(z)] &= \int_z f(z)q(z)dz \\ &\neq \int_z f(z) \left(\left| \det \frac{\partial \epsilon}{\partial z} \right| q_\epsilon(\epsilon) \right) dz \\ \implies \mathbb{E}_{q(z)}[f(z)] &\neq \int_\epsilon f(g(\epsilon))q_\epsilon(\epsilon)d\epsilon \end{aligned} \quad (79)$$

For example, remember in Normalizing Flow, knowing the relationship between $p(\mathbf{z}_k) = \left| \det \frac{\partial \mathbf{z}_k}{\partial \mathbf{z}_0} \right| p(\mathbf{z}_0)$ we can compute the expectation using samples from \mathbf{z}_0 only:

$$\begin{aligned} \mathbb{E}_{p(\mathbf{z}_K)}[h(\mathbf{z}_K)] &= \int_{\mathbf{z}_K} h(\mathbf{z}_K)p(\mathbf{z}_K)d\mathbf{z}_K \\ &= \int_{\mathbf{z}_0} h(f_K \circ \dots \circ f_2 \circ f_1(\mathbf{z}_0))p(\mathbf{z}_0)d\mathbf{z}_0 \\ &= \mathbb{E}_{p(\mathbf{z}_0)}[h(f_K \circ \dots \circ f_2 \circ f_1(\mathbf{z}_0))] \end{aligned} \quad (80)$$

Another example is the well known Gaussian case, we let:

$$p_\epsilon(\epsilon) = \frac{1}{\sqrt{2\pi}} \exp^{-\epsilon^2/2} \equiv \mathcal{N}(\epsilon; 0, 1) \quad (81)$$

and letting:

$$\begin{aligned} g(\theta, \epsilon) &= \mu + \sigma \epsilon \\ \implies \epsilon &= g^{-1}(z, \theta) = \frac{z - \mu}{\sigma} \\ \implies \left| \frac{d\epsilon}{dz} \right| &= \sigma^{-1} \end{aligned} \quad (82)$$

substituting formula:

$$\begin{aligned}
p_Z(z) &= \left| \frac{d\epsilon}{dz} \right| p_\epsilon(\epsilon) \\
&= \sigma^{-1} p_\epsilon(g^{-1}(\theta, z)) \\
&= \sigma^{-1} \frac{1}{\sqrt{2\pi}} \exp^{-\left(\frac{z-\mu}{\sigma}\right)^2/2} \\
&= \frac{1}{\sqrt{2\pi}\sigma} \exp^{-\frac{1}{2}\left(\frac{z-\mu}{\sigma}\right)^2}
\end{aligned} \tag{83}$$

7.1.2 discrete case

$$p_Y(Y = y) = \sum_{x \in g^{-1}(y)} p_X(x) \tag{84}$$

apply to Gumbel-max case:

$$\begin{aligned}
p_K(K = k) &= \int_{\{z_i\} : \arg \max_i (\{z_i + \phi_i\}) = k} \text{ga}_Z(\{z_i\}) \\
&= \frac{\exp(\phi_k)}{\sum_i \exp(\phi_i)}
\end{aligned} \tag{85}$$

it's actually simpler, as there is no Jacobian stuff

7.2 independent case

coming back to our focus of re-parameterization in ∇_θ :

$$\begin{aligned}
\nabla_\theta \mathbb{E}_{q(z;\theta)}[f(z)] &= \nabla_\theta \int_z \textcolor{red}{f}(z) \textcolor{blue}{q}(z; \theta) dz \\
&= \nabla_\theta \int_z \textcolor{red}{f}(z) \left| \det \frac{\partial \epsilon}{\partial z} \right| \textcolor{blue}{q}_\epsilon(\epsilon) dz \\
&= \int_z \textcolor{blue}{q}_\epsilon(\epsilon) \nabla_\theta \textcolor{brown}{f}(z) \left| \det \frac{\partial \epsilon}{\partial z} \right| dz \\
&= \int_z \underbrace{\textcolor{blue}{q}_\epsilon(g^{-1}(z)) \nabla_\theta \textcolor{brown}{f}(z)}_{\phi(z)} \left| \det \frac{\partial \epsilon}{\partial z} \right| dz
\end{aligned} \tag{86}$$

using generic change of variable in integral, and assuming:

$$\int_{z \in \mathcal{Z}} \phi(T(z)) \left| \det \frac{\partial \epsilon}{\partial z} \right| dz = \int_{\epsilon \in T(\mathcal{Z})} \phi(\epsilon) d\epsilon \tag{87}$$

let $T(z) \rightarrow g^{-1}(z)$, the RHS of Eq.(86):

$$\begin{aligned}
\nabla_\theta \mathbb{E}_{q(z;\theta)}[f(z)] &= \nabla_\theta \int_z \textcolor{red}{f}(z) \textcolor{blue}{q}(z; \theta) dz \\
&= \int_\epsilon \underbrace{\textcolor{blue}{q}_\epsilon(\epsilon) \nabla_\theta \textcolor{brown}{f}(g(\epsilon, \theta))}_{\phi(\epsilon)} d\epsilon \quad \text{RHS of Eq.(87)}
\end{aligned} \tag{88}$$

it does **not** work if one was attempt to change $q_\epsilon(\epsilon) \rightarrow q(z; \theta)$ back, i.e.,:

$$\int_z q(z; \theta) \nabla_\theta f(z) dz \quad (89)$$

if you can do this, then ∇_θ has to be put in front of $q(z; \theta)$ to look like Eq.(86). Otherwise, conducting this “virtual” change of variable will make:

$$\int_z q(z; \theta) \nabla_\theta f(z) dz \neq \int_z \nabla_\theta f(z) q(z; \theta) dz \quad (90)$$

challenging conventional thinking, one usually think re-parameterization works by $\epsilon \sim q_\epsilon(\epsilon)$ first; however, one can replace this by:

$$z \sim q(z; \theta) \quad \epsilon = g^{-1}(z, \theta) \quad (91)$$

one way to think about this is that the θ effect by $q(z; \theta)$ is “counteracted” by $g^{-1}(z, \theta)$, such as ϵ independent of θ

7.3 “weakly” dependent case

now look at a situation where ϵ and θ are “weakly” dependent, this means both ϵ and z are dependant on θ :

$$q_\epsilon(\epsilon; \theta) \longleftrightarrow q(z; \theta) \quad (92)$$

$$\begin{aligned} \nabla_\theta \mathbb{E}_{q(z; \theta)}[f(z)] &= \nabla_\theta \int_z f(z) q(z; \theta) dz \\ &= \int_\epsilon \nabla_\theta \left(\underbrace{f(g(\epsilon, \theta))}_u \underbrace{q(\epsilon; \theta)}_v \right) d\epsilon \quad \text{re-param } \epsilon \text{ using Eq.(88), but placing everything after } \nabla_\theta \\ &= \int_\epsilon q_\epsilon(\epsilon; \theta) \nabla_\theta f(g(\epsilon, \theta)) d\epsilon + \int_\epsilon f(g(\epsilon, \theta)) \nabla_\theta q_\epsilon(\epsilon; \theta) d\epsilon \\ &= \underbrace{\int_\epsilon q_\epsilon(\epsilon; \theta) \nabla_\theta f(g(\epsilon, \theta)) d\epsilon}_{g^{\text{rep}}} + \underbrace{\int_\epsilon f(g(\epsilon, \theta)) \underline{q_\epsilon(\epsilon; \theta)} \nabla_\theta \log q_\epsilon(\epsilon; \theta) d\epsilon}_{g^{\text{corr}}} \quad \text{Reinforce trick} \end{aligned} \quad (93)$$

7.3.1 look at g^{rep} :

$$\begin{aligned}
g^{\text{rep}} &= \int_{\epsilon} q_{\epsilon}(\epsilon, \theta) \nabla_{\theta} f(g(\epsilon, \theta)) d\epsilon \\
&= \int_{\epsilon} \left| \det \frac{\partial z}{\partial \epsilon} \right| q(z; \theta) \nabla_{\theta} f(g(\epsilon, \theta)) d\epsilon \quad \text{using Eq.(74)} \\
&= \int_{\epsilon} p \left| \det \frac{\partial z}{\partial \epsilon} \right| q(z; \theta) \nabla_z f(z) \Big|_{z=g(\epsilon, \theta)} \nabla_{\theta} g(\epsilon, \theta) d\epsilon \quad \text{using chain rule} \\
&= \int_{\epsilon} q(z; \theta) \nabla_z f(z) h(\epsilon, \theta) \left| \det \frac{\partial z}{\partial \epsilon} \right| d\epsilon \\
&= \int_z q(z; \theta) \nabla_z f(z) \nabla_{\theta} g(g^{-1}(z, \theta), \theta) dz \quad \text{re-parameterization to } z \text{ via Eq.(87)} \\
&= \mathbb{E}_{q(z; \theta)} \left[\nabla_z f(z) \nabla_{\theta} g(g^{-1}(z, \theta), \theta) \right]
\end{aligned} \tag{94}$$

7.3.2 look at g^{corr} :

$$\text{let } u(\epsilon, \theta) \equiv \nabla_{\theta} \log \left| \det \frac{\partial z}{\partial \epsilon} \right| \tag{95}$$

$$\begin{aligned}
g^{\text{corr}} &= \int_{\epsilon} q_{\epsilon}(\epsilon; \theta) f(g(\epsilon, \theta)) \nabla_{\theta} \log q_{\epsilon}(\epsilon; \theta) d\epsilon \\
&= \int_{\epsilon} \left| \det \frac{\partial z}{\partial \epsilon} \right| \underbrace{q(z; \theta) f(z)} \nabla_{\theta} \log \left(\left| \det \frac{\partial z}{\partial \epsilon} \right| q(z; \theta) \right) d\epsilon \quad \text{using Eq.(74)} \tag{96} \\
&= \int_{\epsilon} \underbrace{q(z; \theta) f(z)} \left(\nabla_{\theta} \log \left| \det \frac{\partial z}{\partial \epsilon} \right| + \nabla_{\theta} \log q(g(\epsilon, \theta); \theta) \right) \left| \det \frac{\partial z}{\partial \epsilon} \right| d\epsilon
\end{aligned}$$

looking at $\nabla_{\theta} \log q(g(\epsilon, \theta); \theta)$, there are two path of derivatives for computing ∇_{θ} :

$$\begin{aligned}
\nabla_{\theta} \log q(\underbrace{g(\epsilon, \theta)}_1; \underbrace{\theta}_2) &= \underbrace{\nabla_z \log q(z; \theta) \nabla_{\theta} g(\epsilon; \theta)}_1 + \underbrace{\nabla_{\theta} \log q(z; \theta)}_2 \tag{97} \\
&= \int_z q(z; \theta) f(z) \left(\nabla_{\theta} \log \left| \det \frac{\partial z}{\partial \epsilon} \right| + \nabla_z \log q(z; \theta) \nabla_{\theta} g(\epsilon; \theta) + \nabla_{\theta} \log q(z; \theta) \right) dz \quad \text{re-parameterize to } z \\
&= \int_z q(z; \theta) f(z) \left(u(\epsilon, \theta) + \nabla_z \log q(z; \theta) \nabla_{\theta} g(\epsilon, \theta) + \nabla_{\theta} \log q(z; \theta) \right) dz \\
&= \int_z q(z; \theta) f(z) \left(u(g^{-1}(z, \theta), \theta) + \nabla_z \log q(z; \theta) \nabla_{\theta} g(g^{-1}(z, \theta), \theta) + \nabla_{\theta} \log q(z; \theta) \right) dz \quad \text{replace all } \epsilon \rightarrow g^{-1}(\epsilon, \theta) \\
&= \mathbb{E}_{q(z; \theta)} \left[f(z) \left(u(g^{-1}(z, \theta), \theta) + \nabla_z \log q(z; \theta) \nabla_{\theta} g(g^{-1}(z, \theta), \theta) + \nabla_{\theta} \log q(z; \theta) \right) \right] \tag{98}
\end{aligned}$$

7.3.3 to compute weakly-dependent re-parameterization

to summarize:

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{q(z;\theta)}[f(z)] &= g^{\text{rep}} + g^{\text{corr}} \\ &= \mathbb{E}_{q(z;\theta)} \left[\nabla_z f(z) \nabla_{\theta} g(g^{-1}(z, \theta), \theta) \right] \\ &\quad + \mathbb{E}_{q(z;\theta)} \left[f(z) \left(u(g^{-1}(z, \theta), \theta) + \nabla_z \log q(z; \theta) \nabla_{\theta} g(g^{-1}(z, \theta), \theta) + \nabla_{\theta} \log q(z; \theta) \right) \right]\end{aligned}\tag{99}$$

therefore, it seems that to re-parameterize any distribution $q(z; \theta)$, one needs to first compute:

1. $\nabla_{\theta} g(\epsilon, \theta) = \nabla_{\theta} g(g^{-1}(z, \theta), \theta)$
2. $u(\epsilon, \theta) \equiv \nabla_{\theta} \log \left| \det \frac{\partial z}{\partial \epsilon} \right|$

7.4 Example: Dirichlet Distribution

looking at definition of Dirichlet distribution:

$$q(\mathbf{z}; \boldsymbol{\alpha}) \equiv \text{Dir}(\boldsymbol{\alpha}) \text{ with } \boldsymbol{\alpha} \equiv [\alpha_1, \dots, \alpha_k], \text{ and } \boldsymbol{\alpha}_0 = \sum_i \alpha_i$$

$$\begin{aligned}\mathbf{z} &= g(\boldsymbol{\epsilon}, \boldsymbol{\alpha}) = \exp \left(\boldsymbol{\Sigma}^{\frac{1}{2}}(\boldsymbol{\alpha}) \boldsymbol{\epsilon} + \boldsymbol{\mu}(\boldsymbol{\alpha}) \right) \\ &= g(\boldsymbol{\epsilon}, \boldsymbol{\alpha}) = \exp \left(\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\epsilon} + \boldsymbol{\mu} \right) \quad \text{for short}\end{aligned}\tag{100}$$

and

$$\boldsymbol{\mu} = \mathbb{E}_{q(\mathbf{z}; \boldsymbol{\alpha})}[\log(\mathbf{z})] = \begin{bmatrix} \psi(\alpha_1) - \psi(\boldsymbol{\alpha}_0) \\ \vdots \\ \psi(\alpha_k) - \psi(\boldsymbol{\alpha}_0) \end{bmatrix}\tag{101}$$

looking at individual element of co-variance matrix $\Sigma_{i,j}$

$$\Sigma_{i,j} = \text{Cov}(\log(z_i), \log(z_j)) = \begin{cases} \psi(\alpha_i) - \psi'(\boldsymbol{\alpha}_0) & \text{if } i = j \\ -\psi'(\boldsymbol{\alpha}_0) & \text{otherwise} \end{cases}\tag{102}$$

the whole $\boldsymbol{\Sigma}$:

$$\boldsymbol{\Sigma} = \begin{bmatrix} \psi'(\alpha_1) & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & 0 & \psi'(\alpha_k) \end{bmatrix} - \psi'(\boldsymbol{\alpha}_0) \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}\tag{103}$$

by construction, $\boldsymbol{\Sigma}$ is positive definite, hence:

$$\begin{aligned}\boldsymbol{\Sigma} &= \mathbf{V} \mathbf{D} \mathbf{V}^{\top} \\ \implies \boldsymbol{\Sigma}^{\frac{1}{2}} &= \mathbf{V} \mathbf{D}^{\frac{1}{2}} \mathbf{V}^{\top}\end{aligned}\tag{104}$$

let's look at:

$$\begin{aligned}
J &= \frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} = \frac{\partial \exp(\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\epsilon} + \boldsymbol{\mu})}{\partial \boldsymbol{\epsilon}} \\
&= \frac{\partial \left(\begin{bmatrix} z_1 \\ \vdots \\ z_k \end{bmatrix} \right)}{\partial \boldsymbol{\epsilon}} = \frac{\partial \exp \left(\begin{bmatrix} \boldsymbol{\Sigma}^{\frac{1}{2}}_{1,:} \boldsymbol{\epsilon} + \mu_1 \\ \vdots \\ \boldsymbol{\Sigma}^{\frac{1}{2}}_{k,:} \boldsymbol{\epsilon} + \mu_k \end{bmatrix} \right)}{\partial \boldsymbol{\epsilon}} \\
\Rightarrow J_{i,j} &= \Sigma_{i,j}^{1/2} \exp(\boldsymbol{\Sigma}^{\frac{1}{2}}_{i,:} \boldsymbol{\epsilon} + \mu_i) \\
&= \Sigma_{i,j}^{1/2} z_i
\end{aligned} \tag{105}$$

it is then obvious that the Jacobian J looks like:

$$\begin{aligned}
\mathbf{J} &= \begin{bmatrix} \Sigma_{1,1}^{1/2} z_1 & \dots & \Sigma_{1,k}^{1/2} z_1 \\ \vdots & & \ddots \\ \Sigma_{k,1}^{1/2} z_k & \dots & \Sigma_{k,k}^{1/2} z_k \end{bmatrix} = \underbrace{\begin{bmatrix} \Sigma_{1,1}^{1/2} z_1 & \dots & \Sigma_{1,k}^{1/2} z_1 \\ \vdots & & \ddots \\ \Sigma_{1,k}^{1/2} z_k & \dots & \Sigma_{k,k}^{1/2} z_k \end{bmatrix}}_{\text{co-variance matrix is symmetric}} = \boldsymbol{\Sigma}^{\frac{1}{2}} \begin{bmatrix} z_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & z_k \end{bmatrix} \\
\end{aligned} \tag{106}$$

substitute J in:

$$\begin{aligned}
\left| \det \frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right| &= \left| \det J \right| = \left| \det \left(\boldsymbol{\Sigma}^{\frac{1}{2}} \begin{bmatrix} z_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & z_k \end{bmatrix} \right) \right| \\
&= \det(\boldsymbol{\Sigma}^{\frac{1}{2}}) \prod_{i=1}^k z_i \quad \text{drop absolute operator as all positive}
\end{aligned} \tag{107}$$

from definition of Dirichlet distribution:

$$\begin{aligned}
\text{Dir}(\mathbf{z}; \boldsymbol{\alpha}) &= \frac{\Gamma(\boldsymbol{\alpha}_0)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K z_i^{\alpha_i - 1} \\
\Rightarrow \log(\text{Dir}(\mathbf{z}; \boldsymbol{\alpha})) &= \log(\Gamma(\boldsymbol{\alpha}_0)) - \sum_i \log \Gamma(\alpha_i) + \sum_{i=1}^k (\alpha_i - 1) \log(z_i)
\end{aligned} \tag{108}$$

derivative of $\log q(\mathbf{z}; \boldsymbol{\alpha})$ is:

$$\begin{aligned}
\frac{\partial}{\partial z_i} \log q(\mathbf{z}; \boldsymbol{\alpha}) &= \frac{\alpha_i - 1}{z_i} \\
\frac{\partial}{\partial \alpha_i} \log q(\mathbf{z}; \boldsymbol{\alpha}) &= \frac{\partial \log \Gamma(\boldsymbol{\alpha}_0)}{\partial \alpha_i} + \frac{\partial \log \Gamma(\alpha_i)}{\partial \alpha_i} + \frac{\partial \alpha_i \log(z_i)}{\partial \alpha_i} \\
&= \psi(\boldsymbol{\alpha}_0) \frac{\partial \boldsymbol{\alpha}_0}{\partial \alpha_i} + \psi(\alpha_i) + \log(z_i) \\
&= \psi(\boldsymbol{\alpha}_0) + \psi(\alpha_i) + \log(z_i)
\end{aligned} \tag{109}$$

7.4.1 Finding $\nabla_{\theta} g(\epsilon, \theta)$ and $u(\epsilon, \theta) \equiv \nabla_{\theta} \log \left| \det \frac{\partial \mathbf{z}}{\partial \epsilon} \right|$ for Dirichlet

and be reminded again that:

$$\mathbf{z} = g(\epsilon, \alpha) = \exp \left(\Sigma^{\frac{1}{2}}(\alpha) \epsilon + \mu(\alpha) \right) \quad (110)$$

for $\nabla_{\theta} g(\epsilon, \theta)$, looking at individual Jacobian matrix: $\frac{\partial z_i}{\alpha_j}$:

$$\begin{aligned} \frac{\partial z_i}{\alpha_j} &= \frac{\exp \left(\Sigma^{\frac{1}{2}}_{i,:} \epsilon + \mu_i \right)}{\partial \left(\Sigma^{\frac{1}{2}}_{i,:} \epsilon + \mu_i \right)} \frac{\partial \left(\Sigma^{\frac{1}{2}}_{i,:} \epsilon + \mu_i \right)}{\partial \alpha_j} \\ &= z_i \frac{\partial \left(\Sigma^{\frac{1}{2}}_{i,:} \epsilon + \mu_i \right)}{\partial \alpha_j} \\ &= z_i(\epsilon, \alpha) \frac{\partial \left(\Sigma^{\frac{1}{2}}_{i,:}(\alpha) \epsilon + \mu_i(\alpha) \right)}{\partial \alpha_j} \quad \text{add } \alpha \\ &= z_i(\epsilon, \alpha) \frac{\partial \Sigma^{\frac{1}{2}}_{i,:}(\alpha)}{\partial \alpha_j} \epsilon + \frac{\partial \mu_i(\alpha)}{\partial \alpha_j} \end{aligned} \quad (111)$$

fill the whole Jacobian: $\frac{\partial \mathbf{z}}{\alpha}$:

$$\begin{aligned} \nabla_{\theta} g(\epsilon, \theta) &= \begin{bmatrix} \frac{\partial z_1}{\partial \alpha_1} & \cdots & \frac{\partial z_1}{\partial \alpha_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_k}{\partial \alpha_1} & \cdots & \frac{\partial z_k}{\partial \alpha_k} \end{bmatrix} \\ &= \begin{bmatrix} z_1(\epsilon, \alpha) \frac{\partial \Sigma^{\frac{1}{2}}_{1,:}(\alpha)}{\partial \alpha_1} \epsilon + \frac{\partial \mu_1(\alpha)}{\partial \alpha_1} & \cdots & z_1(\epsilon, \alpha) \frac{\partial \Sigma^{\frac{1}{2}}_{1,:}(\alpha)}{\partial \alpha_k} \epsilon + \frac{\partial \mu_1(\alpha)}{\partial \alpha_k} \\ \vdots & \ddots & \vdots \\ z_k(\epsilon, \alpha) \frac{\partial \Sigma^{\frac{1}{2}}_{k,:}(\alpha)}{\partial \alpha_1} \epsilon + \frac{\partial \mu_k(\alpha)}{\partial \alpha_1} & \cdots & z_k(\epsilon, \alpha) \frac{\partial \Sigma^{\frac{1}{2}}_{k,:}(\alpha)}{\partial \alpha_k} \epsilon + \frac{\partial \mu_k(\alpha)}{\partial \alpha_k} \end{bmatrix} \end{aligned} \quad (112)$$

to find $u(\epsilon, \theta) \equiv \nabla_{\theta} \log \left| \det \frac{\partial \mathbf{z}}{\partial \epsilon} \right|$, since we know

$$\begin{aligned} \left| \det \frac{\partial \mathbf{z}}{\partial \epsilon} \right| &= \det \left(\Sigma^{\frac{1}{2}} \right) \prod_{i=1}^k z_i \\ \implies \nabla_{\alpha} \log \left| \det \frac{\partial \mathbf{z}}{\partial \epsilon} \right| &= \nabla_{\alpha} \log \left(\det \left(\Sigma^{\frac{1}{2}} \right) \prod_{i=1}^k \exp \left(\Sigma^{\frac{1}{2}}_{i,:} \epsilon + \mu_i \right) \right) \\ &= \nabla_{\alpha} \left[\log \det \left(\Sigma^{\frac{1}{2}} \right) + \sum_{i=1}^k \Sigma^{\frac{1}{2}}_{i,:} \epsilon + \mu_i \right] \end{aligned} \quad (113)$$

note that $\nabla_{\alpha} \log \left| \det \frac{\partial \mathbf{z}}{\partial \epsilon} \right|$ is a vector, not a matrix:

$$\nabla_{\alpha} \log \left| \det \frac{\partial \mathbf{z}}{\partial \epsilon} \right| = \begin{bmatrix} \frac{\partial \log \det(\Sigma^{\frac{1}{2}})}{\partial \alpha_1} + \frac{\partial \sum_{i=1}^k \Sigma^{\frac{1}{2}}_{i,:}(\alpha)_{1,:} \epsilon + \mu_i(\alpha)}{\partial \alpha_1} \\ \vdots \\ \frac{\partial \log \det(\Sigma^{\frac{1}{2}})}{\partial \alpha_k} + \frac{\partial \sum_{i=1}^k \Sigma^{\frac{1}{2}}_{i,:}(\alpha)_{k,:} \epsilon + \mu_i(\alpha)}{\partial \alpha_k} \end{bmatrix} \quad (114)$$

so the only thing remain is to compute in $\nabla_{\theta} g(\epsilon, \theta)$ and $\nabla_{\alpha} \log \left| \det \frac{\partial \mathbf{z}}{\partial \epsilon} \right|$ is:

$$1. \frac{\partial \boldsymbol{\mu}}{\partial \alpha_i}: \text{ since } \boldsymbol{\mu} = \begin{bmatrix} \psi(\alpha_1) - \psi(\boldsymbol{\alpha}_0) \\ \vdots \\ \psi(\alpha_k) - \psi(\boldsymbol{\alpha}_0) \end{bmatrix}, \text{ therefore:}$$

$$\frac{\partial \boldsymbol{\mu}}{\partial \alpha_i} = \begin{bmatrix} -\frac{\partial \psi(\boldsymbol{\alpha}_0)}{\partial \alpha_i} \\ \vdots \\ \frac{\psi(\alpha_i)}{\partial \alpha_i} - \frac{\partial \psi(\boldsymbol{\alpha}_0)}{\partial \alpha_i} \\ \vdots \\ -\frac{\partial \psi(\boldsymbol{\alpha}_0)}{\partial \alpha_i} \end{bmatrix} = \begin{bmatrix} -\frac{\partial \psi(\boldsymbol{\alpha}_0)}{\partial \boldsymbol{\alpha}_0} \frac{\partial \boldsymbol{\alpha}_0}{\partial \alpha_i} \\ \vdots \\ \frac{\psi(\alpha_i)}{\partial \alpha_i} - \frac{\partial \psi(\boldsymbol{\alpha}_0)}{\partial \boldsymbol{\alpha}_0} \frac{\partial \boldsymbol{\alpha}_0}{\partial \alpha_i} \\ \vdots \\ -\frac{\partial \psi(\boldsymbol{\alpha}_0)}{\partial \boldsymbol{\alpha}_0} \frac{\partial \boldsymbol{\alpha}_0}{\partial \alpha_i} \end{bmatrix} = \begin{bmatrix} -\psi'(\boldsymbol{\alpha}_0) \\ \vdots \\ \psi'(\alpha_i) - \psi'(\boldsymbol{\alpha}_0) \\ \vdots \\ -\psi'(\boldsymbol{\alpha}_0) \end{bmatrix} \quad (115)$$

2. $\frac{\partial \boldsymbol{\Sigma}^{\frac{1}{2}}}{\partial \alpha_i}$, this is hard, and from the original paper stating, $\frac{\partial \boldsymbol{\Sigma}^{\frac{1}{2}}}{\partial \alpha_i}$ is the solution of Lyapunov equation:

$$\frac{\partial \boldsymbol{\Sigma}}{\partial \alpha_i} = \frac{\partial \boldsymbol{\Sigma}^{\frac{1}{2}}}{\partial \alpha_i} \boldsymbol{\Sigma}^{\frac{1}{2}} + \boldsymbol{\Sigma}^{\frac{1}{2}} \frac{\partial \boldsymbol{\Sigma}^{\frac{1}{2}}}{\partial \alpha_i} \quad (116)$$

where, by change $\boldsymbol{\Sigma} \rightarrow \frac{\partial \boldsymbol{\Sigma}}{\partial \alpha_i}$ in Eq.(103):

$$\frac{\partial \boldsymbol{\Sigma}}{\partial \alpha_i} = \begin{bmatrix} \psi''(\alpha_1) & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & 0 & \psi''(\alpha_k) \end{bmatrix} - \psi''(\boldsymbol{\alpha}_0) \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix} \quad (117)$$

3. $\frac{\partial \log \det(\boldsymbol{\Sigma}^{\frac{1}{2}})}{\partial \alpha_i}$:

this is easy and follow matrix cookbook, one may obtain:

$$\frac{\partial \log \det(\boldsymbol{\Sigma}^{\frac{1}{2}})}{\partial \alpha_i} = \text{trace} \left(\boldsymbol{\Sigma}^{\frac{1}{2}} \frac{\partial \boldsymbol{\Sigma}^{\frac{1}{2}}}{\partial \alpha_i} \right) \quad (118)$$

7.5 Apply to Variational equation

$$\begin{aligned} \mathcal{L}(\theta) &= \mathbb{E}_{q(z|\theta)} [\log p(x, z) - \log q(z; \theta)] \\ &= \mathbb{E}_{q(z|\theta)} [\underbrace{\log p(x, z)}_{f(z)}] - \underbrace{\mathbb{E}_{q(z;\theta)} [\log q(z; \theta)]}_{H[q(z;\theta)]} \\ &= \mathbb{E}_{q(z|\theta)} [f(z)] + H[q(z; \theta)] \\ \nabla_{\theta} \mathcal{L}(\theta) &= \nabla_{\theta} \mathbb{E}_{q(z;\theta)} [f(z)] + \nabla_{\theta} H[q(z; \theta)] \end{aligned} \quad (119)$$

$$\begin{aligned}
\nabla_{\theta} \mathcal{L}(\theta) &= \nabla_{\theta} \mathbb{E}_{q(z|\theta)} [f(z)] + \nabla_{\theta} H[q(z|\theta)] \\
&= g^{\text{rep}} + g^{\text{corr}} + \nabla_{\theta} H[q(z|\theta)] \\
&= \mathbb{E}_{q(z|\theta)} [\nabla_z f(z) h(g^{-1}(z, \theta), \theta)] \\
&\quad + \mathbb{E}_{q(z|\theta)} [f(z) (u(g^{-1}(z, \theta), \theta) + \nabla_z \log q(z|\theta) h(g^{-1}(z, \theta), \theta))] \\
&\quad + \nabla_{\theta} H[q(z|\theta)]
\end{aligned} \tag{120}$$

special case 1 : if distribution $q(\epsilon, \theta)$ does not depend on the variational parameters θ , then:

$$\begin{aligned}
\nabla_{\theta} \log q_{\epsilon}(\epsilon, \theta) &= 0 \\
\implies \nabla_{\theta} \mathcal{L}(\theta) &= \int q_{\epsilon}(\epsilon, \theta) \nabla_{\theta} f(g(\epsilon, \theta)) d\epsilon + \nabla_{\theta} H[q(z|\theta)]
\end{aligned} \tag{121}$$

special case 2 : if $g(\epsilon, \theta) = \epsilon$:

$$\begin{aligned}
g(\epsilon, \theta) = \epsilon &\implies \nabla_{\theta} g(\epsilon, \theta) = \nabla_{\theta} \epsilon = 0 \\
\implies u(\epsilon, \theta) &\equiv \nabla_{\theta} \log |\det \nabla_{\epsilon} g(\epsilon, \theta)| = \nabla_{\theta} \log |\det \nabla_{\epsilon} \epsilon| = 0
\end{aligned} \tag{122}$$

$$\begin{aligned}
\implies \nabla_{\theta} \mathcal{L}(\theta) &= g^{\text{rep}} + g^{\text{corr}} + H[q(z|\theta)] \\
&= \mathbb{E}_{q(z|\theta)} [\nabla_z f(z) \underbrace{\nabla_{\theta} g(g^{-1}(z, \theta), \theta)}_{=0}] \\
&\quad + \mathbb{E}_{q(z;\theta)} [f(z) (\underbrace{u(g^{-1}(z, \theta), \theta)}_{=0} + \nabla_z \log q(z; \theta) \underbrace{\nabla_{\theta} g(g^{-1}(z, \theta), \theta)}_{=0} + \nabla_{\theta} \log q(z; \theta))] \\
&\quad + \nabla_{\theta} H[q(z|\theta)] \\
&= \mathbb{E}_{q(z|\theta)} [f(z) \nabla_{\theta} \log q(z; \theta)] + \nabla_{\theta} H[q(z|\theta)]
\end{aligned} \tag{123}$$