

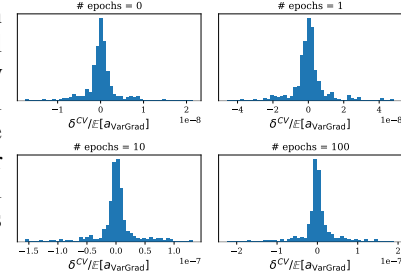
1 We thank all the reviewers for their stimulating comments and engaging questions. We are happy that they appreciated
 2 our theoretical analyses and argued that those, as well as VarGrad’s practical usefulness, are of benefit to the NeurIPS
 3 community. We now address some questions and comments in the sequel.

4 *R1, R2: The gradient estimator was already known.* We agree and did not claim otherwise. We also appreciate that most
 5 reviewers did not find this a concern. We emphasise that the novelty in our work is: (i) the derivation of this estimator
 6 in terms of the log-variance loss, and (ii) the theoretical analysis of its variance. (i) This new perspective is simpler and
 7 allows for a natural interpretation in terms of divergences (which is beneficial in other ML areas, e.g., reinforcement
 8 learning). We also note that the connection with the log-variance loss is of practical interest, as it enables a simple
 9 implementation algorithm based on automatic differentiation. Moreover, we believe this connection opens the door for
 10 further research. (ii) Our work is the first to provide a theoretical analysis of the variance of this estimator. This analysis
 11 shows that VarGrad’s control variate coefficients are close to the optimal ones for the score function control variate.

12 *R1: What class of distributions satisfies the kurtosis condition of Proposition 2?* In general, the kurtosis will not
 13 negatively affect the bound as long as the tails of the variational distribution do not become too heavy during the
 14 optimisation. More precisely, in Lemma 2 we translated the kurtosis condition to a condition that is more easily
 15 verifiable in the case of exponential families and in Lemma 3 we explicitly proved that the condition is fulfilled in the
 16 Gaussian case (where the 1-dimensional case carries over to the multidimensional case). In practice, q is often modelled
 17 by a Gaussian (potentially with a neural network outputting its mean and covariance). The analysis for the Gaussian
 18 case can be found in Section C. Note that in higher dimensions (usually implying that $\text{KL}(q||p)$ is large, as elaborated
 19 in Lemma 4) the kurtosis is allowed to be large (but must be bounded—see also Corollary 1). More importantly, our
 20 numerical experiments suggest that the kurtosis condition is satisfied for more complicated cases (see Figures 1, 2).

21 *R1: The order notation might hide large constants; elaborate the significance of Corollary 1.* Numerical evidence
 22 suggests that, in practice, the ratio $\delta^{\text{CV}}/\mathbb{E}[a_{\text{VarGrad}}]$ is usually very small (e.g., Figure B.4 shows that the order of
 23 magnitude is 10^{-8} , and we have observed similar results in all other experiments). Corollary 1 offers a rigorous
 24 comparison of the variances of Reinforce and VarGrad that does not depend on assumptions that might be hard to check.

25 *R1, R2, R3, R4: Empirical evaluation.* Please see the included figure for an
 26 empirical validation of $\delta^{\text{CV}}/\mathbb{E}[a_{\text{VarGrad}}]$ on a non-linear (deep) DVAE model
 27 used in Tucker, et. al. [2017] on Omniglot. In the attached plot this quantity
 28 is typically very small and fluctuates around zero, validating our analytical
 29 results in the challenging non-linear setting (a similar result is shown in Figure
 30 B.4 for the two-layer linear case with an in-depth discussion that follows for
 31 the non-linear case too). For the non-linear case, we also observe practical
 32 performance gains similar to the linear case in the trace plots. The new results
 33 have now been added to the manuscript.



34 *R1: Is the log-variance well-defined?* Since p , q and r are assumed to admit densities, it follows (under the support
 35 condition stated in the paper) that measure-zero sets of r are necessarily measure-zero sets of p and q , implying the
 36 divergence property. This has now been clarified in the paper.

37 *R3: VarGrad does not account for model structure.* The advantage of VarGrad is that it is easily implementable and
 38 applicable to a large array of models (black-box). In Figures 3, B.7 and B.8 we demonstrate that it is still competitive
 39 w.r.t. to alternative (tuned and structured) estimators.

40 *R3: Deeper analysis w.r.t. the choice of r .* We agree; however, as our work was focused on the theoretical analysis of
 41 the induced ELBO gradient estimator rather than a new variational objective, we left this question out intentionally for
 42 subsequent studies, in order to not confuse the presentation. Taking $r = q$ after differentiation reproduces the gradient
 43 of the KL divergence/ELBO (Proposition 1) and this choice was shown to be effective in the empirical evaluation
 44 (Nüsken & Richter [2020] provide further arguments for this choice).

45 *R4: Comparison to reparametrisation.* The reparametrisation trick is often not applicable (e.g., in discrete models),
 46 whereas VarGrad is general-purpose. That said, in Figures 3, B.7 and B.8 we compared against RELAX + REBAR,
 47 which use reparametrisation estimators as control variates.

48 *R4: Application to mixtures of non-diagonal Gaussians.* Indeed this is a potential application, since VarGrad can be
 49 applied easily to many variational families including non-diagonal mixtures of Gaussians. We will explore this further.

50 *R1: Figure 1 error bars and clarification of the oracle estimator.* We have now added the error bars. The oracle
 51 estimator in Figure 2 takes 1000 extra samples. We have now clarified this in the figure caption.

52 *R1: Why did you focus on the ELBO in the presentation?* As per your suggestion, we have now shifted the presentation
 53 on the ELBO to streamline the background section. Extra details were moved to the appendix.