
Improved Marginal Unbiased Score Expansion (MUSE) via Implicit Differentiation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We apply the technique of implicit differentiation to boost performance, reduce
2 numerical error, and remove required user-tuning in the Marginal Unbiased Score
3 Expansion (MUSE) algorithm for hierarchical Bayesian inference. We demon-
4 strate these improvements on three representative inference problems: 1) an ex-
5 tended Neal’s funnel 2) Bayesian neural networks, and 3) probabilistic principal
6 component analysis. On our particular test cases, MUSE with implicit differen-
7 tiation is faster than Hamiltonian Monte Carlo by factors of 155, 397, and 5, re-
8 spectively, or factors of 65, 278, and 1 without implicit differentiation, and yields
9 good approximate marginal posteriors. The Julia and Python MUSE packages
10 have been updated to use implicit differentiation, and can solve problems defined
11 by hand or with any of a number of popular probabilistic programming languages
12 and automatic differentiation backends.

13 1 Introduction

14 MUSE is an algorithm for fast approximate hierarchical Bayesian inference, recently proposed by
15 [1, 2]. The user denotes some subset of model parameters as the “parameters of interest,” and the
16 algorithm will approximate their marginal posterior while integrating out remaining “latent” param-
17 eters. MUSE is efficient for very high-dimensional latent spaces and can often provide near-exact
18 inference at orders of magnitude lower computational cost than other methods such as Hamiltonian
19 Monte Carlo (HMC) or variational inference (VI) [2].

20 The requirements for using MUSE on a given problem are that 1) samples can be generated from
21 the prior and 2) gradients of the joint posterior probability distribution can be calculated. The latter
22 requirement is the same as for HMC, VI, and many other tools. The former requirement is not
23 strictly a requirement for some of these, but is generally even easier. All problems defined via a
24 probabilistic programming language satisfy the requirements automatically. Owing to its reliance
25 on prior samples, MUSE can be considered a form of simulation-based inference, extended to use
26 readily available joint posterior gradients, similar to the proposal by [3].

27 At its core, MUSE is based on an approximation to the marginal score formed from solutions to a
28 series of optimization problems. As part of the algorithm, we must compute derivatives of these
29 solutions, and, in this work, we improve MUSE by making use of implicit differentiation (ID) to
30 perform this calculation. While ID is not a new development, it has recently been shown to be
31 particularly powerful in conjunction with automatic differentiation (AD) [4, 5]. We follow this
32 approach, and demonstrate that it leads to significant improvements in both speed and usability for
33 MUSE, strengthening its case as a generic inference tool.

34 **2 Summary of the MUSE method**

35 Here we give a brief and practical summary of MUSE to help understand where ID fits in (for a
 36 comprehensive introduction, see [2]). MUSE is applicable to inference problems where the posterior
 37 probability of some parameters of interest, θ , given data, x , requires marginalization over a high-
 38 dimensional latent space parameterized by z ,

$$\mathcal{P}(\theta | x) = \int d^n z \mathcal{P}(x, z | \theta) \mathcal{P}(\theta). \quad (1)$$

39 The algorithm provides a fast estimate of the marginal posterior mean and covariance, which is
 40 computed under an approximation to the integral over z . This approximation involves solving a
 41 series of optimization problems wherein we maximize the joint likelihood, $\mathcal{P}(x, z | \theta)$, over the
 42 latent parameters z , given fixed x and θ ,

$$\hat{z}(\theta, x) \equiv \underset{z}{\operatorname{argmax}} \log \mathcal{P}(x, z | \theta). \quad (2)$$

43 These correspond to maximum a posteriori (MAP) estimates of z , and they are used to define the
 44 score at the MAP,

$$s_i^{\text{MAP}}(\theta, x) \equiv \frac{d}{d\theta_i} \log \mathcal{P}(x, \hat{z}(\theta, x) | \theta). \quad (3)$$

45 The MUSE estimate of the posterior mean, $\bar{\theta}$, is then implicitly defined as the solution to

$$s_i^{\text{MAP}}(\bar{\theta}, x) = \left\langle s_i^{\text{MAP}}(\bar{\theta}, x) \right\rangle_{x \sim \mathcal{P}(x | \bar{\theta})}, \quad (4)$$

46 and the posterior covariance is $\Sigma = H^{-1} J H^{-\dagger}$, with

$$J_{ij} = \left\langle s_i^{\text{MAP}}(\bar{\theta}, x) s_j^{\text{MAP}}(\bar{\theta}, x) \right\rangle_{x \sim \mathcal{P}(x | \bar{\theta})} - \left\langle s_i^{\text{MAP}}(\bar{\theta}, x) \right\rangle_{x \sim \mathcal{P}(x | \bar{\theta})} \left\langle s_j^{\text{MAP}}(\bar{\theta}, x) \right\rangle_{x \sim \mathcal{P}(x | \bar{\theta})} \quad (5)$$

$$H_{ij} = \frac{d}{d\theta_j} \left[\left\langle s_i^{\text{MAP}}(\bar{\theta}, x) \right\rangle_{x \sim \mathcal{P}(x | \theta)} \right] \Big|_{\theta = \bar{\theta}}. \quad (6)$$

47 This definition gives MUSE a number of useful properties (see [2] for proofs): 1) it is an asymptotically unbiased estimate of θ irregardless of any non-Gaussianity in the likelihood, 2) it is asymptotically optimal for a Gaussian likelihood, where it becomes equivalent to the marginal maximum likelihood estimate and the covariance becomes the inverse Fisher information matrix, 3) no dense operators of the dimensionality of z ever need to be computed, meaning it is well-suited for high-dimensional problems and 4) it requires few tuning parameters, setting it apart from HMC, VI, or many other simulation-based inference methods, which need user-provided mass matrices, surrogate distributions, or neural network architectures to work or to achieve optimal performance on complicated latent spaces. MUSE is approximate, so it does not aim to generically replace exact algorithms like HMC, but in many cases, its speed and aforementioned properties make it a very advantageous alternative.

58 In practice, the optimization problem in Eq. (2) is performed with LBFGS using user-provided
 59 or AD gradients. An existing challenge for MUSE is that naively computing Eq. (6) with AD
 60 would require propagating second-order derivatives through the optimizer, since a chain rule term
 61 involving $d\hat{z}/d\theta$ arises. With few or no AD libraries robustly supporting second-order AD through
 62 an optimizer, we have previously resorted to computing this term with finite differences (FD). This
 63 has not been completely prohibitive as FD are needed only over the low-dimensional θ and not over
 64 the high-dimensional z , so the solution remains tractable despite a linear computational scaling with
 65 the dimensionality of θ . However, it requires tuning the FD step size for each dimension of θ , and
 66 can at times incur large numerical errors. The main development of this paper is to demonstrate that
 67 this term can instead be computed more simply and exactly with ID.

68 **3 Using implicit differentiation**

69 To compute H with ID, first note that Eq. (6) can be written as $H_{ij} = \frac{1}{N} \sum_{\alpha=1}^N h_{ij}(\Omega_\alpha)$, where Ω_α
 70 are some independent random states, and

$$h_{ij}(\Omega) = \frac{d}{d\theta'_j} \frac{d}{d\theta_i} \log \mathcal{P} \left(x(\Omega, \theta'), \hat{z}(x(\Omega, \theta'), \bar{\theta}) \mid \theta \right) \Big|_{\theta = \theta' = \bar{\theta}}. \quad (7)$$

71 Here, we consider a single realization of x as dependent on θ in the sense that any simulated x can
 72 be written as a deterministic function of θ and a random state (think of Ω as the machine’s pseudo
 73 random number generator). Expanding the chain rule once and omitting Ω and the final evaluation
 74 at $\bar{\theta}$ for brevity yields

$$\frac{d}{d\theta'_j} \frac{d}{d\theta_i} \log \mathcal{P}(x(\theta'), \hat{z}(x(\bar{\theta}), \bar{\theta}) \mid \theta) + \frac{d}{dz_n} \frac{d}{d\theta_i} \log \mathcal{P}(x(\bar{\theta}), z \mid \theta) \Big|_{z=\hat{z}} \frac{d\hat{z}_n(x(\theta'), \bar{\theta})}{d\theta'_j}. \quad (8)$$

75 The first term can be computed with second-order AD through the likelihood and through the prior
 76 samples of x . In practice, this means simply using the same random state on the forwards and/or
 77 backwards AD passes and otherwise considering random number generation constant (this is the
 78 default in most AD libraries). The second term, where we will use ID, involves a derivative of the
 79 MAP solution, $d\hat{z}/d\theta$. The MAP solution by definition obeys

$$\frac{d}{dz} \log \mathcal{P}(x(\theta'), z \mid \bar{\theta}) \Big|_{z=\hat{z}(x(\theta'), \bar{\theta})} = 0. \quad (9)$$

80 Taking a θ' derivative of this equation and solving the resulting equation for $d\hat{z}/d\theta'$ yields

$$\frac{d\hat{z}_n}{d\theta'_j} = \left[\frac{d^2}{dz_m dz_n} \log \mathcal{P}(x(\bar{\theta}), z \mid \bar{\theta}) \right]^{-1} \frac{d}{d\theta'_j} \frac{d}{dz_m} \log \mathcal{P}(x(\theta'), z \mid \bar{\theta}) \Big|_{z=\hat{z}}. \quad (10)$$

81 This quantity now only requires derivatives through the likelihood rather than through an optimizer;
 82 in fact, it is independent of the particular optimizer used to obtain \hat{z} . Computing it involves solving
 83 a linear problem with the same dimensionality as z . Because z is assumed high-dimensional where
 84 forming an explicit matrix is impossible, we solve the system iteratively, with the action of the quan-
 85 tity in brackets above given by a jacobian-vector product. Note, however, that the linear operator
 86 is symmetric since it is a Hessian, and, by definition if the MAP exists (which is a requirement for
 87 MUSE anyway), it is positive definite. Thus, we can use an efficient conjugate gradient solver which
 88 exploits this structure, as opposed to generic linear solvers which must be used in more general ID
 89 problems.

90 4 Results

91 We compare HMC and MUSE with or without ID on three representative inference problems:

92 **Funnel problem** We consider an embedding of several Neal’s funnels into a toy hierarchical
 93 problem [6, 2]. The model is:

$$\theta_i \sim \text{Normal}(0, 3) \quad z_{ij} \sim \text{Normal}(0, \exp(\theta_i/2)) \quad x_{ij} \sim \text{Normal}(\tanh(z_{ij}), 1) \quad (11)$$

94 with $i \in 1:10$ parameters and $j \in 1:500$ latent dimensions per parameter. Although the embedded
 95 funnels are independent, for demonstration, we solve the entire problem as one large system when
 96 running either NUTS or MUSE.

97 **Bayesian Neural Network** Following the example given in [7], we consider a Bayesian neural
 98 network (BNN) analysis, where we interpolate some noisy one-dimensional data with a three-layer
 99 neural network. The model is:

$$\begin{aligned} \sigma_i &\sim \text{LogNormal}(0, 1) & W_i &\sim \text{Normal}(0, \sigma_i) \\ \tau &\sim \text{Gamma}(3, 1) & Y_j &\sim \text{Normal}(\text{NN}(W_i), 1/\tau) \end{aligned} \quad (12)$$

100 where $i \in 1:3$ layers, the layer weights, W_i , which parameterize the network, NN, contain 45 latent
 101 dimensions and map the data coordinates to 5 hidden units and finally to the data space, and the data
 102 Y_j consists of $j \in 1:500$ data points. The goal is to infer the σ_i and τ . We note that each internal
 103 optimization solution in Eq. (2) involves training the network given some prior on the weights. For
 104 this simple example we use our standard LBFGS solver, but other more machine learning oriented
 105 solvers can readily be used for the internal MUSE optimization step as well.

106 **Probabilistic Principal Component Analysis** Finally, we consider a probabilistic principal com-
 107 ponent analysis (PPCA) with automatic relevance determination [8]. The model is:

$$\begin{aligned} \alpha_i &\sim \text{InverseGamma}(1, 1) & W_{ki} &\sim \text{Normal}(0, \sqrt{\alpha_i}) \\ Z_{ij} &\sim \text{Normal}(0, 1) & X_{kjl} &\sim \text{Normal}(W_{ki} Z_{ij}, 1) \end{aligned} \quad (13)$$

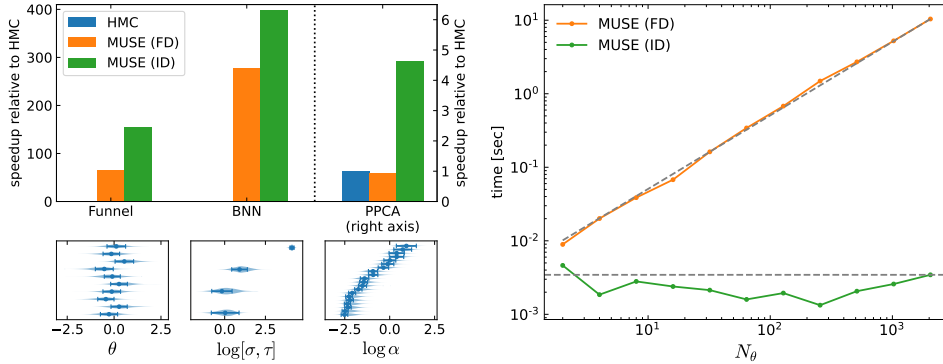


Figure 1: (Top left) Speedups which are possible with MUSE both with or without ID as compared to HMC on a variety of hierarchical Bayesian inference problems (described in Sec. 4). (Bottom left) HMC posteriors as violin plots, compared to MUSE results as error bars. (Right) Empirical check of the asymptotical scaling of the H computation with FD or ID.

108 with $i \in 1:10$ principal components, $j, k \in 1:100$ observations, and $l \in 5$ batches. The goal is to find
 109 the largest principal component amplitudes, α , given observations of X , while marginalizing over
 110 the entries in the Z and W matrices.

111 Our benchmarks compare the number of posterior gradient evaluations needed such that for all
 112 parameters of interest, we reach 1) a 10% error on the mean relative to the standard deviation and
 113 2) a 10% relative error on the standard deviation. Given a Gaussian sampling distribution, these
 114 criteria impose the same constraint. For HMC, this corresponds to achieving an effective sample
 115 size of 100 for all parameters. We use NumPyro to implement each model [9, 10], and sample with
 116 NumPyro’s NUTS implementation with default parameters. For MUSE this corresponds to running
 117 MUSE with 100 simulations and setting the θ tolerance to 10%. We use the existing Jax [11] MUSE
 118 implementation to run MUSE on the same NumPyro model.

119 The results are summarized in the left panels in Fig. 1. We see that for each of the three problems,
 120 ID outperforms the previous FD approach by as much as a factor of 5. In all cases, MUSE with ID
 121 significantly outperforms HMC, including by a factor of 391 in the most dramatic case (the BNN).
 122 The bottom panel shows a comparison of the inferred values of the parameters of interest, confirming
 123 the quality of the MUSE approximation.

124 We also expect a more favorable computational scaling for ID over FD as we increase the dimen-
 125 sionality of θ . This is because computing h_{ij} with FD requires perturbing each element of θ and
 126 recomputing a MAP each time, whereas ID requires just one MAP that is then used evaluating all
 127 terms in Eqn. (8), with the tradeoff of also needing to solve a linear problem. To confirm this trade-
 128 off is beneficial, we modify our funnel problem to increase the dimensionality of θ (while keeping
 129 the latent dimensionality the same), and plot resulting timings for the H computation in the right
 130 panel of Fig. 1. We find that FD scales linearly with the dimensionality of θ as expected, but that
 131 ID is nearly constant, meaning the cost of the linear solver is subdominant. For the configurations
 132 considered, we reach multiple orders of magnitude speedups over FD.

133 5 Conclusions

134 In this work, we have shown that ID makes the MUSE algorithm faster and removes reliance on
 135 numerically-noisy FD. It requires second-order derivatives through the joint likelihood, but not
 136 through any optimizer, and never fully with respect to the latent space, meaning MUSE with ID
 137 is still well-suited for very high-dimensional problems. We have also provided examples of MUSE
 138 applied to BNNs and PPCA, demonstrating the extended applicability of the algorithm, which had
 139 previously been tested only on simpler toy problems or more complex but less general problems in
 140 cosmology [12, 2]. Beyond speed and accuracy improvements, removing the need to verify or tweak
 141 FD step-sizes represents a significant usability advancement for the algorithm.

References

- 142
- 143 [1] Uros Seljak, Grigor Aslanyan, Yu Feng, and Chirag Modi. Towards optimal extraction of cos-
144 mological information from nonlinear data. *Journal of Cosmology and Astroparticle Physics*,
145 2017(12):009–009, December 2017.
- 146 [2] Marius Millea and Uroš Seljak. Marginal unbiased score expansion and application to CMB
147 lensing. *Physical Review D*, 105:103531, May 2022.
- 148 [3] Kyle Cranmer, Johann Brehmer, and Gilles Louppe. The frontier of simulation-based infer-
149 ence. *Proceedings of the National Academy of Sciences*, 117(48):30055–30062, December
150 2020.
- 151 [4] D Duvenaud, J. Z. Kolter, and M. Johnson. Deep Implicit Layers: Neural ODEs, Equilibrium
152 Models, and Differentiable Optimization. <https://implicit-layers-tutorial.org>.
- 153 [5] Mathieu Blondel, Quentin Berthet, Marco Cuturi, Roy Frostig, Stephan Hoyer, Felipe Llinares-
154 López, Fabian Pedregosa, and Jean-Philippe Vert. Efficient and Modular Implicit Differentia-
155 tion, May 2022.
- 156 [6] Radford M. Neal. Slice sampling. *The Annals of Statistics*, 31(3):705–767, June 2003.
- 157 [7] NumPyro Documentation: Bayesian Neural Network.
158 <https://num.pyro.ai/en/stable/examples/bnn.html>.
- 159 [8] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and*
160 *Statistics)*. Springer-Verlag, Berlin, Heidelberg, 2006.
- 161 [9] Du Phan, Neeraj Pradhan, and Martin Jankowiak. Composable Effects for Flexible and Accel-
162 erated Probabilistic Programming in NumPyro, December 2019.
- 163 [10] Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofa-
164 nis Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall, and Noah D. Goodman. Pyro: Deep
165 Universal Probabilistic Programming. *Journal of Machine Learning Research*, 20(28):1–6,
166 2019.
- 167 [11] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
168 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao
169 Zhang. JAX: Composable transformations of Python+NumPy programs, 2018.
- 170 [12] Benjamin Horowitz, Uros Seljak, and Grigor Aslanyan. Efficient Optimal Reconstruction of
171 Linear Fields and Band-powers from Cosmological Data. *Journal of Cosmology and Astropar-*
172 *ticle Physics*, 2019(10):035–035, October 2019.