
On Multilevel Monte Carlo Unbiased Gradient Estimation For Deep Latent Variable Models: Supplementary Material

1 Preliminary Lemmas

In this section we first establish some preliminary results which we use to prove Theorems 2, 3 and 7.

Lemma 1. *The following inequalities hold.*

$$(i) \forall 0 < p \leq 1, \exists c_p > 0 \text{ such that } \forall x \geq 0, \log(1+x) \leq c_p x^p.$$

$$(ii) \forall p \geq 1, \forall j, k > 0, \exists c_{p,j,k} > 0 \text{ such that } \forall x > 0, |\log x|^p \leq c_{p,j,k} \left(x^j + \frac{1}{x^k}\right).$$

Proof. (i) For $p = 1$ we have the standard inequality $\log(1+x) \leq x$. For $0 < p < 1$, we have $\lim_{x \rightarrow \infty} \frac{\log(1+x)}{x^p} = 0$ and $\lim_{x \rightarrow 0^+} \frac{\log(1+x)}{x^p} = \lim_{x \rightarrow 0^+} \frac{1}{1+x} \frac{1}{px^{p-1}} = 0$. By continuity, $\frac{\log(1+x)}{x^p}$ attains its maximum and is bounded above on $(0, \infty)$.

(ii) Similarly, this is due to $|\log x|^p$ increasing slower than x^j as $x \rightarrow \infty$ and slower than x^{-k} as $x \rightarrow 0^+$. More rigorously, we have by standard results $0 \leq \lim_{x \rightarrow \infty} \frac{(\log x)^p}{x^j + x^{-k}} \leq \lim_{x \rightarrow \infty} \frac{(\log x)^p}{x^j} = 0$ and $0 \leq \lim_{x \rightarrow 0^+} \frac{(-\log x)^p}{x^j + x^{-k}} \leq \lim_{x \rightarrow 0^+} \frac{(-\log x)^p}{x^{-k}} = 0$. Hence by continuity $\frac{|\log x|^p}{x^j + x^{-k}}$ attains its maximum and is bounded above on $(0, \infty)$. \square

Lemma 2. *Let X_1, \dots, X_n be i.i.d. zero mean random variables, whose p -th moments are finite for some $p \geq 2$. Then there exists a constant C_p , independent of n and the distribution of X_i , such that*

$$\mathbb{E} \left[\left| \frac{1}{n} \sum_{i=1}^n X_i \right|^p \right] \leq C_p n^{-p/2} \mathbb{E} [|X_i|^p] \quad (1)$$

Proof. This follows from a direct application of Marcinkiewicz-Zygmund inequality then Jensen inequality. \square

Lemma 3. *Let X_1, \dots, X_n be positive i.i.d. random variables, then $\mathbb{E} \left[\left(\frac{1}{n} \sum_{i=1}^n X_i \right)^{-p} \right] \leq \mathbb{E} [X_1^{-p}]$ for all $p \geq 0$.*

Proof. By Jensen's inequality,

$$\left(\frac{1}{n} \sum_{i=1}^n X_i \right)^{-p} \leq \frac{1}{n} \sum_{i=1}^n X_i^{-p}.$$

Taking expectation over both sides yields the result. \square

Lemma 4. *Suppose $\{X_k\}_{k=1}^\infty, \{Y_k\}_{k=1}^\infty$ are two sequences of random variables satisfying $\mathbb{E} [|X_k|^p] = O(2^{-ck})$, $\mathbb{E} [|Y_k|^p] = O(2^{-ck})$, where $p \geq 1$, then $\mathbb{E} [|X_k + Y_k|^p] = O(2^{-ck})$.*

Proof. By the Cp-inequality, we have for $p \geq 1$

$$\begin{aligned} \mathbb{E} [|X_k + Y_k|^p] &\leq 2^{p-1} \cdot (\mathbb{E} [|X_k|^p] + \mathbb{E} [|Y_k|^p]). \\ &= O(2^{-ck}) \end{aligned}$$

\square

2 Proof of Theorem 2

We now give a proof of Theorem 2. We note that concurrently with our work, Goda and Ishikawa (2019); Ishikawa and Goda (2020) independently proposed the estimator $\hat{\ell}^{\text{ML-SS}}(\boldsymbol{\theta})$ and a related theoretical result under slightly different assumptions (requiring finite moments for $\log w(\mathbf{z})$ when $\mathbf{z} \sim q_\phi$).

Theorem. *Assume there exists $\epsilon, \delta > 0$ such that $\mathbb{E}_{q_\phi} [w(\mathbf{z})^{2+\epsilon} + w(\mathbf{z})^{-\delta}] < \infty$. Then $\hat{\ell}^{\text{ML}}(\boldsymbol{\theta})$ satisfies Theorem 1 for $r \in (\frac{1}{2}, 1 - \frac{1}{2^{1+\alpha}})$, where $\alpha = \min(\frac{\epsilon}{2}, 1) > 0$.*

Proof. To prove this Theorem, we are going to check that the assumptions of Theorem 1 hold. We first start proving Assumption (a) and Assumption (c) and then finally Assumption (b) whose proof is more involved.

Assumption (a): We first show that $\hat{\ell}^{(k)}(\boldsymbol{\theta})$ is uniformly bounded in L^2 . Observe that by Lemma 1, $\mathbb{E} [|\hat{\ell}^{(k)}(\boldsymbol{\theta})|^2] = \mathbb{E} \left[\left| \log \left(\frac{1}{k} \sum_{i=1}^k w_i \right) \right|^2 \right] \leq c \cdot \mathbb{E} \left[\left(\frac{1}{k} \sum_{i=1}^k w_i \right) + \left(\frac{1}{k} \sum_{i=1}^k w_i \right)^{-\delta} \right]$. But now $\mathbb{E}[|w_i|] < \infty$, and $\mathbb{E} \left[\left(\frac{1}{k} \sum_{i=1}^k w_i \right)^{-\delta} \right]$ is bounded for all k by Lemma 3. It follows that $\sup_k \mathbb{E} [|\hat{\ell}^{(k)}(\boldsymbol{\theta})|^2] < \infty$. Therefore

$$\begin{aligned} \mathbb{E}[I_0] + \sum_{k=0}^{\infty} \mathbb{E}[\Delta_k] &= \mathbb{E}[I_0] + \sum_{k=0}^{\infty} \mathbb{E} \left[\log \left(\frac{O_k + E_k}{2} \right) - \frac{1}{2} (\log O_k + \log E_k) \right] \\ &= \mathbb{E}[I_0] + \sum_{k=0}^{\infty} \mathbb{E} \left[\log \left(\frac{O_k + E_k}{2} \right) \right] - \frac{1}{2} (\mathbb{E}[\log O_k] + \mathbb{E}[\log E_k]) \\ &= \mathbb{E} [\hat{\ell}^{(1)}(\boldsymbol{\theta})] + \sum_{k=0}^{\infty} (\mathbb{E} [\hat{\ell}^{(2^{k+1})}(\boldsymbol{\theta})] - \mathbb{E} [\hat{\ell}^{(2^k)}(\boldsymbol{\theta})]) \\ &= \lim_{k \rightarrow \infty} \mathbb{E} [\hat{\ell}^{(2^{k+1})}(\boldsymbol{\theta})] \\ &= \mathbb{E} \left[\lim_{k \rightarrow \infty} \hat{\ell}^{(2^{k+1})}(\boldsymbol{\theta}) \right] \\ &= \log p_\theta(\mathbf{x}). \end{aligned}$$

Here the integrability requirements of the second equality follow from boundedness of $\hat{\ell}^{(k)}(\boldsymbol{\theta})$ in L^2 . The second last equality holds due to uniform integrability, which likewise follows from boundedness in L^2 . Finally, the last equality holds because $\hat{\ell}^{(k)}(\boldsymbol{\theta}) = \log \left(\frac{1}{k} \sum_{i=1}^k w_i \right) \rightarrow \log p_\theta(\mathbf{x})$ almost surely by the SLLN (which is applicable due to our moment conditions, which entail $\mathbb{E}[|w_i|] < \infty$) and the continuous mapping theorem. This shows that Assumption (a) in Theorem 1 is satisfied.

Assumption (c): This trivially holds.

Assumption (b): We divide our analysis into two cases where $O_k E_k$ is large and $O_k E_k$ is small, specifically into $O_k E_k \geq 2^{-2\kappa} Z^2$ and $O_k E_k < 2^{-2\kappa} Z^2$, where $\mathbb{E}[O_k] = \mathbb{E}[E_k] = Z = p_\theta(\mathbf{x})$, and $\kappa > 0$ is a constant, independent of k , taken such that $2^{-\kappa} Z < 1$.

1. Case $O_k E_k$ large, i.e. $O_k E_k \geq 2^{-2\kappa} Z^2$:

We have $\Delta_k = \log \left(\frac{O_k + E_k}{2} \right) - \frac{1}{2} (\log O_k + \log E_k) = \log \left(\frac{O_k + E_k}{2\sqrt{O_k E_k}} \right) = \log \left(1 + \frac{(O_k - E_k)^2}{2\sqrt{O_k E_k}(\sqrt{O_k} + \sqrt{E_k})^2} \right)$. By Lemma 1, $\forall 0 < p \leq 1, \exists c_p > 0$ such that

$$\Delta_k \leq c_p \left(\frac{(O_k - E_k)^2}{2\sqrt{O_k E_k}(\sqrt{O_k} + \sqrt{E_k})^2} \right)^p \quad \text{and} \quad \Delta_k^2 \leq c_p^2 \left(\frac{(O_k - E_k)^4}{4O_k E_k (\sqrt{O_k} + \sqrt{E_k})^4} \right)^p. \quad (2)$$

The denominator in (2) is clearly bounded from below in this case, since $(\sqrt{O_k} + \sqrt{E_k})^4 \geq O_k E_k \geq 2^{-2\kappa} Z^2$.

For the numerator, we can take $p = \min(\frac{2+\epsilon}{4}, 1)$. Hence by Lemma 4,

$$\mathbb{E} [|O_k - E_k|^{4p}] = O(2^{-2pk}),$$

since $\mathbb{E} [|O_k - Z|^{4p}], \mathbb{E} [|E_k - Z|^{4p}]$ are both $O(2^{-2pk})$ by Lemma 2. Hence $\mathbb{E} [\Delta_k^2 \mathbb{I}(O_k E_k \geq 2^{-2\kappa} Z^2)] = O(2^{-(1+\alpha)k})$, where $\alpha = \min(\frac{\epsilon}{2}, 1)$.

2. Case $O_k E_k$ small, i.e. $O_k E_k < 2^{-2\kappa} Z^2$: In this case, we must have either $O_k < 2^{-\kappa} Z < 1$ or $E_k < 2^{-\kappa} Z < 1$. Suppose that $O_k < 2^{-\kappa} Z < 1$, then since $\log(1+x) \leq \max(1+\log x, 1)$, we have

$$\begin{aligned} \Delta_k &= \log \left(1 + \frac{(\sqrt{O_k} - \sqrt{E_k})^2}{2\sqrt{O_k E_k}} \right) \\ &\leq \max \left(1 + \log \left(\frac{(\sqrt{O_k} - \sqrt{E_k})^2}{2\sqrt{O_k E_k}} \right), 1 \right) \\ &= \max \left(1 + 2 \log |\sqrt{O_k} - \sqrt{E_k}| - \frac{1}{2} (\log O_k + \log E_k) - \log 2, 1 \right). \end{aligned}$$

For $2 \log |\sqrt{O_k} - \sqrt{E_k}|$, we know that if $O_k \leq E_k$, then $2 \log |\sqrt{O_k} - \sqrt{E_k}| \leq \log E_k \leq \log E_k \mathbb{I}(E_k \geq 1)$. If $O_k > E_k$, then $2 \log |\sqrt{O_k} - \sqrt{E_k}| \leq \log O_k < 0 \leq \log E_k \mathbb{I}(E_k \geq 1)$. Hence we have

$$\Delta_k \leq \max \left(1 + \log E_k \mathbb{I}(E_k \geq 1) - \frac{1}{2} (\log O_k + \log E_k), 1 \right).$$

We also have $-\frac{1}{2} \log E_k \leq -\frac{1}{2} \log E_k \mathbb{I}(E_k < 1)$, and $-\frac{1}{2} \log O_k > 0$. Overall, we thus have

$$\begin{aligned} \Delta_k &\leq \max \left(1 + \log E_k \mathbb{I}(E_k \geq 1) - \frac{1}{2} \log E_k \mathbb{I}(E_k < 1) - \frac{1}{2} \log O_k, 1 \right) \\ &= 1 + \log E_k \mathbb{I}(E_k \geq 1) - \frac{1}{2} \log E_k \mathbb{I}(E_k < 1) - \frac{1}{2} \log O_k. \end{aligned}$$

Let $Y_k = 1 + \log E_k \mathbb{I}(E_k \geq 1) - \frac{1}{2} \log E_k \mathbb{I}(E_k < 1)$. Then we have

$$\begin{aligned} \Delta_k^2 &\leq Y_k^2 - Y_k \log O_k + \frac{1}{4} (\log O_k)^2, \\ Y_k^2 &= 1 + \left((\log E_k)^2 + 2 \log E_k \right) \mathbb{I}(E_k \geq 1) + \left(\frac{1}{4} (\log E_k)^2 - \log E_k \right) \mathbb{I}(E_k < 1). \end{aligned}$$

By independence of O_k and E_k , we have

$$\mathbb{E} [\Delta_k^2 \mathbb{I}(O_k < 2^{-\kappa} Z)] \leq \mathbb{E} [Y_k^2] \cdot \mathbb{P}(O_k < 2^{-\kappa} Z) - \mathbb{E} [Y_k] \mathbb{E} [\log O_k \mathbb{I}(O_k < 2^{-\kappa} Z)] + \frac{1}{4} \mathbb{E} [(\log O_k)^2 \mathbb{I}(O_k < 2^{-\kappa} Z)].$$

We first prove that $\mathbb{E} [Y_k]$ and $\mathbb{E} [Y_k^2]$ are $O(1)$. By Hölder's inequality, we have for any event A $\mathbb{E} [|\log E_k|^p \mathbb{I}(A)] \leq \mathbb{E} [|\log E_k|^{2p}]^{\frac{1}{2}}$. However, $\forall p \geq 1, \exists c_p > 0$ s.t. $|\log x|^p \leq c_p (x^{-\delta} + x)$ by Lemma 1. By Lemma 3, $\mathbb{E} [E_k^{-\delta}]$ is bounded above. Therefore, $\mathbb{E} [Y_k]$ and $\mathbb{E} [Y_k^2]$ are both $O(1)$.

Next, we prove that $\mathbb{P}(O_k < 2^{-\kappa} Z)$ and $\mathbb{E} [|\log O_k|^p \mathbb{I}(O_k < 2^{-\kappa} Z)]$ are both arbitrarily close to the convergence speed $O(2^{-(1+\alpha)k})$. By Hölder's inequality, $\mathbb{E} [|\log O_k|^p \mathbb{I}(O_k < 2^{-\kappa} Z)] \leq \mathbb{E} [|\log O_k|^{pq}]^{\frac{1}{q}} \mathbb{P}(O_k < 2^{-\kappa} Z)^{(1-\frac{1}{q})}$. We have shown above that $\mathbb{E} [|\log O_k|^{pq}]^{\frac{1}{q}}$ is $O(1)$ for fixed p, q , and $\mathbb{P}(O_k < 2^{-\kappa} Z) \leq \mathbb{P}(|O_k - Z| > (1 - 2^{-\kappa})Z) \leq \frac{\mathbb{E}[|O_k - Z|^{2+\epsilon}]}{((1-2^{-\kappa})Z)^{2+\epsilon}} = O(2^{-(1+\epsilon/2)k})$, by Markov's inequality and Lemma 2. Take q to be sufficiently large to obtain $\mathbb{P}(O_k < 2^{-\kappa} Z)^{(1-\frac{1}{q})} = O(2^{-(1+\alpha)k})$, where α can be arbitrarily close to $\frac{\epsilon}{2}$.

Hence $\mathbb{P}(O_k < 2^{-\kappa} Z)$ and $\mathbb{E} [|\log O_k|^p \mathbb{I}(O_k < 2^{-\kappa} Z)]$ are both $O(2^{-(1+\alpha)k})$, where α can be arbitrarily close to $\frac{\epsilon}{2}$. Therefore, $\mathbb{E} [\Delta_k^2 \mathbb{I}(O_k < 2^{-\kappa} Z)] = O(2^{-(1+\alpha)k})$, and by symmetry the case is similar for $E_k < 2^{-\kappa} Z$.

Combining all the results, since $\mathbb{E} [\Delta_k^2] \leq \mathbb{E} [\Delta_k^2 \mathbb{I}(O_k E_k \geq 2^{-2\kappa} Z^2)] + \mathbb{E} [\Delta_k^2 \mathbb{I}(O_k < 2^{-\kappa} Z)] + \mathbb{E} [\Delta_k^2 \mathbb{I}(E_k < 2^{-\kappa} Z)]$, we obtain $\mathbb{E} [\Delta_k^2] = O(2^{-(1+\alpha)k})$, where α can be taken arbitrarily close to $\min(\frac{\epsilon}{2}, 1)$. Therefore, all values of r strictly in $(\frac{1}{2}, 1 - \frac{1}{2^{1+\alpha}})$ admit finite variance for $\hat{\ell}^{\text{ML}}(\theta)$ by Theorem 1, where $\alpha = \min(\frac{\epsilon}{2}, 1)$. \square

3 Proof of Theorem 3

We first provide here a general result for the SNIS estimator

$$\hat{\pi}^{(k)}[\psi] = \sum_{i=1}^k \bar{w}_i \psi(\mathbf{z}_i), \quad \bar{w}_i = \frac{w_i}{\sum_{j=1}^k w_j}$$

introduced in Section 2.1. We will use the following notation for the standard Monte Carlo estimator of the expectation of a function ψ with respect to a measure μ , $\hat{\mu}_{\text{MC}}^{(k)}[\psi] := \frac{1}{k} \sum_{i=1}^k \psi(\mathbf{z}_i)$ for $\mathbf{z}_i \stackrel{\text{i.i.d.}}{\sim} \mu$.

Lemma 5. *Assume that the SNIS estimator $\hat{\pi}^{(k)}[\psi]$, with unnormalized target density $\gamma(\mathbf{z})$, normalizing constant Z , normalized density $\pi(\mathbf{z})$, proposal density $q(\mathbf{z})$ and importance weight $w(\mathbf{z}) = \frac{\gamma(\mathbf{z})}{q(\mathbf{z})}$, satisfies $\mathbb{E}_q[w(\mathbf{z})^{2p}] < \infty$ for some $p \geq 1$. Then for any bounded test function ψ , $\mathbb{E} \left[\left| \hat{\pi}^{(k)}[\psi] - \pi[\psi] \right|^{2p} \right] = O(k^{-p})$.*

Proof. We note that

$$\pi[\psi] = \frac{q[w\psi]}{q[w]} = \frac{q[w\psi]}{Z}, \quad \hat{\pi}^{(k)}[\psi] = \frac{\hat{q}_{\text{MC}}^{(k)}[w\psi]}{\hat{q}_{\text{MC}}^{(k)}[w]}.$$

We have the following decomposition

$$\begin{aligned} \hat{\pi}^{(k)}[\psi] - \pi[\psi] &= \hat{\pi}^{(k)}[\psi] - \frac{1}{Z} \hat{q}_{\text{MC}}^{(k)}[w\psi] + \frac{1}{Z} \hat{q}_{\text{MC}}^{(k)}[w\psi] - \pi[\psi] \\ &= \frac{1}{Z} \left[\left(q[w] - \hat{q}_{\text{MC}}^{(k)}[w] \right) \hat{\pi}^{(k)}[\psi] + \left(\hat{q}_{\text{MC}}^{(k)}[w\psi] - q[w\psi] \right) \right] \end{aligned}$$

By the Cp-inequality, we have for $p \geq 1/2$

$$\begin{aligned} \left| \hat{\pi}^{(k)}[\psi] - \pi[\psi] \right|^{2p} &\leq \frac{1}{Z^{2p}} \left(\left| \left(q[w] - \hat{q}_{\text{MC}}^{(k)}[w] \right) \hat{\pi}^{(k)}[\psi] \right| + \left| \hat{q}_{\text{MC}}^{(k)}[w\psi] - q[w\psi] \right| \right)^{2p} \\ &\leq \frac{2^{2p-1}}{Z^{2p}} \left(\left| \left(q[w] - \hat{q}_{\text{MC}}^{(k)}[w] \right) \hat{\pi}^{(k)}[\psi] \right|^{2p} + \left| \hat{q}_{\text{MC}}^{(k)}[w\psi] - q[w\psi] \right|^{2p} \right). \end{aligned} \quad (3)$$

For $|\psi| \leq B$,

$$\left| \hat{\pi}^{(k)}[\psi] - \pi[\psi] \right|^{2p} \leq \frac{2^{2p-1}}{Z^{2p}} \left(B^{2p} \left| q[w] - \hat{q}_{\text{MC}}^{(k)}[w] \right|^{2p} + \left| \hat{q}_{\text{MC}}^{(k)}[w\psi] - q[w\psi] \right|^{2p} \right).$$

Hence we can conclude by Lemma 2. \square

Lemma 6. *With the same notation as in Lemma 5, suppose instead $\mathbb{E}_q[w(\mathbf{z})^{2ps}] < \infty$ and $\mathbb{E}_q[|\psi(\mathbf{z})|^{2pt}] < \infty$, where $s, t \geq 1$ and $\frac{1}{s} + \frac{1}{t} = 1$, then $\mathbb{E} \left[\left| \hat{\pi}^{(k)}[\psi] - \pi[\psi] \right|^{2p} \right] = O(k^{-p+1/t})$.*

Proof. Following (3), it is sufficient to study the term $\mathbb{E} \left[\left| \left(q[w] - \hat{q}_{\text{MC}}^{(k)}[w] \right) \hat{\pi}^{(k)}[\psi] \right|^{2p} \right]$. By Hölder's inequality, we have

$$\mathbb{E} \left[\left| \left(q[w] - \hat{q}_{\text{MC}}^{(k)}[w] \right) \hat{\pi}^{(k)}[\psi] \right|^{2p} \right] \leq \mathbb{E} \left[\left| q[w] - \hat{q}_{\text{MC}}^{(k)}[w] \right|^{2ps} \right]^{1/s} \mathbb{E} \left[\left| \hat{\pi}^{(k)}[\psi] \right|^{2pt} \right]^{1/t}$$

where $s, t \geq 1$ such that $\frac{1}{s} + \frac{1}{t} = 1$. The first term $\mathbb{E} \left[\left| q[w] - \hat{q}_{\text{MC}}^{(k)}[w] \right|^{2ps} \right]^{1/s}$ decays at $O(k^{-p})$ by Lemma 2.

For the second term $\mathbb{E} \left[\left| \hat{\pi}^{(k)}[\psi] \right|^{2pt} \right]^{1/t} = \mathbb{E} \left[\left| \sum_{i=1}^k \bar{w}_i \psi(\mathbf{z}_i) \right|^{2pt} \right]^{1/t}$, we have $\mathbb{E} \left[\left| \sum_{i=1}^k \bar{w}_i \psi(\mathbf{z}_i) \right|^{2pt} \right]^{1/t} \leq \mathbb{E} \left[\left| \sum_{i=1}^k \bar{w}_i |\psi(\mathbf{z}_i)| \right|^{2pt} \right]^{1/t} \leq \mathbb{E} \left[\sum_{i=1}^k |\psi(\mathbf{z}_i)|^{2pt} \right]^{1/t} \leq k^{1/t} \mathbb{E} \left[|\psi(\mathbf{z}_1)|^{2pt} \right]^{1/t}$ where the second inequality follows from $\bar{w}_i \geq 0, \sum_{i=1}^k \bar{w}_i = 1$. Hence $\mathbb{E} \left[\left| \left(q[w] - \hat{q}_{\text{MC}}^{(k)}[w] \right) \hat{\pi}^{(k)}[\psi] \right|^{2p} \right]$ converges to 0 at a $O(k^{-p+1/t})$ rate. \square

We recall below the statement of Theorem 3 and give its proof.

Theorem. Assume there exists $a > 4$, $b \geq \frac{2a}{a-4}$ such that $\mathbb{E}_{q\phi} [w(\mathbf{z})^a + |\psi(\mathbf{z})|^b] < \infty$. Then $\hat{\pi}^{\text{ML}}[\psi]$ satisfies Theorem 1 for $r \in (\frac{1}{2}, 1 - \frac{1}{2^{1+\alpha}})$, where $\alpha = 1 - \frac{2}{b} > 0$.

Alternatively, if $2 < a \leq 4$, $b > \frac{2a+4}{a-2}$, or $a > 4$, $\frac{2a+4}{a-2} < b \leq \frac{2a}{a-4}$, then $\hat{\pi}^{\text{ML}}[\psi]$ satisfies Theorem 1 for $r \in (\frac{1}{2}, 1 - \frac{1}{2^{1+\alpha}})$, where $\alpha = \frac{a}{2} - \frac{a+2}{b} - 1 > 0$.

Proof. We check that the assumptions of Theorem 1 are satisfied.

Assumption (a): $\hat{\pi}^{(k)}[\psi] \rightarrow \pi[\psi]$ in L^{2q} by Lemma 6 (choice of q, s, t established later), so $\hat{\pi}^{(k)}[\psi]$ is uniformly integrable, hence similarly to the proof of Theorem 2 $\hat{\pi}^{\text{ML}}[\psi]$ is an unbiased estimator of $\pi[\psi]$.

Assumption (c): This trivially holds.

Assumption (b): We have

$$\begin{aligned} \Delta_k &= \hat{\pi}_{O \cup E}^{(2^{k+1})}[\psi] - \frac{1}{2} \left(\hat{\pi}_O^{(2^k)}[\psi] + \hat{\pi}_E^{(2^k)}[\psi] \right) \\ &= \sum_{i=1}^{2^k} \bar{w}_i^{O/O \cup E} \psi(\mathbf{z}_i^O) + \sum_{i=1}^{2^k} \bar{w}_i^{E/O \cup E} \psi(\mathbf{z}_i^E) - \frac{1}{2} \left(\sum_{i=1}^{2^k} \bar{w}_i^O \psi(\mathbf{z}_i^O) + \sum_{i=1}^{2^k} \bar{w}_i^E \psi(\mathbf{z}_i^E) \right) \end{aligned}$$

where

$$\bar{w}_i^{O/O \cup E} = \frac{w_i^O}{\sum_{j=1}^{2^k} w_j^O + \sum_{j=1}^{2^k} w_j^E}, \quad \bar{w}_i^{E/O \cup E} = \frac{w_i^E}{\sum_{j=1}^{2^k} w_j^O + \sum_{j=1}^{2^k} w_j^E}.$$

Consider the coefficients c_i^O, c_i^E for each $\psi(\mathbf{z}_i^O)$ and $\psi(\mathbf{z}_i^E)$. Rearranging the terms gives that

$$\begin{aligned} c_i^O &= \bar{w}_i^{O/O \cup E} - \frac{1}{2} \bar{w}_i^O = \frac{w_i^O}{\sum_{j=1}^{2^k} w_j^O + \sum_{j=1}^{2^k} w_j^E} - \frac{1}{2} \frac{w_i^O}{\sum_{j=1}^{2^k} w_j^O} \\ &= \frac{w_i^O \left(\sum_{j=1}^{2^k} w_j^O - \sum_{j=1}^{2^k} w_j^E \right)}{2 \left(\sum_{j=1}^{2^k} w_j^O \right) \left(\sum_{j=1}^{2^k} w_j^O + \sum_{j=1}^{2^k} w_j^E \right)} = \frac{w_i^O (O_k - E_k)}{2 \sum_{j=1}^{2^k} w_j^O (O_k + E_k)} \end{aligned}$$

and similarly $c_i^E = \frac{w_i^E (E_k - O_k)}{2 \sum_{j=1}^{2^k} w_j^E (O_k + E_k)}$. Hence,

$$\Delta_k = \sum_{i=1}^{2^k} c_i^O \psi(\mathbf{z}_i^O) + \sum_{i=1}^{2^k} c_i^E \psi(\mathbf{z}_i^E) = \frac{O_k - E_k}{2(O_k + E_k)} \left(\sum_{i=1}^{2^k} \bar{w}_i^O \psi(\mathbf{z}_i^O) - \sum_{i=1}^{2^k} \bar{w}_i^E \psi(\mathbf{z}_i^E) \right).$$

Consequently, we have

$$\mathbb{E} [\Delta_k^2] = \mathbb{E} \left[\frac{(O_k - E_k)^2}{4(O_k + E_k)^2} \left(\hat{\pi}_O^{(2^k)}[\psi] - \hat{\pi}_E^{(2^k)}[\psi] \right)^2 \right] \leq \mathbb{E} \left[\frac{|O_k - E_k|^{2p}}{4^p (O_k + E_k)^{2p}} \right]^{1/p} \mathbb{E} \left[\left| \hat{\pi}_O^{(2^k)}[\psi] - \hat{\pi}_E^{(2^k)}[\psi] \right|^{2q} \right]^{1/q}$$

for $\frac{1}{p} + \frac{1}{q} = 1$ by Hölder's inequality.

By Lemma 6, $\mathbb{E} \left[\left| \hat{\pi}_O^{(2^k)}[\psi] - \pi[\psi] \right|^{2q} \right]$ decays at rate $O(2^{-(q-1/t)k})$ if $\mathbb{E}_q[w(\mathbf{z})^{2qs}] < \infty$, $\mathbb{E}_q[|\psi(\mathbf{z})|^{2qt}] < \infty$, and

$\frac{1}{s} + \frac{1}{t} = 1$. By Lemma 4, we conclude that if $a \geq 2qs$, $b \geq 2qt$, $\mathbb{E} \left[\left| \hat{\pi}_O^{(2^k)}[\psi] - \hat{\pi}_E^{(2^k)}[\psi] \right|^{2q} \right]^{1/q} = O \left(2^{-(1-\frac{1}{qt})k} \right)$.

For $\mathbb{E} \left[\frac{|O_k - E_k|^{2p}}{4^p (O_k + E_k)^{2p}} \right]$, we can similarly consider two cases.

1. Case O_k, E_k large, i.e. $O_k \geq 2^{-\kappa}Z$ and $E_k \geq 2^{-\kappa}Z$, where $\kappa > 0$ is a constant:

Since w_i has finite a moments, if we choose $2p \leq a$, we have for the numerator $\mathbb{E} \left[|O_k - E_k|^{2p} \right] = O(2^{-pk})$ using Lemma 2 and 4, and the denominator is bounded below. Hence $\mathbb{E} \left[\frac{|O_k - E_k|^{2p}}{(O_k + E_k)^{2p}} \mathbb{I} \{O_k \geq 2^{-\kappa} Z, E_k \geq 2^{-\kappa} Z\} \right] = O(2^{-pk})$.

If we choose $2p \geq a$, we have that

$$\frac{|O_k - E_k|^{2p}}{(O_k + E_k)^{2p}} \leq \frac{|O_k - E_k|^a}{(O_k + E_k)^a},$$

since $\frac{|O_k - E_k|^{2p}}{(O_k + E_k)^{2p}}$ is bounded above by 1, and similarly $\mathbb{E} \left[\frac{|O_k - E_k|^{2p}}{(O_k + E_k)^{2p}} \mathbb{I} \{O_k \geq 2^{-\kappa} Z, E_k \geq 2^{-\kappa} Z\} \right] = O(2^{-ak/2})$.

2. Case O_k or E_k small, i.e. $O_k < 2^{-\kappa} Z$ or $E_k < 2^{-\kappa} Z$:

We have that $\frac{|O_k - E_k|^{2p}}{(O_k + E_k)^{2p}}$ is bounded above by 1. Hence, we have

$$\begin{aligned} \mathbb{E} \left[\frac{|O_k - E_k|^{2p}}{(O_k + E_k)^{2p}} \mathbb{I} (O_k < 2^{-\kappa} Z) \right] &\leq \mathbb{P} (O_k < 2^{-\kappa} Z) \\ &= O(2^{-ak/2}), \end{aligned}$$

where the last line is by the proof of Theorem 2. The case is similar for $E_k < 2^{-\kappa} Z$.

Therefore, we have that if we choose $2p \leq a$, $\mathbb{E} \left[\frac{|O_k - E_k|^{2p}}{4^p (O_k + E_k)^{2p}} \right]$ is $O(2^{-pk})$, hence $\mathbb{E} \left[\frac{|O_k - E_k|^{2p}}{4^p (O_k + E_k)^{2p}} \right]^{1/p}$ is $O(2^{-k})$. Aggregating the results, we have that $\mathbb{E} [\Delta_k^2] = O \left(2^{-k} \cdot 2^{-(1-\frac{1}{qt})k} \right) = O \left(2^{-(2-\frac{1}{qt})k} \right)$. Collecting together the conditions, we must have that

$$\begin{aligned} 2qs &\leq a \\ 2qt &\leq b \\ 2p &\leq a \end{aligned}$$

as well as the Hölder conjugate conditions $\frac{1}{p} + \frac{1}{q} = 1$, $\frac{1}{s} + \frac{1}{t} = 1$, $p, q, s, t \geq 1$.

Let $u = qt$, which we would like to maximize. Then solving the equation we have equivalently that

$$\begin{aligned} u &\leq \frac{a}{2}t - \frac{a}{2} \\ u &\leq \frac{b}{2} \\ u &\geq \frac{a}{a-2}t. \end{aligned}$$

In order for this set to be valid, one can check that we must have $\frac{a+b}{a} \leq \frac{b(a-2)}{2a}$, i.e. $a > 4$, $b \geq \frac{2a}{a-4}$. Solving for the maximum u , we have that $u = \frac{b}{2}$, and the rest can be taken as $t = \frac{b(a-2)}{2a}$, $s = \frac{t}{t-1}$, $p = \frac{u}{u-t} = \frac{a}{2}$ and $q = \frac{u}{t} = \frac{a}{a-2}$. Thus we have overall a convergence rate $O \left(2^{-(2-\frac{2}{b})k} \right)$. For multilevel Monte Carlo to work, we must have the strict inequality that $\frac{2}{b} < 1$, i.e. $b > 2$, so that $\alpha = 1 - \frac{2}{b} > 0$, but this is implied by $b \geq \frac{2a}{a-4}$.

If we choose $2p \geq a$, $\mathbb{E} \left[\frac{|O_k - E_k|^{2p}}{4^p (O_k + E_k)^{2p}} \right]$ is $O(2^{-ak/2})$, hence $\mathbb{E} \left[\frac{|O_k - E_k|^{2p}}{4^p (O_k + E_k)^{2p}} \right]^{1/p}$ is $O(2^{-\frac{ak}{2p}})$. Aggregating the results, we have that $\mathbb{E} [\Delta_k^2] = O \left(2^{-\frac{ak}{2p}} \cdot 2^{-(1-\frac{1}{qt})k} \right) = O \left(2^{-(1+\frac{a}{2p}-\frac{1}{qt})k} \right)$. Collecting together the conditions, we must have that

$$\begin{aligned} 2qs &\leq a \\ 2qt &\leq b \\ 2p &> a \end{aligned}$$

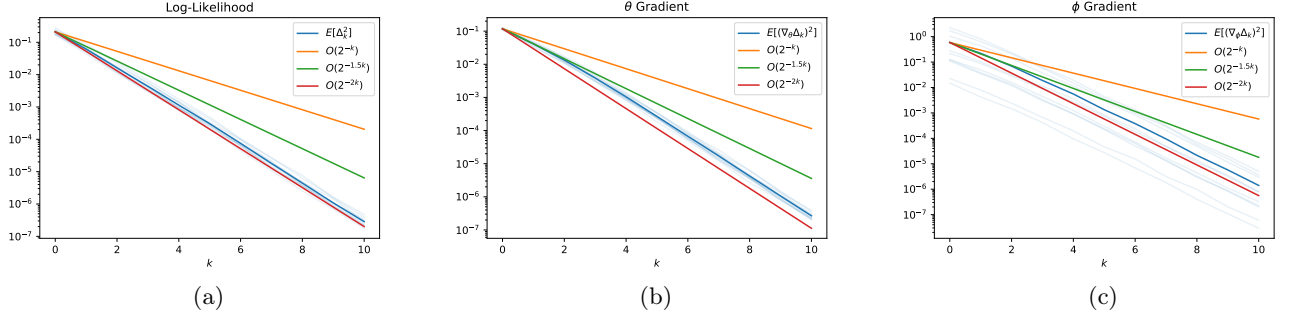


Figure 1: Empirical confirmation of the convergence rate of $\mathbb{E}[\Delta_k^2]$ for (a) $\hat{\ell}^{\text{ML}}(\boldsymbol{\theta})$; (b) $\hat{\pi}^{\text{ML}}[\psi]$ with $\psi(\mathbf{z}) := \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})$; and (c) $\hat{\pi}^{\text{ML}}[\psi]$ with $\psi(\mathbf{z}) := -\nabla_{\boldsymbol{\phi}} \log q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})$. For all three estimators, the convergence rate of $\mathbb{E}[\Delta_k^2]$ is close to the theoretical upper bound $O(2^{-2k})$.

as well as the Hölder conjugate conditions $\frac{1}{p} + \frac{1}{q} = 1$, $\frac{1}{s} + \frac{1}{t} = 1$, $p, q, s, t \geq 1$. Notice here that $\frac{a}{2p} \leq 1$ by our assumption, and that $2qt \leq b$, so the overall convergence rate $O\left(2^{-\left(1 + \frac{a}{2p} - \frac{1}{qt}\right)k}\right)$ is slower than $O\left(2^{-\left(2 - \frac{2}{b}\right)k}\right)$. Therefore, when $a > 4$, $b \geq \frac{2a}{a-4}$ holds, we already have the best convergence rate $O\left(2^{-\left(2 - \frac{2}{b}\right)k}\right)$.

We would like to maximize $1 + \frac{a}{2p} - \frac{1}{qt} = 1 + \frac{a}{2} - \frac{a}{2q} - \frac{1}{qt}$, so we would like to minimize $\frac{a}{2q} + \frac{1}{qt} = \frac{at+2}{2qt}$. Again letting $u = qt$, we have equivalently that

$$\begin{aligned} u &\leq \frac{a}{2}t - \frac{a}{2} \\ u &\leq \frac{b}{2} \\ u &\leq \frac{a}{a-2}t \\ u &\geq t. \end{aligned}$$

(In the case $a \leq 4$ or $a > 4, b \leq \frac{2a}{a-4}$, the third inequality is slack.) Solving again, the minimum $\frac{at+2}{2qt} = 1 + \frac{a+2}{b}$ is attained at $u = \frac{b}{2}$, $t = \frac{a+b}{a}$, $s = \frac{a+b}{b}$, $p = \frac{u}{u-t} = \frac{ab}{ab-2a-2b}$ and $q = \frac{u}{t} = \frac{ab}{2(a+b)}$. Hence the overall convergence rate $O\left(2^{-\left(1 + \frac{a}{2p} - \frac{1}{qt}\right)k}\right) = O\left(2^{-\left(\frac{a}{2} - \frac{a+2}{b}\right)k}\right)$. For multilevel Monte Carlo to work, we must have the strict inequality that $\frac{a}{2} - \frac{a+2}{b} > 1$, i.e. $a > 2$, $b > \frac{2a+4}{a-2}$, so that $\alpha = \frac{a}{2} - \frac{a+2}{b} - 1 > 0$. \square

Corollary 7. Assume there exists $\epsilon, \delta > 0$ such that $\mathbb{E}_{q_{\boldsymbol{\phi}}}[w(\mathbf{z})^{2+\epsilon} + w(\mathbf{z})^{-\delta}] < \infty$. Then $\hat{\pi}^{\text{ML}}[\log w]$ satisfies Theorem 1 for $r \in \left(\frac{1}{2}, 1 - \frac{1}{2(1+\alpha)}\right)$, where $\alpha = \min\left(\frac{\epsilon}{2}, 1\right) > 0$.

Proof. By Lemma 1, all moments of $\log w$ are finite under $q_{\boldsymbol{\phi}}$. Therefore, taking b arbitrarily large in Theorem 3 gives us the result. \square

4 Empirical Confirmation of the Convergence Rate of $\mathbb{E}[\Delta_k^2]$

We confirm our theoretical results and verify the convergence rate of $\mathbb{E}[\Delta_k^2]$ using the same example as in Section 7.1 for the three estimators: $\hat{\ell}^{\text{ML}}(\boldsymbol{\theta})$, $\hat{\pi}^{\text{ML}}[\psi]$ with $\psi(\mathbf{z}) := \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})$, and $\hat{\pi}^{\text{ML}}[\psi]$ with $\psi(\mathbf{z}) := -\nabla_{\boldsymbol{\phi}} \log q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})$. The experiments are conducted with 10 random perturbations and 1000 estimator samples per each perturbation, and the values of $\mathbb{E}[\Delta_k^2]$ are taken as the mean of the 10 trials. We observe that $\mathbb{E}[\Delta_k^2]$ for all estimators converges much faster than $O(2^{-1.5k})$ (corresponding to $\alpha = 0.5$) and very close to our theoretical limit $O(2^{-2k})$ (corresponding to $\alpha = 1$) when k is large. Therefore, Theorem 1 justifies the use of our approach with approximately $r \in (0.5, 0.75)$.

5 Comparing Single Sample (ss) and Russian Roulette (RR) Estimators

For simplicity, we consider the first term I_0 and the telescoping part $\tilde{s}\tilde{s} = \frac{\Delta_K}{p(K)}$, $\tilde{r}\tilde{r} = \sum_{k=0}^K \frac{\Delta_k}{\mathbb{P}(K \geq k)}$ in SS and RR separately. If I_0 is sampled independently, the variance of SS or RR is simply the sum of the variance of the two parts.

5.1 Variance of Single Sample (ss) Estimator

It is straightforward to show that

$$\mathbb{E}[\tilde{s}\tilde{s}] = \sum_{k=0}^{\infty} \mathbb{E}[\Delta_k], \quad \mathbb{E}[\tilde{s}\tilde{s}^2] = \sum_{k=0}^{\infty} \frac{\mathbb{E}[\Delta_k^2]}{p(k)}. \quad (4)$$

Hence we have

$$\text{Var}(\tilde{s}\tilde{s}) = \mathbb{E}[\tilde{s}\tilde{s}^2] - \mathbb{E}[\tilde{s}\tilde{s}]^2 \quad (5)$$

$$= \sum_{k=0}^{\infty} \frac{\mathbb{E}[\Delta_k^2]}{p(k)} - \left(\sum_{k=0}^{\infty} \mathbb{E}[\Delta_k] \right)^2 \quad (6)$$

$$= \underbrace{\sum_{k=0}^{\infty} \frac{\text{Var}(\Delta_k)}{p(k)}}_{\textcircled{1}} + \underbrace{\left(\sum_{k=0}^{\infty} \frac{\mathbb{E}[\Delta_k]^2}{p(k)} - \sum_{k=0}^{\infty} \mathbb{E}[\Delta_k]^2 - 2 \sum_{i,j:i < j} \mathbb{E}[\Delta_i] \mathbb{E}[\Delta_j] \right)}_{\textcircled{2}}. \quad (7)$$

5.2 Variance of Russian Roulette (RR) Estimator

$$\begin{aligned} \text{Var}(\tilde{r}\tilde{r}) &= \mathbb{E}[\tilde{r}\tilde{r}^2] - \mathbb{E}[\tilde{r}\tilde{r}]^2 \\ &= \mathbb{E} \left[\mathbb{E} \left[\left(\sum_{k=0}^K \frac{\Delta_k}{\mathbb{P}(K \geq k)} \right)^2 \middle| K \right] \right] - \left(\sum_{k=0}^{\infty} \mathbb{E}[\Delta_k] \right)^2 \\ &= \mathbb{E} \left[\sum_{k=0}^K \frac{\mathbb{E}[\Delta_k^2]}{\mathbb{P}(K \geq k)^2} + 2 \sum_{i,j:i < j}^K \frac{\mathbb{E}[\Delta_i \Delta_j]}{\mathbb{P}(K \geq i) \mathbb{P}(K \geq j)} \right] - \sum_{k=0}^{\infty} \mathbb{E}[\Delta_k]^2 - 2 \sum_{i,j:i < j} \mathbb{E}[\Delta_i] \mathbb{E}[\Delta_j] \\ &= \sum_{k=0}^{\infty} \frac{\mathbb{E}[\Delta_k^2]}{\mathbb{P}(K \geq k)} + 2 \sum_{i,j:i < j} \frac{\mathbb{E}[\Delta_i \Delta_j]}{\mathbb{P}(K \geq i)} - \sum_{k=0}^{\infty} \mathbb{E}[\Delta_k]^2 - 2 \sum_{i,j:i < j} \mathbb{E}[\Delta_i] \mathbb{E}[\Delta_j] \\ &= \underbrace{\sum_{k=0}^{\infty} \frac{\text{Var}(\Delta_k)}{\mathbb{P}(K \geq k)} + 2 \sum_{i,j:i < j} \frac{\text{Cov}(\Delta_i, \Delta_j)}{\mathbb{P}(K \geq i)}}_{\textcircled{1}} \\ &\quad + \underbrace{\left(\sum_{k=0}^{\infty} \frac{\mathbb{E}[\Delta_k]^2}{\mathbb{P}(K \geq k)} + 2 \sum_{i,j:i < j} \frac{\mathbb{E}[\Delta_i] \mathbb{E}[\Delta_j]}{\mathbb{P}(K \geq i)} - \sum_{k=0}^{\infty} \mathbb{E}[\Delta_k]^2 - 2 \sum_{i,j:i < j} \mathbb{E}[\Delta_i] \mathbb{E}[\Delta_j] \right)}_{\textcircled{2}} \end{aligned}$$

5.3 Proof of Theorem 1 for Russian Roulette (RR)

In Blanchet et al. (2019), only Theorem 1 for SS is mentioned. Here we show that the theorem similarly holds for RR.

Theorem. *If the following conditions hold:*

$$(a) \mathbb{E}[I_0] + \sum_{k=0}^{\infty} \mathbb{E}[\Delta_k] = I_{\infty};$$

(b) I_0 has finite variance, and there exists $\alpha, c > 0$ such that $\mathbb{E}[\Delta_k^2] \leq c \cdot 2^{-(1+\alpha)k}$ for all k ;

(c) $\mathbb{E}[C_k] \leq c' \cdot 2^k$ for some $c' > 0$, where C_k is the sampling cost of Δ_k ;

then for $K \sim \text{Geom}(r)$ where $r \in (\frac{1}{2}, 1 - \frac{1}{2^{1+\alpha}})$, RR is an unbiased estimator of I_∞ , whose variance and expected sampling cost are both finite.

Proof. The unbiasedness of RR is standard. Finite expected sampling cost follows directly from Theorem 1 for SS , since the sampling cost of SS and RR are the same. To show finite variance, we show that

$$\mathbb{E}[\tilde{\text{RR}}^2] = \sum_{k=0}^{\infty} \frac{\mathbb{E}[\Delta_k^2]}{\mathbb{P}(K \geq k)} + 2 \sum_{i,j:i < j} \frac{\mathbb{E}[\Delta_i \Delta_j]}{\mathbb{P}(K \geq i)}$$

is finite. Since $\mathbb{E}[\Delta_k^2] \leq c \cdot 2^{-(1+\alpha)k}$, for fixed i

$$\begin{aligned} \sum_{j>i} |\mathbb{E}[\Delta_i \Delta_j]| &\leq \sqrt{c} \cdot 2^{-\frac{1}{2}(1+\alpha)i} \cdot \sqrt{c} \cdot 2^{-\frac{1}{2}(1+\alpha)(i+1)} \cdot \frac{1}{1 - 2^{-\frac{1}{2}(1+\alpha)}} \\ &= O(2^{-(1+\alpha)i}). \end{aligned}$$

Now $\mathbb{P}(K \geq k) = (1-r)^k$, where $1-r \in (2^{-(1+\alpha)}, \frac{1}{2})$. Therefore $\mathbb{E}[\tilde{\text{RR}}^2] < \infty$, and RR has finite variance. \square

6 Using Thermodynamic Integration

6.1 Proof of Proposition 5

Proposition. *The following identity holds*

$$\hat{\ell}^{(k)}(\boldsymbol{\theta}) = \int_0^1 \hat{\pi}_\beta^{(k)}[\log w] d\beta. \quad (8)$$

It follows that $\hat{\ell}^{\text{ML}}(\boldsymbol{\theta}) = \int_0^1 \hat{\ell}_{\text{TI}}^{\text{ML}}(\boldsymbol{\theta}) d\beta$, i.e. $\hat{\ell}^{\text{ML}}(\boldsymbol{\theta})$ is a Rao-Blackwellized version of $\hat{\ell}_{\text{TI}}^{\text{ML}}(\boldsymbol{\theta})$ and thus $\text{Var}(\hat{\ell}_{\text{TI}}^{\text{ML}}(\boldsymbol{\theta})) \geq \text{Var}(\hat{\ell}^{\text{ML}}(\boldsymbol{\theta}))$.

Proof. We have that

$$\hat{\pi}_\beta^{(k)}[\log w] = \sum_{i=1}^k \frac{w_i^\beta}{\sum_{j=1}^k w_j^\beta} \log w_i. \quad (9)$$

By directly differentiating $\beta \mapsto F(\beta) = \log\left(\frac{1}{k} \sum_{i=1}^k w_i^\beta\right)$ with respect to β , we have that $F'(\beta) = \sum_{i=1}^k \frac{w_i^\beta}{\sum_{j=1}^k w_j^\beta} \log w_i$. Hence $\int_0^1 \hat{\pi}_\beta^{(k)}[\log w] d\beta = F(1) - F(0) = \log\left(\frac{1}{k} \sum_{i=1}^k w_i\right) = \ell^{(k)}(\boldsymbol{\theta})$.

For $\hat{\ell}_{\text{TI}}^{\text{ML}}(\boldsymbol{\theta})$,

$$\begin{aligned} \Delta_k &= \sum_{i=1}^{2^k} \frac{(w_i^O)^\beta}{\sum_{j=1}^{2^k} (w_j^O)^\beta + \sum_{j=1}^{2^k} (w_j^E)^\beta} \log w_i^O + \sum_{i=1}^{2^k} \frac{(w_i^E)^\beta}{\sum_{j=1}^{2^k} (w_j^O)^\beta + \sum_{j=1}^{2^k} (w_j^E)^\beta} \log w_i^E \\ &\quad - \frac{1}{2} \left(\sum_{i=1}^{2^k} \frac{(w_i^O)^\beta}{\sum_{j=1}^{2^k} (w_j^O)^\beta} \log w_i^O + \sum_{i=1}^{2^k} \frac{(w_i^E)^\beta}{\sum_{j=1}^{2^k} (w_j^E)^\beta} \log w_i^E \right). \end{aligned}$$

Directly integrating gives that $\int_0^1 \Delta_k(\beta) d\beta = \log\left(\frac{1}{2}(O_K + E_K)\right) - \frac{1}{2}(\log O_K + \log E_K)$, which corresponds to Δ_k for $\hat{\ell}^{\text{ML}}(\boldsymbol{\theta})$. Hence $\hat{\ell}^{\text{ML}}(\boldsymbol{\theta}) = \int_0^1 \hat{\ell}_{\text{TI}}^{\text{ML}}(\boldsymbol{\theta}) d\beta$. \square

6.2 Proof of Proposition 6

Proposition. $\mathbb{E}_{\mathbf{z}_{1:k}}[\hat{\ell}_{\text{TVO}}^{(k)}(\boldsymbol{\theta})] \leq \mathbb{E}_{\mathbf{z}_{1:k}}[\hat{\ell}^{(k)}(\boldsymbol{\theta})]$, regardless of the placement of all β_t , i.e. $\hat{\ell}_{\text{TVO}}^{(k)}(\boldsymbol{\theta})$ is less tight than IWAE due to the SNIS bias.

Proof. $\hat{\ell}_{\text{TVO}}^{(k)}(\boldsymbol{\theta}) = \sum_{t=0}^{T-1} (\beta_{t+1} - \beta_t) \hat{\pi}_{\beta_t}^{(k)}[\log w]$ can be viewed as a left Riemann sum approximation of the integral $\int_0^1 \hat{\pi}_{\beta}^{(k)}[\log w] d\beta = \hat{\ell}^{(k)}(\boldsymbol{\theta})$, where the equality is due to Proposition 5. However,

$$\frac{\partial}{\partial \beta} \hat{\pi}_{\beta}^{(k)}[\log w] = \frac{\left(\sum_{i=1}^k w_i^{\beta} (\log w_i)^2 \right) \left(\sum_{j=1}^k w_j^{\beta} \right) - \left(\sum_{i=1}^k w_i^{\beta} \log w_i \right) \left(\sum_{j=1}^k w_j^{\beta} \log w_j \right)}{\left(\sum_{j=1}^k w_j^{\beta} \right)^2} \geq 0$$

by the Cauchy-Schwarz inequality. Therefore, the integrand $\hat{\pi}_{\beta}^{(k)}[\log w]$ is non-decreasing in β , and the left Riemann sum approximation $\hat{\ell}_{\text{TVO}}^{(k)}(\boldsymbol{\theta}) \leq \hat{\ell}^{(k)}(\boldsymbol{\theta})$, hence $\mathbb{E}[\hat{\ell}_{\text{TVO}}^{(k)}(\boldsymbol{\theta})] \leq \mathbb{E}[\hat{\ell}^{(k)}(\boldsymbol{\theta})] = \ell_{\text{IWAE}}^{(k)}(\boldsymbol{\theta})$. \square

6.3 Proof of Theorem 7

By the propositions, it is thus not favorable to unbiasedly estimate the log-likelihood via the thermodynamic identity, or use SNIS to biasedly estimate $\ell_{\text{TVO}}(\boldsymbol{\theta})$. However, it is possible to unbiasedly estimate $\ell_{\text{TVO}}(\boldsymbol{\theta})$, and an advantage is that the moment conditions in Theorem 1 can be relaxed for multilevel Monte Carlo to have finite variance:

Theorem. Assume there exists $\epsilon, \delta > 0$ such that $\mathbb{E}_{q_{\phi}} [w(\mathbf{z})^{\beta_{T-1}(2+\epsilon)} + w(\mathbf{z})^{-\delta}] < \infty$. Then $\hat{\ell}_{\text{TVO}}^{\text{ML}}(\boldsymbol{\theta})$ satisfies Theorem 1 for $r \in (\frac{1}{2}, 1 - \frac{1}{2^{1+\alpha}})$, where $\alpha = \min(\frac{\epsilon}{2}, 1)$.

Proof. This is a direct application of Corollary 7. Unbiasedness and finite expected computation follows directly from the unbiasedness and finite expected computation of $\hat{\pi}_{\beta_t}^{(k)}[\log w]$ for each β_t .

Finite variance follows directly from Corollary 7 as well, since in order for the tempered importance weights $w(\mathbf{z})^{\beta_t}$ in $\hat{\pi}_{\beta_t}^{(k)}[\log w]$ to have finite $2 + \epsilon$ moments, equivalently $w(\mathbf{z})$ only needs to have finite $\beta_t(2 + \epsilon)$ moments. Since we require this for all β_t in $0 = \beta_0 < \beta_1 < \dots < \beta_{T-1} < 1$, we require that $\mathbb{E}_{q_{\phi}} [w(\mathbf{z})^{\beta_{T-1}(2+\epsilon)}] < \infty$. \square

6.4 SNIS Bias in the Covariance Gradient Estimator for $\boldsymbol{\theta}$

We note that Masrani et al. (2019) directly consider the gradient via the covariance gradient estimator

$$\nabla_{\boldsymbol{\lambda}} \pi_{\beta}[\log w] = \pi_{\beta}[\nabla_{\boldsymbol{\lambda}} \log w] + \text{Cov}_{\pi_{\beta}}[\nabla_{\boldsymbol{\lambda}} \log \tilde{\pi}_{\beta}, \log w]$$

for $\boldsymbol{\lambda} = \{\boldsymbol{\theta}, \phi\}$, where

$$\text{Cov}_{\pi_{\beta}}[\nabla_{\boldsymbol{\lambda}} \log \tilde{\pi}_{\beta}, \log w] = \mathbb{E}_{\pi_{\beta}} [(\nabla_{\boldsymbol{\lambda}} \log \tilde{\pi}_{\beta} - \mathbb{E}_{\pi_{\beta}}[\nabla_{\boldsymbol{\lambda}} \log \tilde{\pi}_{\beta}]) (\log w - \mathbb{E}_{\pi_{\beta}}[\log w])] .$$

However, Masrani et al. (2019) use SNIS for the RHS, and in particular for the $\boldsymbol{\theta}$ gradients the SNIS estimator for the covariance gradient estimator is precisely

$$\begin{aligned} & \sum_{i=1}^k \overline{w_i^{\beta}} \nabla_{\boldsymbol{\theta}} \log w_i + \sum_{i=1}^k \overline{w_i^{\beta}} \left(\nabla_{\boldsymbol{\theta}} \log \tilde{\pi}_{\beta}(\mathbf{z}_i) - \sum_{j=1}^k \overline{w_j^{\beta}} \nabla_{\boldsymbol{\theta}} \log \tilde{\pi}_{\beta}(\mathbf{z}_j) \right) \left(\log w_i - \sum_{j=1}^k \overline{w_j^{\beta}} \log w_j \right) \\ &= \sum_{i=1}^k \overline{w_i^{\beta}} \nabla_{\boldsymbol{\theta}} \log w_i + \sum_{i=1}^k \overline{w_i^{\beta}} \log w_i \nabla_{\boldsymbol{\theta}} \log \tilde{\pi}_{\beta}(\mathbf{z}_i) - \left(\sum_{i=1}^k \overline{w_i^{\beta}} \log w_i \right) \left(\sum_{i=1}^k \overline{w_i^{\beta}} \nabla_{\boldsymbol{\theta}} \log \tilde{\pi}_{\beta}(\mathbf{z}_i) \right) \\ &= \sum_{i=1}^k \overline{w_i^{\beta}} \nabla_{\boldsymbol{\theta}} \log w_i + \beta \sum_{i=1}^k \overline{w_i^{\beta}} \log w_i \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}_i) - \left(\sum_{i=1}^k \overline{w_i^{\beta}} \log w_i \right) \left(\beta \sum_{i=1}^k \overline{w_i^{\beta}} \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}_i) \right) \end{aligned}$$

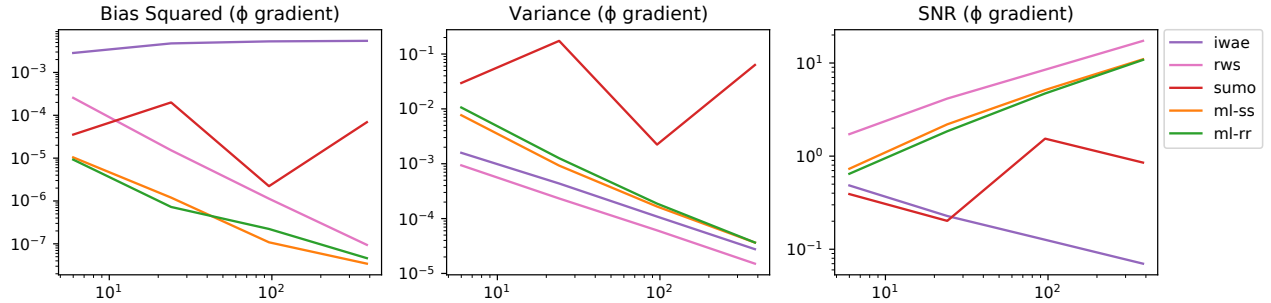


Figure 2: Empirical bias squared, variance and signal-to-noise ratio (SNR) of estimators of $\nabla_{\phi} \text{KL}(p_{\theta}(\mathbf{z}|\mathbf{x})||q_{\phi}(\mathbf{z}|\mathbf{x}))$, plotted against expected computational cost.

This is identical to differentiating the SNIS estimator $\hat{\pi}_{\beta}^{(k)}[\log w]$ directly, since

$$\begin{aligned} \nabla_{\theta} \hat{\pi}_{\beta}^{(k)}[\log w] &= \nabla_{\theta} \sum_{i=1}^k \frac{w_i^{\beta}}{\sum_{j=1}^k w_j^{\beta}} \log w_i \\ &= \sum_{i=1}^k \overline{w_i^{\beta}} \nabla_{\theta} \log w_i + \beta \sum_{i=1}^k \overline{w_i^{\beta}} \log w_i \nabla_{\theta} \log w_i - \frac{\left(\sum_{i=1}^k w_i^{\beta} \log w_i\right) \left(\beta \sum_{j=1}^k w_j^{\beta} \nabla_{\theta} \log w_j\right)}{\left(\sum_{j=1}^k w_j^{\beta}\right)^2} \\ &= \sum_{i=1}^k \overline{w_i^{\beta}} \nabla_{\theta} \log w_i + \beta \sum_{i=1}^k \overline{w_i^{\beta}} \log w_i \nabla_{\theta} \log p_{\theta}(\mathbf{x}, \mathbf{z}_i) - \left(\sum_{i=1}^k \overline{w_i^{\beta}} \log w_i\right) \left(\beta \sum_{i=1}^k \overline{w_i^{\beta}} \nabla_{\theta} \log p_{\theta}(\mathbf{x}, \mathbf{z}_i)\right). \end{aligned}$$

Therefore, using the covariance gradient estimator for θ is the same as $\nabla_{\theta} \hat{\pi}_{\beta}^{(k)}[\log w]$, so the SNIS bias in $\hat{\ell}_{\text{TVO}}^{(k)}(\theta)$ carries over to the covariance estimator.

We note that this analysis applies to the θ objective and its gradients, when viewing TVO as an estimator of $\ell(\theta)$, but not the use of TVO for the ϕ gradients.

7 Further Experiment Details

We provide some further experiment details and results in this section.

7.1 Linear Gaussian Experiment

With the same linear Gaussian example, we further analyze different estimators of the gradient of the forward KL $\nabla_{\phi} \text{KL}(p_{\theta}(\mathbf{z}|\mathbf{x})||q_{\phi}(\mathbf{z}|\mathbf{x}))$. We compare the ML-SS, ML-RR and SUMO estimators of $\hat{\pi}[\psi]$, $\psi(\mathbf{z}) := -\nabla_{\phi} \log q_{\phi}(\mathbf{z}|\mathbf{x})$ against the ϕ gradient in Reweighted Wake-Sleep (RWS), which uses a biased SNIS estimator.

From Figure 2, we observe that SS and RR achieve the lowest bias squared again compared to RWS and SUMO at the same computational cost, while obtaining higher variance than RWS but lower variance than SUMO. As IWAE suffers from the vanishing gradient problem, SS and RR also achieve higher SNR than the IWAE estimator, even at a low computational cost.

7.2 2D Density Modeling Experiment

For the second example, we compare VAE training using different estimators with the Figure-8 dataset proposed by Yacoby et al. (2020). We use the same flexible VAE architecture, with the encoder and decoder both parameterized by a neural network with 3 hidden layers and 50 leaky ReLU units in each layer, so that the model is flexible enough to learn the ground truth model. We use the Adam optimizer with a learning rate of $2.5 \cdot 10^{-4}$ and train for a maximum of 2000 epochs with early stopping applied based on the validation set.

Table 1: Test negative log-likelihood of the trained models for the Figure-8 dataset when trained using different estimators and q_ϕ objectives. All estimators have an expected sampling cost of 5 terms (except for SUMO which has a expected sampling cost of 6 terms).

(a)		(b)				
		q_ϕ objectives				
Baselines		Estimators	ELBO	IWAE	Unbiased RWS-STL	Var
ELBO	0.9659±0.0046	SUMO	0.9836±0.0107	0.9278±0.0058	0.9222±0.0036	0.9226±0.0091
IWAE	0.9325±0.0054	ML-SS	0.9843±0.0029	0.9263±0.0019	0.9216±0.0041	0.9228±0.0076
RWS-STL	0.9387±0.0074	ML-RR	0.9727±0.0057	0.9258±0.0027	0.9202±0.0048	0.9222±0.0071

Table 2: Other test metrics of the trained models for the Figure-8 dataset when trained with the ML-SS log-likelihood estimator and different ϕ objectives.

q_ϕ objectives				
	ELBO	IWAE	Unbiased RWS-STL	Var
$\text{KL}(q_\phi \parallel p_\theta)$	0.1570±0.0114	0.7390±0.2292	1.0066±0.2469	1.3684±0.0663
$\mathbb{E}[w^4]^{1/4}$	3.6667±0.4088	2.1612±0.7074	1.7557±0.1921	1.7534±0.2612
\hat{k}	1.0122±0.2337	0.2475±0.1783	0.3259±0.1668	-0.7768±0.1823

We provide a further study of suitable training objectives for q_ϕ here. Although the SUMO, ML-SS, and ML-RR log-likelihood estimators are unbiased estimators under reasonable conditions, we observe from Table 1 that the choice of the q_ϕ objective can still have a large impact on training. Maximizing the ELBO (or equivalently minimizing the reverse KL) resulted in the worst test log-likelihood among all tested q_ϕ training objectives due to the mode-seeking property of the reverse KL. In particular, an important quantity we considered is $\mathbb{E}[w^4]^{1/4}$, which is an essential quantity in our proofs of finite variance. We also compute the Pareto-smoothed importance sampling diagnostic \hat{k} (Vehtari et al., 2016), whose inverse gives an estimate of the number of existing moments of w when $\hat{k} > 0$. Finite variance of $\hat{\ell}^{\text{ML}}(\theta)$ can be achieved approximately when $\hat{k} < 0.5$. From Table 2, we see that maximizing the ELBO for q_ϕ resulted in the lowest reverse KL but the highest sample 4th moment and $\hat{k} > 1$, which confirms that $q_\phi(\mathbf{z}|\mathbf{x})$ does not cover the true posterior $p_\theta(\mathbf{z}|\mathbf{x})$ well and can lead to infinite variance even biased estimators. On the contrary, using the variance or the unbiased RWS-STL objective (both reuse samples and do not require generating new samples) achieved the best test log-likelihood, and in particular the variance objective achieves the lowest value of $\hat{k} < 0$, which justifies our approach and indicates finite variance of $\hat{\ell}^{\text{ML}}(\theta)$ for all $r \in (0.5, 0.75)$.

7.3 Image Modeling Experiment

For the image modeling experiment, we follow the training scheme proposed by Luo et al. (2020, Appendix A.8) with batch size 100, the Amsgrad optimizer (Reddi et al., 2018) and the same gradient norm clipping scheme for p_θ . However, in order to accurately compare different estimators under same budget, we fix the number of training epochs to 3280 for all estimators and use a linear learning rate decay from 10^{-3} to 10^{-4} . We also only modify the p_θ objective for this experiment, with different estimators in place of IWAE for p_θ , which provides a fair comparison to the IWAE baseline. IWAE is used to train q_ϕ and new samples are drawn for all estimators.

For this example, we find that the convergence rate of $\mathbb{E}[\Delta_k^2]$ at initialization is $O(2^{-(1+\alpha)k})$ with $\alpha \approx 0.5$, which justifies the use of $r \in (0.5, 0.6464)$, but the value of α decays to less than 0 as training progresses. To solve this, similar to Luo et al. (2020), we propose to modify the tail of the sampling distribution of K as described in Section 5.4. We set k_{max} so that ML-RR is an unbiased estimator of IWAE with $k = 128$, while having a much smaller expected computational cost (5 or 15 terms in our experiments)¹. This can also lower the

¹Luo et al. (2020) propose to softly truncates the tail of K after $k = 80$ terms, so our modification is in line with

computational cost and limit the computation/memory usage to be less than $2^{k_{\max}+1}$ samples. For ML-TVO-RR, we use 5 intermediate β_t values and follow Masrani et al. (2019) for their placement. By applying thermodynamic integration, faster convergence rate of $\mathbb{E}[\Delta_k^2]$ can be established. This shows that the relaxation of the moment condition from $\mathbb{E}_{q_\phi} [w(\mathbf{z})^{2+\epsilon}] < \infty$ in Theorem 2 of main text to $\mathbb{E}_{q_\phi} [w(\mathbf{z})^{\beta_{T-1}(2+\epsilon)}] < \infty$ in Theorem 7 is crucial to the convergence speed of $\mathbb{E}[\Delta_k^2]$, at the cost of an approximation error in the numerical integration of the 1D integral.

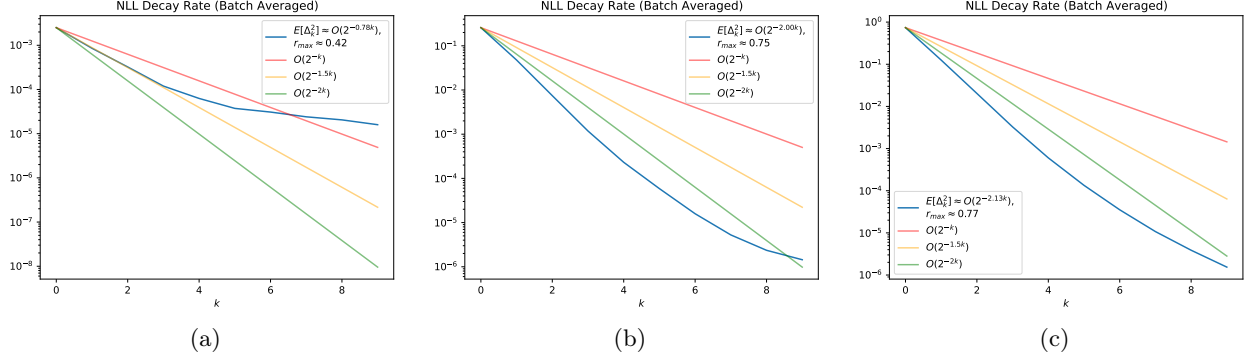


Figure 3: Empirical confirmation of the convergence rate of $\mathbb{E}[\Delta_k^2]$ for $\hat{\ell}^{\text{ML}}(\theta)$, evaluated on the models trained with p_θ objective $\hat{\ell}^{\text{ML-SS}}(\theta)$ and q_ϕ objective (a) ELBO; and (b) unbiased RWS-STL; (c) $\text{Var} \hat{\ell}^{\text{ML-SS}}(\theta)$ on the Figure-8 dataset. Clearly training with ELBO as the q_ϕ objective results in much slower convergence speed of $\mathbb{E}[\Delta_k^2]$.

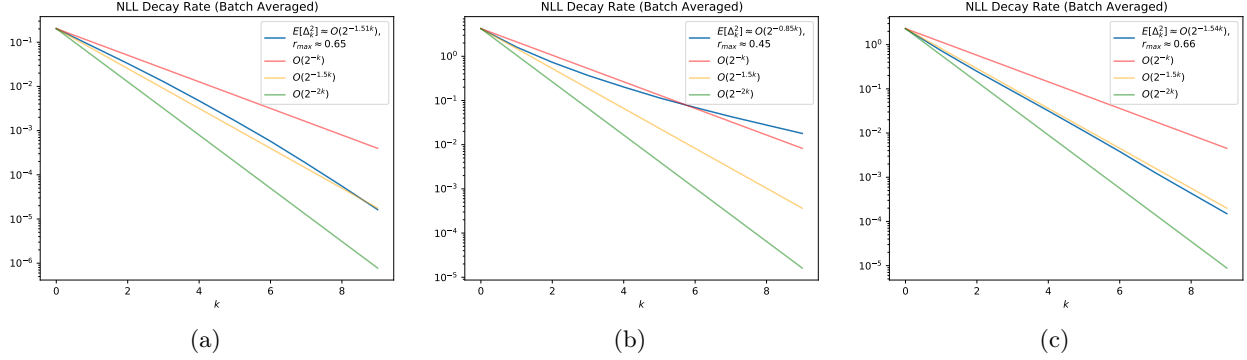


Figure 4: Empirical confirmation of the convergence rate of $\mathbb{E}[\Delta_k^2]$ for $\hat{\ell}^{\text{ML}}(\theta)$, evaluated on (a) model at initialization; (b) model trained with $\hat{\ell}^{\text{ML}}(\theta)$ on the MNIST dataset. (c) Convergence rate of $\mathbb{E}[\Delta_k^2]$ for $\hat{\ell}_{\text{TVO}}^{\text{ML}}(\theta)$, evaluated on model trained with $\hat{\ell}_{\text{TVO}}^{\text{ML}}(\theta)$. $T = 5$ and $\beta_1 = 0.01, \beta_{T-1} \approx 0.3162$ is taken using the log-uniform spacing as in Masrani et al. (2019).

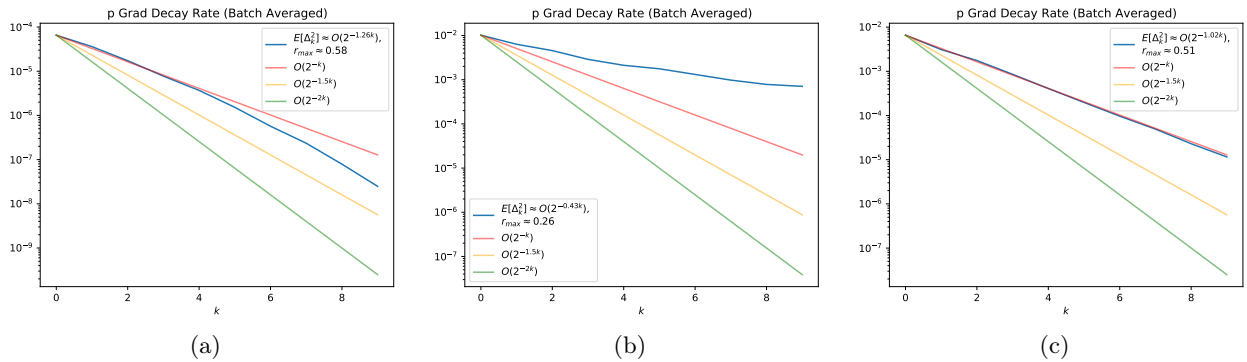


Figure 5: Empirical confirmation of the convergence rate of $\mathbb{E}[\Delta_k^2]$ for $\nabla_\theta \hat{\ell}^{\text{ML}}(\theta)$, evaluated on (a) model at initialization; (b) model trained with $\hat{\ell}^{\text{ML}}(\theta)$ on the MNIST dataset. (c) Convergence rate of $\mathbb{E}[\Delta_k^2]$ for $\nabla_\theta \hat{\ell}_{\text{TVO}}^{\text{ML}}(\theta)$, evaluated on model trained with $\hat{\ell}_{\text{TVO}}^{\text{ML}}(\theta)$.

References

- Jose H Blanchet, Peter W Glynn, and Yanan Pei. Unbiased multilevel Monte Carlo: Stochastic optimization, steady-state simulation, quantiles, and other applications. *arXiv preprint arXiv:1904.09929*, 2019.
- Takashi Goda and Kei Ishikawa. Multilevel Monte Carlo estimation of log marginal likelihood. *arXiv preprint arXiv:1912.10636*, 2019.
- Kei Ishikawa and Takashi Goda. Efficient debiased variational Bayes by multilevel Monte Carlo methods. *arXiv preprint arXiv:2001.04676*, 2020.
- Yucen Luo, Alex Beatson, Mohammad Norouzi, Jun Zhu, David Duvenaud, Ryan P. Adams, and Ricky T. Q. Chen. SUMO: Unbiased estimation of log marginal probability for latent variable models. In *International Conference on Learning Representations*, 2020.
- Vaden Masrani, Tuan Anh Le, and Frank Wood. The thermodynamic variational objective. In *Advances in Neural Information Processing Systems*, pages 11525–11534, 2019.
- Sashank Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018.
- Aki Vehtari, Andrew Gelman, and Jonah Gabry. Practical bayesian model evaluation using leave-one-out cross-validation and WAIC. *Statistics and Computing*, 27(5):1413–1432, August 2016.
- Yaniv Yacoby, Weiwei Pan, and Finale Doshi-Velez. Failure modes of variational autoencoders and their effects on downstream tasks. *arXiv preprint arXiv:2007.07124*, 2020.