
Variational Inference with Gaussian Score Matching

Anonymous Author(s)

Affiliation

Address

email

Abstract

Variational inference (VI) is a method to approximate the computationally intractable posterior distributions that arise in Bayesian statistics. Typically, VI fits a simple parametric distribution to be close to the target posterior, minimizing an appropriate objective such as the evidence lower bound (ELBO). In this work, we present a new approach to VI. Our method is based on the principle of score matching, that if two distributions are equal then their score functions (i.e., gradients of the log density) are equal at every point on their support. With this principle, we develop score matching VI, an iterative algorithm that seeks to match the scores between the variational approximation and the exact posterior. At each iteration, score matching VI solves an inner optimization, one that minimally adjusts the current variational estimate to match the scores at a newly sampled value of the latent variables. We show that when the variational family is a Gaussian, this inner optimization enjoys a closed form solution, which we call Gaussian score matching VI (GSM-VI). GSM-VI is also a “black box” variational algorithm in that it only requires a differentiable joint distribution, and as such it can be applied to a wide class of models. We compare GSM-VI to black box variational inference (BBVI), which has similar requirements but instead optimizes the ELBO. We first study how GSM-VI behaves as a function of the problem dimensionality, the condition number of the target covariance matrix (when the target is Gaussian), and the degree of mismatch between the approximating and exact posterior distribution. We then study GSM-VI on a collection of real-world Bayesian inference problems from the posteriorDB database of datasets and models. In all of our studies we find that GSM-VI is faster than BBVI, but without sacrificing accuracy. It requires 10-100x fewer gradient evaluations to obtain a comparable quality of approximation.

1 Introduction

This paper is about variational inference for approximate Bayesian computation. Consider a statistical model $p(\theta, \mathbf{x})$ of parameters $\theta \in \mathbb{R}^d$ and observations \mathbf{x} . Bayesian inference aims to infer the posterior distribution $p(\theta | \mathbf{x})$, which is often intractable to compute. Variational inference is an optimization-based approach to approximate the posterior [4, 16].

The idea behind VI is to approximate the posterior with a member of a *variational family* of distributions $q_w(\theta)$, parameterized by *variational parameters* w [4, 16]. Specifically, VI methods establish a measure of closeness between $q_w(\theta)$ and the posterior, and then minimize it with an optimization algorithm. Researchers have explored many aspects of VI, including different objectives [7, 8, 18, 22, 23, 25, 30] and optimization strategies [1, 13, 24].

In its modern form, VI typically minimizes $\text{KL}(q_w(\theta) || p(\theta | \mathbf{x}))$ with stochastic optimization, and further satisfies the so-called “black-box” criteria [1, 24, 29]. The resulting black-box VI (BBVI) only

requires the practitioner to specify the log joint $\log p(\theta, x)$ and (often) its gradient $\nabla_{\theta} \log p(\theta, x)$, which for many models can be obtained by automatic differentiation. For these reasons, BBVI has been widely implemented, and it is available in many probabilistic programming systems [3, 19, 27].

In this paper, we propose a new approach to VI. We begin with the principle of *score matching* [14], that when two densities are equal then their gradients are equal as well, and we use this principle to derive a new way to fit a variational distribution to be close to the exact posterior. The result is *score-matching VI*. Rather than explicitly minimize a divergence, score-matching VI iteratively projects the variational distribution onto the exact score matching constraint. This strategy enables a new black-box VI algorithm.

Score-matching VI relies on the same ingredients as reparameterization BBVI [19]—a differentiable variational family and a differentiable log joint—and so it can be as easily incorporated into probabilistic programming systems as well. Further, when the variational family is a Gaussian, score-matching VI is particularly efficient: each iteration is computable in closed form. We call the resulting algorithm Gaussian score matching VI (GSM-VI).

Unlike BBVI, GSM-VI does not rely on stochastic gradient descent (SGD) for its core optimization. Though SGD has the appeal of simplicity, it is also known to require the careful tuning of learning rates. GSM-VI was inspired by a different tradition of constraint-based algorithms for online learning [2, 6, 11, 12, 20]. These algorithms have been extensively developed and analyzed for problems in classification, and under the right conditions, they have been observed to outperform SGD. This paper shows how to extend this constraint-based framework—and the powerful machinery behind it—from the problem of classification to the workhorse of Gaussian VI. The key insight is that score-matching (unlike ELBO maximization) lends itself naturally to a constraint-based formulation.

We empirically compared GSM-VI to reparameterization BBVI on several classes of models, and with both synthetic and real-world data. In general, we found that GSM requires 10-100x fewer gradient evaluations to converge to an equally good approximation. When the exact posterior is Gaussian, we found GSM-VI scales significantly better with respect to dimensionality and is insensitive to the condition number of the target covariance. When the exact posterior is non-Gaussian, we found GSM-VI enjoys faster convergence without sacrificing the quality of the final approximation.

This paper makes the following contributions:

- We introduce *score matching variational inference*, a new black-box approach to fitting $q_w(\theta)$ to be close to $p(\theta | x)$. Score matching VI requires no tunable optimization hyperparameters, to which BBVI can be sensitive.
- When the variational family is Gaussian, we develop *Gaussian score matching variational inference* (GSM-VI). It establishes efficient closed-form iterates for score matching VI.
- We empirically compare GSM-VI to reparameterization BBVI. Across many models and datasets, we found that GSM-VI enjoys faster convergence to an equally good approximation.

We develop score matching VI in Section 2 and study its performance in Section 3.

Related work. Our work introduces a new method for black-box variational inference that relies only on having access to the gradients of the variational distribution and the log joint. GSM-VI has similar goals to automatic-differentiation variational inference (ADVI) [19] and Pathfinder [31], which also fit multivariate Gaussian variational families, but do so by maximizing the ELBO using stochastic optimization. Similar to GSM-VI, the algorithm of ref. [28] also seeks to match the scores of the variational and the target posterior, but it does so by minimizing the L2 loss between them.

A novel aspect of this work is how GSM-VI fits the variational parameters. Rather than minimize a loss function, it aims to solve a set of nonlinear equations. Similar ideas have been pursued in the context of fitting a model to data using empirical risk minimization (ERM). For example, passive aggressive (PA) methods [6] and the stochastic polyak stepsize (SPS) are also derived via projections onto sampled nonlinear equations [2, 12, 20]. A probabilistic extension of PA methods is known as confidence-weighted (CW) learning [11]. In this framework, the learner maintains a multivariate Gaussian distribution over the weight vector of a linear classifier. Like CW learning, the second

87 step of GSM-VI also minimizes a KL divergence between multivariate Gaussians. But it involves a
 88 different projection, one of score-matching versus linear classification.

89 2 Score Matching Variational Inference

90 Suppose for the moment that the variational family $q_w(\theta)$ is rich enough to perfectly capture the
 91 posterior $p(\theta | x)$. That is, there exists a w^* such that

$$\log q_{w^*}(\theta) = \log p(\theta | x), \quad \forall \theta. \quad (1)$$

92 If we could solve Eq. 1 for w^* , the resulting variational distribution would be a perfect fit. The
 93 challenge is that the posterior on the right side is intractable to compute.

94 To help, we appeal to score matching [14]. Define the score of a distribution to be the gradient of its
 95 log with respect to the variable¹, e.g., $\nabla_\theta \log q_w(\theta)$. The principle of score matching is that if two
 96 distributions are equal at each point in their support then their score functions are also equal.

97 To use score matching for VI, we first write the log posterior as the log joint minus the normalizing
 98 constant, i.e., the marginal distribution of x ,

$$\log p(\theta | x) = \log p(\theta, x) - \log p(x). \quad (2)$$

99 With this expression, the principle of score matching leads to the following Lemma.

Lemma 2.1. The parameter w^* satisfies

$$\nabla_\theta \log q_{w^*}(\theta) = \nabla_\theta \log p(\theta, x), \quad \forall \theta, \quad (3)$$

if and only if w^* also satisfies Eq. 1.

100 What is notable about Eq. 3 is that the right side is the gradient of the log joint. Unlike the posterior,
 101 the gradient of the log joint is tractable to compute for a large class of probabilistic models. (The
 102 proof is in the appendix.)

103 This lemma motivates a new algorithm, *score matching VI*. The idea is to iteratively refine the
 104 variational parameters w to try to solve the system of equations in Eq. 3 as well as possible. At each
 105 iteration t , it first samples a new θ_t from the current variational approximation and then minimally
 106 adjusts w to satisfy Eq. 3 for that value of θ_t .

Score matching variational inference

At iteration t :

1. Sample $\theta_t \sim q_{w_t}(\theta)$.
2. Update the variational parameters:

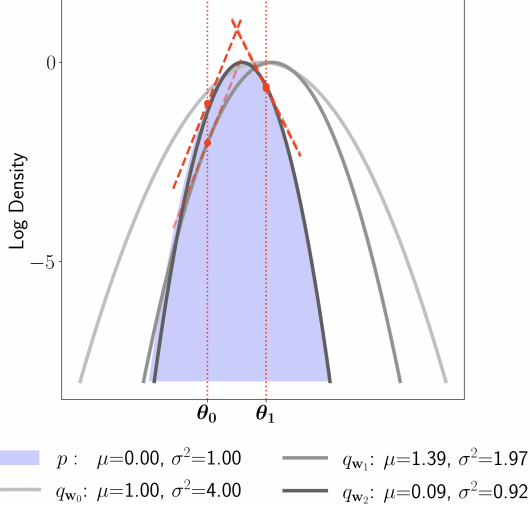
$$w_{t+1} = \arg \min_w \text{KL}(q_{w_t}(\theta) || q_w(\theta)) \quad (4)$$

such that $\nabla_\theta \log q_w(\theta_t) = \nabla_\theta \log p(\theta_t, x)$.

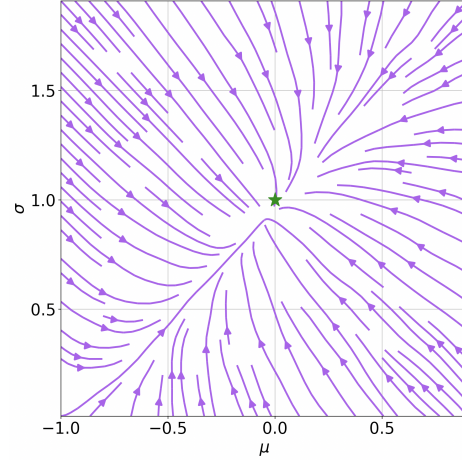
107 This algorithm for score matching VI was inspired by earlier online algorithms for learning a classifier
 108 from a stream of labeled examples. One particularly elegant algorithm in this setting is known as
 109 passive-aggressive (PA) learning [6], in which a model is incrementally updated by the minimal
 110 amount to classify each example correctly by a large margin. This approach was subsequently
 111 extended to a probabilistic setting, known as confidence-weighted (CW) learning [11] in which one
 112 minimally updates a *distribution* over classifiers. Our algorithm is similar in that it minimally updates
 113 an approximating distribution for VI, but it is different in that enforces constraints for score matching
 114 instead of large margin classification.

115 At a high level, what makes this approach to VI likely to succeed or fail? Certainly it is necessary
 116 that there are more variational parameters than elements of the latent variable θ ; when this is not the

¹We make this clear because, in some literature, the score is the gradient with respect to the parameter.



(a) Two iterations of GSM-VI. The log density of the target posterior p is shaded blue; the initial distribution q_{w_0} is light grey; the first update q_{w_1} is medium grey; and the second update q_{w_2} is dark grey.



(b) The vector field of Eq. 4, averaged over 5 independent samples, where $p(\theta | \mathbf{x}) = \mathcal{N}(0, 1)$. The solution $(\mu, \sigma) = (0, 1)$ is the green star.

case, it may be impossible to satisfy a *single* score matching constraint in Eq. 4. That said, setting the number of variational parameters to be at least as large as the latent variable is standard, as in, e.g., a factorized (or mean-field) variational family. It is also apparent that the algorithm may never converge if the target posterior is not contained in the variational family, or that the variational approximation collapses to a point mass, which stalls the updates altogether. While we cannot dismiss these possibilities out of hand, we did not see either issue in any of the empirical studies of Section 3.

For more intuition, Figure 1a illustrates the effect of the update in Eq. 4 when both the target and approximating distribution are 1d Gaussian. The target posterior $p(\theta | \mathbf{x})$ is shaded blue. The plot shows the initial variational distribution q_{w_0} (light grey curve) and its update to q_{w_1} (medium grey) so that the gradient of the updated distribution matches the gradient of the target at the sampled θ_0 (dotted red tangent line). It also shows the update from w_1 to w_2 , now matching the gradient at θ_1 . With these two updates, q_{w_2} (dark grey) is very close to the target $p(\theta | \mathbf{x})$. With this picture in mind, we now develop the details of this algorithm for a widely applicable setting.

Gaussian Score Matching VI. Suppose the variational distribution belongs to a multivariate Gaussian family $q_w(\theta) := \mathcal{N}(\mu, \Sigma)$, which is a common setting especially in systems for automated approximate inference [1, 19]. One of our main contributions is to show that in this case Eq. 4 has a closed form solution. The solution $w_{t+1} = (\mu_{t+1}, \Sigma_{t+1})$ has the following form:

$$\mu_{t+1} = \mu_t + A_t (\nabla_{\theta} \log p(\theta_t, \mathbf{x}) - \nabla_{\theta} \log q_{w_t}(\theta_t)) \quad (5)$$

$$\Sigma_{t+1} = \Sigma_t + (\mu_t - \theta_t)(\mu_t - \theta_t)^{\top} - (\mu_{t+1} - \theta_t)(\mu_{t+1} - \theta_t)^{\top} \quad (6)$$

where $A_t \in \mathbb{R}^{d \times d}$ is a matrix defined in the theorem below. These exact updates only require the score of the log joint $\nabla_{\theta} \log p(\theta, \mathbf{x})$ and the score of the variational distribution $\nabla_{\theta} \log q_w(\theta)$.

Eqs. 5 and 6 also provide intuition. Consider the approximation at the t th iteration q_{w_t} and the current sample θ_t . First suppose the scores already match at this sample, that is $\nabla_{\theta} \log p(\theta_t, \mathbf{x}) = \nabla_{\theta} \log q_{w_t}(\theta_t)$. Then the mean does not change $\mu_{t+1} = \mu_t$ and, similarly, the two rank-one terms in the covariance update in Eq. 6 cancel out so $\Sigma_{t+1} = \Sigma_t$. This shows that when $q_{w_t}(\theta) = p(\theta, \mathbf{x})$ for all θ , the method stops. On the other hand, if the scores do not match, then the mean is updated proportionally to the difference between the scores, and the covariance is updated by a rank-two correction. For a one dimensional target $p(\theta, \mathbf{x}) = \mathcal{N}(0, 1)$, Figure 1b illustrates the vector field of updates. The vector field points to the solution (green star) and, once there, the method stops.

We now formalize this result and give the exact expression for A_t .

Algorithm 1: Gaussian Score Matching VI

Input : Initial mean estimate μ_0 , initial covariance estimate Σ_0 , target distribution $p(\theta|x)$,
number of iterations $N \in \mathbb{N}$, batch size $B \in \mathbb{N}$.

Output : Multivariate normal variational distribution $q_w(\theta) := \mathcal{N}(\mu, \Sigma)$

for $i = 0, \dots, N - 1$ \triangleright iteration loop

do

for $j = 0, \dots, B - 1$ \triangleright batch loop

do

Sample $\theta^{(j)} \sim \mathcal{N}(\mu_i, \Sigma_i)$

$g \leftarrow \nabla_{\theta} \log p(\theta^{(j)}|x)$

$\varepsilon \leftarrow \Sigma_i g - \mu_i + \theta$

Solve $\rho(1+\rho) = g^{\top} \Sigma_i g + [(\mu_i - \theta)^{\top} g]^2$ for $\rho > 0$

$\delta \mu^{(j)} \leftarrow \frac{1}{1+\rho} \left[\mathbf{I} - \frac{(\mu_i - \theta) g^{\top}}{1+\rho + (\mu_i - \theta)^{\top} g} \right] \varepsilon$

$\mu_i^{(j)} \leftarrow \mu_i + \delta \mu^{(j)}$

$\delta \Sigma^{(j)} \leftarrow (\mu_i - \theta)(\mu_i - \theta)^{\top} - (\mu_i^{(j)} - \theta)(\mu_i^{(j)} - \theta)^{\top}$

end

Update $\mu_{i+1} \leftarrow \mu_i + \sum_j \delta \mu^{(j)} / B$

Update $\Sigma_{i+1} \leftarrow \Sigma_i + \sum_j \delta \Sigma^{(j)} / B$

end

$q_w(\theta) \leftarrow \mathcal{N}(\mu_N, \Sigma_N)$

Theorem 2.2. (GSM-VI updates) Let $p(\theta, x)$ be given for some $\theta \in \mathbb{R}^d$, and let $q_{w^t}(\theta)$ and $q_w(\theta)$ be multivariate normal distributions with means μ_t and μ and covariance matrices Σ_t and Σ , respectively. As shorthand, let $g_t := \nabla_{\theta} \log p(\theta_t, x)$ and let

$$\mu_{t+1}, \Sigma_{t+1} = \underset{\mu, \Sigma \succeq 0}{\operatorname{argmin}} \left[\text{KL}(q_t, q) \right] \quad \text{such that} \quad \nabla_{\theta} \log q(\theta_t) = \nabla_{\theta} \log p(\theta_t, x). \quad (7)$$

The solution to eq. (7) is given by Eqs. 5 and 6 where

$$\mathbf{A}_t := \frac{1}{1+\rho} \left[\mathbf{I} - \frac{(\mu_t - \theta_t) g_t^{\top}}{1+\rho + (\mu_t - \theta_t)^{\top} g_t} \right] \Sigma_t, \quad (8)$$

and ρ is the positive root of the quadratic equation

$$\rho(1+\rho) = g_t^{\top} \Sigma_t g_t + [(\mu_t - \theta_t)^{\top} g_t]^2. \quad (9)$$

146 With the definition of \mathbf{A}_t in Eq. 8 we can see that the computational complexity of updating μ and Σ
147 via Eqs. 5 and 6 is $\mathcal{O}(d^2)$, where $\theta \in \mathbb{R}^d$ and we assume the cost of computing the gradients is $\mathcal{O}(d)$.
148 Note this is the best possible iteration complexity we can hope for, since we store and maintain the
149 full covariance matrix of d^2 elements. (The proof is in the appendix.)

150 Algorithm 1 presents the full GSM-VI algorithm. Here we also use mini-batching, where we average
151 over $B \in \mathbb{N}$ independently sampled updates of Eqs. 5 and 6 before updating the mean and covariance.

152 3 Empirical Studies

153 We evaluate the performance of GSM-VI in different settings. GSM-VI uses a multivariate Gaussian
154 distribution as its variational family. We separately investigate when the target posterior is in this
155 family and when it is not.

156 **Algorithmic details and comparisons.** We compare GSM-VI with a reparameterization variant of
157 BBVI as the baseline, similar to [19]. BBVI uses the same multivariate Gaussian variational family,
158 which we fit by maximizing the ELBO. (Maximizing the ELBO is equivalent to minimizing KL). We

Algorithm 2: Black-box variational inference

Input : Initial mean estimate μ_0 , Initial covariance estimate Σ_0 , target distribution $p(\theta|x)$, number of iterations N , batch size B , learning rate ϵ

Output : Multivariate normal variational distribution $q_w(\theta) := \mathcal{N}(\mu, \Sigma)$

$q_w \leftarrow \mathcal{N}(\mu_0, \Sigma_0)$;

for $i = 0, \dots, N - 1$ \triangleright iteration loop

do

$\{\theta^{(0)}, \theta^{(1)}, \dots, \theta^{(B)}\} \sim q_w(\theta)$ \triangleright Sample a batch of B points;

$\text{ELBO} = \sum_j \log(p(\theta^{(j)}, x) - \log q_w(\theta^{(j)})$;

$w \leftarrow w - \epsilon \nabla_w \text{ELBO}$ \triangleright Optimization step, we use ADAM;

end

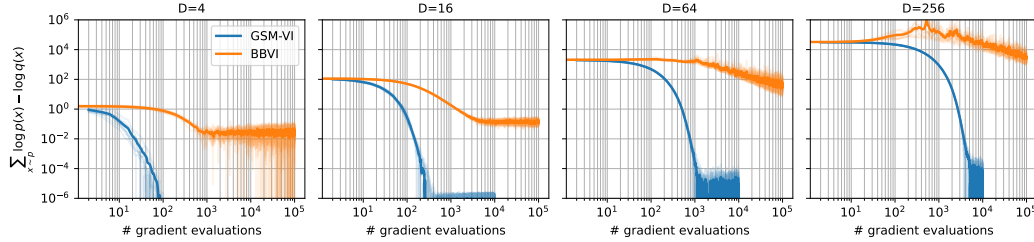


Figure 2: Scaling with dimensions: evolution of FKL with the number of gradient evaluations of the target distribution, which is a Gaussian distribution with dense covariance. Different panels show the results for different dimensions D of the distribution, specified in the title. Translucent lines show the scatter of 10 different runs and the solid line shows the average.

159 use the ADAM optimizer [17] with default settings but vary the learning rate between 10^{-1} and 10^{-3} .
160 We report results only for the best performing setting. The full BBVI algorithm is in Algorithm 2.

161 The only free parameter in GSM-VI is the batch size B . We find that $B = 2$ is better than $B = 1$, but
162 there is no improvement beyond that. In all studies, we report results for $B = 2$.

163 Both algorithms require an initial variational distribution. Unless specified otherwise, we initialize
164 the variational distribution with zero mean and identity covariance matrix.

165 **Evaluation metric.** GSM-VI does not explicitly minimize any loss function. Hence to compare its
166 performance against BBVI, we estimate empirical divergences between the variational and the target
167 distribution and show their evolution with the number of gradient evaluations. In the experiments with
168 synthetic models in Sections 3, 3.1, and 3.1 we have access to the true distribution; so we measure
169 the forward KL divergence (FKL) empirically ($\text{FKL} = \sum_{\theta_i \sim p(\theta)} \log p(\theta_i) - \log q(\theta_i)$). To reduce
170 stochasticity, we always use the same pre-generated set of 1000 samples from the target distribution.
171 In Section 3.3, we do not have access to the samples from the target distribution; so we monitor the
172 negative ELBO. In all experiments, we show the results for 10 independent runs.

173 3.1 GSM-VI for Gaussian approximation

174 We begin by studying GSM-VI where the target distribution is also a multivariate Gaussian.

175 **Scaling with dimensions.** How does GSM-VI scale with respect to the dimensions of the sample
176 space? Figure 2 shows the convergence of FKL for GSM-VI and BBVI as the dimension (D) of the
177 sample space increases. Empirically, we find that the number of iterations required for convergence
178 increases almost linearly with dimensions for GSM. The scaling for BBVI is worse, and it requires
179 100 times more iterations even for small problems ($D < 64$), while also converging to a sub-optimal
180 solution as measured by the FKL metric.

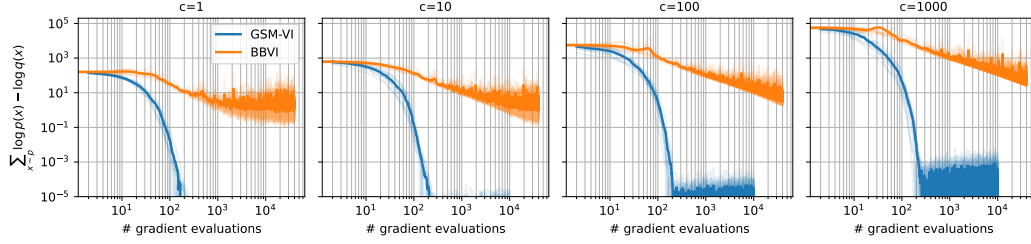


Figure 3: Impact of condition number: evolution of FKL with the number of gradient evaluations of the target distribution. The target is a 10-dimensional Gaussian albeit with a dense covariance matrix of different condition numbers c specified in the title of different panels. Translucent lines show the scatter of 10 different runs and the solid line shows the average.

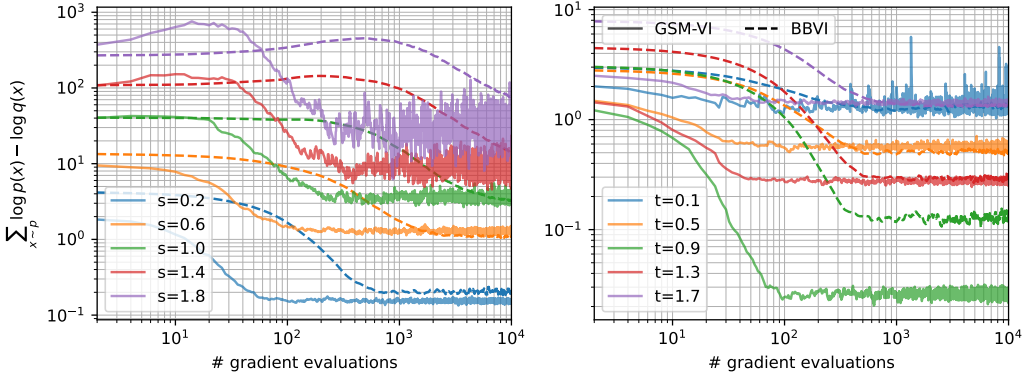


Figure 4: Impact of non-Gaussianity: evolution of FKL with the number of gradient evaluations for Sinh-arcsinh distributions with 10-dimensional dense Gaussian as the base distribution. Gaussian distribution has $s = 0$, $t = 1$. In the left panel, we vary skewness s while fixing $t = 1$, and in the right panel we vary the tail-weight t with skewness fixed to $s = 0$. Solid lines are the results for GSM, dashed for BBVI.

Impact of condition number. What is the impact of the conditioning of the target distribution? We again consider a Gaussian target distribution, but vary the condition number of its covariance matrix by fixing its smallest eigenvalue to 0.1, and scaling the largest eigenvalue to $0.1 \times c$. Figure 3 shows the results for a 10 dimensional Gaussian where we vary the condition number c from 1 to 1000. Convergence of GSM-VI seems to be largely insensitive to the condition number of the covariance matrix. BBVI on the other hand struggles with poorly conditioned problems, and it does not converge for $c > 100$ even with 100 times more iterations than GSM.

3.2 GSM-VI for non-Gaussian target distributions

GSM-VI was designed to solve the exact score-matching equations Eq. 3, which only have a solution when the family of variational distributions contains the target distribution (see Lemma 2.1). Here we investigate the sensitivity of GSM-VI to this assumption by fitting non-Gaussian target distributions with varying degrees of non-Gaussianity. Specifically, we suppose that the target has a multivariate Sinh-arcsinh normal distribution [15]

$$z \sim \mathcal{N}(\mu, \Sigma); \quad x = \sinh \left(\frac{1}{t} [\sinh^{-1}(z) + s] \right) \quad (10)$$

where the scalar parameters s and t control, respectively, the skewness and the heaviness of the tails, and the choices $s = 0$ and $t = 1$ reduce a Gaussian distribution as a special case.

Figure 4 shows the result for fitting the variational Gaussian to a 10-dimensional Sinh-arcsinh normal distribution for different values of s and t . As the target departs further from Gaussianity, the quality

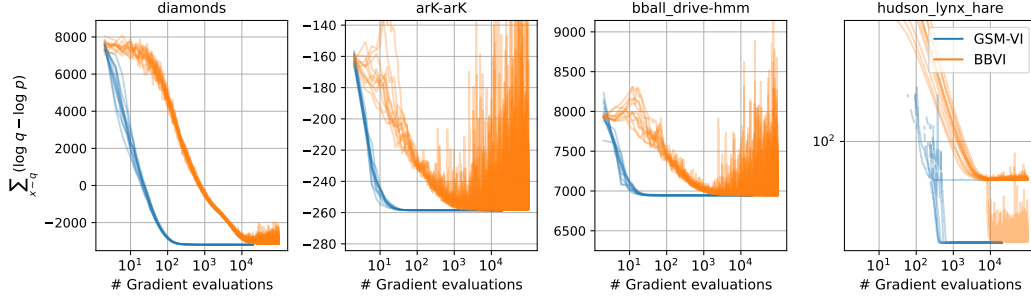


Figure 5: Models from posteriordb: Convergence of the ELBO for four models with multivariate normal posteriors. We show results for 10 runs.

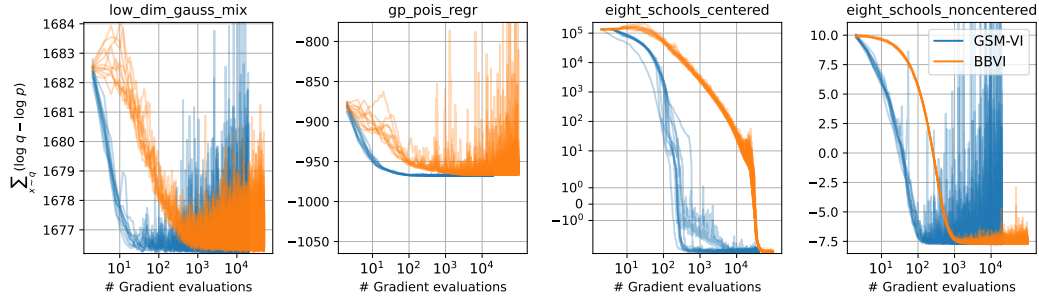


Figure 6: Models from posteriordb: Coverage of ELBO for four models with non-Gaussian posteriors. We show results for 10 runs.

of variational fit worsens for both GSM-VI (solid lines) and BBVI (dashed lines), but they converge to a fit of similar quality in terms of average FKL. GSM converges to this solution at least 10 times faster than BBVI. For highly non-Gaussian targets ($s \geq 1$ or $|t - 1| \geq 0.8$), we have found that GSM-VI does not converge to a fixed point, and it can experience oscillations that are larger in amplitude than BBVI, see for instance $s = 1.8$ and $t = 0.1$ on the left and right of Figure 4, respectively.

3.3 GSM-VI on real-world data

We evaluate GSM-VI for approximate on real-world data with 8 models from the posteriordb database [21]. The database provides the Stan code, data and reference posterior samples, and we use bridgestan to access the gradients of these models [5, 26]. We study the following models: diamonds (generalized linear models), hudson-lynx-hare (differential equation dynamics), bball-drive (hidden Markov models) and arK (time-series), eight-schools-centered and non-centered (hierarchical meta-analysis), gp-pois-regr (Gaussian processes), low-dim-gauss-mix (Gaussian mixture).

For each model (except hudson-lynx-hare), we initialize the variational parameter μ_0 at the mode of the distribution, and we set $\Sigma_0 = 0.1 \mathbf{I}_d$ where \mathbf{I}_d is the identity matrix of dimension d . For hudson-lynx-hare, we initialize the variational distribution as standard normal. We also experimented with other initializations. We find that they do not qualitatively change the conclusions, but can have larger variance between different runs.

We show the evolution of the ELBO for 10 runs of these models. Four of the models have posteriors that can be fit with multivariate normal distribution: diamonds, hudson-lynx-hare, bball-drive, and arK. Figure 5 shows the result for these models. The other models have non-Gaussian posteriors: eight-schools-centered, eight-schools-non-centered, gp-pois-regr,, and low-dim-gauss-mix. Figure 6 shows the results.

Overall, GSM-VI outperforms BBVI by a factor of 10-100x. When the target posterior is Gaussian, GSM-VI leads to more stable solutions. When the target is non-Gaussian, it converges to the same quality of variational approximation as BBVI. Further, though the ELBO estimate is noisy at the convergence, the 1-D marginals and moments of parameters remain stable.

4 Conclusion and Future Work

In this paper we proposed Gaussian score matching VI (GSM-VI), a new approach for VI when the variational family is multivariate Gaussian. GSM-VI is not based on minimizing a divergence or loss function between the variational and target distribution; instead, it repeatedly solves the exact score matching equations with closed-form updates for the mean and covariance matrix of the variational distribution.

Unlike approaches that are rooted in stochastic gradient descent, GSM-VI does not require the tuning of step-size hyper-parameters. It has only one free parameter, the batch size, and we found a batch-size of 2 to perform competitively across all experiments. Another choice is how to initialize the variational distribution. For the experiments in this paper, we initialized the covariance matrix as the identity matrix, but additional gains could potentially be made with more informed choices derived from a Laplace approximation or L-BFGS Hessian approximation [31].

We evaluated the performance of GSM-VI on synthetic targets and real-world models from posteriordb. In general, we found that it requires 10-100x fewer gradient evaluations than BBVI for the target distribution to converge. When the target distribution is itself multivariate Gaussian, we observed that GSM-VI scales almost *linearly* with dimensionality, which is significantly better than BBVI, and that GSM-VI is almost insensitive to the condition number of the target covariance matrix. Compared to BBVI, we also found that GSM-VI converges more quickly to a solution with a larger ELBO, which is surprising given that BBVI explicitly maximizes the ELBO.

GSM-VI is derived from a principled goal and justification, and the empirical studies indicate that it is a promising method. An important avenue for future work is to provide a proof that GSM-VI converges. We note that good convergence results have been obtained for analogous methods that project onto interpolation equations for empirical risk minimization. For instance the Stochastic Polyak Step achieves the min-max optimal rates of convergence for SGD [20]. Note that convergence of VI is a generally challenging problem, with no known rates of convergence even for BBVI [9, 10].

In another avenue of future work, the score-matching VI idea can potentially be used to design other methods for VI. As one example, we can consider non-Gaussian variational approximations, such as those in the exponential family. As another example, if the variational family is a mixture of Gaussians, we can employ GSM-VI to update the individual components of the mixture.

References

- [1] A. Agrawal, D. R. Sheldon, and J. Domke. Advances in black-box vi: normalizing flows, importance weighting, and optimization. *Neural Information Processing Systems*, 2020.
- [2] L. Berrada, A. Zisserman, and M. P. Kumar. Training neural networks for and by interpolation. In *Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, 2020.
- [3] E. Bingham, J. Chen, M. Jankowiak, F. Obermeyer, N. Pradhan, T. Karaletsos, R. Singh, P. Szerlip, P. Horsfall, and N. Goodman. Pyro: Deep universal probabilistic programming. *arXiv:1810.09538*, 2018.
- [4] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- [5] B. Carpenter, A. Gelman, M. D. Hoffman, D. Lee, B. Goodrich, M. Betancourt, M. Brubaker, J. Guo, P. Li, and A. Riddell. Stan: A probabilistic programming language. *Journal of Statistical Software*, 76(1):1–32, 2017.
- [6] K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer. Online passive-aggressive algorithms. *Journal of Machine Learning Research*, 7:551–585, 2006.
- [7] A. K. Dhaka, A. Catalina, M. Welandawe, M. Riis Andersen, J. Huggins, and A. Vehtari. Challenges and Opportunities in High-dimensional Variational Inference. art. arXiv:2103.01085, March 2021.
- [8] A. Dieng, D. Tran, R. Ranganath, J. Paisley, and D. Blei. Variational inference via χ upper bound minimization. In *Neural Information Processing Systems*, 2017.
- [9] J. Domke. Provable gradient variance guarantees for black-box variational inference. In *Neural Information Processing Systems*, volume 32, 2019.
- [10] J. Domke. Provable smoothness guarantees for black-box variational inference. In *International Conference on Machine Learning*, volume 119, pages 2587–2596. PMLR, 2020.
- [11] M. Dredze, K. Crammer, and F. Pereira. Confidence-weighted linear classification. In *International Conference on Machine Learning*, volume 307, pages 264–271, 2008.
- [12] R. M. Gower, A. Defazio, and M. Rabbat. Stochastic polyak stepsize with a moving target. *arXiv:2106.11851*, 2021.
- [13] M. Hoffman, D. Blei, C. Wang, and J. Paisley. Stochastic variational inference. *Journal of Machine Learning Research*, 14:1303–1347, 2013.
- [14] A. Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6:695–709, 2005.
- [15] C. Jones and A. Pewsey. The Sinh-Arcsinh Normal Distribution. *Significance*, 16(2):6–7, 04 2019. ISSN 1740-9705.
- [16] M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. An introduction to variational methods for graphical models. *Machine learning*, 37:183–233, 1999.
- [17] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [18] J. Knoblauch, J. Jewson, and T. Damoulas. An optimization-centric view on Bayes’ rule: Reviewing and generalizing variational inference. *Journal of Machine Learning Research*, 23: 1–109, 2022.
- [19] A. Kucukelbir, D. Tran, R. Ranganath, A. Gelman, and D. M. Blei. Automatic Differentiation Variational Inference. art. arXiv:1603.00788, March 2016.
- [20] N. Loizou, S. Vaswani, I. Laradji, and S. Lacoste-Julien. Stochastic polyak step-size for sgd: An adaptive learning rate for fast convergence. *International Conference on Artificial Intelligence and Statistics*, 2021.
- [21] M. Magnusson, P. Bürkner, and A. Vehtari. posteriordb: a set of posteriors for Bayesian inference and probabilistic programming, November 2022. URL <https://github.com/stan-dev/posteriordb>.
- [22] T. Minka. Expectation propagation for approximate Bayesian inference. In *Uncertainty in Artificial Intelligence (UAI)*, pages 362–369, 2001.

- 304 [23] C. Naesseth, F. Lindsten, and D. Blei. Markovian score climbing: Variational inference with
305 $KL(p \parallel q)$. In *Neural Information Processing Systems*, 2020.
- 306 [24] R. Ranganath, S. Gerrish, and D. M. Blei. Black Box Variational Inference. art. arXiv:1401.0118,
307 December 2013.
- 308 [25] R. Ranganath, J. Alotaib, D. Tran, and D. Blei. Operator variational inference. In *Neural*
309 *Information Processing Systems*. 2016.
- 310 [26] E. Roualdes, B. Ward, S. Axen, and B. Carpenter. BridgeStan: Efficient in-memory access
311 to Stan programs through Python, Julia, and R, March 2023. URL [https://github.com/](https://github.com/roualdes/bridgestan)
312 [roualdes/bridgestan](https://github.com/roualdes/bridgestan).
- 313 [27] J. Salvatier, Wiecki T., and C. Fonnesbeck. Probabilistic programming in Python using PyMC3.
314 *PeerJ Computer Science*, 2016.
- 315 [28] U. Seljak and B. Yu. Posterior inference unchained with EL_2O. art. arXiv:1901.04454, January
316 2019.
- 317 [29] M. W., M. R. Andersen, A. Vehtari, and J. H. Huggins. Robust, automated, and accurate
318 black-box variational inference. *arXiv:2203.15945*, 2022.
- 319 [30] Y. Yang, R. Martin, and H. Bondell. Variational approximations using Fisher divergence. *arXiv*
320 *e-prints*, art. arXiv:1905.05284, May 2019.
- 321 [31] L. Zhang, B. Carpenter, A. Gelman, and A. Vehtari. Pathfinder: Parallel quasi-newton variational
322 inference. *Journal of Machine Learning Research*, 23(306):1–49, 2022.

323 A Proof of Lemma 2.1

Lemma 2.1. The parameter w^* satisfies

$$\nabla_{\theta} \log q_{w^*}(\theta) = \nabla_{\theta} \log p(\theta, x), \quad \forall \theta, \quad (3)$$

if and only if w^* also satisfies Eq. 1.

324 *Proof.* (1) \implies (3): Differentiating both sides of (1) in θ gives

$$\begin{aligned} \nabla_{\theta} \log q_w(\theta) &= \nabla_{\theta} \log p(\theta|x) = \nabla_{\theta} (\log p(\theta, x) - \log p(x)) \\ &= \nabla_{\theta} \log p(\theta, x), \quad \forall \theta. \end{aligned}$$

325 (3) \implies (1): The inverse implication follows by using that $\nabla_{\theta} \log p(\theta|x) = \nabla_{\theta} \log p(\theta, x)$, as we
326 did in the above, and then integrating both sides of (3) in θ , which gives

$$\log q_w(\theta) = \log p(\theta|x) + C, \quad \forall \theta,$$

327 where C is some unknown constant. By exponentiating both sides and integrating in θ we have that

$$1 = \int_{\theta} q_w(\theta) d\theta = e^C \int_{\theta} p(\theta|x) d\theta = e^C.$$

328 Consequently $C = 0$, which gives our result. \square

329 B Proof of Theorem 2.2

330 Here we give the proof for Theorem 2.2. We also re-introduce the theorem with a simplified notation,
331 where we use (μ_0, Σ_0) to denote the mean and covariance at the previous time step of the method,
332 thus dropping the iteration counter t .

Theorem B.1. (GSM updates) Let $p(\theta, x)$ be given for some $\theta \in \mathbb{R}^d$, and let $q_0(\theta)$ and $q(\theta)$ be the multivariate normal distributions, respectively, with means μ_0 and μ and covariance matrices Σ_0 and Σ . We seek the distribution

$$\arg \min_{\mu, \Sigma = \Sigma^\top} \left[\text{KL}(q_0, q) \right] \quad \text{such that} \quad \nabla_{\theta} \log q(\theta) = \nabla_{\theta} \log p(\theta, x). \quad (11)$$

As shorthand, let $g := \nabla_{\theta} \log p(\theta, x)$, and let ρ be the positive root of the quadratic equation

$$\rho(1+\rho) = g^\top \Sigma_0 g + [(\mu_0 - \theta)^\top g]^2. \quad (12)$$

Then the solution is given by the following closed-form updates:

$$\mu = \mu_0 + \frac{1}{1+\rho} \left[\mathbf{I} - \frac{(\mu_0 - \theta)g^\top}{1+\rho + (\mu_0 - \theta)^\top g} \right] \Sigma_0 (g - \nabla_{\theta} \log q_0(\theta_0)), \quad (13)$$

$$\Sigma = \Sigma_0 + (\mu_0 - \theta)(\mu_0 - \theta)^\top - (\mu - \theta)(\mu - \theta)^\top. \quad (14)$$

Furthermore, if Σ_0 is symmetric positive definite then so is Σ .

333 *Proof.* The constraint in this optimization is given by

$$g = \nabla_{\theta} \log q(\theta) \quad (15)$$

$$= \nabla_{\theta} \left[-\frac{1}{2}(\theta - \mu)^\top \Sigma^{-1}(\theta - \mu) - \frac{1}{2} \log((2\pi)^d |\Sigma|) \right] \quad (16)$$

$$= -\Sigma^{-1}(\theta - \mu). \quad (17)$$

334 The KL divergence is given by

$$\text{KL}(q_0, q) = \frac{1}{2} \left\{ \text{tr}[\Sigma^{-1} \Sigma_0] + \log \frac{|\Sigma|}{|\Sigma_0|} + (\mu - \mu_0)^\top \Sigma^{-1}(\mu - \mu_0) - d \right\}. \quad (18)$$

335 Dropping irrelevant terms from the optimization, we obtain the Lagrangian

$$\mathcal{L}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\lambda}) = \frac{1}{2} \{ \text{tr}[\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}_0] - \log |\boldsymbol{\Sigma}^{-1}| + (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_0) \} + \boldsymbol{\lambda}^\top (\boldsymbol{g} - \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \boldsymbol{\theta})). \quad (19)$$

336 It is easier to optimize the matrix $\boldsymbol{\Sigma}^{-1}$ instead of $\boldsymbol{\Sigma}$. We can enforce the symmetry of $\boldsymbol{\Sigma}^{-1}$ by writing
337

$$\boldsymbol{\Sigma}^{-1} = \frac{1}{2} (\boldsymbol{\Phi} + \boldsymbol{\Phi}^\top) \quad (20)$$

338 and performing an unconstrained optimization over $\boldsymbol{\Phi}$. With respect to the latter, the gradients of the
339 Lagrangian are given by

$$\frac{\partial \mathcal{L}}{\partial \Phi_{ij}} = \sum_{kl} \left(\frac{\partial \mathcal{L}}{\partial \Sigma_{kl}^{-1}} \right) \left(\frac{\partial \Sigma_{kl}^{-1}}{\partial \Phi_{ij}} \right) = \sum_{kl} \left(\frac{\partial \mathcal{L}}{\partial \Sigma_{kl}^{-1}} \right) \left(\frac{1}{2} \delta_{ki} \delta_{lj} + \frac{1}{2} \delta_{kj} \delta_{li} \right) = \frac{1}{2} \left(\frac{\partial \mathcal{L}}{\partial \Sigma_{ij}^{-1}} + \frac{\partial \mathcal{L}}{\partial \Sigma_{ji}^{-1}} \right). \quad (21)$$

340 Next we examine where the gradients of the Lagrangian vanish:

$$\mathbf{0} = \frac{\partial \mathcal{L}}{\partial \boldsymbol{\mu}} \implies \mathbf{0} = \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_0) - \boldsymbol{\Sigma}^{-1} \boldsymbol{\lambda} \implies \boxed{\boldsymbol{\lambda} = \boldsymbol{\mu} - \boldsymbol{\mu}_0} \quad (22)$$

$$\mathbf{0} = \frac{\partial \mathcal{L}}{\partial \boldsymbol{\lambda}} \implies \mathbf{0} = \boldsymbol{g} - \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \boldsymbol{\theta}) \implies \boxed{\boldsymbol{\mu} - \boldsymbol{\theta} = \boldsymbol{\Sigma} \boldsymbol{g}} \quad (23)$$

$$\mathbf{0} = \frac{\partial \mathcal{L}}{\partial \boldsymbol{\Phi}} \implies \mathbf{0} = \boldsymbol{\Sigma}_0 - \boldsymbol{\Sigma} - (\boldsymbol{\mu} - \boldsymbol{\mu}_0)(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^\top - [\boldsymbol{\lambda}(\boldsymbol{\mu} - \boldsymbol{\theta})^\top + (\boldsymbol{\mu} - \boldsymbol{\theta})\boldsymbol{\lambda}^\top], \quad (24)$$

$$\implies \boxed{\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_0 + (\boldsymbol{\mu} - \boldsymbol{\mu}_0)(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^\top - \boldsymbol{\lambda}(\boldsymbol{\mu} - \boldsymbol{\theta})^\top - (\boldsymbol{\mu} - \boldsymbol{\theta})\boldsymbol{\lambda}^\top} \quad (25)$$

341 We claim that these equations (though nonlinear) can be solved in closed form. The first step is to
342 eliminate $\boldsymbol{\lambda}$ from eq. (25) using eq. (22). In this way we find

$$\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_0 + (\boldsymbol{\mu} - \boldsymbol{\mu}_0)(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^\top - (\boldsymbol{\mu} - \boldsymbol{\mu}_0)(\boldsymbol{\mu} - \boldsymbol{\theta})^\top - (\boldsymbol{\mu} - \boldsymbol{\theta})(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^\top \quad (26)$$

$$= \boldsymbol{\Sigma}_0 - \boldsymbol{\mu}\boldsymbol{\mu}^\top + \boldsymbol{\mu}\boldsymbol{\theta}^\top + \boldsymbol{\theta}\boldsymbol{\mu}^\top + \boldsymbol{\mu}_0\boldsymbol{\mu}_0^\top - \boldsymbol{\mu}_0\boldsymbol{\theta}^\top - \boldsymbol{\theta}\boldsymbol{\mu}_0^\top \quad (27)$$

$$= \boldsymbol{\Sigma}_0 + (\boldsymbol{\mu}_0 - \boldsymbol{\theta})(\boldsymbol{\mu}_0 - \boldsymbol{\theta})^\top - (\boldsymbol{\mu} - \boldsymbol{\theta})(\boldsymbol{\mu} - \boldsymbol{\theta})^\top. \quad (28)$$

It is worth highlighting the form of this equation:

$$\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_0 + (\boldsymbol{\mu}_0 - \boldsymbol{\theta})(\boldsymbol{\mu}_0 - \boldsymbol{\theta})^\top - (\boldsymbol{\mu} - \boldsymbol{\theta})(\boldsymbol{\mu} - \boldsymbol{\theta})^\top$$

343 This is a simple rank-two update for $\boldsymbol{\Sigma}$. Note that $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_0$ if $\boldsymbol{\mu} = \boldsymbol{\mu}_0$; also, the solution for $\boldsymbol{\Sigma}$ is
344 determined by the solution for $\boldsymbol{\mu}$.

345 Ultimately we must solve for $\boldsymbol{\mu}$, but first it is useful to solve for the intermediate quantity $\boldsymbol{g}^\top \boldsymbol{\Sigma} \boldsymbol{g} > 0$.
346 From eq. (28), we obtain

$$\boldsymbol{g}^\top \boldsymbol{\Sigma} \boldsymbol{g} = \boldsymbol{g}^\top \boldsymbol{\Sigma}_0 \boldsymbol{g} + [(\boldsymbol{\mu}_0 - \boldsymbol{\theta})^\top \boldsymbol{g}]^2 - [(\boldsymbol{\mu} - \boldsymbol{\theta})^\top \boldsymbol{g}]^2, \quad (29)$$

347 and from eq. (23), we obtain

$$\boldsymbol{g}^\top \boldsymbol{\Sigma} \boldsymbol{g} = \boldsymbol{g}^\top \boldsymbol{\Sigma}_0 \boldsymbol{g} + [(\boldsymbol{\mu}_0 - \boldsymbol{\theta})^\top \boldsymbol{g}]^2 - (\boldsymbol{g}^\top \boldsymbol{\Sigma} \boldsymbol{g})^2. \quad (30)$$

As shorthand, let $\rho = \boldsymbol{g}^\top \boldsymbol{\Sigma} \boldsymbol{g}$. Then from eq. (30) we see that ρ satisfies the quadratic equation

$$\rho(1 + \rho) = \boldsymbol{g}^\top \boldsymbol{\Sigma}_0 \boldsymbol{g} + [(\boldsymbol{\mu}_0 - \boldsymbol{\theta})^\top \boldsymbol{g}]^2.$$

348 Note that there are no unknowns on the right side of this equation. The correct solution is given by
349 the positive root since $\rho = \boldsymbol{g}^\top \boldsymbol{\Sigma} \boldsymbol{g} > 0$. Also note that $\rho = (\boldsymbol{\mu} - \boldsymbol{\theta})^\top \boldsymbol{g}$ from eq. (23).

It is useful to define one final intermediate quantity before solving for $\boldsymbol{\mu}$. Let

$$\boldsymbol{\varepsilon}_0 = \boldsymbol{\Sigma}_0 \boldsymbol{g} - \boldsymbol{\mu}_0 + \boldsymbol{\theta}.$$

350 Note that $\boldsymbol{\varepsilon}_0$ simply measures the degree to which the parameters of $q_0(\boldsymbol{\theta})$ violate the desired
351 constraint $\nabla_{\boldsymbol{w}} \log q(\boldsymbol{\theta}) = \nabla_{\boldsymbol{w}} \log p(\boldsymbol{\theta}, \boldsymbol{y})$. Put another way, if $\boldsymbol{\varepsilon}_0 = \mathbf{0}$, then we have the trivial
352 solution $\boldsymbol{\mu} = \boldsymbol{\mu}_0$ and $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_0$.

Now we have everything to express the solution for μ in a highly intuitive form; in particular, it will be immediately evident that $\mu \rightarrow \mu_0$ as $\varepsilon_0 \rightarrow \mathbf{0}$. Starting from eqs. (23) and (28), we find

$$\mu - \mu_0 = \theta - \mu_0 + \Sigma g, \quad (31)$$

$$= \theta - \mu_0 + [\Sigma_0 + (\mu_0 - \theta)(\mu_0 - \theta)^\top - (\mu - \theta)(\mu - \theta)^\top] g, \quad (32)$$

$$= \varepsilon_0 + (\mu_0 - \theta)(\mu_0 - \theta)^\top g - (\mu - \theta)(\mu - \theta)^\top g, \quad (33)$$

$$= \varepsilon_0 + (\mu_0 - \theta)(\mu_0 - \theta)^\top g - (\mu - \mu_0 + \mu_0 - \theta)(\mu - \theta)^\top g, \quad (34)$$

$$= \varepsilon_0 - \rho(\mu - \mu_0) + (\mu_0 - \theta)[(\mu_0 - \theta) - (\mu - \theta)^\top] g, \quad (35)$$

$$= \varepsilon_0 - \rho(\mu - \mu_0) + (\mu_0 - \theta)(\mu_0 - \mu)^\top g, \quad (36)$$

$$= \varepsilon_0 - (\rho \mathbf{I} + (\mu_0 - \theta)g^\top)(\mu_0 - \mu). \quad (37)$$

Note what has happened here: eq. (32) is a system of *nonlinear* equations for μ , but in eq. (37), all the nonlinearity has been expressed in terms of ρ . Since ρ can be determined via eq. (30), we arrive effectively at a system of *linear* equations for μ . Collecting terms, we obtain

$$[(1 + \rho)\mathbf{I} + (\mu_0 - \theta)g^\top](\mu - \mu_0) = \varepsilon_0. \quad (38)$$

We thus arrive at the closed-form update

$$\mu = \mu_0 + [(1 + \rho)\mathbf{I} + (\mu_0 - \theta)g^\top]^{-1} \varepsilon_0 \quad (39)$$

It is evident from this update that $\mu \rightarrow \mu_0$ as $\varepsilon_0 \rightarrow \mathbf{0}$. The matrix inverse in this update can also be computed efficiently from the Woodbury matrix identity.

In sum, the joint update for μ and Σ can be efficiently computed as follows:

1. Set $g = \nabla_w \log p(\theta, y)$ and $\varepsilon_0 = \Sigma_0 g - \mu_0 + \theta$.
2. Solve $\rho(1 + \rho) = g^\top \Sigma_0 g + [(\mu_0 - \theta)^\top g]^2$ for $\rho > 0$.
3. Compute $\mu = \mu_0 + [(1 + \rho)\mathbf{I} + (\mu_0 - \theta)g^\top]^{-1} \varepsilon_0$.
4. Compute $\Sigma = \Sigma_0 + (\mu_0 - \theta)(\mu_0 - \theta)^\top - (\mu - \theta)(\mu - \theta)^\top$.

Solving the quadratic in 2. for ρ we have the positive

$$\rho = \frac{\sqrt{1 + 4(g^\top \Sigma_0 g + [(\mu_0 - \theta)^\top g]^2)} - 1}{2}$$

Solving the above linear equation for μ and using the Sherman Morrison formula

$$[a\mathbf{I} + \mathbf{u}g^\top]^{-1} = \frac{1}{a} \left(\mathbf{I} - \frac{\mathbf{u}g^\top}{a + \mathbf{u}^\top g} \right), \quad \text{for every } \mathbf{u}, g, a$$

gives

$$\mu = \mu_0 + \frac{1}{1 + \rho} \left[\mathbf{I} - \frac{(\mu_0 - \theta)g^\top}{1 + \rho + (\mu_0 - \theta)^\top g} \right] \varepsilon_0. \quad (40)$$

Using that $\nabla_\theta \log q_0(\theta_0) = -\Sigma_0^{-1}(\theta - \mu_0)$ we have that

$$\varepsilon_0 = \Sigma_0 (g - \Sigma_0^{-1}(\mu_0 + \theta)) = \Sigma_0 (g - \nabla_\theta \log q_0(\theta_0)).$$

Finally substituting out ε_0 in (40) the result

$$\mu = \mu_0 + \frac{1}{1 + \rho} \left[\mathbf{I} - \frac{(\mu_0 - \theta)g^\top}{1 + \rho + (\mu_0 - \theta)^\top g} \right] \Sigma_0 (g - \nabla_\theta \log q_0(\theta_0)).$$

372 **Proof that Σ_0 p.s.d. $\implies \Sigma$ p.s.d.** It remains to prove that our solution for Σ is positive-definite,
 373 or equivalently, that all of its eigenvalues are positive. We begin by rewriting our results for Σ in
 374 eq. (28) and ρ in eq. (30) in a more convenient form. As shorthand, let

$$\mathbf{M}_0 = \Sigma_0 + (\mu_0 - \theta)(\mu_0 - \theta)^\top, \quad (41)$$

375 so that \mathbf{M}_0 captures the first two terms on the right side of eq. (28). Note that \mathbf{M}_0 is positive-definite,
 376 a fact that we will exploit repeatedly in what follows. In addition, recall that $\mu - \theta = \Sigma g$ from
 377 eq. (23). Thus with this notation we can rewrite eqs. (28) and (30) as

$$\Sigma = \mathbf{M}_0 - (\Sigma g)(\Sigma g)^\top, \quad (42)$$

$$\rho(1 + \rho) = g^\top \mathbf{M}_0 g. \quad (43)$$

378 Now let e be any normalized eigenvector of Σ ; we want to show that its corresponding eigenvalue λ_e
 379 is positive. From eq. (42), it follows that

$$\lambda_e = e^\top \Sigma e \quad (44)$$

$$= e^\top [\mathbf{M}_0 - (\Sigma g)(\Sigma g)^\top] e \quad (45)$$

$$= e^\top \mathbf{M}_0 e - \lambda_e^2 (e^\top g)^2. \quad (46)$$

380 Note that if $e^\top g = 0$, then it follows trivially that $\lambda_e = e^\top \mathbf{M}_0 e > 0$. So we only need to consider
 381 the non-trivial case $e^\top g \neq 0$. To proceed, we note that

$$(e^\top \mathbf{M}_0 g)^2 = (e^\top \mathbf{M}_0^{\frac{1}{2}} \mathbf{M}_0^{\frac{1}{2}} g)^2 \leq (e^\top \mathbf{M}_0 e)(g^\top \mathbf{M}_0 g), \quad (47)$$

382 where we have used the Cauchy-Schwartz inequality to bound $(e^\top \mathbf{M}_0 g)$ in terms of $(e^\top \mathbf{M}_0 e)$, the
 383 latter of which appears in eq. (46). Substituting this inequality into eq. (46), we find that

$$\lambda_e \geq \frac{(e^\top \mathbf{M}_0 g)^2}{g^\top \mathbf{M}_0 g} - \lambda_e^2 (e^\top g)^2. \quad (48)$$

384 To prove that $\lambda_e > 0$ we need one more intermediate result. Focusing on the rightmost term in this
 385 equality, we note that

$$\lambda_e (e^\top g) = e^\top \Sigma g = e^\top [\mathbf{M}_0 - (\Sigma g)(\Sigma g)^\top] g = e^\top \mathbf{M}_0 g - \lambda_e (e^\top g)(g^\top \Sigma g), \quad (49)$$

386 and rearranging the terms in this equation, we find

$$e^\top \mathbf{M}_0 g = \lambda_e (e^\top g)(1 + g^\top \Sigma g). \quad (50)$$

387 This intermediate result is useful because it relates the two terms on the right side of eq. (48). In
 388 particular, using eq. (50) to eliminate the term $e^\top \mathbf{M}_0 g$ in eq. (48), we find:

$$\begin{aligned} \lambda_e &\geq \frac{[\lambda_e (e^\top g)(1 + g^\top \Sigma g)]^2}{g^\top \mathbf{M}_0 g} - \lambda_e^2 (e^\top g)^2 \\ &= \lambda_e^2 (e^\top g)^2 \left[\frac{(1 + g^\top \Sigma g)^2}{g^\top \mathbf{M}_0 g} - 1 \right] \\ &= \lambda_e^2 (e^\top g)^2 \left[\frac{(1 + \rho)^2}{\rho(1 + \rho)} - 1 \right] \\ &= \frac{\lambda_e^2 (e^\top g)^2}{\rho}, \\ &> 0, \end{aligned}$$

389 where the final inequality follows because the individual terms λ_e^2 , $(e^\top g)^2$, and ρ are all strictly
 390 positive; note that λ_e cannot be equal to zero because this contradicts the equality in eq. (46). This
 391 completes the proof. Perhaps it is useful that this derivation also gives upper bounds on λ_e , namely

$$\frac{1}{\lambda_e} \geq \frac{(e^\top g)^2}{\rho} \implies \lambda_e \leq \frac{\rho}{(e^\top g)^2} = \frac{g^\top \Sigma g}{(e^\top g)^2}. \quad (51)$$

392

□