

---

# Dual Control Variate for Faster Black-box Variational Inference

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Black-box variational inference is a widely-used framework for Bayesian posterior  
2 inference, but in some cases suffers from high variance in gradient estimates, harm-  
3 ing accuracy and efficiency. This variance comes from two sources of randomness:  
4 Data subsampling and Monte Carlo sampling. Whereas existing control variates  
5 only address Monte Carlo noise and incremental gradient methods typically only  
6 address data subsampling, we propose a new "dual" control variate capable of  
7 *jointly* reducing variance from both sources of noise. We confirm that this leads to  
8 reduced variance and improved optimization in several real-world applications.

## 9 1 Introduction

10 Black-box variational inference (BBVI) [12, 22, 16, 2] has become a popular alternative to Markov  
11 Chain Monte Carlo (MCMC) methods. The idea is to posit a variational family and optimize it to be  
12 close to the posterior, using only "black-box" access to the target model (evaluations of the density  
13 or gradient). This is done by estimating a stochastic gradient of the KL-divergence and deploying  
14 it in stochastic optimization. A key advantage of this procedure is that it allows the use of data  
15 subsampling in each iteration, which can greatly speed-up optimization with large datasets.

16 The optimization of BBVI is often described as a doubly-stochastic optimization problem [30, 27] in  
17 that BBVI's gradient estimation involves two sources of randomness: Monte Carlo sampling from  
18 the variational posterior and data subsampling from the full dataset. Because of the doubly-stochastic  
19 nature, one common challenge for BBVI is the variance of the gradient estimates: If this is very high,  
20 it forces very small stepsizes, leading to slow optimization convergence [20, 3].

21 Numerous works have been devoted to reducing the "Monte Carlo" noise that results from drawing  
22 samples from the current variational distribution [19, 25, 8, 9, 4]. These methods can typically be  
23 seen as creating an approximation of the objective for which the Monte Carlo noise can be integrated  
24 exactly, and using this to define a zero-mean random variable, i.e. a control variate, that is negatively  
25 correlated with the original gradient estimator. These methods can be used with subsampling by  
26 creating approximations for each datum. However, they are only able to reduce Monte Carlo noise  
27 for each datum—they do not reduce subsampling noise. This is critical, as subsampling noise is often  
28 the dominant source of gradient variance (Sec. 3).

29 At the same time, for (non-BBVI) optimization problems with *only* subsampling noise, the opti-  
30 mization community has developed incremental gradient methods that "recycle" previous gradient  
31 evaluations [26, 28, 13, 6, 7], leading to faster convergence. These methods do not address Monte  
32 Carlo noise. In fact, due to the way these methods rely on efficiently maintaining running averages,  
33 they cannot typically be applied to doubly-stochastic problems at all.

34 In this paper, we present a method that *jointly* controls Monte Carlo and subsampling noise in BBVI.  
35 The idea is to create approximations of the target for each datum where the Monte Carlo noise

36 can be integrated exactly. Then, we maintain running averages of the *approximate* gradients, with  
 37 noise integrated, overcoming the issue of applying incremental gradient ideas to doubly-stochastic  
 38 problems. The resulting method not only addresses both forms of noise but *interactions* between  
 39 them as well. We demonstrate through a series of experiments with diagonal Gaussian variational  
 40 inference on a range of probabilistic models that the method leads to lower variance and significantly  
 41 faster convergence than existing methods.

## 42 2 Background: Black-box variational inference

43 Given a probabilistic model  $p(x, z) = \prod_{n=1}^N p(x_n | z)p(z)$  and observed data  $\{x_1, \dots, x_N\}$ ,  
 44 variational inference’s goal is to find a tractable distribution  $q_w(z)$  with parameters  $w$  to approximate  
 45 the (often intractable) posterior  $p(z | x)$  over the latent variable  $z \in \mathbb{R}^d$ . BBVI achieves this by  
 46 finding the set of parameters  $w$  that minimize the KL-divergence from  $q_w(z)$  to  $p(z | x)$ , which is  
 47 equivalent to minimizing the negative Evidence Lower Bound (ELBO)

$$f(w) = -\mathbb{E}_n \mathbb{E}_{q_w(z)} \left[ N \log p(x_n | z) + \log p(z) \right] - \mathbb{H}(w), \quad (1)$$

48 where  $\mathbb{H}(w)$  denotes the entropy of  $q_w$ .

49 Since the inner expectation with respect to  $z$  is typically intractable, BBVI methods rely on stochastic  
 50 optimization with unbiased gradient estimates. These gradient estimates are typically obtained using  
 51 the score function method [33] or the reparameterization trick [15, 23, 30]. The latter is often the  
 52 method of choice, as it usually seems to yield estimators with lower variance. The idea is to define a  
 53 fixed base distribution  $s(\epsilon)$  and a deterministic transformation  $\mathcal{T}_w(\epsilon)$  such that for  $\epsilon \sim s$ ,  $\mathcal{T}_w(\epsilon)$  is  
 54 equal in distribution to  $q_w$ . Then, the objective from Equation (1) can be re-written as

$$f(w) = \mathbb{E}_n \mathbb{E}_\epsilon f(w; n, \epsilon), \quad \text{where} \quad f(w; n, \epsilon) = -N \log p(x_n | \mathcal{T}_w(\epsilon)) - \log p(\mathcal{T}_w(\epsilon)) - \mathbb{H}(w), \quad (2)$$

55 and its gradient can be estimated "naively" by drawing a random  $n$  and  $\epsilon$ , and evaluating

$$g_{\text{naive}}(w; n, \epsilon) = \nabla f(w; n, \epsilon). \quad (3)$$

56 BBVI has two advantages. First, since it only evaluates  $\log p$  (and its gradient) at various points,  
 57 it can be applied to a diverse range of models, including those with complex and non-conjugate  
 58 likelihoods. Second, by subsampling data it can be applied to large datasets that might be impractical  
 59 for traditional methods like MCMC [12, 16].

## 60 3 Sources of gradient variance in BBVI

61 Let  $\mathbb{V}_{n, \epsilon}[\nabla f(w; n, \epsilon)]$  denote the variance<sup>1</sup> of the naive estimator from Eq. 3. The two sources for  
 62 this variance correspond to data subsampling ( $n$ ) and Monte Carlo noise ( $\epsilon$ ). It is natural to ask how  
 63 much variance each of these sources contributes. We study this by (numerically) integrating out each  
 64 of these random variables individually and comparing the variances of the resulting estimators.

65 Let  $f(w; n) = \mathbb{E}_\epsilon f(w; n, \epsilon)$  be the objective for a single datum  $n$  with Monte Carlo noise integrated  
 66 out. This can be thought of as an estimator for datum  $n$  with a "perfect" control variate. Similarly,  
 67 let  $f(w; \epsilon) = \mathbb{E}_n f(w; n, \epsilon)$  be the objective for a fixed  $\epsilon$  evaluated on the full dataset. In Fig. 1 we  
 68 generate a single optimization trace using our gradient estimator (described below). Then, for each  
 69 iteration, we estimate the variance of  $\nabla f(w; n, \epsilon)$ ,  $\nabla f(w; \epsilon)$ , and  $\nabla f(w; n)$ <sup>2</sup> using a large number  
 70 of samples. In Table 1 we show the variance at the final iterate on a variety of datasets. (For large  
 71 datasets, it is too expensive to compute the variance this way at each iteration.)

72 Our empirical findings suggest that, despite the exact mix of the two sources being task dependent,  
 73 subsampling noise is usually larger than MC noise. They also show the limits of reducing a single  
 74 source of noise: No control variate applied to each datum could do better than  $\nabla f(w; n)$ , while no  
 75 incremental-gradient-type method could do better than  $\nabla f(w; \epsilon)$ .

<sup>1</sup>For a vector-valued random variable  $z$ , we let  $\mathbb{V}[z] = \text{tr } \mathbb{C}[z]$

<sup>2</sup>Aligned with the experiments in Sec. 7, our evaluation of subsampling variance uses mini-batches, i.e.  $\mathbb{V}_B[\mathbb{E}_{n \in B} \nabla f(w; n)]$ , where  $B$  are mini batches sampled without replacement from  $\{1, \dots, N\}$ .

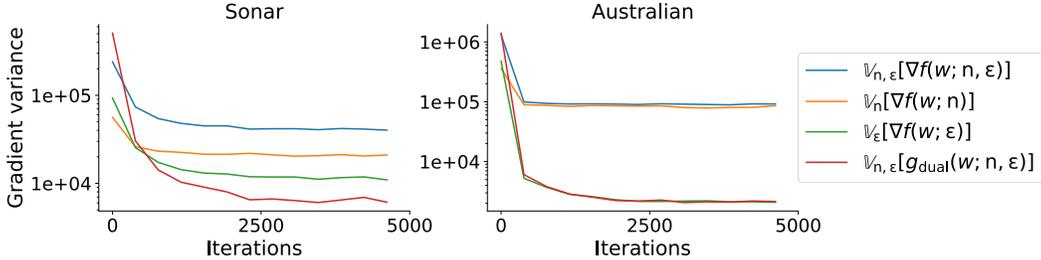


Figure 1: **Gradient Variance Decomposition in Bayesian Logistic Regression using Mean-field BBVI.** The orange line denotes variance from data subsampling ( $n$ ), and the green line denotes Monte Carlo (MC) noise variance ( $\epsilon$ ). For Sonar, both noise sources exhibit similar scales with a batch size of 5. However, for Australian, subsampling noise dominates. Regardless, our proposed gradient estimator  $g_{\text{dual}}$  (red line, Eq. (4)) mitigates subsampling noise and controls MC noise, aligning closely with or below the green line (i.e. the variance without data subsampling) in both datasets.

Task	$\mathbb{V}_{n,\epsilon}[\nabla f(w; n, \epsilon)]$	$\mathbb{V}_n[\nabla f(w; n)]$	$\mathbb{V}_\epsilon[\nabla f(w; \epsilon)]$
Sonar	$4.04 \times 10^4$	$2.02 \times 10^4$	$1.16 \times 10^4$
Australian	$9.16 \times 10^4$	$8.61 \times 10^4$	$2.07 \times 10^3$
MNIST	$4.21 \times 10^8$	$3.21 \times 10^8$	$1.75 \times 10^4$
PPCA	$1.69 \times 10^{10}$	$1.68 \times 10^{10}$	$3.73 \times 10^7$
Tennis	$9.96 \times 10^7$	$9.59 \times 10^7$	$8.56 \times 10^4$

Table 1: BBVI gradient variance decomposition across various tasks, computed at the optimization endpoint. Using a batch size of 100, step size of  $1e-2$  for MNIST, PPCA, and Tennis, and a batch size of 5, step size of  $5e-4$  for Sonar and Australian. We generally observe subsampling noise  $\mathbb{V}_n[\nabla f(w; n)]$  surpassing MC noise  $\mathbb{V}_\epsilon[\nabla f(w; \epsilon)]$ .

## 76 4 Dual Control Variate

77 We now introduce the *dual control variate*, a new approach for controlling the variance of gradient  
 78 estimators for BBVI. Control variates [24] can reduce the variance of a gradient estimator by adding  
 79 a zero-mean random variable that is negatively correlated with the gradient estimator. Considering  
 80 that the objective of BBVI is a function of both  $n$  and  $\epsilon$ , an ideal control variate should also be a  
 81 function of these variables. We take two steps to construct such a control variate.

- 82 1. Inspired by existing control variates for BBVI [19, 9], we create an approximation  $\tilde{f}(w; n, \epsilon)$   
 83 of the true objective  $f(w; n, \epsilon)$ , designed so that the expectation  $\mathbb{E}_\epsilon \nabla \tilde{f}(w; n, \epsilon)$  can easily be  
 84 computed for any datum  $n$ . A common strategy for this is a Taylor-approximation—to replace  
 85  $f$  with a low-order polynomial. Then, if the base distribution  $s(\epsilon)$  is simple, the expectation  
 86  $\mathbb{E}_\epsilon[\nabla \tilde{f}(w; n, \epsilon)]$  is often available in closed-form.
- 87 2. Inspired by SAGA [6], we maintain a table  $W = \{w^1, \dots, w^N\}$  that stores the variational  
 88 parameters at the last iteration each of the data points  $x_1, \dots, x_N$  were accessed. We also  
 89 maintain a running average of gradient estimates evaluated at the stored parameters, denoted by  
 90  $M$ . Unlike SAGA, however, this running average is for the gradients of the *approximation*  $\tilde{f}$ , with  
 91 the Monte Carlo noise  $\epsilon$  integrated out, i.e.  $M = \mathbb{E}_n \mathbb{E}_\epsilon \nabla \tilde{f}(w^n; n, \epsilon)$ .

92 Intuitively, as optimization nears the solution, the weights  $w$  tend to change slowly. This means  
 93 that the entries  $w^n$  in  $W$  will tend to become close to the current iterate  $w$ . Thus, if  $\tilde{f}$  is a good  
 94 approximation of the true objective, we can expect  $\nabla f(w; n, \epsilon)$  to be close to  $\nabla \tilde{f}(w^n; n, \epsilon)$ , meaning  
 95 the two will be strongly correlated. However, thanks to the running average  $M$ , the full expectation  
 96 of  $\nabla \tilde{f}(w^n; n, \epsilon)$  is available in closed-form. This leads to our proposed gradient estimator

$$g_{\text{dual}}(w; n, \epsilon) = \nabla f(w; n, \epsilon) + \underbrace{\mathbb{E}_{\mathbf{m}, \boldsymbol{\eta}} \nabla \tilde{f}(w^{\mathbf{m}}; \mathbf{m}, \boldsymbol{\eta}) - \nabla \tilde{f}(w^n; n, \epsilon)}_{\text{zero mean control variate } c_{\text{dual}}(w; n, \epsilon)}. \quad (4)$$

---

**Algorithm 1** Black-box variational inference with the dual control variate.

---

**Require:** Learning rate  $\lambda$ , variational family  $q_w(z)$ , target  $p(z, x)$

**Require:** Estimator  $f(w; n, \epsilon)$  whose expectation over  $n$  and  $\epsilon$  is the negative ELBO from  $q_w$  and  $p$  (Eq. 2)

**Require:** Approximate estimator  $\tilde{f}(w; n, \epsilon)$  that has an expectation over  $\epsilon$  in closed form

Initialize the parameter  $w_0$ , the parameter table  $W = \{w^1, \dots, w^N\}$

Initialize running mean  $M = \mathbb{E}_m \mathbb{E}_\eta \nabla \tilde{f}(w_0; m, \eta)$   $\triangleright$  Closed-form expectation over  $\eta$ , explicit sum over  $m$

**for**  $k = 1, 2, \dots$  **do**

  Sample  $n$  and  $\epsilon$

  Extract the value of  $w^n$  from the table  $W$

  Compute the base gradient  $g \leftarrow \nabla f(w_k; n, \epsilon)$

  Compute the control variate  $c \leftarrow \mathbb{E}_m \mathbb{E}_\eta \nabla \tilde{f}(w^m; m, \eta) - \nabla \tilde{f}(w^n; n, \epsilon)$  using  $\mathbb{E}_m \mathbb{E}_\eta \nabla \tilde{f}(w^m; m, \eta) = M$

  Update the running mean  $M \leftarrow M + \frac{1}{N} (\mathbb{E}_\eta \nabla \tilde{f}(w_k; n, \eta) - \mathbb{E}_\eta \nabla \tilde{f}(w^n; n, \eta))$   $\triangleright$  Closed-form over  $\eta$

  Update the table  $w^n \leftarrow w_k$

  Update the parameter  $w_{k+1} \leftarrow w_k - \lambda(g + c)$ .  $\triangleright$  Or use  $g + c$  in any stochastic optimization algorithm

**end for**

---

97 The running average  $M = \mathbb{E}_n \mathbb{E}_\epsilon \nabla \tilde{f}(w^n; n, \epsilon)$  can be cheaply maintained through optimization,  
98 since a single value  $w^n$  changes per iteration and  $\mathbb{E}_\epsilon \nabla \tilde{f}(w; n, \epsilon)$  is known in closed form. The  
99 variance of the proposed gradient estimator is given by

$$\mathbb{V}[g_{\text{dual}}] = \mathbb{V}_{\epsilon, n}[\nabla f(w; n, \epsilon) - \nabla \tilde{f}(w^n; n, \epsilon)]. \quad (5)$$

100 Critically, this expression illustrates that the variance of  $g_{\text{dual}}$  can be arbitrarily small, only limited by  
101 how close  $\tilde{f}$  is to  $f$  and how close the stored values  $w^n$  are to the current parameters  $w$ .

102 We illustrate how this gradient estimator can be used for black-box variational inference in Alg. 1.  
103 The same idea could be applied more generally to doubly-stochastic objectives in other domains,  
104 using the more generic version of the algorithm given in Appendix. D.

## 105 5 Variance reduction for stochastic optimization

106 This section considers existing variance reduction techniques and how they compare to the proposed  
107 dual estimator.

### 108 5.1 Monte Carlo sampling and approximation-based control variates

109 Consider the variational objective from Eq. 2 where we sum over the full dataset in each iteration to  
110 define the objective  $f(w) = \mathbb{E}_\epsilon f(w; \epsilon)$ . The gradient estimator obtained by sampling  $\epsilon$  has been  
111 observed to sometimes have problematic variance. Previous work [21, 31, 11, 4] proposed to reduce  
112 this variance by constructing a (zero-mean) control variate  $c(w; \epsilon)$  and defining the new estimator

$$g(w; \epsilon) = \nabla f(w; \epsilon) + c(w; \epsilon). \quad (6)$$

113 The hope is that  $c(w; \epsilon) \approx \nabla f(w) - \nabla f(w; \epsilon)$  approximates the noise of the original estimator,  
114 which can lead to large reductions in variance and thus more efficient and reliable inference.

115 A general way to construct control variates involves using an approximation function  $\tilde{f} \approx f$  for  
116 which the expectation  $\mathbb{E}_\epsilon \tilde{f}(w; \epsilon)$  is available in closed-form [19, 9]. Then, the control variate is  
117 defined as  $c(w; \epsilon) = \mathbb{E}_\eta \nabla \tilde{f}(w; \eta) - \nabla \tilde{f}(w; \epsilon)$ , and the estimator from Eq. (6) becomes

$$g(w; \epsilon) = \nabla f(w; \epsilon) + \mathbb{E}_\eta \nabla \tilde{f}(w; \eta) - \nabla \tilde{f}(w; \epsilon). \quad (7)$$

118 Intuitively, the better  $\tilde{f}$  approximates  $f$ , the lower the variance of this estimator tends to be (for a  
119 perfect approximation, the variance is fully removed). A popular choice for  $\tilde{f}$  involves a quadratic  
120 function, either learned [9] or obtained through a second order Taylor expansion [19], since their  
121 expectation under general Gaussian variational distributions is tractable.

122 In doubly-stochastic problems with objectives of the form  $f(w; n, \epsilon)$ , data  $n$  is subsampled as well as  
 123  $\epsilon$ . While the above control variate has most commonly been used without subsampling, it can also be  
 124 used with subsampling, by developing an approximation  $\tilde{f}(w; n, \epsilon)$  to  $f(w; n, \epsilon)$  for each datum  $n$ .  
 125 This leads to the control variate  $\mathbb{E}_\eta \nabla \tilde{f}(w; n, \eta) - \nabla \tilde{f}(w; n, \epsilon)$  and gradient estimator

$$g_{\text{cv}}(w; n, \epsilon) = \nabla f(w; n, \epsilon) + \underbrace{\mathbb{E}_\eta \nabla \tilde{f}(w; n, \eta) - \nabla \tilde{f}(w; n, \epsilon)}_{\text{zero mean control variate } c_{\text{cv}}(w; n, \epsilon)}. \quad (8)$$

126 It is important to note that this control variate is unable to reduce variance coming from data  
 127 subsampling. Even if  $\tilde{f}(w; n, \epsilon)$  were a *perfect* approximation there would still be gradient variance  
 128 due to  $n$  being sampled randomly. This can be shown by noting that the variance of this estimator is  
 129 given by (see Appendix B.1 for a full derivation using the law of total variance)

$$\mathbb{V}[g_{\text{cv}}] = \mathbb{E}_n \mathbb{V}_\epsilon [\nabla f(w; n, \epsilon) - \nabla \tilde{f}(w; n, \epsilon)] + \mathbb{V}_n [\nabla f(w; n)] \geq \mathbb{V}_n [\nabla f(w; n)]. \quad (9)$$

130 While the first term of the expression above can be made arbitrarily small in the ideal case of a perfect  
 131 approximation  $\tilde{f} \approx f$ , the second term is irreducible, regardless of the quality of the approximation  
 132 used. Therefore, this approach cannot reduce subsampling variance. As shown in Fig. 2 and Table 1,  
 133 subsampling variance is typically substantial, and often several orders of magnitude larger than  
 134 Monte-Carlo variance. When this is true, this control variate, which is only able to reduce variance  
 135 coming from Monte Carlo sampling, will have minimal effect on the overall gradient variance.

## 136 5.2 Data subsampling and incremental gradient methods

137 We now consider a stochastic optimization problem with objective  $f(w) = \mathbb{E}_n f(w; n)$ , where  $n$  is  
 138 uniformly distributed on  $\{1, \dots, N\}$ , representing data indices, but no other stochasticity (i.e. no  
 139 Monte Carlo sampling). While one could compute  $f$  or its gradient exactly, this is expensive when  $N$   
 140 is large. A popular alternative involves drawing a random  $n$  and using the estimator  $\nabla f(w; n)$  with a  
 141 stochastic optimization method, such as stochastic gradient descent. Alternatively, for such problems,  
 142 *incremental gradient* methods [26, 28, 13, 7, 10] often lead to faster convergence.

143 While details vary by algorithm, the basic idea of incremental gradient methods is to "recycle"  
 144 previous gradient evaluations to reduce randomness. For example, SAGA [6] stores the parameters  
 145  $w^n$  of the most recent iteration where  $f(w; n)$  was evaluated and takes a step as

$$w \leftarrow w - \lambda \left( \nabla f(w; n) + \mathbb{E}_m \nabla f(w^m; m) - \nabla f(w^n; n) \right), \quad (10)$$

146 where  $\lambda$  is a step size and the expectation over  $m$  is tracked efficiently using a running average,  
 147 meaning the cost per iteration is independent of  $N$ . The update rule above can be interpreted as using  
 148 a control variate to reduce the variance of the naive estimator  $\nabla f(w; n)$  as

$$g(w; n) = \nabla f(w; n) + \underbrace{\mathbb{E}_m \nabla f(w^m; m) - \nabla f(w^n; n)}_{\text{zero mean control variate}}. \quad (11)$$

149 When  $w^m \approx w$ , the first and last terms in Eq. (11) will approximately cancel, leading to a gradient  
 150 estimator with significantly lower variance.

151 We now consider a doubly-stochastic objective  $f(w; n, \epsilon)$ . In principle, one might consider computing  
 152 the estimator from Eq. (11) for each value of  $\epsilon$ , i.e. using the gradient estimator

$$g_{\text{inc}}(w; n, \epsilon) = \nabla f_n(w; n, \epsilon) + \underbrace{\mathbb{E}_m \nabla f(w^m; m, \epsilon) - \nabla f(w^n; n, \epsilon)}_{\text{zero mean control variate } c_{\text{inc}}(w; n, \epsilon)}. \quad (12)$$

153 This has two issues. First, the resulting method does not address Monte Carlo noise due to sampling  
 154  $\epsilon$ . This can be shown by noting that the variance of this estimator is given by (see Appendix B.2)

$$\mathbb{V}[g_{\text{inc}}] = \mathbb{E}_\epsilon \mathbb{V}_n [\nabla f(w; n, \epsilon) - \nabla f(w^n; n, \epsilon)] + \mathbb{V}_\epsilon [\nabla f(w; \epsilon)] \geq \mathbb{V}_\epsilon [\nabla f(w; \epsilon)]. \quad (13)$$

155 Since the second term in the variance expression above is irreducible, the variance cannot be expected  
 156 to go to zero, no matter how close all the stored vectors  $w^n$  are to the current parameters. Intuitively,  
 157 this approach cannot do better than simply evaluating the objective on the full dataset for a random  $\epsilon$ .

158 The second issue is more critical:  $g_{\text{inc}}$  *cannot be implemented efficiently*. The value of  $\nabla f(w^n; n, \epsilon)$   
 159 is dependent on  $\epsilon$ , which is resampled at each iteration. Therefore, it is not possible to efficiently  
 160 maintain  $\mathbb{E}_m \nabla f(w^m; m, \epsilon)$  needed by Eq. (12) as a running average. In general, this can only  
 161 be computed by looping over the full dataset in each iteration. While possible, this destroys the  
 162 computational advantage of subsampling. For some models with special structure [32, 34] it is  
 163 possible to efficiently maintain the needed running gradient. However, this can only be done in  
 164 special cases with model-specific derivations, breaking the universality of BBVI.

165 It may seem odd that  $g_{\text{inc}}$  has these computational issues, while  $g_{\text{dual}}$ —an estimator intended to  
 166 reduce variance even further—does not. The fundamental reason is that the dual estimator only stores  
 167 (approximate) gradients after integrating over the Monte Carlo variable  $\epsilon$ , so the needed running  
 168 average is independent of  $\epsilon$ .

### 169 5.3 Ensembles of control variate

170 It is possible to combine multiple control variates. For example, [8] combined control variates that  
 171 reduced Monte Carlo noise [19] with one that reduced subsampling noise [32] (for a special case  
 172 where  $g_{\text{inc}}$  is tractable). While this approach can be better than either control variate alone, it still does  
 173 not reduce *joint* variance. To see this, consider a gradient estimator that uses a convex combination of  
 174 the two above control variates. For any  $\beta \in (0, 1)$  write

$$g_{\text{combo}}(w; n, \epsilon) = \nabla f(w; n, \epsilon) + \underbrace{\beta c_{\text{cv}}(w; n, \epsilon) + (1 - \beta) c_{\text{inc}}(w; n, \epsilon)}_{c_{\text{combo}}(w; n, \epsilon)}. \quad (14)$$

175 It can be shown (Appendix B.3) that if both  $c_{\text{cv}}$  and  $c_{\text{inc}}$  are "perfect", that is, if  $\tilde{f}(w; n, \epsilon) =$   
 176  $f(w; n, \epsilon)$  and  $w^n = w$  for all  $n$ , then

$$\mathbb{V}[g_{\text{combo}}] = \beta^2 \mathbb{V}_n[\nabla f(w; n)] + (1 - \beta)^2 \mathbb{V}_\epsilon[\nabla f(w; \epsilon)]. \quad (15)$$

177 Even in this idealized scenario, such an estimator cannot reduce variance to zero, because each of  
 178 the individual control variates leaves one source of noise uncontrolled. The dual control variate  
 179 overcomes this because it models interactions between  $\epsilon$  and  $n$ .

## 180 6 Related work

181 Recent work proposed to approximate the optimal batch-dependent control variate for BBVI using  
 182 a recognition network [4]. Similar to our work, they take into account the usage of subsampling  
 183 when designing their variance reduction techniques for BBVI. However, like  $g_{\text{cv}}$ , their control variate  
 184 reduces the *conditional* variance of MC noise (conditioned on  $n$ ) but is unable to reduce subsampling  
 185 noise (like  $g_{\text{cv}}$ ).

186 It is also worth discussing a special incremental gradient method called SMISO [1], designed  
 187 for doubly-stochastic problems. Intuitively, SMISO uses exponential averaging to approximately  
 188 marginalize out  $\epsilon$ , and then runs MISO/Finito [7, 18] (an incremental gradient method similar to  
 189 SAGA) to control the subsampling noise. While the method is similar to running SGD with an  
 190 incremental control variate, it is not obvious how to separate the control variate from the algorithm,  
 191 meaning we cannot use the SMISO idea as a control variate to get a gradient estimator that can be  
 192 used with other optimizers like Adam, we include a detailed discussion on this issue in Appendix A.  
 193 Nevertheless, we still include SMISO as one of our baselines.

## 194 7 Experiments

195 This section empirically demonstrates the effectiveness of the dual control variate for BBVI. We  
 196 focus on mean-field Gaussian BBVI, where the variational posterior follows a multivariate Gaussian  
 197 with diagonal covariance  $q_w(z) = \mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$ , with parameters  $w = (\boldsymbol{\mu}, \log(\boldsymbol{\sigma}))$ .

198 The gradient estimators  $g_{\text{cv}}(w; n, \epsilon)$  and  $g_{\text{dual}}(w; n, \epsilon)$  require an approximation function with  
 199 expectation over  $\epsilon$  available in closed form. Inspired by previous work [19], we get an approximation  
 200 for  $f(w; n, \epsilon)$  using a second order Taylor expansion for the negative total likelihood  $k_n(z) =$

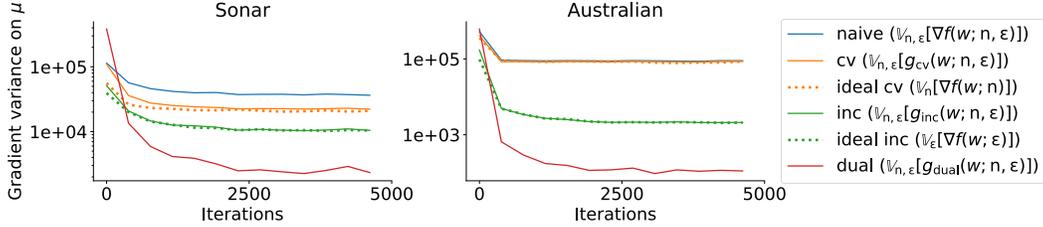


Figure 2: **Dual control variate helps reduce gradient variance.** The naive gradient estimator (Eq. (3)) is the baseline, while the cv estimator (Eq. (8)) controls the Monte Carlo noise, the inc estimator (Eq. (12)) controls for subsampling noise, and the proposed dual estimator (Eq. (4)) controls for both. The variance of cv and inc, as is shown in Eq. (9) and Eq. (13) are lower-bounded by the dotted lines, while dual is capable of reducing the variance to significantly lower values, leading to better and faster convergence (Fig. 3).

201  $N \log p(x_n | z) + \log p(z)$  around  $z_0 = \mathcal{T}_w(0)$ <sup>3</sup>, which yields

$$\tilde{f}(w; n, \epsilon) = k_n(z_0) + (\mathcal{T}_w(\epsilon) - z_0)^\top \nabla k_n(z_0) + \frac{1}{2} (\mathcal{T}_w(\epsilon) - z_0)^\top \nabla^2 k_n(z_0) (\mathcal{T}_w(\epsilon) - z_0) + \mathbb{H}(w), \quad (16)$$

202 where we assume the entropy can be computed in closed-form. For a mean-field Gaussian variational  
 203 distribution, the expected gradient of the approximation Eq. (16) can only be computed efficiently  
 204 (via Hessian-vector products) with respect to the mean parameter  $\mu$  but not for the scale parameter  
 205  $\sigma$ , which means  $g_{cv}(w; n, \epsilon)$  and  $g_{dual}(w; n, \epsilon)$  can only be used as the gradient estimator for  $\mu$ .  
 206 Fortunately, controlling only the gradient variance on  $\mu$  often means controlling most of the variance,  
 207 as, with mean-field Gaussians, the total gradient variance is often dominated by variance from  $\mu$  [9].

## 208 7.1 Experiment setup

209 We evaluate our methods by performing BBVI on a range of tasks: binary Bayesian logistic regression  
 210 on two datasets, Sonar (number of samples  $N = 208$ , dimensionality  $D = 60$ ) and Australian ( $N =$   
 211  $690$ ,  $D = 14$ ); multi-class Bayesian logistic regression on MNIST [17] ( $N = 60000$ ,  $D = 7840$ );  
 212 probabilistic principal component analysis [29] (PPCA,  $N = 60000$ ,  $D = 12544$ ); and Bradley-  
 213 Terry model [5] for tennis player ranking (Tennis,  $N = 169405$ ,  $D = 5525$ ). We give full model  
 214 descriptions in Sec. 7.3.

215 **Baselines.** We compare  $g_{dual}$  (Eq. (4)) with  $g_{naive}$  (Eq. (3)) and  $g_{cv}$  (Eq. (8)). For Sonar and  
 216 Australian (small datasets) we include  $g_{inc}$  (Eq. (12)) as well, which requires a full pass through the  
 217 full dataset at each iteration. For larger-scale problems,  $g_{inc}$  becomes intractable, so we use SMISO  
 218 instead.

219 **Optimization details.** We optimize using Adam [14] for the larger-scale MNIST, PPCA, and  
 220 Tennis datasets and SGD without momentum for the small-scale Sonar and Australian dataset for  
 221 transparency. The optimizer for SMISO is pre-determined by its algorithmic structure and cannot  
 222 be changed. For all estimators, we perform a step-size search (see Appendix C) to ensure a fair  
 223 comparison and use a single shared  $\epsilon$  for all samples in the batch.

224 **Mini-batching.** In practice, for efficient implementation on GPUs, we draw a mini-batch  $B$  of data at  
 225 each iteration (reshuffling for each epoch). For inc, dual, and SMISO, we update multiple entities  
 226 in the parameter table per iteration and adjust the running mean accordingly. For the Sonar and  
 227 Australian datasets, due to their small sizes, we use  $|B| = 5$ . For other datasets we use  $|B| = 100$ .

228 **Evaluation metrics.** We track the ELBO on the full dataset, explicitly computing  $\mathbb{E}_n$  (summing  
 229 over the full dataset) and approximating  $\mathbb{E}_\epsilon$  with 5000 Monte Carlo samples. We present ELBO  
 230 vs. iterations plots for a single example learning rate as well as ELBO values for the best learning  
 231 rate chosen retrospectively for each iteration. In addition, we present the final ELBO after training  
 232 vs. step size at different iterations. For the Sonar and Australian datasets, given the small size, we  
 233 include a detailed trace of gradient variance on  $\mu$  across different estimators. This enables empirical  
 234 validation of the lower bounds derived in Eq. (9) and Eq. (13).

<sup>3</sup>We use  $z_0 = \text{stop\_gradient}(\mathcal{T}_w(0))$  so that the gradient does not backpropagate from  $z_0$  to  $w$ .

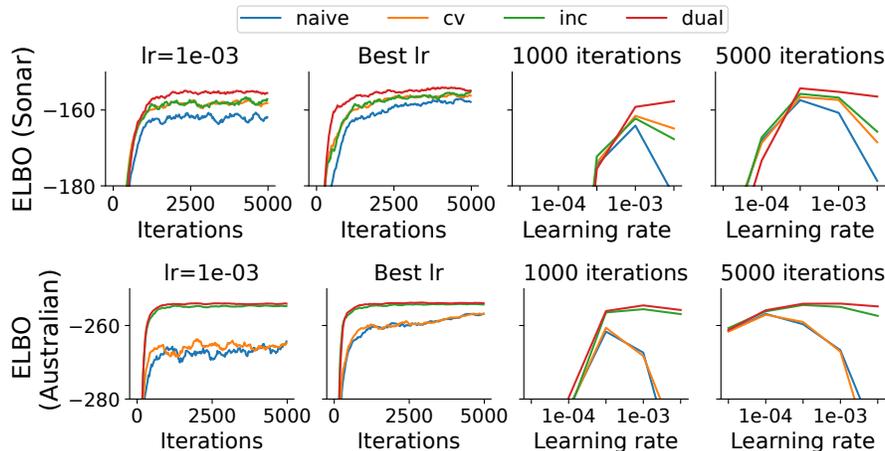


Figure 3: **With reduced variance (Fig. 2), the dual estimator provides better convergence at a larger step size.** On Sonar, Monte Carlo noise and subsampling noise are of similar scale, therefore jointly controlling them shows better performance than methods that only control one source of noise. On Australian, where the subsampling noise dominates, dual shows similar performance compared with inc, which controls subsampling noise but *cannot* be efficiently computed (requires pass over the full dataset at each iteration).

235 **Initialization.** The variational parameters are randomly initialized using a standard Gaussian and  
 236 all results reported are averages over multiple independent trials: We run 10 trials for Sonar and  
 237 Australian, and 5 for the larger scale problems due to resource constraint.

## 238 7.2 Results

239 The experiment results for Sonar and Australian are presented in Fig. 2 and Fig. 3. Both the inc and  
 240 cv estimators have lower variance than the naive estimator, but the improvement varies by the dataset.  
 241 The excellent performance of the (impractical) inc estimator on Australian shows the importance  
 242 of reducing subsampling noise. Overall, the dual estimator has the lowest variance, which enables  
 243 larger learning rates and thus faster optimization.

244 Similar results can be observed on MNIST, PPCA, and Tennis in Fig. 4 (for these datasets inc  
 245 is intractable, so we include SMISO as a baseline instead). Again, dual yields faster and better  
 246 convergence than naive and cv. Whereas SMISO, which does not adopt momentum nor adaptive step  
 247 size, suffers from slow convergence speed in that it has to utilize a small step size to prevent diverging  
 248 during optimization. We provide comparisons of different estimators using SGD in Appendix. E.

## 249 7.3 Model descriptions

250 **Binary/Multi-class Bayesian logistic regression.** A standard logistic regression model with standard  
 251 Gaussian prior.

252 **Probabilistic principal component analysis (PPCA).** Given a centered dataset  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^D$ ,  
 253 PPCA [29] seeks to extract its principal axes  $\mathbf{W} \in \mathbb{R}^{D \times K}$  by assuming  $\mathbf{x}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{W}\mathbf{W}^\top +$   
 254  $\text{diag}(\lambda^2))$ . In our experiments, we employ a standard Gaussian prior on  $\mathbf{W}$  and use BBVI to  
 255 approximate the posterior over  $\mathbf{W}$ . We then test PPCA on the standardized training set of MNIST  
 256 with  $K = 16$  and  $\lambda = 1$ .

257 **Bradley Terry model (Tennis).** This is a model used to rank players from pair-wise matches.  
 258 Each player is represented by a score  $\theta_i$ , and each score is assigned a standard Gaussian prior. The  
 259 result of a match between two players is modeled by the inverse logit of their score difference  
 260  $y_n \sim \text{Bernoulli}(\text{logit}^{-1}(\theta_i - \theta_j))$  where  $y_n = 1$  denotes a win by player  $n$ . We subsample over  
 261 matches and perform inference over the score of each player. We evaluate the model on men’s tennis  
 262 matches log starting from 1960, which contains the results of 169405 matches among 5525 players.

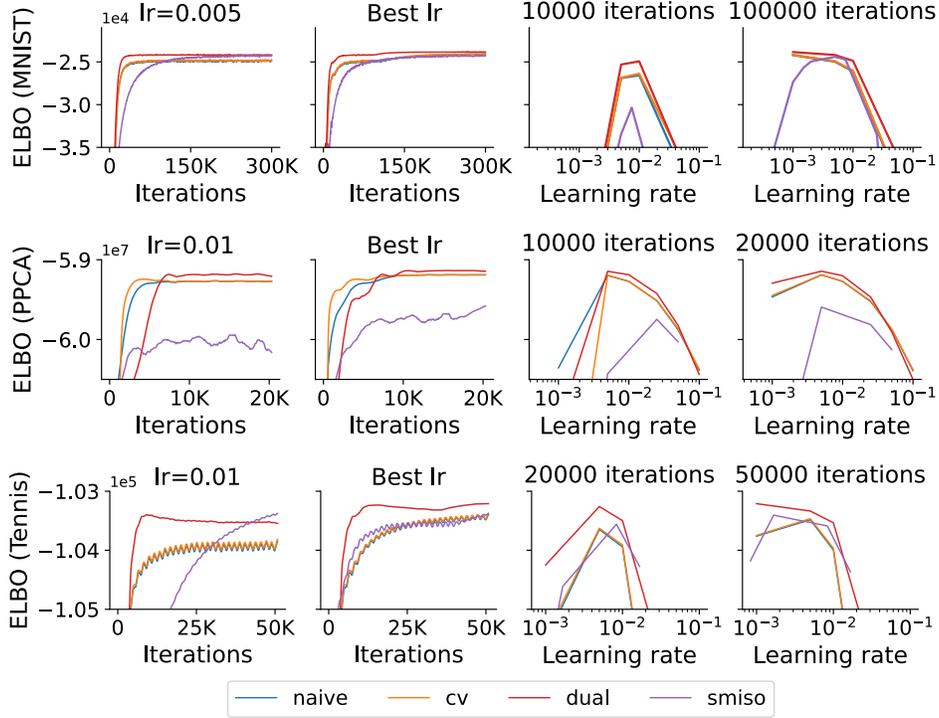


Figure 4: **On larger scale problems, the dual estimator leads to improved convergence.** In large-scale problems, cv shows little or no improvement upon naive while dual converges faster. We suspect that most of the improvement in the dual estimator comes from reducing subsampling variance. SMISO shows slow convergence. We suspect that is because it is an “SGD-type” algorithm while all others use Adam. Note that the step size for SMISO is rescaled for visualization. The loss shows periodic structure in Tennis, this happens because gradients have correlated noise that cancels out at the end of each epoch.

Estimator	Variance lower bound	$\nabla f$ evals per iteration	Wall-clock time per iteration		
			MNIST	PPCA	Tennis
naive	$\mathbb{V}_{n,\epsilon}[\nabla f(w; n, \epsilon)]$	1	10.4ms	12.8ms	10.2ms
cv	$\mathbb{V}_n[\nabla f(w; n)]$	2	12.8ms	18.5ms	14.6ms
inc	$\mathbb{V}_\epsilon[\nabla f(w; \epsilon)]$	N+2	328ms	897ms	588ms
dual	0	3	17.6ms	31.2ms	29.6ms
Fullbatch-naive	$\mathbb{V}_\epsilon[\nabla f(w; \epsilon)]$	N	201ms	740ms	203ms
Fullbatch- $c_{cv}$	0	2N	360ms	1606ms	246ms

Table 2: Variance, oracle complexity, and wall-clock time for different estimators. Notice that inc is more expensive than Fullbatch-naive. We hypothesize this is because inc uses separate  $w^n$  for different data points, which is less efficient for parallelism.

## 263 7.4 Efficiency analysis

264 We now study the computational cost of different estimators. In terms of the number of "oracle"  
 265 evaluations (i.e. evaluations of  $f(w; n, \epsilon)$  or its gradient), the naive estimator is the most efficient,  
 266 requiring a single oracle evaluation per iteration. The cv estimator requires one gradient and also  
 267 one Hessian-vector product, while the dual estimator needs one gradient and two Hessian-vector  
 268 products, one for the control variate and one for updating the running mean  $M$ .

269 Additionally, Table 2 shows measured runtimes based on a JAX implementation on an Nvidia 2080ti  
 270 GPU. All numbers are for a single optimization step, averaged over 200 steps. Overall, each iteration  
 271 with the dual estimator is between 1.5 to 2.5 times slower than naive, and around 1.2 times slower  
 272 than cv. Lastly, given that dual achieves a given performance in an order of magnitude fewer  
 273 iterations (Figs. 3 and 4), it is the fastest in terms of wall-clock time. The exact wall-clock time v.s.  
 274 ELBO results are presented in Appendix. F.

275 **References**

- 276 [1] Alberto Bietti and Julien Mairal. Stochastic optimization with variance reduction for infinite  
277 datasets with finite sum structure. *Advances in Neural Information Processing Systems*, 30:1623–  
278 1633, 2017.
- 279 [2] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for  
280 statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- 281 [3] Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine  
282 learning. *Siam Review*, 60(2):223–311, 2018.
- 283 [4] Ayman Boustati, Sattar Vakili, James Hensman, and ST John. Amortized variance reduction  
284 for doubly stochastic objective. In *Conference on Uncertainty in Artificial Intelligence*, pages  
285 61–70. PMLR, 2020.
- 286 [5] Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the  
287 method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- 288 [6] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. Saga: A fast incremental gradient  
289 method with support for non-strongly convex composite objectives. In *Advances in neural  
290 information processing systems*, pages 1646–1654, 2014.
- 291 [7] Aaron Defazio, Justin Domke, et al. Finito: A faster, permutable incremental gradient method  
292 for big data problems. In *International Conference on Machine Learning*, pages 1125–1133.  
293 PMLR, 2014.
- 294 [8] Tomas Geffner and Justin Domke. Using large ensembles of control variates for variational  
295 inference. In *Advances in Neural Information Processing Systems*, pages 9982–9992, 2018.
- 296 [9] Tomas Geffner and Justin Domke. Approximation based variance reduction for reparameteriza-  
297 tion gradients. *Advances in Neural Information Processing Systems*, 33, 2020.
- 298 [10] Robert M Gower, Mark Schmidt, Francis Bach, and Peter Richtárik. Variance-reduced methods  
299 for machine learning. *Proceedings of the IEEE*, 108(11):1968–1983, 2020.
- 300 [11] Will Grathwohl, Dami Choi, Yuhuai Wu, Geoff Roeder, and David Duvenaud. Backpropagation  
301 through the void: Optimizing control variates for black-box gradient estimation. In *International  
302 Conference on Learning Representations*, 2018.
- 303 [12] Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley. Stochastic variational  
304 inference. *Journal of Machine Learning Research*, 2013.
- 305 [13] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance  
306 reduction. *Advances in neural information processing systems*, 26:315–323, 2013.
- 307 [14] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint  
308 arXiv:1412.6980*, 2014.
- 309 [15] Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. In *International  
310 Conference on Learning Representations*, 2014.
- 311 [16] Alp Kucukelbir, Dustin Tran, Rajesh Ranganath, Andrew Gelman, and David M Blei. Automatic  
312 differentiation variational inference. *The Journal of Machine Learning Research*, 18(1):430–474,  
313 2017.
- 314 [17] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning  
315 applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 316 [18] Julien Mairal. Incremental majorization-minimization optimization with application to large-  
317 scale machine learning. *SIAM Journal on Optimization*, 25(2):829–855, 2015.
- 318 [19] Andrew C Miller, Nicholas J Foti, Alexander D’Amour, and Ryan P Adams. Reducing  
319 reparameterization gradient variance. *Advances in Neural Information Processing Systems*,  
320 2017:3709–3719, 2017.

- 321 [20] Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic  
322 approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–  
323 1609, 2009.
- 324 [21] John Paisley, David M Blei, and Michael I Jordan. Variational bayesian inference with stochastic  
325 search. In *Proceedings of the 29th International Conference on International Conference on*  
326 *Machine Learning*, pages 1363–1370, 2012.
- 327 [22] Rajesh Ranganath, Sean Gerrish, and David Blei. Black box variational inference. In *Artificial*  
328 *intelligence and statistics*, pages 814–822. PMLR, 2014.
- 329 [23] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation  
330 and approximate inference in deep generative models. In *International conference on machine*  
331 *learning*, pages 1278–1286. PMLR, 2014.
- 332 [24] Christian P Robert, George Casella, and George Casella. *Monte Carlo statistical methods*,  
333 volume 2. Springer, 1999.
- 334 [25] Geoffrey Roeder, Yuhuai Wu, and David K Duvenaud. Sticking the landing: Simple, lower-  
335 variance gradient estimators for variational inference. In I. Guyon, U. Von Luxburg, S. Bengio,  
336 H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Informa-*  
337 *tion Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 338 [26] Nicolas Roux, Mark Schmidt, and Francis Bach. A stochastic gradient method with an expo-  
339 nential convergence \_rate for finite training sets. *Advances in neural information processing*  
340 *systems*, 25, 2012.
- 341 [27] Hugh Salimbeni and Marc Deisenroth. Doubly stochastic variational inference for deep gaussian  
342 processes. *Advances in neural information processing systems*, 30, 2017.
- 343 [28] Shai Shalev-Shwartz and Tong Zhang. Stochastic dual coordinate ascent methods for regularized  
344 loss minimization. *Journal of Machine Learning Research*, 14(2), 2013.
- 345 [29] Michael E Tipping and Christopher M Bishop. Probabilistic principal component analysis.  
346 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 61(3):611–622,  
347 1999.
- 348 [30] Michalis Titsias and Miguel Lázaro-Gredilla. Doubly stochastic variational bayes for non-  
349 conjugate inference. In *International conference on machine learning*, pages 1971–1979. PMLR,  
350 2014.
- 351 [31] George Tucker, Andriy Mnih, Chris J Maddison, Dieterich Lawson, and Jascha Sohl-Dickstein.  
352 Rebar: low-variance, unbiased gradient estimates for discrete latent variable models. In  
353 *Proceedings of the 31st International Conference on Neural Information Processing Systems*,  
354 pages 2624–2633, 2017.
- 355 [32] Chong Wang, Xi Chen, Alexander J Smola, and Eric P Xing. Variance reduction for stochastic  
356 gradient optimization. *Advances in neural information processing systems*, 26, 2013.
- 357 [33] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforce-  
358 ment learning. *Reinforcement learning*, pages 5–32, 1992.
- 359 [34] Shuai Zheng and James Tin-Yau Kwok. Lightweight stochastic optimization for minimizing  
360 finite sums with infinite data. In *International Conference on Machine Learning*, pages 5932–  
361 5940. PMLR, 2018.