

Variational bagging: a robust approach for Bayesian uncertainty quantification

Shitao Fan, Ilsang Ohn, David Dunson and Lizhen Lin

Abstract: Variational Bayes methods are popular due to computational efficiency and adaptivity to diverse applications. In specifying the variational family, mean-field classes are commonly chosen, which enables efficient algorithms such as coordinate ascent variational inference (CAVI), but fails to capture parameter dependence and typically underestimates uncertainty. In this work, we introduce a variational bagging approach that integrates a bagging procedure with variational Bayes, resulting in a *bagged variational posterior* for improved inference. We establish strong theoretical guarantees, including posterior contraction rates for general models and a Bernstein–von Mises (BVM)-type theorem that ensures valid uncertainty quantification. Notably, our results show that even when using a mean-field variational family, our approach can recover off-diagonal elements of the limiting covariance structure and, crucially, provide proper uncertainty quantification. In addition, variational bagging is robust to model misspecification, with the covariance structures matching that of the target covariance. We illustrate our variational bagging in numerical studies through applications to parametric models, finite mixture models, deep neural networks, and variational autoencoders (VAEs).

Keywords: Bagging; Bagged VAE; Deep neural networks; Model misspecification; Posterior contraction rates; Robustness; Uncertainty quantification; Variational Bayes.

1. Introduction

Variational Bayes has emerged as a powerful framework for scalable Bayesian inference by approximating posterior distributions in complex and high-dimensional models with simpler variational families bypassing the need for MCMC sampling. There have been wide applications including for topic modeling [26], graphical models [13, 2], Bayesian nonparametric modeling [3], and high-dimensional sparse models [24]. Variational Bayes also plays an important role in modern generative AI through the variational autoencoder (VAE) [16, 6]. In specifying the variational family, a commonly used approximation within this framework is the mean-field family, which simplifies computation by assuming independence among model parameters. Although this assumption facilitates efficient algorithms like coordinate ascent varia-

tional inference (CAVI), it can lead to poor uncertainty quantification by ignoring parameter dependencies and systematically underestimating posterior variance. Although there have been attempts to mitigate problems with mean-field approaches [10, 27, 14, 20], most existing approaches are tailored to specific model classes, for example Gaussian or sparse Gaussian process models.

In this work, we propose *variational bagging*, a method that combines variational Bayes with a bootstrap aggregation (bagging) procedure to produce a bagged variational posterior with improved inferential properties, especially uncertainty quantification. We will show that variational bagging, which yields a *bagged variational posterior*, comes with strong theoretical guarantees, including posterior contraction rates for general models and a Bernstein–von Mises (BvM)-type result that ensures valid uncertainty quantification asymptotically. Remarkably, even when restricting to a mean-field variational family, our method can recover aspects of the full covariance structure, including off-diagonal elements.

Beyond the well-known undercoverage of variational Bayes, there is a broader concern about robustness to model misspecification. Bayesian statistics, including variational Bayes, is a model-based approach, which comes with the implicit assumption that the likelihood is correctly specified. In reality, however, knowledge about the true model or model class is rarely given, and model misspecification is more often than not. For many model classes, the posterior distribution (or their variational approximations) and our corresponding inferences and predictions are sensitive to model misspecification. As sample size increases in the misspecified case, the posterior typically concentrates around the pseudo-true parameter value corresponding to the minimal Kullback-Leibler (KL) divergence from the true data-generating model. Kleijn and van der Vaart [17] provide a Bernstein von Mises (BvM) theory characterizing the limiting form of the posterior under misspecification, showing that in general Bayesian credible sets are not valid confidence sets when the model is misspecified even asymptotically. We will show that variational bagging is robust to model misspecification, and yield a covariance that matches the target covariance.

‘BayesBag’ is a simple and general approach to obtain a robustified posterior by averaging over posteriors defined conditionally on different bootstrap-replicated datasets [5, 12]. This is an application of the bootstrap aggregation (bagging) approach of [4]. Let $X_{1:M}^* = (X_1^*, \dots, X_M^*)$ be a bootstrapped sample of size M from the data $X = (X_1, \dots, X_n)$ with sample size n . Let $\pi(\theta|X^*)$ be the posterior distribution conditioned on the bootstrapped copy

of the data. The bagged posterior [12] is defined as

$$\tilde{\pi}(\theta|X_{1:n}) = \frac{1}{n^M} \sum_{X_{1:M}^*} \pi(\theta|X_{1:M}^*), \quad (1)$$

which is defined over all possible bootstrap datasets. To reduce computational burden, we may use B bootstrap samples, with $B \ll n^M$, to obtain

$$\tilde{\pi}(\theta|X_{1:n}) \approx \frac{1}{B} \sum_{b=1}^B \pi\left(\theta|X_{(b)}^*\right),$$

where each $X_{(b)}^*$ is a bootstrap sample of size M . Huggins and Miller [12] studied asymptotic properties of the bagged posterior, including a BvM type theorem, and showed good predictive and uncertainty quantification (UQ) performance under model misspecification in simulation studies.

One potential drawback with BayesBag is its computation cost in the typical case in which $\pi(\theta|X_{1:M}^*)$ is intractable. Running MCMC sampling for each bootstrap replicate is often infeasible. We propose a *variational bagging approach* by first providing a variational approximation of each $\pi(\theta|X_{1:M}^*)$, the collection of which is then aggregated to produce the final *bagged variational posterior distribution* for inference. Compared to BayesBag, variational bagging massively speeds up the computation while yielding new and interesting theoretical results that rely on the unique properties of variational Bayes.

To summarize, we outline our theoretical results:

1. First, although variational approximations are well known to underestimate posterior variance in many cases, we provide a BvM theorem showing that bootstrap aggregation not only accommodates model misspecification but also appropriately inflates the variance so that the bagged posterior has theoretical guarantees of accurate uncertainty quantification. Remarkably, the limiting covariance structure (with simple adjustment) matches the target covariance even if the popular mean-field variational family is adopted. When variational inference is not conducted over the latent variable Z , the covariance coincides with the asymptotic variance of the MLE under model misspecification.
2. Second, we provide theoretical results on posterior contraction rates of the resulting bagged VB posterior and conditions on when such rates are minimax optimal. Our results encompass complex and nonparametric models and required the development of new technical tools.

3. Finally, we show that bagged variational posterior distributions lack overconfident credible sets.

From the practical side, a major advantage of variational bagging is in providing a natural mechanism for better approximating complex multi-modal posteriors by combining many local variational approximations from different bootstrapped data. In many modern complex models, ranging from intricate mixture models to deep neural networks, MCMC algorithms can fail to adequately explore the complex and multi-modal posterior landscape, while simple variational approximations without bagging only capture local features of this landscape. The variational bagging procedure can naturally overcome the above issues and we demonstrate this in a rich variety of examples including finite mixture models, deep neural networks, and variational auto-encoders.

The paper is organized as follows. Section 2 introduces notation and background on MLE, Bayes and variational Bayes under model misspecification. In Section 3, we describe the variational bagging algorithm in general and provide associated theory, while discussing practical aspects. Section 3.4 shows applications of variational bagging to a variety of models, providing algorithmic details and theory support in these specific contexts. Section 4 provides an empirical evaluation and demonstration through a simulation study and real data analyses.

2. Preliminaries

2.1. Model setup and notation

We consider a probability triple $(\mathcal{X}, \mathcal{F}, P_0)$ with corresponding density p_0 . Let $X_{1:n} = (X_1, \dots, X_n)$ be an i.i.d. sample from P_0 . We model the data using a parametric family $\mathcal{P}_\Theta = \{P_\theta : \theta \in \Theta\}$, where each P_θ has density $p(\cdot | \theta)$ indexed by a parameter θ .

The true distribution P_0 is not assumed to belong to the parametric family \mathcal{P}_Θ ; that is, we explicitly allow for model misspecification. We denote by θ_0 the pseudo-true parameter, defined as the minimizer of the Kullback–Leibler (KL) divergence:

$$\theta_0 = \arg \min_{\theta \in \Theta} \text{KL}(P_\theta, P_0),$$

where $\text{KL}(\cdot, \cdot)$ denotes the KL divergence

$$\text{KL}(P_1, P_2) = \int \log\left(\frac{dP_1}{dP_2}\right) dP_1,$$

whenever P_1 is absolutely continuous with respect to P_2 . When the corresponding densities p_1 and p_2 exist, we write $\text{KL}(p_1, p_2)$ for the induced KL divergence. When the model is correctly specified, we have $P_{\theta_0} = P_0$.

The maximum likelihood estimator (MLE),

$$\hat{\theta}_{\text{mle}} = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \log p(X_i | \theta),$$

is asymptotically centered at θ_0 and, under standard regularity conditions, satisfies

$$\sqrt{n}(\hat{\theta}_{\text{mle}} - \theta_0) \xrightarrow{d} N(0, V(\theta_0)^{-1} D(\theta_0) V(\theta_0)^{-1}),$$

where

$$\begin{aligned} V(\theta) &= -E_{P_0}[\nabla^2 \log p(X | \theta)], \\ D(\theta) &= E_{P_0}[\nabla \log p(X | \theta) \nabla \log p(X | \theta)^\top]. \end{aligned}$$

When the model is correctly specified, the above ‘‘sandwich’’ covariance reduces to the inverse Fisher information matrix. Under misspecification, any unbiased estimator $\hat{\theta}$ of θ_0 satisfies

$$\text{Var}(\hat{\theta}) \geq V(\theta_0)^{-1} D(\theta_0) V(\theta_0)^{-1},$$

so that the MLE is asymptotically efficient in the usual sandwich sense.

2.2. Bernstein–von Mises (BvM) theorem for a Bayesian model

We first introduce a *local asymptotic normality* (LAN) condition for misspecified models, which will be useful in our later discussions. A model is said to satisfy a stochastic LAN condition around $\theta_0 \in \Theta$ relative to a rate $\delta_n \rightarrow 0$ if there exists a random vector Δ_{n,θ_0} , bounded in P_0 -probability, such that for every compact set $K \subset \mathbb{R}^d$,

$$\sup_{h \in K} \left| \log \prod_{i=1}^n \frac{p(X_i | \theta_0 + \delta_n h)}{p(X_i | \theta_0)} - h^\top V(\theta_0) \Delta_{n,\theta_0} - \frac{1}{2} h^\top V(\theta_0) h \right| \xrightarrow{P_0} 0.$$

Under the above stochastic LAN condition, Kleijn and van der Vaart [17] prove an asymptotic normality result for the posterior distribution. Writing $\vartheta \sim \pi(\theta | X_{1:n})$ for a draw from the posterior, one has

$$\sqrt{n}(\vartheta - \theta_0) - \Delta_{n,\theta_0} \xrightarrow{d} N(0, V(\theta_0)^{-1}).$$

For a misspecified model, the covariance of the limiting normal distribution differs from the sandwich covariance $V(\theta_0)^{-1}D(\theta_0)V(\theta_0)^{-1}$ introduced above. This discrepancy implies that, under misspecification, credible sets derived from the usual Bayesian posterior are in general not expected to have asymptotically valid frequentist coverage.

2.3. Variational Bayes and asymptotic properties

Variational Bayes approximates the posterior distribution by a member of a prespecified parametric family (the variational family) by minimizing the KL divergence between the posterior distribution and distributions in this family [7]. Variational Bayes has become one of the most popular approaches for posterior approximation due to its simplicity, generality, and computational efficiency. There is also an emerging literature providing theoretical guarantees for variational methods [35, 23, 1, 21].

We consider a setting in which the unknowns consist of latent variables $Z_{1:n} = (Z_1, \dots, Z_n)$ and a global parameter $\theta = (\theta_1, \dots, \theta_d)$. Variational Bayes aims to approximate the joint posterior distribution $\pi(\theta, Z_{1:n} | X_{1:n})$ by solving the optimization problem

$$q^*(\theta, Z_{1:n}) = \arg \min_{q \in \mathcal{Q}} \text{KL}(q, \pi(\cdot | X_{1:n})),$$

where \mathcal{Q} denotes the variational family. In this paper, we focus on the mean-field variational family

$$\mathcal{Q} = \left\{ q : q(\theta, Z_{1:n}) = \prod_{j=1}^d q_{\theta_j}(\theta_j) \prod_{i=1}^n q_{Z_i}(Z_i) \right\}. \quad (2)$$

Recently, frequentist asymptotic properties of variational Bayes approximations have been established, including consistency, contraction rates, and BVM [35, 23, 34, 1, 30, 29]. We briefly review the relevant results.

Definition 2.1. The *variational log-likelihood* is defined as

$$\log p_{VB}(X | \theta) = \max_{q_Z} E_{q_Z} [\log p(X, Z | \theta) - \log q_Z(Z)]. \quad (3)$$

Wang and Blei [29] introduced a LAN-type condition for the variational log-likelihood. A model is said to satisfy a stochastic *variational local asymptotic normality (VLAN)* condition around an interior point $\theta_0 \in \Theta$ relative

to a rate $\delta_n \rightarrow 0$ if the following holds: there exists a random vector Δ_{n,θ_0} , bounded in P_0 -probability, such that for every compact set $K \subset \mathbb{R}^d$,

$$\sup_{h \in K} \left| \log \prod_{i=1}^n \frac{p_{\text{VB}}(X_i | \theta_0 + \delta_n h)}{p_{\text{VB}}(X_i | \theta_0)} - h^\top V_{\text{VB}}(\theta_0) \Delta_{n,\theta_0} - \frac{1}{2} h^\top V_{\text{VB}}(\theta_0) h \right| \xrightarrow{P_0} 0,$$

where

$$V_{\text{VB}}(\theta) = -E_{P_0} [\nabla^2 \log p_{\text{VB}}(X | \theta)], \quad (4)$$

$$D_{\text{VB}}(\theta) = E_{P_0} [\nabla \log p_{\text{VB}}(X | \theta) \nabla \log p_{\text{VB}}(X | \theta)^\top]. \quad (5)$$

Under the VLAN condition (together with additional regularity conditions), Wang and Blei [29] show that the limiting distribution of θ under the variational posterior is Gaussian:

$$\sqrt{n}(\vartheta - \theta_0) - \Delta_{n,\theta_0} \xrightarrow{d} N(0, (\tilde{V}_{\text{VB}}^0)^{-1}), \quad (6)$$

for $\vartheta \sim \int q^*(\theta, Z) dZ$, where \tilde{V}_{VB}^0 is the diagonal matrix that has the same diagonal entries as $V_{\text{VB}}(\theta_0)$. As is well-known and as shown in (6), the VB covariance with a mean-field class is only diagonal.

2.4. Asymptotic properties of BayesBag under model misspecification

Huggins and Miller [12] show that, under suitable regularity conditions on the log density $\log p(x | \theta)$, the BayesBag posterior $\tilde{\pi}(\theta | X_{1:n})$ in (1) satisfies the following asymptotic normality: for $\vartheta \sim \tilde{\pi}(\theta | X_{1:n})$,

$$\sqrt{n}(\vartheta - \theta_0) - \Delta_n | X_{1:n} \xrightarrow{d} N\left(0, \frac{1}{c} V(\theta_0)^{-1} + \frac{1}{c} V(\theta_0)^{-1} D(\theta_0) V(\theta_0)^{-1}\right),$$

for some random vector Δ_n that is bounded in P_0 -probability, where $c = \lim_{n \rightarrow \infty} M/n$ and M is the bootstrap sample size.

In contrast to the original Bayesian posterior, whose limiting covariance under misspecification is $V(\theta_0)^{-1}$ (see BvM result in Section 2), the BayesBag posterior has asymptotic covariance $\frac{1}{c} V(\theta_0)^{-1} + \frac{1}{c} V(\theta_0)^{-1} D(\theta_0) V(\theta_0)^{-1}$, which is closer to the sandwich form $V(\theta_0)^{-1} D(\theta_0) V(\theta_0)^{-1}$ and therefore yields better-calibrated credible sets under model misspecification. In particular, taking $c = 1$ already leads to conservative uncertainty quantification. Moreover, as discussed in Section 4.1, one can select c in a principled way so that the resulting covariance approximates the sandwich form $V(\theta_0)^{-1} D(\theta_0) V(\theta_0)^{-1}$.

3. Variational bagging: bagged variational posterior

In this section, we introduce our *variational bagging* approach. For a bootstrap sample $X_{1:M}^* = (X_1^*, \dots, X_M^*)$ generated from the original data $X_{1:n} = (X_1, \dots, X_n)$, define

$$q^*(\theta, Z_{1:M}^* | X_{1:M}^*) = \arg \min_{q \in \mathcal{Q}} \text{KL}(q(\theta, Z_{1:M}^*), \pi(\theta, Z_{1:M}^* | X_{1:M}^*))$$

as the joint variational approximation to the posterior given this bootstrap dataset. We focus on inference for θ via the marginal

$$q^*(\theta | X_{1:M}^*) = \int q^*(\theta, Z_{1:M}^* | X_{1:M}^*) dZ_{1:M}^*.$$

When \mathcal{Q} is a mean-field family, $q^*(\theta | X_{1:M}^*)$ is easily obtained using the factorization structure.

We define the *bagged variational posterior* as the average of the variational posteriors obtained from B bootstrap samples:

$$q^{\text{bVB}}(\theta | X_{1:n}) = \frac{1}{B} \sum_{b=1}^B q^*(\theta | X_{(b)}^*),$$

where each $X_{(b)}^*$ denotes a bootstrap sample of size M . In our theoretical analysis of the bagged variational posterior, we focus on the mean-field variational family.

3.1. Robust uncertainty quantification

In this section, we derive the asymptotic distribution of the bagged variational posterior and discuss implications for uncertainty quantification. Recall the definitions of $V_{\text{VB}}(\cdot)$ and $D_{\text{VB}}(\cdot)$ from (4) and (5), that is,

$$\begin{aligned} V_{\text{VB}}(\theta) &= -E_{P_0}[\nabla^2 \log p_{\text{VB}}(X | \theta)], \\ D_{\text{VB}}(\theta) &= E_{P_0}[\nabla \log p_{\text{VB}}(X | \theta) \nabla \log p_{\text{VB}}(X | \theta)^\top]. \end{aligned}$$

Notably, we show that under model misspecification, the asymptotic covariance of the bagged variational posterior contains an additional “sandwich” term

$$(V_{\text{VB}}^0)^{-1} D_{\text{VB}}^0 (V_{\text{VB}}^0)^{-1}, \quad V_{\text{VB}}^0 = V_{\text{VB}}(\theta_0), \quad D_{\text{VB}}^0 = D_{\text{VB}}(\theta_0),$$

on top of the covariance $(\tilde{V}_{\text{VB}}^0)^{-1}$ arising from the usual variational posterior, where \tilde{V}_{VB}^0 is the diagonal matrix formed from the diagonal entries of V_{VB}^0 . In particular, $(\tilde{V}_{\text{VB}}^0)^{-1}$ corresponds to the covariance under the standard mean-field VB procedure, which only captures marginal variances.

Theorem 3.1 (Bernstein–von Mises theorem for the bagged variational posterior). *Let $\ell_\theta(x) = \log p_{\text{VB}}(x | \theta)$, and assume the following conditions hold:*

1. *The map $\theta \mapsto \ell_\theta(X_1)$ is differentiable at θ_0 in probability.*
2. *There exists an open neighborhood U of θ_0 and a function $m_{\theta_0} : \mathcal{X} \rightarrow \mathbb{R}$ such that $E_{P_0}(m_{\theta_0}^3) < \infty$, and for all $\theta, \theta' \in U$,*

$$|\ell_\theta - \ell_{\theta'}| \leq m_{\theta_0} \|\theta - \theta'\|_2 \quad a.s. [P_0].$$

3. *As $\theta \rightarrow \theta_0$,*

$$-E_{P_0}(\ell_\theta - \ell_{\theta_0}) = \frac{1}{2}(\theta - \theta_0)^\top V_{\text{VB}}^0(\theta - \theta_0) + o(\|\theta - \theta_0\|_2^2).$$

4. *V_{VB}^0 is invertible.*

5. *For every $\epsilon > 0$, there exists a sequence of tests ϕ_n such that*

$$\begin{aligned} \int \phi_n(x_1, \dots, x_n) \prod_{i=1}^n p_0(x_i) dx_i &\rightarrow 0, \\ \sup_{\theta: \|\theta - \theta_0\| > \epsilon} \int \{1 - \phi_n(x_1, \dots, x_n)\} \prod_{i=1}^n \frac{p_{\text{VB}}(x_i | \theta)}{p_{\text{VB}}(x_i | \theta_0)} p_0(x_i) dx_i &\rightarrow 0. \end{aligned}$$

6. *$c = \lim_{n \rightarrow \infty} M/n \in (0, \infty)$, where M is the bootstrap sample size.*

Then, for $\vartheta^\dagger \sim q^{\text{bvB}}(\theta | X_{1:n})$, we have

$$\sqrt{n}(\vartheta^\dagger - \theta_0) - \Delta_n \Big| X_{1:n} \xrightarrow{d} N\left(0, \frac{1}{c}(\tilde{V}_{\text{VB}}^0)^{-1} + \frac{1}{c}(V_{\text{VB}}^0)^{-1} D_{\text{VB}}^0 (V_{\text{VB}}^0)^{-1}\right), \quad (7)$$

where

$$\Delta_n = n^{1/2}(V_{\text{VB}}^0)^{-1}(\mathbb{P}_n - P_0)\dot{\ell}_{\theta_0}, \quad \mathbb{P}_n = n^{-1} \sum_{i=1}^n \delta_{X_i},$$

and \tilde{V}_{VB}^0 is the diagonal matrix with the same diagonal entries as V_{VB}^0 .

Remark 1. The sandwich term $(V_{\text{VB}}^0)^{-1}D_{\text{VB}}^0(V_{\text{VB}}^0)^{-1}$ in Theorem 3.1 can be viewed as the *target covariance*: it corresponds to the covariance of the MLE defined via the variational log-likelihood (3). This term is typically *non-diagonal*, even when a mean-field variational family is used. It is precisely this additional sandwich term that drives robustness and more accurate uncertainty quantification of VB bagging: even under model misspecification, variational bagging attempts to recover the “best” covariance structure allowed by the model and the variational approximation.

When the model is specified correctly, we have $\text{diag}((V_{\text{VB}}^0)^{-1}D_{\text{VB}}^0(V_{\text{VB}}^0)^{-1}) = \text{diag}((\tilde{V}_{\text{VB}}^0)^{-1})$. That is, these two terms share the same diagonal covariance. The covariance in (7) can be decomposed as

$$\frac{1}{c} \left((V_{\text{VB}}^0)^{-1}D_{\text{VB}}^0(V_{\text{VB}}^0)^{-1} - (\tilde{V}_{\text{VB}}^0)^{-1} \right) + \frac{2}{c}(\tilde{V}_{\text{VB}}^0)^{-1}.$$

The first term above contains only off-diagonal entries, while the second term is purely diagonal. In this case, one can recover the off-diagonal entries of the target covariance $(V_{\text{VB}}^0)^{-1}D_{\text{VB}}^0(V_{\text{VB}}^0)^{-1}$ from the bagged VB posterior (e.g., with $c = 1$), and for the diagonal term we need to make the simple adjustment by multiplying them by 1/2. Combining these two pieces yields an estimator of the full target covariance.

When the model is misspecified, we can still recover the off-diagonal structure by setting $c = 1$ and extracting only the off-diagonal entries from the bagged VB covariance. For the diagonal entries, we can use the choice of M described in Section 4.1 to match the desired marginal variances.

Remark 2. If we instead use a variational family of Gaussian distributions with general covariance matrices, $\mathcal{Q} = \{q : q(\theta) = N(\mu, \Sigma)\}$, i.e., without a mean-field restriction, then by inspecting the proof of Theorem 3.1 one can see that the limiting distribution of the bagged variational posterior becomes

$$\sqrt{n}(\vartheta^\dagger - \theta_0) - \Delta_n \Big| X_{1:n} \xrightarrow{d} N\left(0, \frac{1}{c}(V_{\text{VB}}^0)^{-1} + \frac{1}{c}(V_{\text{VB}}^0)^{-1}D_{\text{VB}}^0(V_{\text{VB}}^0)^{-1}\right).$$

The next corollary is a special case of Theorem 3.1 when variational inference is not conducted over a latent variable Z (that is, we marginalize Z in the model), so that $p_{\text{VB}}(x | \theta) = p(x | \theta)$.

Corollary 3.2 (BvM theorem without latent variables). *Assume that $p_{\text{VB}}(x | \theta) = p(x | \theta)$. Let $V^0 = V(\theta_0)$ and $D^0 = D(\theta_0)$, and let \tilde{V}^0 be the diagonal matrix with the same diagonal entries as V^0 . Then, under the same condi-*

tions as in Theorem 3.1, for $\vartheta^\dagger \sim q^{\text{bvB}}(\theta | X_{1:n})$,

$$\sqrt{n}(\vartheta^\dagger - \theta_0) - \Delta_n \Big| X_{1:n} \xrightarrow{d} N\left(0, \frac{1}{c}(\tilde{V}^0)^{-1} + \frac{1}{c}(V^0)^{-1}D^0(V^0)^{-1}\right),$$

where

$$\Delta_n = n^{1/2}(V^0)^{-1}(\mathbb{P}_n - P_0)\dot{\ell}_{\theta_0}.$$

In Theorem 3.2, the sandwich term $(V^0)^{-1}D^0(V^0)^{-1}$ does not depend on the choice of variational family and coincides with the usual sandwich covariance for the MLE. Moreover, the full asymptotic covariance of the bagged variational posterior is larger than that of the MLE, with the difference given by $(\tilde{V}^0)^{-1}/c$. Thus, the bagged variational posterior yields conservative uncertainty quantification: asymptotically, it is never overconfident. This is rigorously formalized in the following corollary.

Corollary 3.3 (No overconfident credible sets). *Assume the same conditions as in Theorem 3.2 and take $M = n$ (so that $c = 1$). Consider the ellipsoid*

$$C(r) = \left\{ \theta \in \Theta : n(\theta - \hat{\theta}_{\text{mle}})^\top \widehat{\Sigma}^{-1}(\theta - \hat{\theta}_{\text{mle}}) \leq r^2 \right\}$$

with radius $r > 0$, where $\hat{\theta}_{\text{mle}}$ is the MLE of θ_0 and $\widehat{\Sigma}$ is a consistent estimator of the asymptotic covariance

$$\Sigma_0 = (\tilde{V}^0)^{-1} + (V^0)^{-1}D^0(V^0)^{-1}$$

of the bagged variational posterior. For $\alpha \in (0, 1)$, let $r_{n,1-\alpha}$ be such that $C(r_{n,1-\alpha})$ is a $(1 - \alpha)$ -credible set for θ_0 under the bagged variational posterior, i.e.,

$$Q^{\text{bvB}}(\vartheta^\dagger \in C(r_{n,1-\alpha})) = 1 - \alpha.$$

Then

$$\lim_{n \rightarrow \infty} P_0(\theta_0 \in C(r_{n,1-\alpha})) \geq 1 - \alpha.$$

3.2. Valid variational Bayes uncertainty quantification

Mean-field variational Bayes is well known for providing a fast and tractable approximation of the Bayesian posterior. As shown by Wang and Blei [30] (see Equation (6)), mean-field approximations are asymptotically normal under standard conditions, but they *only approximate the diagonal terms of*

the covariance structure, ignoring dependence among parameters. For this reason, although variational Bayes often delivers accurate first-order inference (e.g., point estimates), it is generally unsuitable for second-order inference such as uncertainty quantification, even when the model is correctly specified.

An interesting consequence of Theorem 3.1 is that the variational bagging procedure can address this limitation by recovering the off-diagonal elements of the covariance structure that are wiped out in the standard mean-field variational approximation. In particular, when the model is correctly specified, variational bagging yields asymptotically valid uncertainty quantification, comparable to that of the standard Bayesian posterior. We briefly discussed this in Remark 1 for the case of a well-specified model with the presence of latent variables, and the following deals with the case when there is no latent variable.

Corollary 3.4 (BvM theorem when the model is correctly specified). *Let $\ell_\theta(x) = \log p_{\text{VB}}(x | \theta)$ and assume Assumptions 1–5 in Theorem 3.1 hold, with $M = n$. In this case θ_0 is the true parameter so that $P_0 = P_{\theta_0}$. Then, for $\vartheta^\dagger \sim q^{\text{bvB}}(\theta | X_{1:n})$, we have*

$$\sqrt{n}(\vartheta^\dagger - \theta_0) - \Delta_n \mid X_{1:n} \xrightarrow{d} N(0, (\tilde{V}_{\text{VB}}^0)^{-1} + (V_{\text{VB}}^0)^{-1}),$$

where $V_{\text{VB}}^0 = V_{\text{VB}}(\theta_0)$ and \tilde{V}_{VB}^0 is the diagonal matrix with the same diagonal entries as V_{VB}^0 .

When the model is correctly specified and the usual regularity conditions hold, the Fisher information structure implies that

$$\text{diag}((V_{\text{VB}}^0)^{-1}) = \text{diag}((\tilde{V}_{\text{VB}}^0)^{-1}),$$

so that the limiting covariance in Corollary 3.4 can be decomposed as

$$(\tilde{V}_{\text{VB}}^0)^{-1} + (V_{\text{VB}}^0)^{-1} = 2 \text{diag}((V_{\text{VB}}^0)^{-1}) + \text{offdiag}((V_{\text{VB}}^0)^{-1}),$$

where $\text{offdiag}(A)$ denotes the matrix obtained from A by zeroing out its diagonal. The first term above only has nonzero diagonal entries, while the second term only contains off-diagonal entries. Thus:

- the *off-diagonal* part of the target covariance $(V_{\text{VB}}^0)^{-1}$ is recovered by the bagged variational posterior, and
- the *diagonal* part is inflated by a factor of 2 relative to the target covariance.

To recover the same covariance as in the standard Bayesian BvM result, one can simply rescale the diagonal entries of the bagged VB covariance by 1/2, keeping the off-diagonal entries unchanged. In this way, variational bagging can serve as a general tool to improve and calibrate uncertainty quantification based on mean-field variational Bayes, restoring both the correct marginal variances (after a simple adjustment) and the dependence structure.

3.2.1. Toy example illustration

We illustrate the implication of Corollary 3.4 with a simple toy example. Consider estimation of the mean vector μ for two-dimensional Gaussian data X_1, \dots, X_{500} , where

$$X_i \sim N(\mu, \Sigma), \quad \mu = (-1, 1)^\top, \quad \Sigma = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}.$$

We place a conjugate prior on μ and compare three posterior approximations: (i) Hamiltonian Monte Carlo (HMC), (ii) a mean-field variational approximation, and (iii) our variational bagging approach.

Figure 1 shows the resulting 95% posterior credible regions. The mean-field variational posterior provides a reasonable approximation of the posterior mean of μ , but it fails to capture the correlation structure between the components of μ and yields an elliptical region aligned with the coordinate axes. In contrast, the bagged variational posterior closely matches the full Bayesian posterior obtained by HMC, successfully recovering the correlation and orientation of the credible region. This illustrates how variational bagging can substantially improve uncertainty quantification over standard mean-field variational Bayes while retaining its computational benefits.

3.3. Contraction rates of the bagged variational posterior

Our BvM theorem (Theorem 3.1) implies that the bagged variational posterior contracts to θ_0 at the parametric rate $n^{-1/2}$ for finite-dimensional parametric models satisfying its assumptions. In this section, we extend this result to a general setup that encompasses nonparametric models, and we consider convergence in Hellinger distance. Let

$$H(P_1, P_2) = \left(\int (\sqrt{dP_1/dP_2} - 1)^2 dP_2 \right)^{1/2}$$

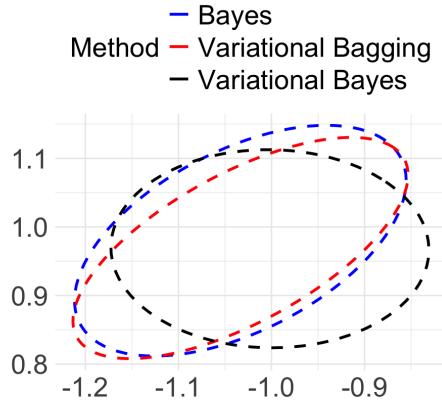


Fig 1: 95% posterior credible regions for a two-dimensional Gaussian mean: comparison of HMC, mean-field VB, and variational bagging.

denote the Hellinger distance between two probability measures P_1 and P_2 . To quantify the corresponding convergence rate, we use the notion of bracketing Hellinger metric entropy.

We say that a finite collection of pairs of functions $\{(f_i^L, f_i^U) : i = 1, \dots, N\}$ is a δ -bracketing of a function space \mathcal{F} if

$$\|(f_i^U)^{1/2} - (f_i^L)^{1/2}\|_2 \leq \delta, \quad i = 1, \dots, N,$$

and for any $f \in \mathcal{F}$ there exists $i \in \{1, \dots, N\}$ such that $f_i^L \leq f \leq f_i^U$. The δ -bracketing Hellinger metric entropy of \mathcal{F} , denoted by $\mathcal{H}_B(\delta, \mathcal{F})$, is defined as the logarithm of the cardinality of a minimal δ -bracketing.

Assumption 3.1. Let $(\epsilon_n)_{n=1,2,\dots}$ be a positive sequence such that $\epsilon_n \rightarrow 0$ and $n\epsilon_n^2 \rightarrow \infty$ as $n \rightarrow \infty$. Assume the following:

(1) **(Prior mass)** There exists a constant $C_1 > 0$ such that

$$\Pi\left(\theta \in \Theta : \text{KL}(P_0, P_\theta) \leq \epsilon_n^2, \text{KLV}(P_0, P_\theta) \leq \epsilon_n^2\right) \geq \exp(-C_1 n \epsilon_n^2),$$

where $\text{KLV}(P_0, P_\theta) = \int (\log(dP_\theta/dP_0))^2 dP_0$.

(2) **(Sieve and complexity)** There exists a subset $\Theta_n \subset \Theta$ such that

$$\Pi(\Theta \setminus \Theta_n) \leq \exp(-(C_1 + 4)n \epsilon_n^2)$$

and, for some constants $C_2 > 0$ and $C_3 > 0$,

$$\int_{\epsilon^2/2^8}^{\sqrt{2}\epsilon} \mathcal{H}_B^{1/2}(u/C_2, \{p(\cdot | \theta) : \theta \in \Theta_n\}) du \leq C_3 \sqrt{n} \epsilon^2$$

for any $\epsilon > \epsilon_n$.

(3) **(Variational family)** There exists a constant $C_4 > 0$ such that

$$\inf_{q \in \mathcal{Q}} \left[n \int \text{KL}(P_\theta, P_0) q(\theta) d\theta + \text{KL}(q, \pi) \right] \leq C_4 n \epsilon_n^2.$$

(4) **(Bootstrap sample size)** $M = n$.

The ‘‘prior mass and testing’’ method, first developed in the seminal paper Ghosal et al. [9], is a powerful and general tool for deriving contraction rates of posterior distributions. We retain a prior mass condition in Assumption 3.1(1), which plays the same role as in the prior mass and testing framework. However, in our setting it is not straightforward to construct suitable tests based on bootstrapped samples. To the best of our knowledge, there is no principled general approach for test construction for Bayesian procedures involving bootstrap-based posteriors.

Instead, we follow the strategy of Shen and Wasserman [25] with a suitable modification to handle bootstrap-weighted likelihoods. Han and Yang [11] adopt a similar idea, but their result is restricted to parametric models. In our nonparametric setting, we use the complexity condition on the Hellinger metric entropy in Assumption 3.1(2), adapted from Shen and Wasserman [25], to control the behavior of the empirical process of the ‘‘bootstrap-weighted’’ log-likelihood ratio. The third assumption on the variational family, which has been employed in several related works such as Ohn and Lin [21] and Zhang and Gao [35], controls the variational approximation error between the variational posterior and the original posterior distribution. We show that this assumption, together with the bootstrap sample size condition in Assumption 3.1(4), can still be used to bound the variational approximation error even when a bootstrapped sample is used.

Theorem 3.5 (Contraction rate). *Suppose that Assumption 3.1 holds. Then the bagged variational posterior satisfies*

$$E \left[Q^{\text{bvB}} (H^2(P_\theta, P_0) \geq M_n \epsilon_n^2 \log n) \right] \rightarrow 0,$$

as $n \rightarrow \infty$ for any diverging sequence (M_n) with $M_n \rightarrow \infty$, where the expectation is taken with respect to $P_0^{\otimes n}$.

3.4. Illustrating examples

In this section, we demonstrate our theory by considering two simple examples: a multivariate Gaussian example and a two-component finite mixture model.

3.4.1. Multivariate Gaussian mean

Consider again the toy example in Section 3.2.1. That is, $X_{1:n} \in \mathbb{R}^2$ with

$$X_i \sim N(\mu_0, \Lambda_0^{-1}),$$

where Λ_0 is known and we are interested in estimating the mean vector μ_0 . For the posterior distribution $\pi(\mu | X_{1:n})$, the Bernstein–von Mises theorem yields

$$\sqrt{n}(\mu - \bar{X}_n) \xrightarrow{d} N(0, \Lambda_0^{-1}),$$

for $\mu \sim \pi(\mu | X_{1:n})$, where \bar{X}_n is the sample mean.

For variational Bayes, we consider the mean-field variational family, which assumes $q(\mu) = q_1(\mu_1) q_2(\mu_2)$. By the BvM theorem for the variational posterior in Wang and Blei [29], we have

$$\sqrt{n}(\mu - \bar{X}_n) \xrightarrow{d} N\left(0, \begin{pmatrix} \Lambda_{022}^{-1} & 0 \\ 0 & \Lambda_{011}^{-1} \end{pmatrix} / \det(\Lambda_0)\right),$$

for μ drawn from the mean-field variational posterior $q^*(\mu)$, where we decompose the true precision matrix as

$$\Lambda_0 = \begin{pmatrix} \Lambda_{011} & \Lambda_{012} \\ \Lambda_{021} & \Lambda_{022} \end{pmatrix}.$$

Compared with the asymptotic distribution of the exact posterior, we see that variational Bayes ignores the correlation between μ_1 and μ_2 .

Now, applying Corollary 3.4, we obtain

$$\sqrt{n}(\mu - \bar{X}_n) | X_{1:n} \xrightarrow{d} N\left(0, \begin{pmatrix} \Lambda_{022}^{-1} & 0 \\ 0 & \Lambda_{011}^{-1} \end{pmatrix} / \det(\Lambda_0) + \Lambda_0^{-1}\right),$$

for μ drawn from the bagged variational posterior $q^{\text{bvB}}(\mu | X_{1:n})$.

Therefore, using variational bagging, we recover the off-diagonal (correlation) structure of the true posterior covariance matrix via the Λ_0^{-1} term,

while the diagonal terms are inflated relative to the target covariance. As discussed in Section 3.2, a simple correction of rescaling the diagonal entries by a factor of 1/2 recovers the full target covariance, providing a concrete illustration of how variational bagging improves uncertainty quantification over standard mean-field variational Bayes in this simple Gaussian setting.

3.4.2. Bayesian mixture models

In this example, we consider model misspecification in a finite mixture model. For technical simplicity, we assume that our working inference model is a symmetric two-component Gaussian mixture with unit variance,

$$p(x | \theta) = \frac{1}{2}N(x; \theta, 1) + \frac{1}{2}N(x; -\theta, 1),$$

with $\theta > 0$, while the true data-generating distribution P_0 is not necessarily in this model class.

The above mixture model can be equivalently written in hierarchical form as

$$\begin{aligned} X_i | Z_i &\sim N((2Z_i - 1)\theta, 1), \\ Z_i &\sim \text{Bernoulli}(1/2), \end{aligned}$$

so we can conduct variational inference jointly over the parameter θ and latent variables Z_1, \dots, Z_n . We consider a mean-field variational family in which each Z_i has a degenerate (point-mass) distribution at either 0 or 1, which is analogous to a hard clustering procedure.

In this case, it is straightforward to see that the variational log-likelihood is given by

$$\begin{aligned} \ell(\theta) &= \log p_{\text{VB}}(X | \theta) = \sup_{q(Z)} \int q(Z) \log \frac{p(X, Z | \theta)}{q(Z)} dZ \\ &= \max_{z \in \{0,1\}} \log p(X, z | \theta) \\ &= \frac{1}{2} \max\{-(X - \theta)^2, -(X + \theta)^2\} - \frac{1}{2} \log(2\pi) \\ &= -\frac{1}{2}(X - \text{sign}(X)\theta)^2 - \frac{1}{2} \log(2\pi). \end{aligned}$$

Differentiating twice in θ gives

$$V_{\text{VB}}(\theta) = -E_{P_0} \left[\frac{d^2}{d\theta^2} \ell(\theta) \right] = 1.$$

It is then easy to verify that the first through fourth conditions of Theorem 3.1 hold for this model. The fifth condition (existence of suitable tests) follows from Theorem 1 of Westling and McCormick [32], which establishes consistency of the maximum variational likelihood estimator in this setting. In our case, the maximum “variational” likelihood estimator is

$$\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i \operatorname{sign}(X_i),$$

and this estimator is consistent for the pseudo-true parameter θ_0 under model misspecification.

Applying Theorem 3.1, we obtain, for $\vartheta^\dagger \sim q^{\text{bvB}}(\theta | X_{1:n})$,

$$\sqrt{n}(\vartheta^\dagger - \hat{\theta}_n) \mid X_{1:n} \xrightarrow{d} N\left(0, \frac{1}{c} \left\{ 1 + E_{P_0}[(X - \operatorname{sign}(X)\theta_0)^2] \right\}\right),$$

where $c = \lim_{n \rightarrow \infty} M/n$ is the limiting ratio of bootstrap sample size to the original sample size. This illustrates how, even under model misspecification in a mixture setting, the bagged variational posterior enjoys a well-defined asymptotic distribution with a variance that incorporates a sandwich-type correction term.

4. Simulation study

In this section, we conduct a simulation study in which we apply variational bagging to several examples, including the multivariate Gaussian model, a finite Gaussian mixture model, sparse linear regression, regression based on feedforward deep neural networks, and a bagged VAE (variational autoencoder). We first discuss some practical aspects of variational bagging, such as the choice of bootstrap sample size and the number of bootstrap replicates.

4.1. Bootstrap sample size and number of bootstrap samples

If we have strong confidence in our model specification, we may set $M = n$ and use mean-field VB to learn the off-diagonal terms of the covariance and take 1/2 of the diagonal terms. Under model misspecification, however, as illustrated in the previous examples, choosing $M = n$ (i.e., $c = 1$) may not yield the desired robust, sandwich-type covariance.

Let \tilde{v}_n and \tilde{v}_n^* denote, respectively, the standard and bagged variational posterior variances of a real-valued functional $f(\theta)$, where the bagged variational posterior is computed with $M = n$. In the asymptotic BvM setting, the bagged variational covariance (with a general c) behaves like

$$\frac{1}{c}(V_{\text{VB}}^0)^{-1} + \frac{1}{c}(V_{\text{VB}}^0)^{-1}D_{\text{VB}}^0(V_{\text{VB}}^0)^{-1},$$

while the “target” (sandwich) covariance is

$$(V_{\text{VB}}^0)^{-1}D_{\text{VB}}^0(V_{\text{VB}}^0)^{-1}.$$

At the level of a scalar functional $f(\theta)$, this corresponds to finding c such that

$$\frac{1}{c}\tilde{v}_n + \frac{1}{c}(\tilde{v}_n^* - \tilde{v}_n) \approx \tilde{v}_n^* - \tilde{v}_n,$$

which yields

$$\frac{1}{c}\tilde{v}_n^* = \tilde{v}_n^* - \tilde{v}_n \implies c = \frac{\tilde{v}_n^*}{\tilde{v}_n^* - \tilde{v}_n}.$$

Hence the corresponding optimal bootstrap sample size is

$$M^* = c \cdot n = \frac{\tilde{v}_n^*}{\tilde{v}_n^* - \tilde{v}_n} n.$$

We therefore suggest the plug-in estimator

$$\hat{M}^* = \frac{\tilde{v}_n^*}{\tilde{v}_n^* - \tilde{v}_n} n. \quad (8)$$

In finite-sample settings, we also need to account for prior influence when choosing a suitable bootstrap sample size. Following Huggins and Miller [12], we define a finite-sample version of the optimal bootstrap size, denoted $\hat{M}_{\text{fs, opt}}^*$.

Let v_0 denote the prior variance of $f(\theta)$ and define

$$\hat{\sigma}_o^2 := n v_0 \tilde{v}_n / (v_0 - \tilde{v}_n),$$

and

$$\hat{s}_o^2 := \frac{v_0^2}{(v_0 - \tilde{v}_n)^2} (\tilde{v}_n^* - \tilde{v}_n) n. \quad (9)$$

Then the finite-sample optimal bootstrap size is estimated by

$$\hat{M}_{\text{fs, opt}} := \frac{n}{2} + \frac{n\hat{\sigma}_o^2}{2\hat{s}_o^2} - \frac{\hat{\sigma}_o^2}{v_0} + \left\{ \left(\frac{n}{2} + \frac{n\hat{\sigma}_o^2}{2\hat{s}_o^2} \right)^2 - \frac{n\hat{\sigma}_o^2}{v_0} \right\}^{1/2}. \quad (10)$$

For the derivation of (10), we refer to Appendix E of Huggins and Miller [12].

Regarding the number of bootstrap samples B , Huggins and Miller [12] recommend $B \approx 50$ to 100 for BayesBag due to the computational cost of repeated MCMC. In variational bagging, by contrast, computing each variational posterior is typically much faster than running MCMC, allowing substantially larger B (e.g., $B \approx 200$). In our experiments, however, we find that even a relatively small number of bootstrap samples, such as $B = 5$, is often sufficient to obtain robust and well-calibrated uncertainty quantification.

4.2. Uncertainty quantification for a multivariate Gaussian

Using the same toy example as in Section 3.2.1, we generate two-dimensional Gaussian data X_1, \dots, X_n , where

$$X_i \sim N(\mu, \Sigma), \quad \mu = (-1, 1)^\top, \quad \Sigma = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}.$$

We vary the number of bootstrap replicates B and the sample size n to assess the effectiveness of uncertainty quantification under variational bagging. Specifically, we take

$$B \in \{5, 10, 20, 30, 50\}, \quad n \in \{50, 100, 200, 300, 500, 1000\}.$$

Because Corollary 3.4 is an asymptotic result, our goal here is to identify practically reasonable choices of B and n for which the asymptotic behavior is already well approximated.

We place a conjugate prior on μ and approximate the posterior using three methods: (i) Hamiltonian Monte Carlo (HMC) with 2 chains and 2000 posterior draws, (ii) mean-field variational Bayes (MFVB), and (iii) variational bagging based on MFVB. For variational bagging, we compute MFVB for each bootstrap dataset and average the corresponding variational posteriors; runs in which the MFVB algorithm fails to converge are discarded.

From Figure 2, we observe that when the number of bootstrap replicates B is at least 30, the credible regions obtained from variational bagging closely match those from the full Bayesian posterior, even for relatively small sample sizes. In particular, the results suggest that our theoretical guarantees are already informative for sample sizes as small as $n = 50$, and that a moderate number of bootstrap replicates (around $B \geq 30$) suffices to capture the improved uncertainty quantification predicted by our theory.

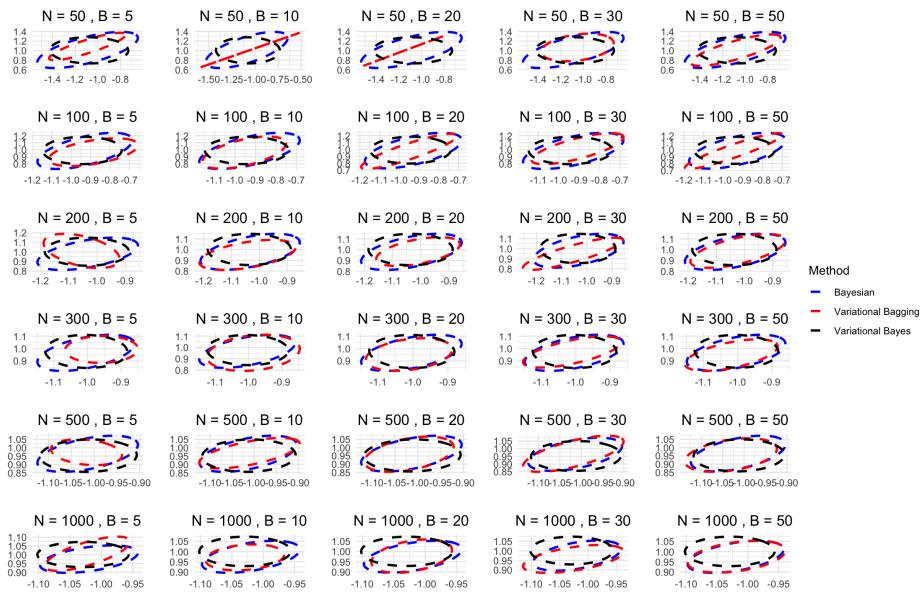


Fig 2: 95% posterior credible regions for the Gaussian mean under HMC, MFVB, and variational bagging, for varying sample size n and number of bootstrap replicates B .

4.3. Gaussian mixture model

This simulation study considers the Gaussian mixture model

$$X_i \sim \sum_{k=1}^K \pi_k N(\mu_k, \sigma_k^2), \quad i = 1, \dots, n.$$

For the prior, we use

$$\pi \sim \text{Dirichlet}(\alpha), \quad \mu_k \mid \sigma_k^2 \sim N(0, \nu_0 \sigma_k^2), \quad \sigma_k^2 \sim IG(a, b).$$

We generate data from heavy-tailed mixture distributions such as t mixtures and double-exponential mixtures, but fit a Gaussian mixture model to these data. We also fit the same data under the correctly specified (true) model to evaluate the performance of variational bagging.

For the full Bayesian posterior, we use *Stan* (R interface), based on Hamiltonian Monte Carlo (HMC), with 2000 iterations and 1000 burn-in per chain. For variational Bayes, we implement a CAVI algorithm. For bagging, we choose $B = 50$ bootstrap replicates and set the bootstrap sample size to be \hat{M}^* from (8). To accelerate computation, variational fits for different bootstrap samples are run in parallel.

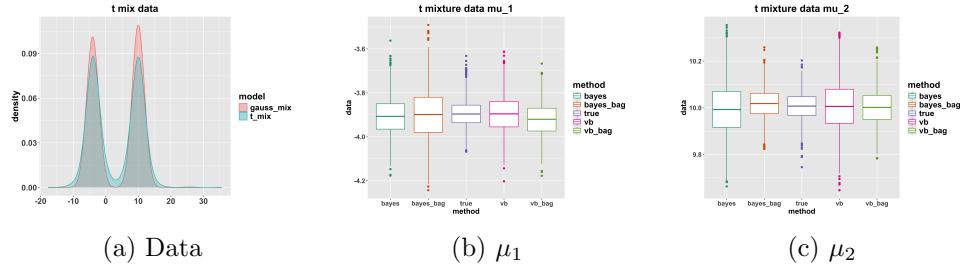


Fig 3: Gaussian mixture fit to t -mixture data.

Figure 3a shows an example where the data are generated from a t -mixture distribution but fitted with a Gaussian mixture model, leading to heavier tails in the data than assumed by the working model. Figures 3b and 3c show results for the component means μ_1 and μ_2 under several approaches, including variational bagging.

When using standard Bayesian methods and mean-field variational Bayes, the posteriors tend to be reasonably centered around the true means but exhibit poorly calibrated uncertainty, with credible intervals that do not match

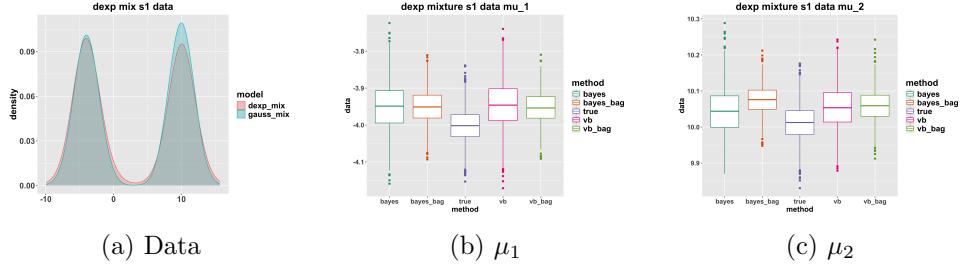


Fig 4: Gaussian mixture fit to data from a double-exponential mixture model.

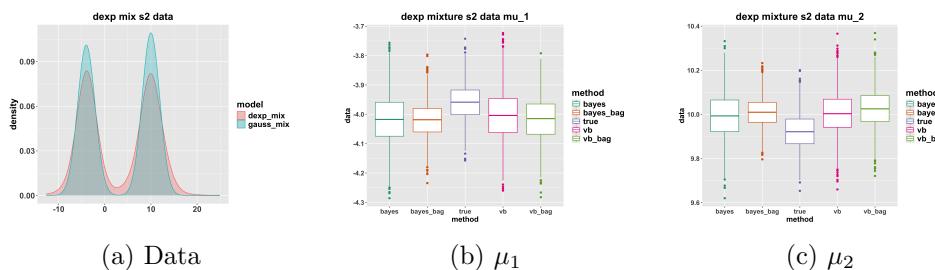


Fig 5: Gaussian mixture fit to double-exponential mixture data with larger variance.

the empirical variability induced by the heavy-tailed data. In contrast, our Variational Bagging procedure corrects for this misspecification effect and yields posterior coverage much closer to the empirical coverage. Here, “robustness” means that, relative to standard Bayesian or VB fits, the bagged procedures produce uncertainty statements that better reflect the true distribution of the estimators under model misspecification.

This is visually summarized in Figures 3b and 3c. The green boxes represent the empirical interquartile range (Q1 to Q3) of the true sampling distribution. The standard methods (olive and blue boxes) produce interquartile ranges that are noticeably misaligned with this target, either too wide or too narrow depending on the setting. By contrast, the bagging procedures (red boxes for BayesBag and purple boxes for variational bagging) yields Q1–Q3 ranges more closely aligned with the green boxes, demonstrating improved accuracy and robustness in the presence of model misspecification.

A similar phenomenon is observed when the data are generated from double-exponential mixture distributions, as shown in Figures 4 and 5. In these settings as well, the bagging-based procedures effectively mitigate the impact of tail misspecification and produce uncertainty quantification that more faithfully reflects the underlying data-generating process.

Consistent with our theoretical results, the outcomes of variational bagging closely mirror those of Bayesian bagging in these simulations. This indicates that variational bagging can inherit the robustness properties of BayesBag while being substantially more computationally efficient, thereby offering a practical and effective tool for robust uncertainty quantification in mixture models and beyond.

4.4. Sparse linear regression model

In this section, we consider a sparse linear regression model with a spike-and-slab prior, a popular and effective prior for variable selection. The model is

$$\begin{aligned} Y | X, \Gamma, \beta, \sigma^2 &\sim N(X\Gamma\beta, \sigma^2 I_n), \\ \beta &\sim N(0, \sigma_\beta^2 I_q), \\ \sigma^2 &\sim \text{InvGamma}(A, B), \\ \gamma_i &\stackrel{i.i.d.}{\sim} \text{Bernoulli}(p), \quad \Gamma = \text{diag}(\gamma_1, \dots, \gamma_q). \end{aligned}$$

We simulate data from several models and fit each using this sparse linear regression specification. Unless otherwise noted, the general settings are:

sample size $n = 1000$, hyperparameters $A = B = 0.1$, $\sigma_\beta^2 = 10$, and $p = 0.5$. For standard Bayesian inference, we obtain 2000 MCMC samples via Gibbs sampling. For variational inference, we use Algorithm 1 of Ormerod et al. [22]. For bagging, we take $B = 50$ bootstrap samples and set the bootstrap sample size to $M = \hat{M}_\infty$ in (8)).

To compare methods, we focus on estimation accuracy for the regression coefficients. Let β_{OLS} be the ordinary least squares estimator (when defined), and let $\hat{\beta}$ denote a point estimate from each method (e.g., posterior mean). We compute the relative squared error (RSE)

$$RSE = \frac{\|\hat{\beta} - \beta_{OLS}\|^2}{\|\beta_{OLS}\|^2},$$

and compare RSEs across standard Bayes, BayesBag, variational Bayes (VB), and Variational Bagging.

We consider the following four scenarios:

- Scenario 1 (S1): 10 covariates, $n = 1000$.
- Scenario 2 (S2): 10 covariates, $n = 2000$.
- Scenario 3 (S3): 20 covariates, $n = 1000$.
- Scenario 4 (S4): 20 covariates, $n = 2000$.

Tables 1 and 2 report RSEs under Gaussian and Student- t errors, respectively.

Scenario	Bayes	BayesBag	VB	VB Bagging
S1	3.16e-4	2.17e-4	1.24e-4	9.67e-5
S2	3.04e-4	2.15e-4	4.09e-4	3.63e-4
S3	8.82e-4	3.33e-4	1.65e-4	1.39e-4
S4	2.50e-5	2.37e-5	3.42e-5	2.15e-5

TABLE 1

RSEs of four methods for sparse linear regression with Gaussian errors.

Scenario	Bayes	BayesBag	VB	VB bagging
S1	6.55e-4	3.34e-4	5.16e-4	4.14e-4
S2	9.73e-4	5.10e-4	1.97e-3	1.54e-3
S3	1.35e-3	7.80e-4	2.81e-4	2.41e-4
S4	6.68e-4	7.48e-4	1.15e-3	7.57e-4

TABLE 2

RSEs of four methods for sparse linear regression with Student- t errors.

Across all scenarios, the bagging regimes (BayesBag and Variational Bagging) generally achieve lower RSEs than their non-bagged counterparts,

highlighting the robustness benefits of bagging. Under Gaussian errors, VB bagging often produces the most accurate results, slightly improving upon both Bayes and standard VB. Under Student- t errors, where the likelihood is misspecified, the gains from bagging are even more pronounced: Bayes-Bag consistently improves on Bayes, and Variational Bagging improves on standard VB in three of the four scenarios.

We also observe that VB can occasionally outperform standard Bayesian inference. This is likely due to Monte Carlo error and limited MCMC exploration in the Gibbs sampler, especially in higher-dimensional settings. Overall, these results indicate that bagging can enhance the accuracy and stability of both Bayesian and variational approaches for sparse linear regression, while retaining the computational advantages of VB.

4.5. Deep learning model for prediction

We implement variational bagging for prediction via

$$q^{\text{bvB}}(X_{\text{new}} \mid X_{1:n}) = \frac{1}{B} \sum_{b=1}^B q(X_{\text{new}} \mid X_{1:M}^{*(b)}),$$

where $X_{1:M}^{*(b)}$ denotes the b -th bootstrap sample. From Section 4.1, the resulting asymptotic covariance for the bagged variational posterior has the form

$$\frac{1}{2} [(\tilde{V}_{\text{VB}}^0)^{-1} + (V_{\text{VB}}^0)^{-1} D_{\text{VB}}^0 (V_{\text{VB}}^0)^{-1}],$$

which can be viewed as a compromise between the model-based covariance $(\tilde{V}_{\text{VB}}^0)^{-1}$ and the sandwich covariance $(V_{\text{VB}}^0)^{-1} D_{\text{VB}}^0 (V_{\text{VB}}^0)^{-1}$. Even with $M = 2n$, we still recover half of the off-diagonal (sandwich) contribution through the term $\frac{1}{2}(V_{\text{VB}}^0)^{-1} D_{\text{VB}}^0 (V_{\text{VB}}^0)^{-1}$.

We apply fully-connected deep neural networks (DNNs) for regression and assess predictive uncertainty. Motivated by the fact that such DNNs can achieve near-optimal nonparametric rates [21, 18], we investigate three regression settings:

- simple linear regression,
- nonlinear regression,
- multivariate linear regression.

For each setting, data is generated with error terms (ϵ) drawn either from a Student- t distribution or a Laplace distribution, but the working regression model assumes Gaussian errors, thereby introducing deliberate model misspecification.

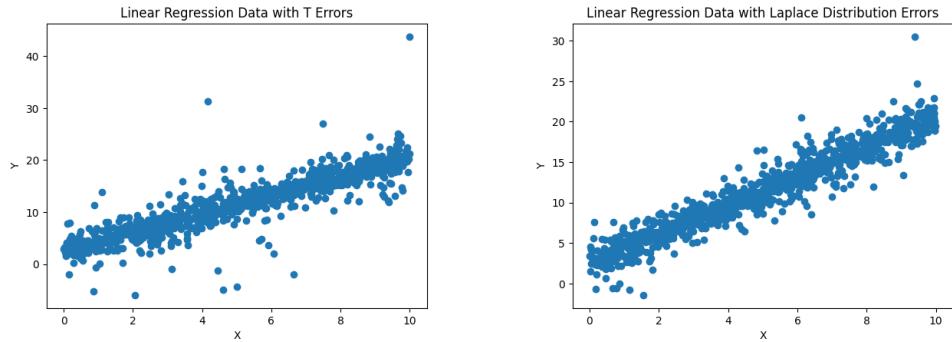


Fig 6: Regression data from a linear model with t and Laplace errors.

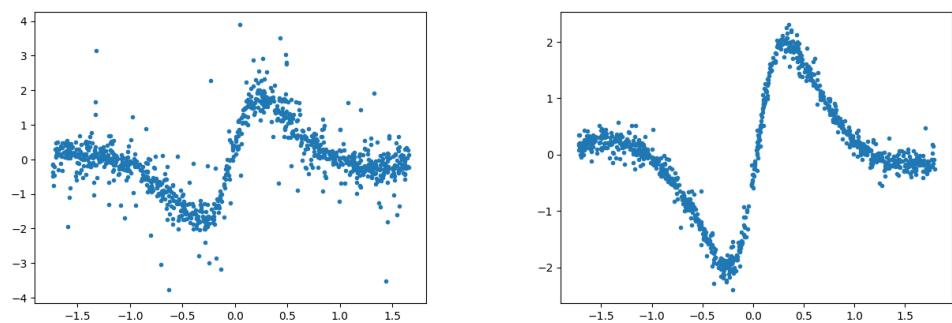


Fig 7: Regression data from a nonlinear model with t and Laplace errors.

- **Simple linear regression.** Data are generated as

$$Y = 2.5 + 1.8X + \epsilon,$$

as illustrated in Figure 6. For this setting, we use a DNN with architecture [1, 10, 1] (one input, one hidden layer with 10 units, and one output), which is sufficient to capture the underlying linear relationship.

- **Nonlinear regression.** Data follow

$$Y = \frac{\sin(X)}{1 + X^2} + \epsilon,$$

as shown in Figure 7. To capture the nonlinear structure, we employ a deeper and wider DNN with architecture [1, 32, 64, 1], which provides increased capacity to approximate the nonlinear regression function.

- **Multivariate linear regression.** We generate responses via

$$Y = X\beta + \epsilon,$$

where $\beta_i \stackrel{i.i.d.}{\sim} N(0, 1)$ and there are 7 predictors (so $X \in \mathbb{R}^{n \times 7}$). For this model, we use a larger DNN with architecture [7, 128, 64, 32, 1] to handle the increased input dimension and capture interactions among predictors.

Our primary evaluation metric is the empirical coverage of 95% predictive intervals on the original data. Specifically, we examine whether the nominal 95% predictive intervals contain approximately 95% of the observed responses, which serves as a diagnostic for the calibration of predictive uncertainty. In the simulation study:

- the number of bootstrap replicates is $B = 10$;
- the bootstrap sample size is set to $M = 2n$.

Model misspecification (Gaussian working errors vs. heavy-tailed true errors) tends to produce predictive intervals that are too narrow, leading to under-coverage of the nominal 95% predictive intervals. This effect is visible in the variational Bayes (VB) results. However, as shown in Table 3, variational bagging (VB bagging) substantially corrects this under-coverage and yields predictive intervals whose empirical coverage is much closer to the nominal level.

Across all regression settings and both error distributions, VB alone yields slightly under-covered predictive intervals (around 93% rather than 95%),

Data setting	VB	VB bagging
linear reg + t errors	93.8	94.6
linear reg + Laplace errors	93.8	94.6
nonlinear reg + t errors	93.2	95.3
nonlinear reg + Laplace errors	92.9	95.7
multivariate reg + t errors	93.0	94.8
multivariate reg + Laplace errors	93.5	95.0

TABLE 3

Empirical coverage (%) of nominal 95% predictive intervals for DNN-based regression under VB and variational bagging.

reflecting the combination of model misspecification and variational under-dispersion. In contrast, variational bagging systematically improves coverage, bringing it close to the nominal 95% level in all cases. These results highlight the robustness of variationla bagging in restoring calibrated predictive uncertainty for deep neural network models, even when the error distribution is misspecified and the underlying regression function is nonlinear or high-dimensional.

4.6. Bagged variational autoencoder (BVAE)

In this simulation study, we consider the application of variational bagging to a variational autoencoder (VAE) model [15], which we refer to as the bagged variational autoencoder. Recall in a deep generative model where a D -dimensional random variable X is modeled by $X = \mathbf{f}(Z) + \epsilon$, where Z is some latent variable of dimension d , ϵ represents the error and \mathbf{f} is the generator typically parametrized by a deep neural network. VAE is one of the primary likelihood based training methods for estimating the deep generative model, where $P(Z | X)$ is parameterized by another deep neural network, the so-called encoder network. For implementation, we adopt the same algorithm described in [6]. To model data as noisy realizations from a distribution on a 1-dimensional manifold, we let $X = \mathbf{f}(Z) + \epsilon$, where $D = 2$, $\sigma_* = 0$ is the true variance of the residual, and Z is a univariate random variable following a standard normal distribution $N(0, 1)$. We examine three different functions for the true generators $\mathbf{f}_* = (f_{*1}, f_{*2})$ as:

- Case 1 (Figure 8a). $f_{*1}(z) = 6(z - 0.5)$, $f_{*2}(z) = 0.5(z - 2)z(z + 2)$.
- Case 2 (Figure 8b). $f_{*1}(z) = \cos(2\pi z)$, $f_{*2}(z) = \sin(2\pi z)$.

- Case 3 (Figure 8c).

$$\begin{cases} f_{*1}(z) = 2 \cos(2\pi z) + 1, & f_{*2}(z) = 2 \sin(2\pi z) + 0.4, & \text{if } z > 0.5 \\ f_{*1}(z) = 2 \cos(2\pi z) - 1, & f_{*2}(z) = 2 \sin(2\pi z) - 0.4, & \text{otherwise} \end{cases}$$

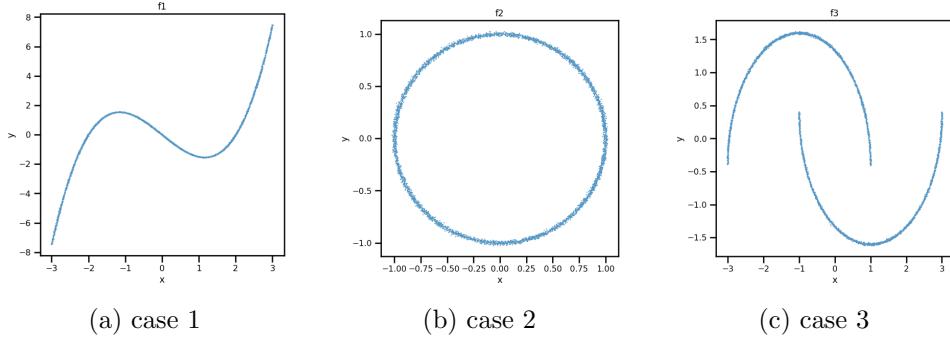


Fig 8: Simulated data concentrated on different 1-dimensional manifolds. These data will be analyzed with VAEs with and without bagging.

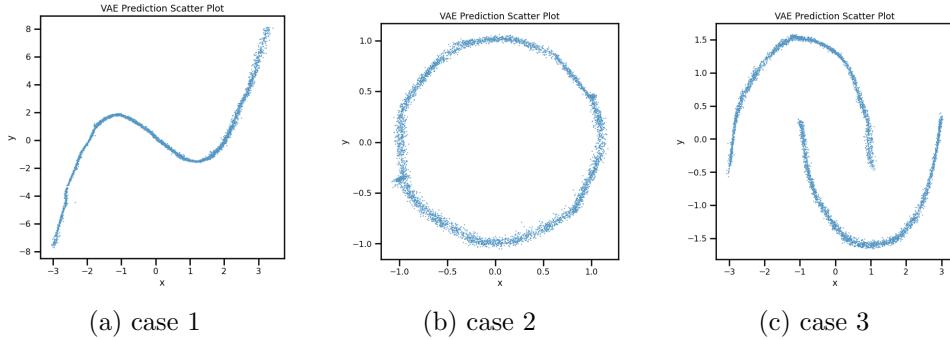


Fig 9: Results from applying a usual (non-bagged) VAE to the one-dimensional manifold data from the previous Figure.

Subsequently, we investigate two additional distributions: one on the Swiss roll in Figure 11 and another featuring a uniform distribution on a sphere in Figure 12. Both distributions are supported on 2-dimensional manifolds within the ambient space \mathbb{R}^3 . The Swiss roll distribution represents the distribution of $\mathbf{f}(Z)$, where Z follows a uniform distribution on $(0, 1)^2$, and the

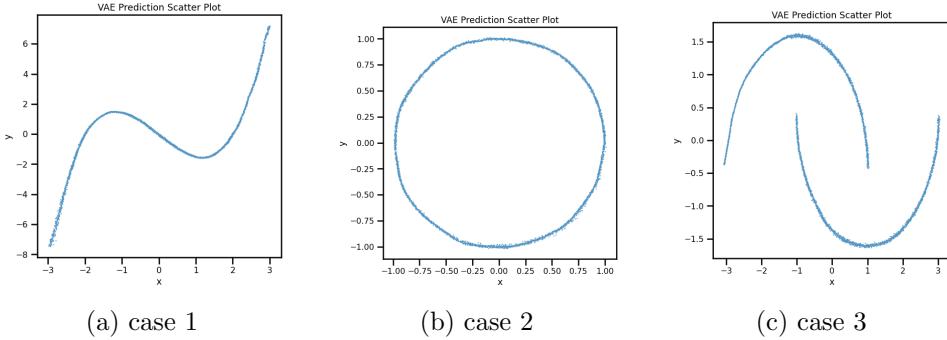


Fig 10: Results of bagged VAE on the simulated one-dimensional manifold data; averaging over bootstrap dataset produces smoother and more accurate estimates

true generator $\mathbf{f}_* = (f_{*1}, f_{*2}, f_{*3}) : (0, 1)^2 \rightarrow \mathbb{R}^3$ is defined as:

$$\begin{aligned} t_1 &= 1.5\pi(1 + 2z_1), t_2 = 21z_2, \\ f_{*1}(z_1, z_2) &= t_1 \cos(t_1), \quad f_{*2}(z_1, z_2) = t_2, \quad f_{*3}(z_1, z_2) = t_1 \sin(t_1). \end{aligned}$$

We first use the standard VAE with 10 epochs for training. We generate the figures by sampling from the posterior predictive distribution. We also apply our proposed bagged VAE approach. In all experiments, we set the validation and test sample sizes to 1000, while the training sample is set to 5000.

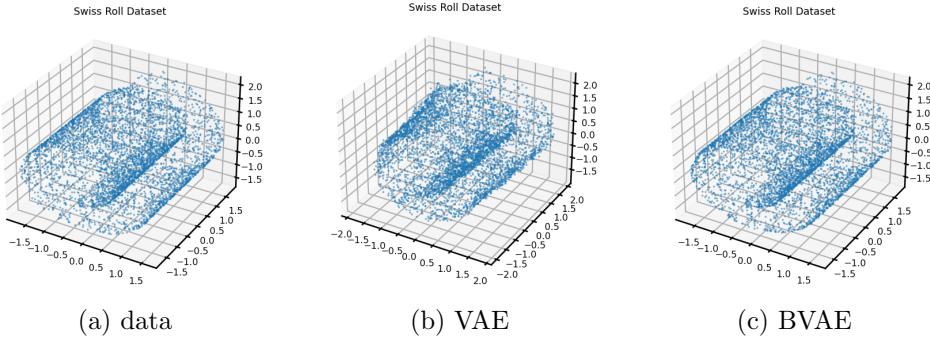


Fig 11: Simulated data on a 2d swiss roll embedded in 3d are shown in panel (a) with the results of VAE in panel (b) and those of BVAE in panel (c).

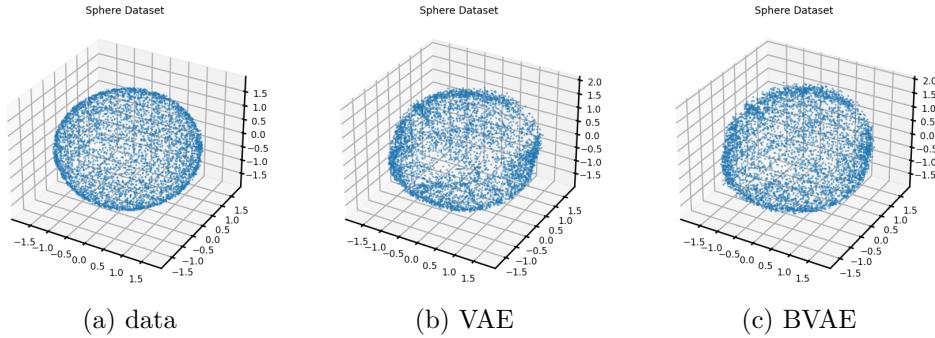


Fig 12: Sphere data

In our study, we have observed distinct differences in the performance of the bagged Variational Autoencoder (BVAE) compared to the standard VAE across various data sets, especially those representing 1-dimensional manifolds. The findings can be summarized as follows:

- General Observation for 1-D Manifolds in Figure 9 and 10: Across all the data sets representing 1-dimensional manifolds, the bagged VAE consistently demonstrates superior capability in reconstructing the generator. This suggests that the bagging technique enhances the VAE's ability to capture and replicate the underlying structure of the data more effectively than the standard VAE.
- Swiss Roll Dataset in Figure 11: In the specific case of the Swiss roll data set, the bagged VAE again shows a notable improvement over the standard VAE. It demonstrates its ability to reconstruct the manifold structure with greater fidelity, indicating its enhanced modeling capability in this more complex data set.
- Sphere Dataset in Figure 12: The sphere example presents a slightly different scenario. Here, the difference in performance between bagged VAE and standard VAE is not as pronounced. While the bagged VAE still performs marginally better in reconstructing the sphere, the improvement is more subtle than in the other cases.

These observations suggest that the bagging technique, when applied to a VAE, generally enhances the model's ability to accurately reconstruct various types of data structures, particularly in the context of 1-dimensional manifolds. The enhanced performance of the bagged VAE in most scenarios indicates its potential as a more robust and effective approach in the field of generative modeling. However, the slight improvement seen in the sphere

data set also highlights that the effectiveness of bagging can vary depending on the specific characteristics of the data set being modeled.

5. Real data analysis: Bagged VAE for the MNIST and Omniglot datasets

We first consider the well-known MNIST dataset [19]. MNIST consists of grayscale images of handwritten digits of size 28×28 , with a training set of 60,000 images and a test set of 10,000 images. We randomly sample 10,000 images from the training set to form a validation set.

In our VAE architecture for images, we utilize convolutional Flipout layers [31] as variational counterparts of standard convolutional layers, and transposed convolutional Flipout layers as variational analogs of transposed convolutional layers in the decoder.

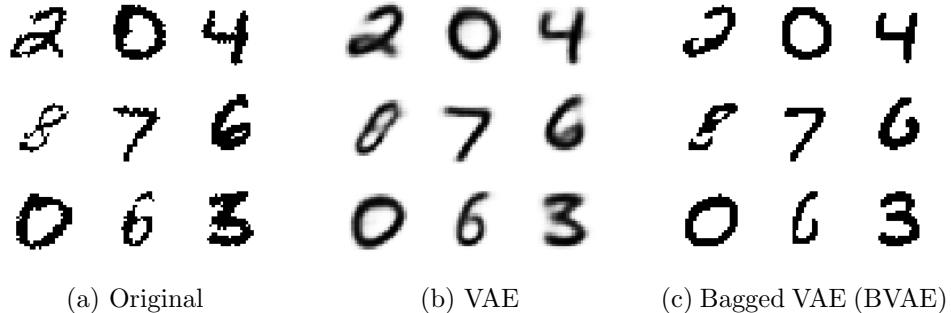


Fig 13: MNIST: original images, standard VAE reconstructions, and bagged VAE (BVAE) reconstructions.

We compare the standard VAE and the bagged VAE (BVAE) using the same protocol as in our earlier VAE simulation study. As illustrated in Figure 13, the bagged VAE reconstructions appear visually cleaner and more representative of the underlying digit structure. In particular, the BVAE tends to denoise the images more effectively and emphasize the main strokes and shapes of each digit, while the standard VAE produces blurrier reconstructions.

Next, we consider the Omniglot dataset, which consists of handwritten character images from 50 different alphabets, each of size 28×28 . The dataset contains 24,345 training samples and 8,070 test samples. As with MNIST, we split the training set into 20,000 images for training and 4,345 images for

validation. We apply the same VAE and BVAE architectures and training procedures as in the MNIST experiment.

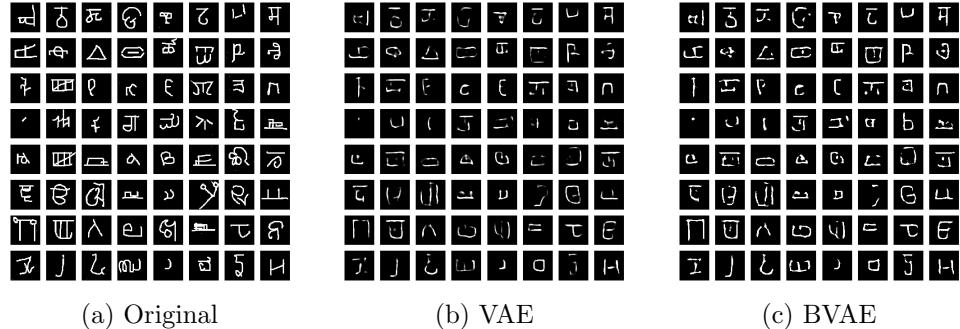


Fig 14: Omniglot: original images, standard VAE reconstructions, and bagged VAE (BVAE) reconstructions.

Figure 14 shows that the qualitative behavior observed on MNIST carries over to Omniglot. The bagged VAE reconstructions more faithfully preserve fine-grained structural details of the characters and better capture their distinctive features compared to the standard VAE. Overall, these experiments suggest that variational bagging can substantially enhance the quality and robustness of VAE-based generative modeling for image data, resulting in reconstructions that are both cleaner and more representative of the original inputs.

Appendix A: Proofs

A.1. Proof of Theorem 3.1

Define $\mathbb{P}_n = n^{-1} \sum_{i=1}^n \delta_{X_i}$ and $\mathbb{G}_n = n^{1/2}(\mathbb{P}_n - P_0)$. Also, define their bootstrapped versions

$$\begin{aligned}\mathbb{P}_n^* &= M^{-1} \sum_{i=1}^n K_i \delta_{X_i} \\ \mathbb{G}_n^* &= M^{1/2}(\mathbb{P}_n^* - \mathbb{P}_n)\end{aligned}$$

where $(K_1, \dots, K_n) \sim \text{Multi}(M, 1/n)$. Under the conditions 1 to 4, we can apply Lemma 19.31 of [28] with $\ell_\theta(x) = \log p_{\text{VB}}(x|\theta)$ to have

$$\mathbb{G}_n \left(\sqrt{n}(\ell_{\theta_0+h_n/\sqrt{n}} - \ell_{\theta_0}) - h_n^\top \dot{\ell}_{\theta_0} \right) \xrightarrow{P_0} 0 \quad (11)$$

for any sequence h_n bounded in P_0 -probability. Moreover, by Theorem 23.7 of [28], conditionally given X_1, X_2, \dots , the bootstrap empirical process \mathbb{G}_n^* and the empirical process \mathbb{G}_n converge weakly to the same limiting random process, thus we also have

$$\mathbb{G}_n^* \left(\sqrt{n}(\ell_{\theta_0+h_n/\sqrt{n}} - \ell_{\theta_0}) - h_n^\top \dot{\ell}_{\theta_0} \right) \xrightarrow{P_0} 0.$$

This implies that

$$(nM)^{1/2} \mathbb{P}_n^*(\ell_{\theta_0+h_n/\sqrt{n}} - \ell_{\theta_0}) - \sqrt{n} h_n^\top (\mathbb{P}_n^* - \mathbb{P}_n) \dot{\ell}_{\theta_0} \\ - n \mathbb{P}_n(\ell_{\theta_0+h_n/\sqrt{n}} - \ell_{\theta_0}) \xrightarrow{P_0} 0,$$

where, from (11), the third term of the above display satisfies

$$n \mathbb{P}_n(\ell_{\theta_0+h_n/\sqrt{n}} - \ell_{\theta_0}) - \sqrt{n} h_n^\top (\mathbb{P}_n - P_0) \dot{\ell}_{\theta_0} \\ - n P_0(\ell_{\theta_0+h_n/\sqrt{n}} - \ell_{\theta_0}) \xrightarrow{P_0} 0$$

Moreover, by the second-order Taylor expansion in the condition 3, we have

$$-P_0(\ell_{\theta_0+h_n/\sqrt{n}} - \ell_{\theta_0}) - \frac{1}{2n} h_n^\top V_{VB}^0 h_n = o(1).$$

Therefore, letting

$$\Delta_n = n^{1/2} (V_{VB}^0)^{-1} (\mathbb{P}_n - P_0) \dot{\ell}_{\theta_0} \\ \Delta_n^* = n^{1/2} (V_{VB}^0)^{-1} (\mathbb{P}_n^* - \mathbb{P}_n) \dot{\ell}_{\theta_0},$$

we have for every compact $K \subset \Theta$,

$$\sup_{h \in K} \left| M \mathbb{P}_n^*(\ell_{\theta_0+h_n/\sqrt{n}} - \ell_{\theta_0}) - h^\top (c V_{VB}^0) (\Delta_n + \Delta_n^*) - \frac{1}{2} h^\top (c V_{VB}^0) h \right| \xrightarrow{P_0} 0. \quad (12)$$

Then we can apply the second result of Theorem 3 of [29] with $X_{1:M}^*$ in place of $X_{1:n}$, \mathbb{P}_n^* in place of \mathbb{P}_n , $c V_{VB}^0$ in place of V_{θ_0} , and $\Delta_n + \Delta_n^*$ in place of Δ_{n,θ_0} , to obtain that

$$\left\| \mathcal{L}(\sqrt{n}(\vartheta' - \theta_0) | X_{1:M}^*) - N(\Delta_n + \Delta_n^*, (c \tilde{V}_{VB}^0)^{-1}) \right\|_{TV} \xrightarrow{P_0} 0 \quad (13)$$

for $\vartheta' \sim q(\theta | X_{1:M}^*)$, where $\mathcal{L}(\sqrt{n}(\vartheta' - \theta_0) | X_{1:M}^*)$ denotes the conditional raw of $\sqrt{n}(\vartheta' - \theta_0)$ given $X_{1:M}^*$. Hence, the characteristic function of $\sqrt{n}(\vartheta' - \theta_0) - \Delta_n | X_{1:n}$ for $\vartheta' \sim q^{\text{bvB}}(\theta | X_{1:n})$ evaluated at $t \in \mathbb{R}^d$ can be written as

$$E \left[\exp \left\{ i(\Delta_n^*)^\top t - \frac{1}{2c} t^\top (\tilde{V}_{VB}^0)^{-1} t \right\} | X_{1:n} \right] + \epsilon_n(t) \quad (14)$$

for some function $\epsilon_n(t)$ such that $\limsup_{n \rightarrow \infty} \sup_t \epsilon_n(t) = 0$. We can expand (14) as

$$\begin{aligned} & E \left[\exp \left\{ in^{1/2} \mathbb{P}_n^* \dot{\ell}_{\theta_0}(V_{\text{VB}}^0)^{-1} t \right\} | X_{1:n} \right] \exp \left\{ -in^{1/2} \mathbb{P}_n \dot{\ell}_{\theta_0}(V_{\text{VB}}^0)^{-1} t \right\} \\ & \times \exp \left\{ -t^\top (\tilde{V}_{\text{VB}}^0)^{-1} t / 2c \right\} + \epsilon_n(t) \end{aligned}$$

The first line can be written as:

$$\begin{aligned} & E \left[\exp \left\{ in^{1/2} M^{-1} \sum_{j=1}^n K_j \dot{\ell}_{\theta_0}(X_j)^\top (V_{\text{VB}}^0)^{-1} t \right\} | X_{1:n} \right] \\ & \times \exp \left\{ -in^{1/2} \mathbb{P}_n \dot{\ell}_{\theta_0}(V_{\text{VB}}^0)^{-1} t \right\} \\ & = \left[\frac{1}{n} \sum_{j=1}^n \exp \left\{ \frac{in^{1/2} \dot{\ell}_{\theta_0}(X_j)^\top (V_{\text{VB}}^0)^{-1} t}{M} \right\} \right]^M \exp \left\{ -in^{1/2} \mathbb{P}_n \dot{\ell}_{\theta_0}(V_{\text{VB}}^0)^{-1} t \right\} \\ & = \left[\frac{1}{n} \sum_{j=1}^n \exp \left\{ \frac{in^{1/2} \delta \dot{\ell}_{\theta_0}(X_j)^\top (V_{\text{VB}}^0)^{-1} t}{M} \right\} \right]^M \end{aligned}$$

where we let $\delta \dot{\ell}_{\theta_0}(X_j) = \dot{\ell}_{\theta_0}(X_j) - \mathbb{P}_n \dot{\ell}_{\theta_0}$. By the second-order Taylor expansion, the last display can be further written as

$$\left[1 - \frac{n \mathbb{P}_n [(\delta \dot{\ell}_{\theta_0})(\delta \dot{\ell}_{\theta_0})^\top] (V_{\text{VB}}^0)^{-1} t}{2M^2} + R_n \right]^M,$$

where $R_n = O_P(1/M^3)$ denotes the remainder term. Then since

$$\mathbb{P}_n [(\delta \dot{\ell}_{\theta_0})(\delta \dot{\ell}_{\theta_0})^\top] \rightarrow D_{\text{VB}}^0$$

almost surely and $M/n \rightarrow c$ by assumption, the last display converges to

$$\exp \left\{ -\frac{1}{2c} t^\top (\tilde{V}_{\text{VB}}^0)^{-1} t - \frac{1}{2c} t^\top (V_{\text{VB}}^0)^{-1} D_{\text{VB}}^0 (V_{\text{VB}}^0)^{-1} t \right\} \quad (15)$$

This implies the desired

$$\sqrt{n}(\vartheta^\dagger - \theta_0) - \Delta_n \mid X_{1:n} \xrightarrow{d} N(0, (\tilde{V}_{\text{VB}}^0)^{-1}/c + (V_{\text{VB}}^0)^{-1} D_{\text{VB}}^0 (V_{\text{VB}}^0)^{-1}/c)$$

by Levy's continuity theorem.

A.2. Proof of Theorem 3.3

Since $\hat{\theta}_{\text{mle}} \rightarrow \theta_0$ and $\Delta_n \rightarrow 0$ in P_0 -probability as well as $\widehat{\Sigma} \rightarrow \Sigma_0$ in P_0 -probability by assumption, by Theorem 3.2, we have

$$\begin{aligned} P_{N(0, \Sigma_0)}(C(r_{n,1-\alpha})) &= Q^{\text{bvB}}(\vartheta^\dagger \in C(r_{n,1-\alpha})) + o_{P_0}(1) \\ &= 1 - \alpha + o_{P_0}(1), \end{aligned}$$

where we denote by $P_{N(0, \Sigma)}$ the probability measure under the normal distribution $N(0, \Sigma)$. Therefore by the continuous mapping theorem, $r_{n,1-\alpha}^2 \rightarrow \chi_{d,1-\alpha}^2$ in P_0 -probability, where $\chi_{d,1-\alpha}^2$ denotes the $1-\alpha$ quantile of the $\chi^2(d)$ distribution and d denotes the dimension of the parameter θ . We use this fact to get, letting $S_0 = (V^0)^{-1}D^0(V^0)^{-1}$ for notational simplicity,

$$\begin{aligned} P_0(\theta_0 \in C(r_{n,1-\alpha})) &= P_0(\sqrt{n}(\hat{\theta}_{\text{mle}} - \theta_0) \in \{u : u^\top \Sigma_0^{-1} u \leq r_{n,1-\alpha}^2\}) + o(1) \\ &= P_0(\sqrt{n}(\hat{\theta}_{\text{mle}} - \theta_0) \in \{u : u^\top \Sigma_0^{-1} u \leq \chi_{d,1-\alpha}^2\}) + o(1) \\ &= P_{U \sim N(0, S_0)}(U^\top \Sigma_0^{-1} U \leq \chi_{d,1-\alpha}^2) + o(1) \end{aligned}$$

where the last equality holds due to

$$\sqrt{n}(\hat{\theta}_{\text{mle}} - \theta_0) \xrightarrow{d} N(0, S_0).$$

Therefore, since $\Sigma_0 - S_0 = (\tilde{V}^0)^{-1}$ is non-negative definite, so is $S_0^{-1} - \Sigma_0^{-1}$, we have

$$\begin{aligned} P_{U \sim N(0, S_0)}(U^\top \Sigma_0^{-1} U \leq \chi_{d,1-\alpha}^2) \\ \geq P_{U \sim N(0, S_0)}(U^\top S_0^{-1} U \leq \chi_{d,1-\alpha}^2) = 1 - \alpha, \end{aligned}$$

which completes the proof.

A.3. Proof of Theorem 3.5

Before giving the proof, we introduce additional notations.

Let $p_{\theta, K_{1:n}}(X_{1:n}) = \prod_{i=1}^n p(X_i | \theta)^{K_i}$ and $p_{0, K_{1:n}}(X_{1:n}) = \prod_{i=1}^n p(X_i | \theta)^{K_i}$ be denote the density of the bootstrapped sample $X_{1:n}^*$ with the “bootstrap weight” $K_{1:n} = (K_1, \dots, K_n) \sim \text{Multi}(n, 1/n)$, under P_θ and P_0 , respectively. Then the posterior given the bootstrapped sample can be written as $\pi(\theta | X_{1:n}^*) = p_{\theta, K_{1:n}}(X_{1:n})\pi(\theta) / p_{\Pi, K_{1:n}}(X_{1:n})$ with $p_{\Pi, K_{1:n}}(X_{1:n}) = \int \prod_{i=1}^n p(X_i | \theta)^{K_i} \pi(\theta) d\theta$. In what follows, the probability measure P and

the expectation operator E are taken on the randomness of both the sample $X_{1:n}$ and the bootstrap weights $K_{1:n}$.

Let Q^* be the variational posterior obtained with a bootstrapped sample $X_{1:n}^*$. The proof is done if we show that the contraction rate of Q^* is given by $\epsilon_n \sqrt{\log n}$. We start with observing that, by Donsker and Varadhan's variational inequality (e.g., Lemma J.1 of [21]),

$$Q(A) \leq \frac{1}{t} \{ \text{KL}(Q, \Pi(\cdot | X_{1:n}^*)) + e^t \Pi(A | X_{1:n}^*) \} \quad (16)$$

for any distribution Q on Θ , any event $A \subset \Theta$ and any positive number t . We apply the above display to $Q = Q^*$ and $A = A_n = \{\theta \in \Theta : H^2(P_\theta, P_0) \geq \eta_n^2\}$ with $\eta_n^2 = M_n \epsilon_n^2 \log n$.

Next, we define the event

$$G_n = \left\{ \int \frac{p_{\theta, K_{1:n}}(X_{1:n})}{p_{0, K_{1:n}}(X_{1:n})} \pi(\theta) d\theta \geq \exp(-(C_1 + 2)n\epsilon_n^2) \right\}.$$

We have

$$E \left[\log \frac{p_{\theta, K_{1:n}}(X_{1:n})}{p_{0, K_{1:n}}(X_{1:n})} \right] = n \text{KL}(P_0, P_\theta)$$

and

$$\begin{aligned} \text{Var} \left(\log \frac{p_{\theta, K_{1:n}}(X_{1:n})}{p_{0, K_{1:n}}(X_{1:n})} \right) &\leq \sum_{i=1}^n E \left[\left(K_i \log \frac{p(X_i | \theta)}{p_0(X_i)} \right)^2 \right] \\ &= \sum_{i=1}^n E(K_i^2) \text{KLV}(P_0, P_\theta) \\ &\leq 2n \text{KLV}(P_0, P_\theta), \end{aligned}$$

where the equality follows from the independence of K_i and X_i . Hence, by a standard argument for bounding the probability of G_n under the prior mass condition (e.g., Lemma 8.10 of [8]), we have $P(G_n^c) \leq 2/(n\epsilon_n^2) \rightarrow 0$. Therefore, it suffices to bound the quantity $E(Q^*(A_n)\mathbb{I}(G_n))$, which, from the inequality in (16), is further bounded by

$$\begin{aligned} &E(Q^*(A_n)\mathbb{I}(G_n)) \\ &\leq \frac{1}{t} \{ E[\text{KL}(Q^*, \Pi(\cdot | X_{1:n}^*))\mathbb{I}(G_n)] + e^t E[\Pi(A_n | X_{1:n}^*)\mathbb{I}(G_n)] \}. \end{aligned} \quad (17)$$

For the first term in (17), we have

$$\begin{aligned}
E[\text{KL}(Q^*, \Pi(\cdot|X_{1:n}^*))\mathbb{I}(G_n)] &\leq E[\text{KL}(Q^*, \Pi(\cdot|X_{1:n}^*))] \\
&= E\left[\int \left\{\log \frac{q^*(\theta)}{\pi(\theta)} + \log \frac{p_{0,K_{1:n}}(X_{1:n})}{p_{\theta,K_{1:n}}(X_{1:n})}\right\} dQ^*(\theta) + \log \frac{p_{\Pi,K_{1:n}}(X_{1:n})}{p_{0,K_{1:n}}(X_{1:n})}\right] \\
&\leq E\left[\int \left\{\log \frac{q(\theta)}{\pi(\theta)} + \log \frac{p_{0,K_{1:n}}(X_{1:n})}{p_{\theta,K_{1:n}}(X_{1:n})}\right\} dQ(\theta) + \log \frac{p_{\Pi,K_{1:n}}(X_{1:n})}{p_{0,K_{1:n}}(X_{1:n})}\right] \\
&= \text{KL}(Q, \Pi) + n \int \text{KL}(P_\theta, P_0) dQ(\theta) + E\left[\log \frac{p_{\Pi,K_{1:n}}(X_{1:n})}{p_{0,K_{1:n}}(X_{1:n})}\right]
\end{aligned}$$

for any $Q \in \mathcal{Q}$, where we use the optimization optimality of Q^* in the third line and use the assumption $\sum_{i=1}^n K_i = n$ in the last line. Moreover, by Jensen's inequality together with that $\sum_{i=1}^n K_i = n$, we have

$$\begin{aligned}
E\left[\log \frac{p_{\Pi,K_{1:n}}(X_{1:n})}{p_{0,K_{1:n}}(X_{1:n})}\right] &= E\left[\log \frac{p_{\Pi,K_{1:n}}(X_{1:n})}{\prod_{i=1}^n p_0(X_i)} + \log \frac{\prod_{i=1}^n p_0(X_i)}{p_{0,K_{1:n}}(X_{1:n})}\right] \\
&\leq E\left[\log \frac{\prod_{i=1}^n p_0(X_i)}{p_{0,K_{1:n}}(X_{1:n})}\right] \\
&= E\left[\sum_{i=1}^n (1 - K_i) \log p_0(X_i)\right] = 0.
\end{aligned}$$

Therefore, by our variational family assumption,

$$E[\text{KL}(Q^*, \Pi(\cdot|X_{1:n}^*))] \leq C'_1 n \epsilon_n^2 \quad (18)$$

for some constant $C'_1 > 0$.

Now we focus on the second term in (17). We can bound it as

$$\Pi(A_n|X_{1:n}^*)\mathbb{I}(G_n) \leq \Pi(A_n \cap \Theta_n|X_{1:n}^*)\mathbb{I}(G_n) + \Pi(\Theta_n^c|X_{1:n}^*)\mathbb{I}(G_n).$$

The expectation of the second term in the above display satisfies

$$E[\Pi(\Theta_n^c|X_{1:n}^*)\mathbb{I}(G_n)] \leq e^{(2+C_1)n\epsilon_n^2} \Pi(\Theta_n^c) \leq e^{-2n\epsilon_n^2} \rightarrow 0$$

by the second sieve assumption. On the other hand, for the first term, we need the following technical result on the bootstrap weight. Since $K_i \sim \text{Binom}(n, 1/n)$, by Markov's inequality

$$\begin{aligned}
P(K_i > 2 \log n) &\leq e^{-2 \log n} E(\exp(K_i)) \\
&= e^{-2 \log n} (1 - n^{-1} + n^{-1}e)^n \\
&\leq e^{-2 \log n} e^{e-1} = e^{e-1}/n^2,
\end{aligned}$$

where the last inequality follows from the inequality $1 + x \leq e^x$ for any $x \in \mathbb{R}$. Thus, by the union bound, we have

$$P\left(\max_{1 \leq i \leq n} K_i > 2 \log n\right) \leq \frac{e^{e-1}}{n} \quad (19)$$

which tends to zero as $n \rightarrow \infty$. Then, following the proof strategy of [11], we define

$$U_\theta(X) = \max \left\{ \log \frac{p(X|\theta)}{p_0(X)}, -\tau \right\}$$

for $\tau > 0$. Then we consider the weighted likelihood ratio empirical process defined as

$$\nu_n(\theta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n [K_i U_\theta(X_i) - E(K_i U_\theta(X_i))]$$

Then by Lemma 12 of [11] together with the fact given in (19) and our complexity assumption, we have

$$P \left(\sup_{\theta \in \Theta_n : H(P_\theta, P_0) \leq \sqrt{2}\epsilon} \nu_n(\theta) \geq \frac{k}{2} \sqrt{n}\epsilon^2 \right) \leq 3 \exp \left(-C'_2 \frac{n\epsilon^2}{1 + 2\log n} \right) \quad (20)$$

for any $\epsilon > \epsilon_n$ for some constants $C'_2 > 0$ and $k \in (1/2, 1)$. By Lemma 4 of [33], $E(\nu_n(\theta)) \leq -(1 - \delta_0)H^2(P_\theta, P_0)$ with $\delta = 2 \exp(-\tau/2)/(1 - \exp(-\tau/2))^2$. Let $\Theta_n(\epsilon; \sqrt{2}\epsilon) = \{\theta \in \Theta_n : \epsilon \leq H(P_\theta, P_0) \leq \sqrt{2}\epsilon\}$. Then we have

$$\begin{aligned} B(\epsilon; \sqrt{2}\epsilon) &= \left\{ \sup_{\theta \in \Theta_n(\epsilon; \sqrt{2}\epsilon)} \frac{p_{\theta, K_{1:n}}(X_{1:n})}{p_{0, K_{1:n}}(X_{1:n})} \geq \exp(-n\epsilon^2(1 - \delta_0 - k/2)) \right\} \\ &\subset \left\{ \sup_{\theta \in \Theta_n(\epsilon; \sqrt{2}\epsilon)} \nu_n(\theta) \geq \frac{k}{2} \sqrt{n}\epsilon^2 \right\} \end{aligned}$$

Hence, by (20), the probability of the event $B(\epsilon; \sqrt{2}\epsilon)$ is bounded above by $3 \exp \left(-C'_1 \frac{n\epsilon^2}{1 + 2\log n} \right)$. Let J be the smallest integer such that $2^J \eta_n^2 \geq 4$.

Then by using a peeling technique, we have

$$\begin{aligned} & P \left(\sup_{\theta \in \Theta_n : H(P_\theta, P_0) \geq \eta_n} \frac{p_{\theta, K_{1:n}}(X_{1:n})}{p_{0, K_{1:n}}(X_{1:n})} \geq \exp(-n\eta_n^2(1 - \delta_0 - k/2)) \right) \\ & \leq \sum_{j=0}^J P \left(B(\sqrt{2^j}\eta_n; \sqrt{2^{j+1}}\eta_n) \right) \leq \sum_{j=0}^J 3 \exp \left(-C'_2 \frac{n2^j\eta_n^2}{1 + 2\log n} \right) \\ & \leq 4 \exp \left(-C'_2 \frac{n\eta_n^2}{1 + 2\log n} \right) \leq 4 \exp \left(-\frac{1}{3} C'_2 M_n n \epsilon_n^2 \right). \end{aligned}$$

Therefore, $E[\Pi(A_n \cap \Theta_n | X_{1:n}^*) \mathbb{I}(G_n)] \leq \exp(-C'_3 n \eta_n^2)$ for some constant $C'_3 > 0$. Combining these derivations together, we have that by taking $t = C'_4 n \eta_n^2$ with C'_4 being a positive constant less than C'_3 , the term (17) converges to zero, which proves that the variational posterior Q^* with a single bootstrapped sample contracts to P_0 at a rate η_n .

References

- [1] Pierre Alquier and James Ridgway. Concentration of tempered posteriors and of their variational approximations. *The Annals of Statistics*, 48(3):1475–1497, 2020.
- [2] Matthew J. Beal and Zoubin Ghahramani. Variational Bayesian learning of directed graphical models with hidden variables. *Bayesian Analysis*, 1(4):793 – 831, 2006. . URL <https://doi.org/10.1214/06-BA126>.
- [3] David M. Blei and Michael I. Jordan. Variational inference for Dirichlet process mixtures. *Bayesian Analysis*, 1(1):121 – 143, 2006. . URL <https://doi.org/10.1214/06-BA104>.
- [4] Leo Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, 1996.
- [5] Peter Bühlmann. Discussion of big bayes stories and bayesbag. *Statistical Science*, 29(1):91–94, 2014. ISSN 08834237, 21688745. URL <http://www.jstor.org/stable/43288454>.
- [6] Minwoo Chae, Dongha Kim, Yongdai Kim, and Lizhen Lin. A likelihood approach to nonparametric estimation of a singular distribution using deep generative models. *J. Mach. Learn. Res.*, 24:1–42, 2023.
- [7] Alp Kucukelbir David M. Blei and Jon D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017. . URL <https://doi.org/10.1080/01621459.2017.1285773>.

- [8] Subhashis Ghosal and Aad Van der Vaart. *Fundamentals of nonparametric Bayesian inference*, volume 44. Cambridge University Press, 2017.
- [9] Subhashis Ghosal, Jayanta K Ghosh, and Aad van der Vaart. Convergence rates of posterior distributions. *The Annals of Statistics*, 28(2):500–531, 2000.
- [10] Ryan Giordano, Tamara Broderick, and Michael I. Jordan. Covariances, robustness, and variational bayes. *J. Mach. Learn. Res.*, 19:51:1–51:49, 2017. URL <https://api.semanticscholar.org/CorpusID:53238793>.
- [11] Wei Han and Yun Yang. Statistical inference in mean-field variational bayes. *arXiv preprint arXiv:1911.01525*, 2019.
- [12] Jonathan H Huggins and Jeffrey W Miller. Robust inference and model criticism using bagged posteriors. *arXiv preprint arXiv:1912.07104*, 2019.
- [13] Michael I. Jordan, Zoubin Ghahramani, Tommi S. Jaakkola, and Lawrence K. Saul. *An introduction to variational methods for graphical models*, page 105–161. MIT Press, Cambridge, MA, USA, 1999. ISBN 0262600323.
- [14] Anya Katsevich and Philippe Rigollet. On the approximation accuracy of gaussian variational inference. *arXiv preprint arXiv:2306.00052*, 2023. URL <https://arxiv.org/abs/2306.00052>.
- [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2014.
- [16] Diederik P Kingma and Max Welling. Auto-Encoding Variational Bayes. *arXiv e-prints*, art. arXiv:1312.6114, December 2013. .
- [17] Bas JK Kleijn and Aad W van der Vaart. The bernstein-von-mises theorem under misspecification. *Electronic Journal of Statistics*, 6:354–381, 2012.
- [18] Michael Kohler and Sophie Langer. On the rate of convergence of fully connected deep neural network regression estimates. *The Annals of Statistics*, 49(4):2231–2249, 2021.
- [19] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [20] Dennis Nieman, Botond Szabo, and Harry van Zanten. Uncertainty quantification for sparse spectral variational approximations in Gaussian process regression. *Electronic Journal of Statistics*, 17(2):2250 – 2288, 2023. . URL <https://doi.org/10.1214/23-EJS2155>.

- [21] Ilsang Ohn and Lizhen Lin. Adaptive variational bayes: Optimality, computation and applications. *The Annals of Statistics*, 52(1):335–363, 2024.
- [22] John T Ormerod, Chong You, and Samuel Müller. A variational bayes approach to variable selection. 2017.
- [23] Debdipati Pati, Anirban Bhattacharya, and Yun Yang. On statistical optimality of variational Bayes. pages 1579–1588, 2018.
- [24] Kolyan Ray and Botond Szabó and. Variational bayes for high-dimensional linear regression with sparse priors. *Journal of the American Statistical Association*, 117(539):1270–1281, 2022. . URL <https://doi.org/10.1080/01621459.2020.1847121>.
- [25] Xiaotong Shen and Larry Wasserman. Rates of convergence of posterior distributions. *Annals of Statistics*, pages 687–714, 2001.
- [26] Yee Teh, David Newman, and Max Welling. A collapsed variational bayesian inference algorithm for latent dirichlet allocation. In B. Schölkopf, J. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems*, volume 19. MIT Press, 2006. URL https://proceedings.neurips.cc/paper_files/paper/2006/file/532b7cbe070a3579f424988a040752f2-Paper.pdf.
- [27] Sattar Vakili, Jonathan Scarlett, Da shan Shiu, and Alberto Bernacchia. Improved convergence rates for sparse approximation methods in kernel-based learning. In *International Conference on Machine Learning*, 2022. URL <https://api.semanticscholar.org/CorpusID:246652434>.
- [28] Aad W Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- [29] Yixin Wang and David Blei. Variational bayes under model misspecification. *Advances in Neural Information Processing Systems*, 32, 2019.
- [30] Yixin Wang and David M Blei. Frequentist consistency of variational bayes. *Journal of the American Statistical Association*, 114(527):1147–1161, 2019.
- [31] Yeming Wen, Paul Vicol, Jimmy Ba, Dustin Tran, and Roger Grosse. Flipout: Efficient pseudo-independent weight perturbations on mini-batches. *arXiv preprint arXiv:1803.04386*, 2018.
- [32] T Westling and TH McCormick. Beyond prediction: A framework for inference with variational approximations in mixture models. *Journal of Computational and Graphical Statistics*, 28(4):778–789, 2019.
- [33] Wing Hung Wong and Xiaotong Shen. Probability inequalities for likelihood ratios and convergence rates of sieve mles. *The Annals of Statistics*, pages 339–362, 1995.

- [34] Yun Yang, Debdeep Pati, and Anirban Bhattacharya. α -variational inference with statistical guarantees. *The Annals of Statistics*, 48(2):886–905, 2020.
- [35] Fengshuo Zhang and Chao Gao. Convergence rates of variational posterior distributions. *The Annals of Statistics*, 48(4):2180 – 2207, 2020.