1 Analysis of Local Expectation Gradients (LEG) [1]

We provide an analysis of the LEG estimator including its computational requirements.

Proposition 1.1. (LEG Trick and Estimator) Given a factorised distribution $p(\mathbf{X}) = \prod_{d=1}^{D} p(X_d | \mathbf{X}_{\leq d})$, Local Expectation Gradients are based on the following equalities: $\mathbb{E}_{\mathbf{X} \sim p(\mathbf{X})} [f(\mathbf{X})\partial_{\lambda} \log p(\mathbf{X})] =$

$$= \sum_{d=1}^{D} \mathbb{E}_{\mathbf{X}_{\neq d} \sim p(\mathbf{X}_{\neq d})} \left\{ \mathbb{E}_{X_{d} \sim p(X_{d} | \mathbf{MB}(X_{d}))} \left[f(\mathbf{X}) \partial_{\lambda} \log p(X_{d} | \mathbf{X}_{< d}) \right] \right\} \quad \text{cf. Eq. 8 in [1]}$$

$$= \sum_{d=1}^{D} \underbrace{\mathbb{E}_{\mathbf{X}_{d} \sim p(\mathbf{X}_{>d} | \mathbf{X}'_{d}, \mathbf{X}_{< d})}}_{\text{Used for sampling pivots}} \left\{ \underbrace{\mathbb{E}_{X_{d} \sim p(X_{d} | \mathbf{X}_{\neq d})}}_{\text{Used for weighted average}} \left[\underbrace{f(\mathbf{X}) \partial_{\lambda} \log p(X_{d} | \mathbf{X}_{< d})}_{\text{Used as score evaluation}} \right] \right\}$$

where $MB(X_d)$ is the Markov blanket of variable X_d , the red terms introduce auxiliary variables for each $d \in \{1, ..., D\}$ and $p(X_d \mid MB(X_d)) = p(X_d \mid \mathbf{X}_{\neq d}) = \frac{p(X_d \mid \mathbf{X}_{< d}) \prod_{j > d} p(X_j \mid \mathbf{X}_{< j})}{\sum_{x_\delta \in \Omega(X_d)} p(x_\delta \mid \mathbf{X}_{< d}) \prod_{j > d} p(X_j \mid \mathbf{X}_{< j})}$ is a weighting distribution.

The Monte Carlo estimate corresponding to the LEG Trick in Eq. 1.1 (cf. Algorithm 1 in [1]) given N pivot samples is given by the following expression:

$$\mathbb{E}_{\mathbf{X} \sim p(\mathbf{X})} \left[f(\mathbf{X}) \partial_{\lambda} \log p(\mathbf{X}) \right] \approx \sum_{d=1}^{D} \frac{1}{N} \sum_{n=1}^{N} \sum_{x_{\delta} \in \Omega(X_d)} p(x_{\delta} \mid \mathbf{x}_{\neq d}^{(n)}) f(\mathbf{x}_{\neq d}^{(n)}, x_{\delta}) \partial_{\lambda} \log p(x_{\delta} \mid \mathbf{x}_{< d}^{(n)})$$

Additionally, the computational complexity is $O(N \cdot D^2 \cdot K)$, with $K = \max_d (|\Omega(X_d)|)$.

Proof. We start by recalling the proof for the first equality in Eq. 1.1.

$$\mathbb{E}_{\mathbf{X}\sim p(\mathbf{X})} \left[f(\mathbf{X})\partial_{\lambda}\log p(\mathbf{X}) \right] = \mathbb{E}_{\mathbf{X}\sim p(\mathbf{X})} \left[f(\mathbf{X})\partial_{\lambda}\log\prod_{d=1}^{D} p(X_{d} \mid \mathbf{X}_{< d}) \right]$$
$$= \sum_{d=1}^{D} \mathbb{E}_{\mathbf{X}\sim p(\mathbf{X})} \left[f(\mathbf{X})\partial_{\lambda}\log p(X_{d} \mid \mathbf{X}_{< d}) \right]$$
$$= \sum_{d=1}^{D} \mathbb{E}_{\mathbf{X}\neq d}\sim p(\mathbf{X}_{\neq d}) \left\{ \underbrace{\mathbb{E}_{X_{d}\sim p(X_{d}\mid\mathbf{X}\neq d)} \left[f(\mathbf{X})\partial_{\lambda}\log p(X_{d}\mid \mathbf{X}_{< d}) \right] \right\}.$$
$$:= g(\mathbf{X}\neq d)$$
(1.2)

By noting that $p(X_d | \mathbf{X}_{\neq d}) = p(X_d | \mathbf{MB}(X_d))$, we obtain the first result.

Regarding the second equality in Eq. 1.1, we observe that $\mathbb{E}_{\mathbf{X}\neq d\sim p(\mathbf{X}\neq d)}[g(\mathbf{X}\neq d)]$ in Eq. 1.2 can be equivalently rewritten in the following way:

$$\mathbb{E}_{\mathbf{X}_{\neq d} \sim p(\mathbf{X}_{\neq d})} \left[g(\mathbf{X}_{\neq d}) \right] = \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X})} \left[g(\mathbf{X}_{\neq d}) \right]$$
$$= \mathbb{E}_{\mathbf{X}_{< d} \sim p(\mathbf{X}_{< d})} \mathbb{E}_{X'_{d} \sim p(X'_{d} | \mathbf{X}_{< d})} \mathbb{E}_{\mathbf{X}_{> d} \sim p(\mathbf{X}_{> d} | X'_{d}, \mathbf{X}_{< d})} \left[g(\mathbf{X}_{\neq d}) \right] \quad (1.3)$$

By plugging this result into Eq. 1.2 we obtain the desired result for the second equality in Eq. 1.1.

The Monte Carlo estimate trivially follows by sampling N pivots from the first three expectations in Eq. 1.1 and by taking the exact weighted average of score evaluations.

Regarding computational complexity, we observe that computing the weighting function in the Monte Carlo estimate requires O(D) evaluations. Also, the estimate requires to iterate over three summations, thus having a computational requirement of $O(N \cdot D \cdot K)$. By combining these two results, we obtain the overall complexity of $O(N \cdot D^2 \cdot K)$.

References

[1] Titsias, Michalis and Lázaro-Gredilla, Miguel. Local Expectation Gradients for Black Box Variational Inference. *NeurIPS*, 2015.