

Location–Scale Calibration for Generalized Posterior

Shu Tamano^{1,2,*} and Yui Tomo²

¹Department of Multidisciplinary Sciences, Graduate School of Arts and Sciences,
The University of Tokyo, 3-8-1 Komaba, Meguro-Ku, Tokyo 153-8902, Japan

²Department of Epidemiology, National Institute of Infectious Diseases, Japan
Institute for Health Security, 1-23-1 Toyama, Shinjuku-Ku, Tokyo 162-0052, Japan

*Email: tamano-shu212@g.ecc.u-tokyo.ac.jp

Abstract

General Bayesian updating replaces the likelihood with a loss scaled by a learning rate, but posterior uncertainty can depend sharply on that scale. We propose a simple post-processing that aligns generalized posterior draws with their asymptotic target, yielding uncertainty quantification that is invariant to the learning rate. We prove total-variation convergence for generalized posteriors with an effective sample size, allowing sample-size-dependent priors, non-i.i.d. observations, and convex penalties under model misspecification. Within this framework, we justify and extend the open-faced sandwich adjustment (Shaby, 2014), provide general theoretical guarantees for its use within generalized Bayes, and extend it from covariance rescaling to a location–scale calibration whose draws converge in total variation to the target for any learning rate. In our empirical illustration, calibrated draws maintain stable coverage, interval width, and bias over orders of magnitude in the learning rate and closely track frequentist benchmarks, whereas uncalibrated posteriors vary markedly.

Keywords: Bayesian inference; Bernstein–von Mises; Generalized Bayes; Learning rate; Open-faced sandwich; Penalized estimating equation; Sandwich variance.

1 Introduction

Bayesian inference provides a coherent probabilistic framework that combines prior information with likelihood-based learning and delivers uncertainty quantification. However, when the assumed likelihood is misspecified, or when inference is based on a non-likelihood objective such as quasi- or composite likelihoods, estimating equations, or other loss-based objectives, posterior uncertainty can be miscalibrated (Kleijn and van der Vaart, 2012; Syring and Martin, 2019; Miller, 2021). General Bayesian inference replaces the likelihood with a loss-based construction (Bissiri et al., 2016), thereby avoiding the need to model the entire data distribution explicitly. In this formulation, a single scaling parameter $\eta \in \mathbb{R}_{>0}$, often called the learning rate or temperature, multiplies the loss and directly controls posterior dispersion. When $\eta = 1$ and the loss function is taken as the negative log-likelihood, the generalized posterior coincides with the usual likelihood-based Bayes posterior. With a fixed prior, decreasing η attenuates the contribution of the data so that, in the limit, the update reverts to the prior, whereas increasing η amplifies the data contribution and, under standard regularity conditions, the posterior concentrates around minimizers of the population loss. Thus, the learning rate governs the trade-off between prior information and loss-based evidence.

Therefore, the choice of learning rate has been widely discussed. Bootstrap-based calibration methods choose η to achieve frequentist targets (Lyddon et al., 2019; Syring and Martin, 2019); in a related approach, Matsubara et al. (2024) compute the loss minimizer on bootstrap resamples, obtain a closed-form estimate of the learning rate from the bootstrap spread, and then run Markov chain

Monte Carlo. SafeBayes provides an alternative data-driven choice that adapts the learning rate for robustness (Grünwald and van Ommen, 2017). Information-matching rules select the learning rate so that the generalized posterior with learning rate η aligns, under criteria of divergence or information, with the generalized posterior obtained by setting $\eta = 1$ for the same loss (Holmes and Walker, 2017). Along a related calibration-to-Bayes line, Altamirano et al. (2023) choose the learning rate by minimizing the Kullback–Leibler divergence between the generalized posterior with learning rate η and the generalized posterior obtained by setting $\eta = 1$ for the same loss, both computed on an initial data window, and then keep the resulting learning rate fixed for the full analysis. These approaches can be computationally intensive and typically require a frequentist point estimator, so the resulting learning rate reflects the plug-in distribution of that estimator rather than the full posterior law. Moreover, several methods explicitly or implicitly calibrate to the generalized posterior obtained by setting $\eta = 1$ for the chosen loss; this is natural when the loss is negative log-likelihood under correct specification, but can be misleading for more general loss-based objectives or under model misspecification. Although McLatchie et al. (2025) suggest that, for prediction, the choice of learning rate matters little in moderate to large samples, this does not resolve the problem of the sensitivity of uncertainty quantification to the learning rate.

In this paper, we revisit the open-faced sandwich adjustment (Shaby, 2014) and place it in a general asymptotic framework for generalized Bayes based on loss functions with an effective sample size, allowing sample-size-dependent priors, possibly non-i.i.d. observations, and convex penalties under model misspecification. First, we show that, under sample-size-dependent priors, the generalized posterior admits a normal limit in total variation and we establish a prior–penalty correspondence that identifies the target curvature and variability. Second, building on this justification, we improve the open-faced sandwich adjustment from covariance rescaling to a location–scale calibration that is implementable via plug-in sandwich estimators, requiring only posterior draws and empirical moments (no bootstrap or learning-rate tuning). Third, we prove that the calibrated draws converge to the target distribution for any learning rate, so that the asymptotic uncertainty quantification is invariant to the choice of learning rate. Our empirical illustration for a random-intercept mixed model with Huber loss and Gaussian prior with sample-size-dependent scale shows that the calibrated intervals track frequentist benchmarks across several orders of magnitude in the learning rate, whereas uncalibrated generalized Bayes posteriors exhibit substantial sensitivity.

2 Problem setup

Fix $p \in \mathbb{N}$ and let $\Theta \subset \mathbb{R}^p$ be an open parameter space. Write $\boldsymbol{\theta} \in \Theta$ for the parameter, and let $\|\cdot\|$ denote the Euclidean norm. Let $M_n : \Theta \rightarrow (-\infty, \infty]$ be an empirical criterion with effective scale $s_n \rightarrow \infty$; typical examples include i.i.d. additive losses ($s_n = n$), m -variate U -type losses ($s_n = \binom{n}{m}$ for non-degenerate kernels), and kernel-smoothed losses with bandwidth $h_n > 0$ ($s_n = nh_n$, $h_n \rightarrow 0$ and $nh_n \rightarrow \infty$). Define the rescaled criterion and its first and second derivatives by

$$\bar{M}_n(\boldsymbol{\theta}) := s_n^{-1} M_n(\boldsymbol{\theta}), \quad \mathbf{U}_n(\boldsymbol{\theta}) := \nabla_{\boldsymbol{\theta}} \bar{M}_n(\boldsymbol{\theta}), \quad \mathbf{J}_n(\boldsymbol{\theta}) := \nabla_{\boldsymbol{\theta}}^2 \bar{M}_n(\boldsymbol{\theta}).$$

We write $R(\boldsymbol{\theta}) := \mathbb{E}_{P^*}[\bar{M}_n(\boldsymbol{\theta})]$ for the corresponding population criterion, where P^* denotes the true data-generating distribution.

For a learning rate $\eta > 0$, we adopt the general Bayesian updating framework of Bissiri et al. (2016). Define the composite loss $L_n(\boldsymbol{\theta}) := M_n(\boldsymbol{\theta}) + \lambda_n s_n \rho(\boldsymbol{\theta})$, where $\{\lambda_n\}_{n \geq 1}$ is a deterministic non-negative sequence with $\lambda_n \rightarrow \lambda \in [0, \infty)$ and $\rho : \mathbb{R}^p \rightarrow \mathbb{R}$ is a convex penalty. We introduce a baseline prior $\tilde{\pi}_n(\boldsymbol{\theta}) \propto \exp\{r_n(\boldsymbol{\theta})\}$ that may depend on n but is independent of the data. The generalized posterior is then given by

$$\Pi_n^\eta(d\boldsymbol{\theta}) \propto \exp\{-\eta L_n(\boldsymbol{\theta})\} \tilde{\pi}_n(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$

Equivalently, absorbing the penalty into the prior, we obtain

$$\Pi_n^\eta(d\boldsymbol{\theta}) \propto \exp\{-\eta M_n(\boldsymbol{\theta})\} \pi_n(\boldsymbol{\theta}) d\boldsymbol{\theta}, \quad \log \pi_n(\boldsymbol{\theta}) = -\eta \lambda_n s_n \rho(\boldsymbol{\theta}) + r_n(\boldsymbol{\theta}).$$

This loss-based representation makes explicit that the learning rate η scales the entire loss, in line with Bissiri et al. (2016). Throughout the main text we assume that ρ is C^2 in a neighborhood of the target point introduced below (see Assumption 2.1); the convex non-smooth case is treated in Section D.

Let $\{\mathcal{D}_i\}_{i=1}^n$ be observations from an unknown distribution P^* and $\Psi(\boldsymbol{\theta}) := \mathbb{E}_{P^*}[\psi(\mathcal{D}_1, \boldsymbol{\theta})]$ denote the population estimating function corresponding to M_n , for some measurable ψ . In the M -estimation

case $M_n(\boldsymbol{\theta}) = \sum_{i=1}^n m(\mathcal{D}_i, \boldsymbol{\theta})$, one may take $\psi(\mathcal{D}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} m(\mathcal{D}, \boldsymbol{\theta})$. The penalized population equation is

$$\mathbf{0} = \Psi(\boldsymbol{\theta}) + \lambda \nabla \rho(\boldsymbol{\theta}), \quad (1)$$

with solution $\boldsymbol{\theta}^\lambda$. At $\boldsymbol{\theta}^\lambda$, define

$$\mathbf{J}^* := \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}^\lambda), \quad \mathbf{K}^* := \text{Var}_{P^*} [\psi(\mathcal{D}_1, \boldsymbol{\theta}^\lambda)], \quad \mathbf{H}_\rho(\boldsymbol{\theta}^\lambda) := \nabla^2 \rho(\boldsymbol{\theta}^\lambda).$$

Then set

$$\mathbf{J}_\lambda^* := \mathbf{J}^* + \lambda \mathbf{H}_\rho(\boldsymbol{\theta}^\lambda), \quad \mathbf{V}_{\text{target}}^* := (\mathbf{J}_\lambda^*)^{-1} \mathbf{K}^* (\mathbf{J}_\lambda^*)^{-1}. \quad (2)$$

The matrix $\mathbf{V}_{\text{target}}^*$ is the usual sandwich covariance for penalized M -estimators; see, e.g., van der Vaart (1998); Kosorok (2008).

To state our theoretical results rigorously, we impose the following regularity conditions. All stochastic limits are taken under the true law P^* ; \rightarrow_p and \rightarrow_d denote convergence in probability and in distribution, respectively.

Assumption 2.1 (Penalty). ρ is convex and C^2 on an open neighborhood of $\boldsymbol{\theta}^\lambda$, and $\mathbf{H}_\rho(\boldsymbol{\theta}^\lambda) = \nabla^2 \rho(\boldsymbol{\theta}^\lambda)$ is positive semi-definite.

Assumption 2.2 (General loss). There exists an open neighborhood \mathfrak{N} of $\boldsymbol{\theta}^\lambda$ such that: (i) the penalized population equation (1) has the unique solution $\boldsymbol{\theta}^\lambda \in \mathfrak{N}$; (ii) $R \in C^3(\mathfrak{N})$ and $\sup_{\boldsymbol{\theta} \in \mathfrak{N}} \|\nabla^3 M_n(\boldsymbol{\theta})\| = O_p(1)$; (iii) $\sup_{\boldsymbol{\theta} \in \mathfrak{N}} \|\mathbf{J}_n(\boldsymbol{\theta}) - \mathbf{J}^*\| \rightarrow_p 0$ and \mathbf{J}^* is nonsingular; and (iv) $\sqrt{s_n}\{\mathbf{U}_n(\boldsymbol{\theta}^\lambda) - \Psi(\boldsymbol{\theta}^\lambda)\} \rightarrow_d \mathcal{N}(\mathbf{0}, \mathbf{K}^*)$.

Assumption 2.3 (Prior remainder). Let $\mathcal{U}_n = \{\boldsymbol{\theta} : \|\boldsymbol{\theta} - \boldsymbol{\theta}^\lambda\| \leq c/\sqrt{s_n}\}$ for some fixed constant $c > 0$. The prior remainder r_n is locally Lipschitz on \mathcal{U}_n with Lipschitz constant $L_n = o_p(\sqrt{s_n})$.

For subsequent results we allow the posterior to be centered at arbitrary data-dependent locations that are root- s_n close to a fixed baseline estimator of $\boldsymbol{\theta}^\lambda$.

Assumption 2.4 (Baseline estimator). There exists a measurable sequence of estimators $\check{\boldsymbol{\theta}}_n \in \Theta$ such that (i) $\check{\boldsymbol{\theta}}_n \rightarrow_p \boldsymbol{\theta}^\lambda$; (ii) the penalized estimating equation is solved up to $s_n^{-1/2}$ order, that is, $\mathbf{U}_n(\check{\boldsymbol{\theta}}_n) + \lambda_n \nabla \rho(\check{\boldsymbol{\theta}}_n) = o_p(s_n^{-1/2})$.

Definition 2.5 (Admissible center). Suppose Assumption 2.4 holds. A measurable sequence of data-dependent centers $\tilde{\boldsymbol{\theta}}_n$ is called an admissible center if $\|\tilde{\boldsymbol{\theta}}_n - \check{\boldsymbol{\theta}}_n\| = o_p(s_n^{-1/2})$. In particular, admissibility implies $\tilde{\boldsymbol{\theta}}_n \rightarrow_p \boldsymbol{\theta}^\lambda$.

Remark 2.6 (Examples of admissible centers). Under Assumptions 2.1–2.4, the following centers are admissible: (i) any maximum a posteriori (MAP) estimator, (ii) the posterior mean $\boldsymbol{\theta}_{\text{GB}} := \mathbb{E}_{\Pi_n^\eta}[\boldsymbol{\theta}]$, and (iii) the one-step Newton update from $\boldsymbol{\theta}_{\text{GB}}$ based on the penalized score. Details are given in Section A.

3 Main results

In this section, we present the main theoretical results. First, we give a total-variation limit for generalized Bayes posteriors based on a general empirical criterion with effective scale s_n , allowing for n -dependent priors and non-i.i.d. data at admissible centers. We then establish the proposed location–scale calibration and show that its limiting law is invariant to the learning rate. Proofs of all results in this section are given in Section B.

We denote by $q_n^{\tilde{\boldsymbol{\theta}}_n}$ the density of $\sqrt{s_n}(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}_n)$ under Π_n^η , and by $\mathcal{N}(\cdot \mid \boldsymbol{\mu}, \boldsymbol{\Sigma})$ the density of the p -variate normal distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. Using local asymptotic normality arguments in the spirit of Miller (2021, Theorem 4), adapted to our loss-based criterion with effective scale s_n , to possibly n -dependent priors and to non-i.i.d. data, and evaluated at admissible centers, we obtain the following total-variation limit.

Proposition 3.1 (Total-variation limit with n -dependent priors at an admissible center). *Under Assumptions 2.1–2.4, let $\tilde{\boldsymbol{\theta}}_n$ be any admissible center. Then*

$$\int_{\mathbb{R}^p} \left| q_n^{\tilde{\boldsymbol{\theta}}_n}(\mathbf{x}) - \mathcal{N}\left(\mathbf{x} \mid \mathbf{0}, (\eta \mathbf{J}_\lambda^*)^{-1}\right) \right| d\mathbf{x} \rightarrow 0, \quad n \rightarrow \infty.$$

We next use Proposition 3.1 to construct an affine transformation of posterior draws that aligns their limiting distribution with the target law. Let $\{\boldsymbol{\theta}^{(d)}\}_{d=1}^D$ be posterior draws from Π_n^η with mean $\boldsymbol{\theta}_{\text{GB}}$, and write the working curvature as $\mathbf{H}_0 := \eta \mathbf{J}_\lambda^*$. For any symmetric positive definite matrix \mathbf{A} , we denote by $\mathbf{A}^{1/2}$ its uniquely determined symmetric positive definite square root, that is, the symmetric matrix satisfying $\mathbf{A}^{1/2} \mathbf{A}^{1/2} = \mathbf{A}$. We also define $\mathbf{A}^{-1/2}$ as $(\mathbf{A}^{1/2})^{-1}$. With this convention, define the location-scale calibration map

$$\boldsymbol{\Omega} := (\mathbf{V}_{\text{target}}^*)^{1/2} \mathbf{H}_0^{1/2}, \quad \boldsymbol{\theta}_{\text{calib}}^{(d)} := \tilde{\boldsymbol{\theta}}_n + \boldsymbol{\Omega} (\boldsymbol{\theta}^{(d)} - \boldsymbol{\theta}_{\text{GB}}). \quad (3)$$

Theorem 3.2 (Location-scale calibration). *Under Assumptions 2.1–2.4, let $\tilde{\boldsymbol{\theta}}_n$ be any admissible center. Then, conditionally on the data,*

$$\sqrt{s_n} (\boldsymbol{\theta}_{\text{calib}}^{(d)} - \tilde{\boldsymbol{\theta}}_n) \longrightarrow_d \mathcal{N}(\mathbf{0}, \mathbf{V}_{\text{target}}^*), \quad n \rightarrow \infty.$$

Remark 3.3 (Learning-rate invariance). By construction,

$$\boldsymbol{\Omega} \mathbf{H}_0^{-1} \boldsymbol{\Omega}^\top = \mathbf{V}_{\text{target}}^*.$$

Hence the asymptotic law of $\sqrt{s_n}(\boldsymbol{\theta}_{\text{calib}}^{(d)} - \tilde{\boldsymbol{\theta}}_n)$ is invariant to the learning rate η .

Remark 3.4 (Extension to convex non-smooth penalties). The main text assumes $\rho \in C^2$ near $\boldsymbol{\theta}^\lambda$. For convex non-smooth penalties, the results extend by restricting the analysis to an appropriate active set and then applying the same arguments with all matrices restricted accordingly. Details are provided in Section D.

Remark 3.5 (Scope and exclusions). Our standing assumptions in the main text require convexity and C^2 -smoothness of ρ near $\boldsymbol{\theta}^\lambda$ and a locally Lipschitz prior remainder. These conditions exclude nonconvex or singular shrinkage specifications, including the horseshoe and spike-and-slab.

4 Practical plug-in calibration

We now describe a fully implementable version of the location-scale calibration that requires only posterior draws and empirical moment estimators. The procedure (i) estimates the working curvature from the posterior sample and (ii) plugs in consistent estimators of the target curvature and variability. Throughout this section, let $\{\boldsymbol{\theta}^{(d)}\}_{d=1}^D$ be draws from Π_n^η , and write the Monte Carlo mean $\hat{\boldsymbol{\theta}}_{\text{GB}} := D^{-1} \sum_{d=1}^D \boldsymbol{\theta}^{(d)}$. Define the sample covariance

$$\hat{\boldsymbol{\Sigma}}_{\text{post}} := \frac{1}{D} \sum_{d=1}^D (\boldsymbol{\theta}^{(d)} - \hat{\boldsymbol{\theta}}_{\text{GB}})(\boldsymbol{\theta}^{(d)} - \hat{\boldsymbol{\theta}}_{\text{GB}})^\top.$$

Proofs of all results in this section are given in Section C.

Lemma 4.1 (Working covariance from posterior draws). *Suppose Assumptions 2.1–2.4 hold. Assume further that, conditionally on the observed data,*

$$s_n \left\| \hat{\boldsymbol{\Sigma}}_{\text{post}} - \boldsymbol{\Sigma}_{\text{post},n} \right\|_F \longrightarrow_p 0, \quad \boldsymbol{\Sigma}_{\text{post},n} := \text{Var}_{\Pi_n^\eta}[\boldsymbol{\theta} \mid \{\mathcal{D}_i\}_{i=1}^n],$$

where $\|\cdot\|_F$ denotes the Frobenius norm on $\mathbb{R}^{p \times p}$. Then, conditionally on the data,

$$s_n \hat{\boldsymbol{\Sigma}}_{\text{post}} \longrightarrow_p \mathbf{H}_0^{-1}, \quad n \rightarrow \infty,$$

where $\mathbf{H}_0 := \eta \mathbf{J}_\lambda^*$.

Next, we estimate the target covariance $\mathbf{V}_{\text{target}}^*$ by plugging in a data-dependent center $\bar{\boldsymbol{\theta}}_n$ such that $\bar{\boldsymbol{\theta}}_n \rightarrow_p \boldsymbol{\theta}^\lambda$:

$$\hat{\mathbf{J}}_\lambda := \mathbf{J}_n(\bar{\boldsymbol{\theta}}_n) + \lambda_n \mathbf{H}_\rho(\bar{\boldsymbol{\theta}}_n), \quad \hat{\mathbf{V}}_{\text{target}} := \hat{\mathbf{J}}_\lambda^{-1} \hat{\mathbf{K}} \hat{\mathbf{J}}_\lambda^{-1},$$

where $\hat{\mathbf{K}}$ is any estimator with $\hat{\mathbf{K}} \rightarrow_p \mathbf{K}^*$.

Lemma 4.2 (Plug-in consistency for the target). *Under Assumptions 2.1 and 2.2, if $\bar{\boldsymbol{\theta}}_n \rightarrow_p \boldsymbol{\theta}^\lambda$, then*

$$\hat{\mathbf{J}}_\lambda \rightarrow_p \mathbf{J}_\lambda^*, \quad \hat{\mathbf{V}}_{\text{target}} \rightarrow_p \mathbf{V}_{\text{target}}^*.$$

Using these estimators, we define the empirical location–scale calibration operator and the corresponding calibrated draws:

$$\hat{H}_0^{-1} := s_n \hat{\Sigma}_{\text{post}}, \quad \hat{\Omega} := \hat{V}_{\text{target}}^{1/2} \hat{H}_0^{1/2}, \quad \hat{\theta}_{\text{calib}}^{(d)} := \bar{\theta}_n + \hat{\Omega}(\theta^{(d)} - \hat{\theta}_{\text{GB}}). \quad (4)$$

Proposition 4.3 (Estimated location–scale calibration). *Suppose the conditions of Lemmas 4.1 and 4.2 hold. Then, conditionally on the data,*

$$\sqrt{s_n}(\hat{\theta}_{\text{calib}}^{(d)} - \bar{\theta}_n) \rightarrow_d \mathcal{N}(\mathbf{0}, \mathbf{V}_{\text{target}}^*), \quad n \rightarrow \infty.$$

Therefore, the limiting law is invariant to the learning rate η .

5 Empirical illustration

5.1 Random-intercept linear mixed model with a Huber loss

We consider a random-intercept linear mixed model with a Huber loss and Gaussian prior with n -dependent scale within the framework of generalized Bayes inference.

Let groups $i = 1, \dots, G$ have sizes n_i and total $n = \sum_{i=1}^G n_i$. Observations are $(y_{ij}, \mathbf{x}_{ij}) \in \mathbb{R} \times \mathbb{R}^p$ with $y_{ij} = \mathbf{x}_{ij}^\top \beta + b_i + \varepsilon_{ij}$, $b_i \sim \mathcal{N}(0, \tau^2)$, $\varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2)$, independent across i . Stack $\mathbf{y}_i = (y_{i1}, \dots, y_{in_i})^\top$ and $\mathbf{X}_i \in \mathbb{R}^{n_i \times p}$ with rows \mathbf{x}_{ij}^\top . Take the working marginal covariance $\Sigma_i = \tau^2 \mathbf{1}_{n_i} \mathbf{1}_{n_i}^\top + \sigma^2 \mathbf{I}_{n_i}$ and its symmetric square root $\mathbf{L}_i \mathbf{L}_i^\top = \Sigma_i$, and define whitened objects

$$\tilde{\mathbf{r}}_i(\beta) := \mathbf{L}_i^{-1}(\mathbf{y}_i - \mathbf{X}_i \beta), \quad \tilde{\mathbf{X}}_i := \mathbf{L}_i^{-1} \mathbf{X}_i.$$

With Huber loss $\rho_c(u) = 2^{-1}u^2 \mathbf{1}\{|u| \leq c\} + \{c|u| - 2^{-1}c^2\} \mathbf{1}\{|u| > c\}$, define

$$M_n(\beta) = \sum_{i=1}^G \sum_{j=1}^{n_i} \rho_c(\tilde{\mathbf{r}}_{ij}(\beta)), \quad s_n = n.$$

Let $\psi_c = \rho'_c$ and $\mathbf{W}_i(\beta) = \text{diag}(\mathbf{1}\{|\tilde{\mathbf{r}}_{ij}(\beta)| \leq c\})$. Then

$$\mathbf{U}_n(\beta) = -\frac{1}{n} \sum_{i=1}^G \tilde{\mathbf{X}}_i^\top \psi_c(\tilde{\mathbf{r}}_i(\beta)), \quad \mathbf{J}_n(\beta) = \frac{1}{n} \sum_{i=1}^G \tilde{\mathbf{X}}_i^\top \mathbf{W}_i(\beta) \tilde{\mathbf{X}}_i.$$

Let $\Psi(\beta) = \mathbb{E}[\mathbf{U}_n(\beta)]$. Let $\boldsymbol{\mu} \in \mathbb{R}^p$ be a fixed ridge center and $\mathbf{Q} \in \mathbb{R}^{p \times p}$ be a given symmetric positive definite penalty matrix. For the ridge penalty $\rho(\beta) = 2^{-1}(\beta - \boldsymbol{\mu})^\top \mathbf{Q}(\beta - \boldsymbol{\mu})$ and $\lambda \in [0, \infty)$, define β^λ by $\mathbf{0} = \Psi(\beta) + \lambda \nabla \rho(\beta)$. Set, at β^λ , $\mathbf{J}^* = \nabla_\beta \Psi(\beta^\lambda)$, $\mathbf{K}^* = \lim_{n \rightarrow \infty} n \text{Var}[\mathbf{U}_n(\beta^\lambda)]$, $\mathbf{J}_\lambda^* = \mathbf{J}^* + \lambda \mathbf{Q}$, and the target sandwich $\mathbf{V}_{\text{target}}^* = (\mathbf{J}_\lambda^*)^{-1} \mathbf{K}^* (\mathbf{J}_\lambda^*)^{-1}$.

Take a Gaussian prior with n -dependent scale $\beta \sim \mathcal{N}(\boldsymbol{\mu}, \{(\lambda s_n) \mathbf{Q}\}^{-1})$. The posterior is $\Pi_n^\eta(d\beta) \propto \exp\{-\eta M_n(\beta)\} \pi_n(\beta) d\beta$ and the local precision per s_n is $\mathbf{H}_0 = \eta \mathbf{J}_\lambda^*$. With an admissible center $\tilde{\beta}_n$, estimate $\hat{\mathbf{J}}_\lambda = \mathbf{J}_n(\tilde{\beta}_n) + \lambda_n \mathbf{Q}$, $\hat{\mathbf{K}} = n^{-1} \sum_{i=1}^G \hat{\mathbf{U}}_i \hat{\mathbf{U}}_i^\top$, $\hat{\mathbf{U}}_i = -\tilde{\mathbf{X}}_i^\top \psi_c(\tilde{\mathbf{r}}_i(\tilde{\beta}_n))$, and $\hat{\mathbf{V}}_{\text{target}} = \hat{\mathbf{J}}_\lambda^{-1} \hat{\mathbf{K}} \hat{\mathbf{J}}_\lambda^{-1}$. From posterior draws $\{\beta^{(d)}\}_{d=1}^D$ with mean $\hat{\beta}_{\text{GB}}$, set

$$\hat{\Sigma}_{\text{post}} = \frac{1}{D} \sum_{d=1}^D (\beta^{(d)} - \hat{\beta}_{\text{GB}})(\beta^{(d)} - \hat{\beta}_{\text{GB}})^\top, \quad \hat{H}_0^{-1} = n \hat{\Sigma}_{\text{post}}, \quad \hat{\Omega} = \hat{V}_{\text{target}}^{1/2} \hat{H}_0^{1/2}.$$

Then, calibrated draws are $\beta_{\text{calib}}^{(d)} = \tilde{\beta}_n + \hat{\Omega}(\beta^{(d)} - \hat{\beta}_{\text{GB}})$. By Proposition 4.3, $\sqrt{s_n}(\beta_{\text{calib}}^{(d)} - \tilde{\beta}_n) \rightarrow_d \mathcal{N}(\mathbf{0}, \mathbf{V}_{\text{target}}^*)$. Details of the augmentation and Markov chain Monte Carlo are deferred to Section E.

5.2 Experiment

We illustrate the finite-sample behavior of the plug-in location–scale calibration in the random-intercept linear mixed model with Huber loss and Gaussian prior with n -dependent scale described in Section 5.1. We fix $G = 100$, $n_i = 5$ (so $n = 500$), and set $p = 1$. We generate covariates from $\mathcal{N}(0, 1)$ and contaminate Gaussian errors to induce model misspecification. Since no closed

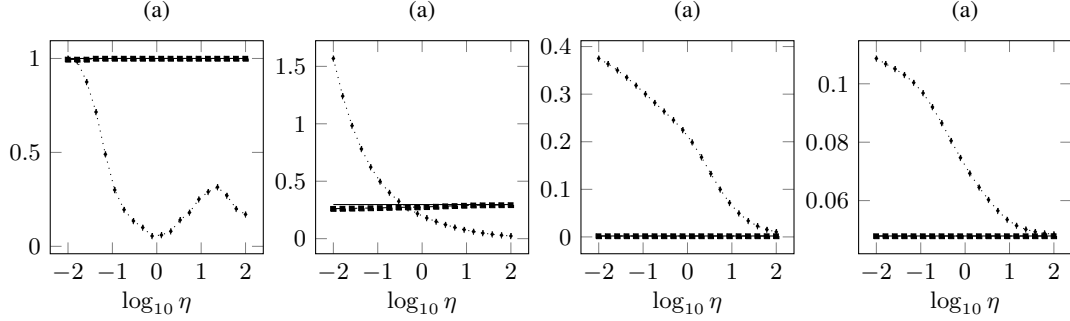


Figure 1: Evaluation metrics for β^λ as the learning rate varies, comparing three procedures. Panels: (a) coverage probability at the nominal 95% level; (b) mean interval width; (c) bias of the point estimator; (d) standard deviation of the bias. Solid line: frequentist benchmark (confidence intervals). Dashed line with square markers: location–scale calibrated posterior (credible intervals). Dotted line with diamond markers: uncalibrated posterior (credible intervals).

form is available for β^λ in this example, we approximate it numerically by computing the penalized estimating equation estimator on large simulated datasets with $G = 5,000$ and averaging over 1,000 replications; see, e.g., Oh and Patton (2013). For each learning rate η on a logarithmic grid over $[0.01, 100]$ we compare three procedures: (i) a frequentist Huber M -estimator with ridge penalty and sandwich-based Wald intervals; (ii) the generalized Bayes posterior based on M_n and the corresponding n -dependent Gaussian prior; and (iii) the location–scale calibration applied to the same posterior draws, using the MAP estimate as an admissible center. For each η we use 200 Monte Carlo replications and record, for all three methods, the empirical coverage, the mean interval width, the mean bias, and the standard deviation of that bias. Further details of the data-generating mechanism, the numerical approximation of β^λ , and the Markov chain Monte Carlo settings are provided in Section E. The Python code for reproducing the experiments is available at <https://github.com/shutech2001/ls-calib-gp>.

Figure 1 shows that the proposed location–scale calibration rendered inference for β^λ essentially invariant to the learning rate η . For the calibrated posterior, coverage probabilities were very close to the frequentist benchmark, mean interval widths are stable, and both the bias and the standard deviation of the bias remain essentially unchanged as η varies. Across the $[0.01, 100.0]$ range of η , the calibrated posterior yielded curves that were very close to the frequentist benchmark. By contrast, the uncalibrated posterior exhibits pronounced sensitivity to η . In this example, the uncalibrated coverage varied substantially with the learning rate. As η increases, the point estimator tracks the loss-based target more closely and its bias decreases, but the associated credible intervals become progressively narrower and eventually exhibit marked undercoverage. This pattern highlights a trade-off between bias and interval width and suggests that procedures which tune the learning rate by optimizing coverage at a single nominal level may be sensitive to local features of this trade-off.

Acknowledgements

Shu Tamano was supported by JSPS KAKENHI Grant Numbers 25K24203.

References

- Altamirano, M., Briol, F.-X., and Knoblauch, J. (2023). Robust and scalable Bayesian online changepoint detection. In *Proceedings of the 40th International Conference on Machine Learning*, pages 642–663.
- Bissiri, P., Holmes, C., and Walker, S. (2016). A general framework for updating belief distributions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(5):1103–1130.

- Clarke, F. H. (1990). *Optimization and Nonsmooth Analysis*. SIAM.
- Grünwald, P. and van Ommen, T. (2017). Inconsistency of bayesian inference for misspecified linear models, and a proposal for repairing it. *Bayesian Analysis*, 12(4):1069–1103.
- Holmes, C. and Walker, S. (2017). Assigning a value to a power likelihood in a general bayesian model. *Biometrika*, 104(2):497–503.
- Kleijn, B. J. K. and van der Vaart, A. W. (2012). The Bernstein-Von-Mises theorem under misspecification. *Electronic Journal of Statistics*, 6:354–381.
- Kosorok, M. R. (2008). *Introduction to empirical processes and semiparametric inference*. Springer.
- Lyddon, S., Holmes, C., and Walker, S. (2019). General Bayesian updating and the loss-likelihood bootstrap. *Biometrika*, 106(2):465–478.
- Matsubara, T., Knoblauch, J., Briol, F.-X., and Oates, C. J. (2024). Generalized Bayesian inference for discrete intractable likelihood. *Journal of the American Statistical Association*, 119(547):2345–2355.
- McLatchie, Y., Fong, E., Frazier, D. T., and Knoblauch, J. (2025). Predictive performance of power posteriors. *Biometrika*, 112(3):asaf034.
- Miller, J. W. (2021). Asymptotic normality, concentration, and coverage of generalized posteriors. *Journal of Machine Learning Research*, 22(168):1–53.
- Oh, D. H. and Patton, A. J. (2013). Simulated method of moments estimation for copula-based multivariate models. *Journal of the American Statistical Association*, 108(502):689–700.
- Shaby, B. (2014). The open-faced sandwich adjustment for MCMC using estimating functions. *Journal of Computational and Graphical Statistics*, 23(3):853–876.
- Syring, N. and Martin, R. (2019). Calibrating general posterior credible regions. *Biometrika*, 106(2):479–486.
- van der Vaart, A. W. (1998). *Asymptotic Statistics*. Cambridge University Press.

A Examples of admissible centers

A.1 MAP centering

We first show that a maximum a posteriori estimator is an admissible center. Write the generalized log-posterior, up to an additive constant, as

$$\ell_n(\boldsymbol{\theta}) := -\eta M_n(\boldsymbol{\theta}) - \eta \lambda_n s_n \rho(\boldsymbol{\theta}) + r_n(\boldsymbol{\theta}) = -\eta s_n \{\bar{M}_n(\boldsymbol{\theta}) + \lambda_n \rho(\boldsymbol{\theta})\} + r_n(\boldsymbol{\theta}).$$

Let $\hat{\boldsymbol{\theta}}_n^{\text{MAP}}$ be a measurable maximizer of ℓ_n over the neighborhood \mathcal{U}_n ; by standard M -estimation arguments, the unique maximizer of the population criterion $R(\boldsymbol{\theta}) + \lambda \rho(\boldsymbol{\theta})$ at $\boldsymbol{\theta}^\lambda$ and Assumption 2.2 imply that such a maximizer exists and that $\hat{\boldsymbol{\theta}}_n^{\text{MAP}} \rightarrow_p \boldsymbol{\theta}^\lambda$. On the event $\{\hat{\boldsymbol{\theta}}_n^{\text{MAP}} \in \mathcal{U}_n\}$, the first-order condition yields

$$\mathbf{0} = \nabla \ell_n(\hat{\boldsymbol{\theta}}_n^{\text{MAP}}) = -\eta s_n \{U_n(\hat{\boldsymbol{\theta}}_n^{\text{MAP}}) + \lambda_n \nabla \rho(\hat{\boldsymbol{\theta}}_n^{\text{MAP}})\} + \nabla r_n(\hat{\boldsymbol{\theta}}_n^{\text{MAP}}),$$

so that

$$U_n(\hat{\boldsymbol{\theta}}_n^{\text{MAP}}) + \lambda_n \nabla \rho(\hat{\boldsymbol{\theta}}_n^{\text{MAP}}) = \eta^{-1} s_n^{-1} \nabla r_n(\hat{\boldsymbol{\theta}}_n^{\text{MAP}}).$$

By Assumption 2.3,

$$\|U_n(\hat{\boldsymbol{\theta}}_n^{\text{MAP}}) + \lambda_n \nabla \rho(\hat{\boldsymbol{\theta}}_n^{\text{MAP}})\| \leq \eta^{-1} s_n^{-1} L_n = o_p(s_n^{-1/2}).$$

Assumption 2.4 (ii) gives

$$U_n(\check{\boldsymbol{\theta}}_n) + \lambda_n \nabla \rho(\check{\boldsymbol{\theta}}_n) = o_p(s_n^{-1/2}).$$

Subtracting these two relations and applying the mean-value theorem to U_n and $\nabla \rho$ along the segment between $\hat{\boldsymbol{\theta}}_n^{\text{MAP}}$ and $\check{\boldsymbol{\theta}}_n$ yields

$$\{J_n(\bar{\boldsymbol{\theta}}_n) + \lambda_n H_\rho(\check{\boldsymbol{\theta}}_n)\}(\hat{\boldsymbol{\theta}}_n^{\text{MAP}} - \check{\boldsymbol{\theta}}_n) = o_p(s_n^{-1/2}), \quad (5)$$

for some random intermediate points $\bar{\boldsymbol{\theta}}_n, \check{\boldsymbol{\theta}}_n$ on the line segment joining $\hat{\boldsymbol{\theta}}_n^{\text{MAP}}$ and $\check{\boldsymbol{\theta}}_n$. By Assumption 2.2 (iii), Assumption 2.1, and the consistency of $\hat{\boldsymbol{\theta}}_n^{\text{MAP}}$ and $\check{\boldsymbol{\theta}}_n$, we have

$$J_n(\bar{\boldsymbol{\theta}}_n) + \lambda_n H_\rho(\check{\boldsymbol{\theta}}_n) \rightarrow_p J_\lambda^*,$$

and J_λ^* is nonsingular. Hence the smallest eigenvalue of $J_n(\bar{\boldsymbol{\theta}}_n) + \lambda_n H_\rho(\check{\boldsymbol{\theta}}_n)$ is bounded away from zero in probability and its inverse is $O_p(1)$. Multiplying (5) by this inverse gives

$$\|\hat{\boldsymbol{\theta}}_n^{\text{MAP}} - \check{\boldsymbol{\theta}}_n\| = o_p(s_n^{-1/2}).$$

A.2 Generalized posterior mean centering

We next consider the generalized Bayes posterior mean $\boldsymbol{\theta}_{\text{GB}} := \mathbb{E}_{\Pi_n^\eta}[\boldsymbol{\theta}]$. To exploit Proposition 3.1 we center at the baseline estimator $\check{\boldsymbol{\theta}}_n$, which is itself an admissible center by Definition 2.5. Let

$$\mathbf{Z}_n := \sqrt{s_n}(\boldsymbol{\theta} - \check{\boldsymbol{\theta}}_n)$$

under Π_n^η , and let $q_n^{\check{\boldsymbol{\theta}}_n}$ denote the density of \mathbf{Z}_n . By Proposition 3.1 applied with $\check{\boldsymbol{\theta}}_n = \check{\boldsymbol{\theta}}_n$,

$$\int_{\mathbb{R}^p} |q_n^{\check{\boldsymbol{\theta}}_n}(\mathbf{x}) - \mathcal{N}(\mathbf{x} \mid \mathbf{0}, \mathbf{H}_0^{-1})| d\mathbf{x} \rightarrow 0, \quad \mathbf{H}_0 := \eta J_\lambda^*.$$

Define

$$\boldsymbol{\delta}_n := \sqrt{s_n}(\boldsymbol{\theta}_{\text{GB}} - \check{\boldsymbol{\theta}}_n) = \int_{\mathbb{R}^p} \mathbf{x} q_n^{\check{\boldsymbol{\theta}}_n}(\mathbf{x}) d\mathbf{x} = \mathbb{E}_{\Pi_n^\eta}[\mathbf{Z}_n].$$

The total-variation convergence together with the fact that the second moments of \mathbf{Z}_n converge to those of $\mathcal{N}(\mathbf{0}, \mathbf{H}_0^{-1})$ (again by Proposition 3.1) implies $\boldsymbol{\delta}_n \rightarrow_p \mathbf{0}$. Thus

$$\|\boldsymbol{\theta}_{\text{GB}} - \check{\boldsymbol{\theta}}_n\| = \frac{\|\boldsymbol{\delta}_n\|}{\sqrt{s_n}} = o_p(s_n^{-1/2}). \quad (6)$$

A.3 One-step Newton centering

Finally, we consider a one-step Newton update from the posterior mean based on the penalized estimating equation. Define the empirical penalized score

$$\mathbf{F}_n(\boldsymbol{\theta}) := \mathbf{U}_n(\boldsymbol{\theta}) + \lambda_n \nabla \rho(\boldsymbol{\theta}), \quad \mathbf{A}_n(\boldsymbol{\theta}) := \mathbf{J}_n(\boldsymbol{\theta}) + \lambda_n \mathbf{H}_\rho(\boldsymbol{\theta}).$$

The one-step Newton estimator from $\boldsymbol{\theta}_{\text{GB}}$ is

$$\tilde{\boldsymbol{\theta}}_n^{(1)} := \boldsymbol{\theta}_{\text{GB}} - \mathbf{A}_n(\boldsymbol{\theta}_{\text{GB}})^{-1} \mathbf{F}_n(\boldsymbol{\theta}_{\text{GB}}). \quad (7)$$

By (6) and Assumption 2.4, $\boldsymbol{\theta}_{\text{GB}} \rightarrow_p \boldsymbol{\theta}^\lambda$ and $\tilde{\boldsymbol{\theta}}_n \rightarrow_p \boldsymbol{\theta}^\lambda$, so both sequences eventually lie in \mathcal{U}_n with probability tending to one. On this event, a Taylor expansion of \mathbf{F}_n about $\boldsymbol{\theta}^\lambda$ yields

$$\mathbf{F}_n(\boldsymbol{\theta}) = \mathbf{F}_n(\boldsymbol{\theta}^\lambda) + \mathbf{A}_n(\boldsymbol{\theta}^\lambda)(\boldsymbol{\theta} - \boldsymbol{\theta}^\lambda) + \mathbf{R}_n(\boldsymbol{\theta}), \quad (8)$$

where, by Assumption 2.2 (ii) and the continuity of \mathbf{H}_ρ , the remainder satisfies

$$\sup_{\boldsymbol{\theta} \in \mathcal{U}_n} \|\mathbf{R}_n(\boldsymbol{\theta})\| = o_p(s_n^{-1/2}).$$

By Assumption 2.2 (iv),

$$\sqrt{s_n} \{ \mathbf{U}_n(\boldsymbol{\theta}^\lambda) - \Psi(\boldsymbol{\theta}^\lambda) \} \rightarrow_d \mathcal{N}(\mathbf{0}, \mathbf{K}^*).$$

Using the penalized population equation $\Psi(\boldsymbol{\theta}^\lambda) + \lambda \nabla \rho(\boldsymbol{\theta}^\lambda) = \mathbf{0}$, we can write

$$\mathbf{F}_n(\boldsymbol{\theta}^\lambda) = \{ \mathbf{U}_n(\boldsymbol{\theta}^\lambda) - \Psi(\boldsymbol{\theta}^\lambda) \} + (\lambda_n - \lambda) \nabla \rho(\boldsymbol{\theta}^\lambda).$$

Hence, if in addition $\sqrt{s_n}(\lambda_n - \lambda) \rightarrow 0$ (e.g. when $\lambda_n \equiv \lambda$), we obtain

$$\sqrt{s_n} \mathbf{F}_n(\boldsymbol{\theta}^\lambda) \rightarrow_d \mathcal{N}(\mathbf{0}, \mathbf{K}^*).$$

Assumptions 2.2 (iii) and 2.1 imply $\mathbf{A}_n(\boldsymbol{\theta}^\lambda) \rightarrow_p \mathbf{J}_\lambda^*$ with nonsingular limit. Applying (8) with $\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}_n$ and using Assumption 2.4 (ii), we obtain the asymptotic linear representation

$$\sqrt{s_n}(\tilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\lambda) = -\mathbf{J}_\lambda^{*-1} \sqrt{s_n} \mathbf{F}_n(\boldsymbol{\theta}^\lambda) + o_p(1). \quad (9)$$

Similarly, applying (8) with $\boldsymbol{\theta} = \boldsymbol{\theta}_{\text{GB}}$ and (6),

$$\mathbf{F}_n(\boldsymbol{\theta}_{\text{GB}}) = \mathbf{F}_n(\boldsymbol{\theta}^\lambda) + \mathbf{A}_n(\boldsymbol{\theta}^\lambda)(\boldsymbol{\theta}_{\text{GB}} - \boldsymbol{\theta}^\lambda) + o_p(s_n^{-1/2}). \quad (10)$$

Moreover, $\mathbf{A}_n(\boldsymbol{\theta}_{\text{GB}}) = \mathbf{A}_n(\boldsymbol{\theta}^\lambda) + o_p(1)$, so its inverse is $\mathbf{A}_n(\boldsymbol{\theta}_{\text{GB}})^{-1} = \mathbf{J}_\lambda^{*-1} + o_p(1)$. Substituting (10) into (7) and simplifying, we find

$$\begin{aligned} \tilde{\boldsymbol{\theta}}_n^{(1)} - \boldsymbol{\theta}^\lambda &= \boldsymbol{\theta}_{\text{GB}} - \boldsymbol{\theta}^\lambda - \mathbf{A}_n(\boldsymbol{\theta}_{\text{GB}})^{-1} \mathbf{F}_n(\boldsymbol{\theta}_{\text{GB}}) \\ &= \boldsymbol{\theta}_{\text{GB}} - \boldsymbol{\theta}^\lambda - \mathbf{J}_\lambda^{*-1} \{ \mathbf{F}_n(\boldsymbol{\theta}^\lambda) + \mathbf{A}_n(\boldsymbol{\theta}^\lambda)(\boldsymbol{\theta}_{\text{GB}} - \boldsymbol{\theta}^\lambda) \} + o_p(s_n^{-1/2}) \\ &= -\mathbf{J}_\lambda^{*-1} \mathbf{F}_n(\boldsymbol{\theta}^\lambda) + o_p(s_n^{-1/2}). \end{aligned}$$

Combining this with (9) gives

$$\sqrt{s_n}(\tilde{\boldsymbol{\theta}}_n^{(1)} - \tilde{\boldsymbol{\theta}}_n) = \sqrt{s_n}(\tilde{\boldsymbol{\theta}}_n^{(1)} - \boldsymbol{\theta}^\lambda) - \sqrt{s_n}(\tilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\lambda) \rightarrow_p \mathbf{0},$$

so that

$$\|\tilde{\boldsymbol{\theta}}_n^{(1)} - \tilde{\boldsymbol{\theta}}_n\| = o_p(s_n^{-1/2}).$$

B Proofs of the main results

B.1 Proof of Proposition 3.1

of Proposition 3.1. Throughout the proof all stochastic limits are taken under the true law P^* . Let $\tilde{\boldsymbol{\theta}}_n$ be any admissible center. We first establish the total-variation limit for a baseline center $\tilde{\boldsymbol{\theta}}_n$, and then apply a translation argument to transfer the result to $\tilde{\boldsymbol{\theta}}_n$.

Introduce local coordinates

$$\mathbf{u} = \sqrt{s_n}(\boldsymbol{\theta} - \check{\boldsymbol{\theta}}_n), \quad \boldsymbol{\theta} = \check{\boldsymbol{\theta}}_n + \frac{\mathbf{u}}{\sqrt{s_n}},$$

and consider the log-density of the generalized posterior $\log \pi_n^\eta$ as

$$\log \pi_n^\eta(\boldsymbol{\theta}) = -\eta M_n(\boldsymbol{\theta}) - \eta \lambda_n s_n \rho(\boldsymbol{\theta}) + r_n(\boldsymbol{\theta})$$

as a function of \mathbf{u} .

Assumption 2.2 (ii) gives a third-order Taylor expansion of M_n at $\check{\boldsymbol{\theta}}_n$,

$$M_n\left(\check{\boldsymbol{\theta}}_n + \frac{\mathbf{u}}{\sqrt{s_n}}\right) = M_n(\check{\boldsymbol{\theta}}_n) + \sqrt{s_n} \mathbf{U}_n(\check{\boldsymbol{\theta}}_n)^\top \mathbf{u} + \frac{1}{2} \mathbf{u}^\top \mathbf{J}_n(\check{\boldsymbol{\theta}}_n) \mathbf{u} + R_{n,1}(\mathbf{u}), \quad (11)$$

where, for each fixed $M < \infty$, $\sup_{\|\mathbf{u}\| \leq M} |R_{n,1}(\mathbf{u})| = o_p(1) \|\mathbf{u}\|^2$.

By Assumption 2.1, ρ is C^2 in a neighborhood of $\boldsymbol{\theta}^\lambda$. Since $\check{\boldsymbol{\theta}}_n \rightarrow_p \boldsymbol{\theta}^\lambda$, a second order Taylor expansion at $\check{\boldsymbol{\theta}}_n$ yields

$$s_n \rho\left(\check{\boldsymbol{\theta}}_n + \frac{\mathbf{u}}{\sqrt{s_n}}\right) = s_n \rho(\check{\boldsymbol{\theta}}_n) + \sqrt{s_n} \nabla \rho(\check{\boldsymbol{\theta}}_n)^\top \mathbf{u} + \frac{1}{2} \mathbf{u}^\top \mathbf{H}_\rho(\check{\boldsymbol{\theta}}_n) \mathbf{u} + R_{n,2}(\mathbf{u}), \quad (12)$$

with $\sup_{\|\mathbf{u}\| \leq M} |R_{n,2}(\mathbf{u})| = o_p(1) \|\mathbf{u}\|^2$.

For the prior remainder r_n , Assumption 2.3 states that r_n is locally Lipschitz on \mathcal{U}_n with Lipschitz constant $L_n = o_p(\sqrt{s_n})$. Lebourg's mean value theorem for locally Lipschitz functions (see, e.g., Clarke (1990, Thm. 2.3.7)) applied along the segment joining $\check{\boldsymbol{\theta}}_n$ and $\check{\boldsymbol{\theta}}_n + \mathbf{u}/\sqrt{s_n}$ gives, for each fixed $M < \infty$ and all $\|\mathbf{u}\| \leq M$,

$$r_n\left(\check{\boldsymbol{\theta}}_n + \frac{\mathbf{u}}{\sqrt{s_n}}\right) = r_n(\check{\boldsymbol{\theta}}_n) + \frac{1}{\sqrt{s_n}} \boldsymbol{\xi}_n(\mathbf{u})^\top \mathbf{u}, \quad (13)$$

where $\boldsymbol{\xi}_n(\mathbf{u})$ is a vector on that segment with $\|\boldsymbol{\xi}_n(\mathbf{u})\| \leq L_n$.

Substituting (11), (12) and (13) into $\log \pi_n^\eta$ at $\check{\boldsymbol{\theta}}_n + \mathbf{u}/\sqrt{s_n}$ yields

$$\log \pi_n^\eta\left(\check{\boldsymbol{\theta}}_n + \frac{\mathbf{u}}{\sqrt{s_n}}\right) = C_n - \frac{1}{2} \mathbf{u}^\top \mathbf{H}_{0,n} \mathbf{u} + R_n(\mathbf{u}),$$

where

$$C_n = -\eta M_n(\check{\boldsymbol{\theta}}_n) - \eta \lambda_n s_n \rho(\check{\boldsymbol{\theta}}_n) + r_n(\check{\boldsymbol{\theta}}_n), \quad \mathbf{H}_{0,n} = \eta \{ \mathbf{J}_n(\check{\boldsymbol{\theta}}_n) + \lambda_n \mathbf{H}_\rho(\check{\boldsymbol{\theta}}_n) \},$$

and $R_n(\mathbf{u})$ collects the quadratic remainders $R_{n,1}$ and $R_{n,2}$.

The linear term in \mathbf{u} arising from M_n and ρ are of order $\sqrt{s_n}$. The admissibility of $\check{\boldsymbol{\theta}}_n$ and the approximate penalized estimating equation imply

$$\mathbf{U}_n(\check{\boldsymbol{\theta}}_n) + \lambda_n \nabla \rho(\check{\boldsymbol{\theta}}_n) = o_p(s_n^{-1/2}),$$

hence

$$-\eta \sqrt{s_n} \mathbf{U}_n(\check{\boldsymbol{\theta}}_n) - \eta \lambda_n \sqrt{s_n} \nabla \rho(\check{\boldsymbol{\theta}}_n) = o_p(1)$$

uniformly on bounded \mathbf{u} .

The linear term coming from r_n is of smaller order: for $\|\mathbf{u}\| \leq M$,

$$\left| \frac{1}{\sqrt{s_n}} \boldsymbol{\xi}_n(\mathbf{u})^\top \mathbf{u} \right| \leq \frac{L_n}{\sqrt{s_n}} M = o_p(1).$$

Combining these bounds with those on $R_{n,1}$ and $R_{n,2}$, we obtain, for each fixed M ,

$$\sup_{\|\mathbf{u}\| \leq M} |R_n(\mathbf{u})| = o_p(1) (1 + \|\mathbf{u}\|^2). \quad (14)$$

after absorbing all linear term into $R_n(\mathbf{u})$. Thus the log posterior admits the quadratic-plus-remainder representation

$$\log \pi_n^\eta \left(\check{\boldsymbol{\theta}}_n + \frac{\mathbf{u}}{\sqrt{s_n}} \right) = C_n - \frac{1}{2} \mathbf{u}^\top \mathbf{H}_{0,n} \mathbf{u} + R_n(\mathbf{u}),$$

with R_n satisfying (14).

Assumption 2.2 (iii), the convergence $\check{\boldsymbol{\theta}}_n \rightarrow_p \boldsymbol{\theta}^\lambda$, and continuity of \mathbf{H}_ρ at $\boldsymbol{\theta}^\lambda$ imply

$$\mathbf{J}_n(\check{\boldsymbol{\theta}}_n) \rightarrow_p \mathbf{J}^*, \quad \mathbf{H}_\rho(\check{\boldsymbol{\theta}}_n) \rightarrow_p \mathbf{H}_\rho(\boldsymbol{\theta}^\lambda),$$

and therefore

$$\mathbf{H}_{0,n} \rightarrow_p \mathbf{H}_0 := \eta \mathbf{J}_\lambda^*,$$

with \mathbf{H}_0 positive definite.

Let P_n^0 denote the Gaussian measure $\mathcal{N}(\mathbf{0}, \mathbf{H}_{0,n}^{-1})$ on \mathbb{R}^p , and let π_n° be the law of \mathbf{u} induced by Π_n^η via the above re-parametrization. The density of π_n° relative to P_n^0 is proportional to $\exp\{R_n(\mathbf{u})\}$. Since $\mathbf{H}_{0,n} \rightarrow_p \mathbf{H}_0 \succ 0$, the eigenvalues of $\mathbf{H}_{0,n}$ are bounded away from zero and infinity with probability tending to one, and Gaussian tails under P_n^0 are uniformly controlled. Choosing $M < \infty$ large enough, we may assume $P_n^0(\|\mathbf{u}\| > M)$ is arbitrarily small uniformly in n . On $\{\|\mathbf{u}\| \leq M\}$, (14) implies $\sup_{\|\mathbf{u}\| \leq M} |R_n(\mathbf{u})| = o_p(1)$, and $\exp\{R_n(\mathbf{u})\}$ converges to 1 uniformly on this set. A truncation and dominated convergence argument then yields

$$\mathbb{E}_{P_n^0}[\exp\{R_n(\mathbf{U})\} - 1] = o_p(1), \quad \mathbb{E}_{P_n^0}[|\exp\{R_n(\mathbf{U})\} - 1|] = o_p(1).$$

By Scheffé's lemma applied to the Radon–Nikodym derivatives with respect to Lebesgue measure,

$$\int_{\mathbb{R}^p} |p_n^\circ(\mathbf{x}) - \mathcal{N}(\mathbf{x} \mid \mathbf{0}, \mathbf{H}_{0,n}^{-1})| d\mathbf{x} \rightarrow_p 0, \quad (15)$$

where p_n° denotes the density of π_n° . Since $\mathbf{H}_{0,n} \rightarrow_p \mathbf{H}_0$ and Gaussian laws depend continuously on the covariance matrix in total variation when eigenvalues are uniformly bounded, we also have

$$\int_{\mathbb{R}^p} |\mathcal{N}(\mathbf{x} \mid \mathbf{0}, \mathbf{H}_{0,n}^{-1}) - \mathcal{N}(\mathbf{x} \mid \mathbf{0}, \mathbf{H}_0^{-1})| d\mathbf{x} = o_p(1). \quad (16)$$

Combining (15) and (16) and using the triangle inequality gives

$$\int_{\mathbb{R}^p} |q_n^{\check{\boldsymbol{\theta}}_n}(\mathbf{x}) - \mathcal{N}(\mathbf{x} \mid \mathbf{0}, \mathbf{H}_0^{-1})| d\mathbf{x} \rightarrow_p 0, \quad n \rightarrow \infty. \quad (17)$$

where $q_n^{\check{\boldsymbol{\theta}}_n}$ is the density of $\sqrt{s_n}(\boldsymbol{\theta} - \check{\boldsymbol{\theta}}_n)$ under Π_n^η .

To obtain the stated result at an arbitrary admissible center $\tilde{\boldsymbol{\theta}}_n$, define

$$\boldsymbol{\delta}_n := \sqrt{s_n}(\tilde{\boldsymbol{\theta}}_n - \check{\boldsymbol{\theta}}_n).$$

Both $\tilde{\boldsymbol{\theta}}_n$ and $\check{\boldsymbol{\theta}}_n$ are admissible, so

$$\|\boldsymbol{\delta}_n\| \leq \sqrt{s_n} \|\tilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\lambda\| + \sqrt{s_n} \|\check{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^\lambda\| = o_p(1),$$

and thus $\boldsymbol{\delta}_n \rightarrow_p \mathbf{0}$. Let $q_n^{\tilde{\boldsymbol{\theta}}_n}$ be the density of $\sqrt{s_n}(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}_n)$ under Π_n^η . For each realization of the data,

$$q_n^{\tilde{\boldsymbol{\theta}}_n}(\mathbf{x}) = q_n^{\check{\boldsymbol{\theta}}_n}(\mathbf{x} + \boldsymbol{\delta}_n), \quad \mathbf{x} \in \mathbb{R}^p.$$

Write $\varphi(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \mathbf{0}, \mathbf{H}_0^{-1})$. Using the change of variables $\mathbf{y} = \mathbf{x} + \boldsymbol{\delta}_n$ and the triangle inequality,

$$\begin{aligned} \int |q_n^{\tilde{\boldsymbol{\theta}}_n}(\mathbf{x}) - \varphi(\mathbf{x})| d\mathbf{x} &= \int |q_n^{\check{\boldsymbol{\theta}}_n}(\mathbf{x}) - \varphi(\mathbf{x} - \boldsymbol{\delta}_n)| d\mathbf{x} \\ &\leq \int |q_n^{\check{\boldsymbol{\theta}}_n}(\mathbf{x}) - \varphi(\mathbf{x})| d\mathbf{x} + \int |\varphi(\mathbf{x} - \boldsymbol{\delta}_n) - \varphi(\mathbf{x})| d\mathbf{x}. \end{aligned} \quad (18)$$

The first term on the right-hand side converges to zero in probability by (17). For the second term, the Gaussian density φ is continuous in L^1 under translations, so for any deterministic sequence $\mathbf{h}_n \rightarrow \mathbf{0}$,

$$\int_{\mathbb{R}^p} |\varphi(\mathbf{x} - \mathbf{h}_n) - \varphi(\mathbf{x})| d\mathbf{x} \rightarrow 0.$$

Since $\delta_n \rightarrow_p \mathbf{0}$, the same convergence holds in probability along the random sequence δ_n . Therefore the second term in (18) is $o_p(1)$, and we obtain

$$\int_{\mathbb{R}^p} \left| q_n^{\tilde{\theta}_n}(\mathbf{x}) - \mathcal{N}(\mathbf{x} | \mathbf{0}, \mathbf{H}_0^{-1}) \right| d\mathbf{x} \rightarrow_p 0.$$

Recalling that $\mathbf{H}_0 = \eta \mathbf{J}_\lambda^*$, this is precisely the claimed total-variation convergence of the law of $\sqrt{s_n}(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}_n)$ under Π_n^η to $\mathcal{N}(\mathbf{0}, (\eta \mathbf{J}_\lambda^*)^{-1})$ at any admissible center $\tilde{\boldsymbol{\theta}}_n$. \square

B.2 Proof of Theorem 3.2

of Theorem 3.2. Recall the working curvature $\mathbf{H}_0 := \eta \mathbf{J}_\lambda^*$ and the target covariance $\mathbf{V}_{\text{target}}^* = (\mathbf{J}_\lambda^*)^{-1} \mathbf{K}^* (\mathbf{J}_\lambda^*)^{-1}$. Under Assumptions 2.1–2.4, the posterior mean $\boldsymbol{\theta}_{\text{GB}}$ is an admissible center. Hence, Proposition 3.1 applied at the center $\boldsymbol{\theta}_{\text{GB}}$ yields, conditionally on the data,

$$\sqrt{s_n}(\boldsymbol{\theta}^{(d)} - \boldsymbol{\theta}_{\text{GB}}) \rightarrow_d \mathcal{N}(\mathbf{0}, \mathbf{H}_0^{-1}), \quad n \rightarrow \infty. \quad (19)$$

Let $\tilde{\boldsymbol{\theta}}_n$ be any admissible center and recall the location–scale calibration map

$$\boldsymbol{\Omega} := (\mathbf{V}_{\text{target}}^*)^{1/2} \mathbf{H}_0^{1/2}, \quad \boldsymbol{\theta}_{\text{calib}}^{(d)} := \tilde{\boldsymbol{\theta}}_n + \boldsymbol{\Omega}(\boldsymbol{\theta}^{(d)} - \boldsymbol{\theta}_{\text{GB}}).$$

The centering at $\tilde{\boldsymbol{\theta}}_n$ cancels:

$$\sqrt{s_n}(\boldsymbol{\theta}_{\text{calib}}^{(d)} - \tilde{\boldsymbol{\theta}}_n) = \boldsymbol{\Omega} \sqrt{s_n}(\boldsymbol{\theta}^{(d)} - \boldsymbol{\theta}_{\text{GB}}), \quad (20)$$

Combining (19) with (20) and applying the continuous mapping theorem to the fixed linear map $\mathbf{x} \mapsto \boldsymbol{\Omega} \mathbf{x}$ gives, conditionally on the data,

$$\sqrt{s_n}(\boldsymbol{\theta}_{\text{calib}}^{(d)} - \tilde{\boldsymbol{\theta}}_n) \rightarrow_d \mathcal{N}(\mathbf{0}, \boldsymbol{\Omega} \mathbf{H}_0^{-1} \boldsymbol{\Omega}^\top), \quad n \rightarrow \infty. \quad (21)$$

By construction of $\boldsymbol{\Omega}$,

$$\boldsymbol{\Omega} \mathbf{H}_0^{-1} \boldsymbol{\Omega}^\top = (\mathbf{V}_{\text{target}}^*)^{1/2} \mathbf{H}_0^{1/2} \mathbf{H}_0^{-1} \mathbf{H}_0^{1/2} (\mathbf{V}_{\text{target}}^*)^{1/2} = \mathbf{V}_{\text{target}}^*. \quad (22)$$

Substituting (22) into (21) yields

$$\sqrt{s_n}(\boldsymbol{\theta}_{\text{calib}}^{(d)} - \tilde{\boldsymbol{\theta}}_n) \rightarrow_d \mathcal{N}(\mathbf{0}, \mathbf{V}_{\text{target}}^*).$$

In particular, since $\mathbf{H}_0 = \eta \mathbf{J}_\lambda^*$ and $\mathbf{V}_{\text{target}}^*$ depends only on $(\mathbf{J}_\lambda^*, \mathbf{K}^*)$, the limiting covariance $\mathbf{V}_{\text{target}}^*$ is independent of the learning rate η . This establishes learning-rate invariance of the limiting law. \square

C Proofs of the plug-in results

C.1 Proof of Lemma 4.1

Let $q_n^{\tilde{\theta}_n}$ be the density of the scaled posterior $\mathbf{X}_n = \sqrt{s_n}(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}_n)$ under Π_n^η for an admissible center $\tilde{\boldsymbol{\theta}}_n$. We first present two generic facts that will be used repeatedly.

Lemma C.1. *Suppose that, for some positive definite matrix \mathbf{H}_0 ,*

$$\int_{\mathbb{R}^p} \left| q_n^{\tilde{\theta}_n}(\mathbf{x}) - \mathcal{N}(\mathbf{x} | \mathbf{0}, \mathbf{H}_0^{-1}) \right| d\mathbf{x} \rightarrow 0, \quad n \rightarrow \infty$$

and that there exist random constants $C_n, c_n > 0$, bounded in probability, such that

$$q_n^{\tilde{\theta}_n}(\mathbf{x}) \leq C_n \exp(-c_n |\mathbf{x}|^2) \quad \text{for all } \mathbf{x} \in \mathbb{R}^p, \text{ all } n.$$

Then, for any function $g : \mathbb{R}^p \rightarrow \mathbb{R}^k$, ($k \in \mathbb{N}$) with polynomial growth,

$$\int g(\mathbf{x}) q_n^{\tilde{\theta}_n}(\mathbf{x}) d\mathbf{x} \rightarrow \int g(\mathbf{x}) \mathcal{N}(\mathbf{x} | \mathbf{0}, \mathbf{H}_0^{-1}) d\mathbf{x}.$$

In particular,

$$\mathbb{E}[\mathbf{X}_n] \rightarrow \mathbf{0}, \quad \mathbb{E}[\mathbf{X}_n \mathbf{X}_n^\top] \rightarrow \mathbf{H}_0^{-1}, \quad \sup_n \mathbb{E}[\|\mathbf{X}_n\|^4] < \infty.$$

of Lemma C.1. Total-variation convergence implies convergence of expectations for bounded measurable function g . For a polynomially growing g , fix $R > 0$ and write

$$\int g(\mathbf{x}) q_n^{\tilde{\theta}_n}(\mathbf{x}) d\mathbf{x} = \int_{\|\mathbf{x}\| \leq R} g(\mathbf{x}) q_n^{\tilde{\theta}_n}(\mathbf{x}) d\mathbf{x} + \int_{\|\mathbf{x}\| > R} g(\mathbf{x}) q_n^{\tilde{\theta}_n}(\mathbf{x}) d\mathbf{x}.$$

On $\{\|\mathbf{x}\| \leq R\}$, g is bounded, so total-variation convergence yields convergence of the first term as $n \rightarrow \infty$, and the limit as $R \rightarrow \infty$ recovers the Gaussian expectation. For the tail term, polynomial growth of g together with the Gaussian domination implies

$$\sup_n C_n \int_{\|\mathbf{x}\| > R} |g(\mathbf{x})| \exp(-c_n \|\mathbf{x}\|^2) d\mathbf{x} \rightarrow 0, \quad R \rightarrow \infty.$$

Indeed, since $(C_n)_n$ and $(c_n)_n$ are bounded in probability, there exists $B < \infty$ such that, with probability tending to one, $C_n \leq B$ and $c_n \geq B^{-1}$ for all n . On this event the integrand is dominated by a fixed Gaussian envelope of the form $(1 + \|\mathbf{x}\|^m) \exp(-B^{-1} \|\mathbf{x}\|^2)$ for some $m \geq 0$, which is integrable, so dominated convergence yields the claimed limit. This gives convergence of expectations for polynomially growing g , and the stated consequences follow by taking $g(\mathbf{x}) = \mathbf{x}$, $g(\mathbf{x}) = \mathbf{x}\mathbf{x}^\top$, and $g(\mathbf{x}) = \|\mathbf{x}\|^4$. \square

The required Gaussian domination is standard in Bernstein–von Mises arguments under local asymptotic normality: in a shrinking neighborhood \mathcal{U}_n of $\boldsymbol{\theta}^\lambda$, the log posterior is a quadratic with positive-definite curvature \mathbf{H}_0 plus an $o_p(1)$ perturbation, while outside \mathcal{U}_n the quadratic term dominates. Under Assumptions 2.2–2.3, this gives sub-Gaussian tails for \mathbf{X}_n , uniformly in n .

Lemma C.2. *For any admissible center $\tilde{\boldsymbol{\theta}}_n$,*

$$\sqrt{s_n}(\boldsymbol{\theta}_{\text{GB}} - \tilde{\boldsymbol{\theta}}_n) \rightarrow_p \mathbf{0}$$

and

$$s_n \text{Var}_{\Pi_n^\eta}[\boldsymbol{\theta}] = \mathbb{E}[\mathbf{X}_n \mathbf{X}_n^\top] + o_p(1),$$

where $\mathbf{X}_n = \sqrt{s_n}(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}_n)$.

of Lemma C.2. By definition,

$$\boldsymbol{\theta}_{\text{GB}} = \mathbb{E}_{\Pi_n^\eta}[\boldsymbol{\theta}] = \tilde{\boldsymbol{\theta}}_n + \frac{1}{\sqrt{s_n}} \mathbb{E}[\mathbf{X}_n],$$

so $\sqrt{s_n}(\boldsymbol{\theta}_{\text{GB}} - \tilde{\boldsymbol{\theta}}_n) = \mathbb{E}[\mathbf{X}_n]$. By Lemma C.1 with $g(\mathbf{x}) = \mathbf{x}$ and the centered Gaussian limit, $\mathbb{E}[\mathbf{X}_n] \rightarrow_p \mathbf{0}$, proving the first claim.

For the variance,

$$s_n \text{Var}_{\Pi_n^\eta}[\boldsymbol{\theta}] = \mathbb{E}[\mathbf{X}_n \mathbf{X}_n^\top] - \mathbb{E}[\mathbf{X}_n] \mathbb{E}[\mathbf{X}_n]^\top = \mathbb{E}[\mathbf{X}_n \mathbf{X}_n^\top] + o_p(1),$$

since $\mathbb{E}[\mathbf{X}_n] \rightarrow_p \mathbf{0}$ and $\sup_n \mathbb{E}[\mathbf{X}_n]^4 < \infty$ by Lemma C.1. \square

of Lemma 4.1. For the asymptotic form of the posterior covariance, fix an admissible center $\tilde{\boldsymbol{\theta}}_n$ and consider $\mathbf{X}_n = \sqrt{s_n}(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}_n)$. By Proposition 3.1, the law of \mathbf{X}_n under Π_n^η converges in total variation to $\mathcal{N}(\mathbf{0}, \mathbf{H}_0^{-1})$. Lemma C.1 with $g(\mathbf{x}) = \mathbf{x}\mathbf{x}^\top$ gives

$$\mathbb{E}[\mathbf{X}_n \mathbf{X}_n^\top] \rightarrow_p \mathbf{H}_0^{-1}.$$

Combining this with Lemma C.2,

$$s_n \text{Var}_{\Pi_n^\eta}[\boldsymbol{\theta}] = \mathbb{E}[\mathbf{X}_n \mathbf{X}_n^\top] - \mathbb{E}[\mathbf{X}_n] \mathbb{E}[\mathbf{X}_n]^\top \rightarrow_p \mathbf{H}_0^{-1},$$

that is,

$$s_n \boldsymbol{\Sigma}_{\text{post},n} \rightarrow_p \mathbf{H}_0^{-1}, \quad n \rightarrow \infty.$$

By the assumption of the Lemma 4.1,

$$s_n \hat{\boldsymbol{\Sigma}}_{\text{post}} = s_n \boldsymbol{\Sigma}_{\text{post},n} + s_n (\hat{\boldsymbol{\Sigma}}_{\text{post}} - \boldsymbol{\Sigma}_{\text{post},n}) \rightarrow_p \mathbf{H}_0^{-1}.$$

This completes the proof. \square

C.2 Proof of Lemma 4.2

of Lemma 4.2. By Assumption 2.2 (iii), there exists a neighborhood \mathfrak{N} of θ^λ such that $\sup_{\theta \in \mathfrak{N}} \|J_n(\theta) - J^*\| \rightarrow_p 0$. Since $\bar{\theta}_n \rightarrow_p \theta^\lambda$, we have $J_n(\bar{\theta}_n) \rightarrow_p J^*$. By Assumption 2.1, ρ is C^2 in a neighborhood of θ^λ , so $H_\rho(\bar{\theta}_n) \rightarrow_p H_\rho(\theta^\lambda)$. With $\lambda_n \rightarrow \lambda$, this yields

$$\hat{J}_\lambda = J_n(\bar{\theta}_n) + \lambda_n H_\rho(\bar{\theta}_n) \rightarrow_p J^* + \lambda H_\rho(\theta^\lambda) = J_\lambda^*.$$

Assumption 2.2 (iii) also guarantees that J_λ^* is nonsingular, so the matrix inverse is continuous in a neighborhood, implying $\hat{J}_\lambda^{-1} \rightarrow_p (J_\lambda^*)^{-1}$.

By assumption, $\hat{K} \rightarrow_p K^*$. A continuous mapping argument then gives

$$\hat{V}_{\text{target}} = \hat{J}_\lambda^{-1} \hat{K} \hat{J}_\lambda^{-1} \rightarrow_p (J_\lambda^*)^{-1} K^* (J_\lambda^*)^{-1} = V_{\text{target}}^*.$$

□

C.3 Proof of Proposition 4.3

of Proposition 4.3. Fix an admissible center $\bar{\theta}_n$ and set $X_n := \sqrt{s_n}(\theta^{(d)} - \bar{\theta}_n)$. By Proposition 3.1, conditionally on the data,

$$X_n \rightarrow_d \mathcal{N}(\mathbf{0}, H_0^{-1}),$$

where $H_0 = \eta J_\lambda^*$.

Lemma C.2 gives $\sqrt{s_n}(\theta_{\text{GB}} - \bar{\theta}_n) \rightarrow_p \mathbf{0}$, and by the Monte Carlo rate assumption, $\sqrt{s_n}(\hat{\theta}_{\text{GB}} - \theta_{\text{GB}}) \rightarrow_p \mathbf{0}$. Thus

$$\sqrt{s_n}(\theta^{(d)} - \hat{\theta}_{\text{GB}}) = X_n - \sqrt{s_n}(\hat{\theta}_{\text{GB}} - \bar{\theta}_n) = X_n + o_p(1),$$

so by Slutsky's theorem,

$$\sqrt{s_n}(\theta^{(d)} - \hat{\theta}_{\text{GB}}) \rightarrow_d \mathcal{N}(\mathbf{0}, H_0^{-1}).$$

By Lemma 4.1, $\hat{H}_0^{-1} \rightarrow_p H_0^{-1}$ and therefore $\hat{H}_0^{1/2} \rightarrow_p H_0^{1/2}$. By Lemma 4.2, $\hat{V}_{\text{target}} \rightarrow_p V_{\text{target}}^*$ and hence $\hat{V}_{\text{target}}^{1/2} \rightarrow_p (V_{\text{target}}^*)^{1/2}$. Consequently,

$$\hat{\Omega} = \hat{V}_{\text{target}}^{1/2} \hat{H}_0^{1/2} \rightarrow_p \Omega := (V_{\text{target}}^*)^{1/2} H_0^{1/2}.$$

From (4),

$$\sqrt{s_n}(\hat{\theta}_{\text{calib}}^{(d)} - \bar{\theta}_n) = \hat{\Omega} \sqrt{s_n}(\theta^{(d)} - \hat{\theta}_{\text{GB}}).$$

Applying Slutsky's theorem again,

$$\sqrt{s_n}(\hat{\theta}_{\text{calib}}^{(d)} - \bar{\theta}_n) \rightarrow_d \mathcal{N}(\mathbf{0}, \Omega H_0^{-1} \Omega^\top) = \mathcal{N}(\mathbf{0}, V_{\text{target}}^*),$$

where the last equality uses

$$\Omega H_0^{-1} \Omega^\top = (V_{\text{target}}^*)^{1/2} H_0^{1/2} H_0^{-1} H_0^{1/2} (V_{\text{target}}^*)^{1/2} = V_{\text{target}}^*.$$

Finally, the definition (2) shows that V_{target}^* does not depend on the learning rate η , so the limiting calibrated law is learning-rate invariant. □

D Non-smooth penalties: active set and subgradient calculus

In the main text we assume that the penalty ρ is twice continuously differentiable in a neighborhood of the target point θ^λ . Here we outline how the arguments extend when ρ is convex but possibly non-smooth. Throughout this section we work with the penalized population inclusion

$$\mathbf{0} \in \Psi(\theta) + \lambda \partial \rho(\theta)$$

with solution θ^λ , where $\partial \rho$ is the convex subdifferential. For an index set $A \subset \{1, \dots, p\}$, we write $\theta_A := (\theta_j)_{j \in A}$ for the subvector of θ with coordinates in A , and, for any matrix M , we write $M_{AA} := P_A M P_A^\top$ for the corresponding principal submatrix, where P_A denotes the coordinate projection onto A .

Assumption D.1 (Non-smooth penalty and active set). The penalty $\rho : \mathbb{R}^p \rightarrow (-\infty, \infty]$ is convex and lower semicontinuous. The penalized population inclusion has a unique solution θ^λ . Moreover, there exists an index set

$$A := \{j : \theta_j^\lambda \neq 0\} \cup \{j : 0 \in \partial \rho_j(\theta^\lambda)\}$$

and a neighborhood \mathfrak{N}_A of θ_A^λ such that the map $\theta_A \mapsto \rho(\theta_A, \theta_{A^c}^\lambda)$ is twice continuously differentiable on \mathfrak{N}_A ; the Clarke generalized Jacobian $\partial^G \nabla \rho(\theta^\lambda)$ is nonempty and contains a symmetric positive semidefinite matrix $\mathbf{H}_\rho(\theta^\lambda)$; and the matrix $\mathbf{J}_\lambda^* := \mathbf{J}^* + \lambda \mathbf{H}_\rho(\theta^\lambda)$ is nonsingular, where $\mathbf{J}^* := -\nabla_\theta \Psi(\theta^\lambda)$ and $\mathbf{K}^* := \text{Var}_{P^*}[\psi(\mathcal{D}_1, \theta^\lambda)]$.

Throughout this section we work on the active set A supplied by Assumption D.1. Under Assumption D.1, we define

$$\mathbf{V}_{\text{target}}^* := (\mathbf{J}_\lambda^*)^{-1} \mathbf{K}^* (\mathbf{J}_\lambda^*)^{-1},$$

exactly as in the smooth case, and we denote by \mathbf{J}_{AA}^* , \mathbf{K}_{AA}^* , and so on, the corresponding AA -blocks.

Lemma D.2 (Active-set smoothness and curvature). *Suppose Assumption D.1 holds and define*

$$\mathbf{H}_{\rho,AA} := \left. \nabla_{\theta_A}^2 \rho(\theta_A, \theta_{A^c}^\lambda) \right|_{\theta_A = \theta_A^\lambda}.$$

Then $\mathbf{H}_{\rho,AA}$ exists, is symmetric positive semidefinite, and coincides with the AA -block of any symmetric selection $\mathbf{H}_\rho(\theta^\lambda) \in \partial^G \nabla \rho(\theta^\lambda)$, that is,

$$\mathbf{H}_{\rho,AA} = P_A \mathbf{H}_\rho(\theta^\lambda) P_A^\top.$$

Consequently,

$$\mathbf{J}_{\lambda,AA}^* := \mathbf{J}_{AA}^* + \lambda \mathbf{H}_{\rho,AA}$$

is nonsingular whenever \mathbf{J}_λ^ is nonsingular. For the ℓ_1 -penalty $\rho(\theta) = \sum_{j=1}^p |\theta_j|$, we have $\mathbf{H}_{\rho,AA} = \mathbf{0}$.*

of Lemma D.2. By Assumption D.1, ρ is convex and lower semicontinuous and is C^2 in the active coordinates on a neighborhood of θ_A^λ , which gives the existence and positive semidefiniteness of $\mathbf{H}_{\rho,AA}$. The Clarke generalized Jacobian $\partial^G \nabla \rho(\theta^\lambda)$ is the convex hull of limits of classical Hessians at differentiability points approaching θ^λ . Since ρ is C^2 in the active coordinates in a neighborhood of θ_A^λ , every such limit has the same AA -block, which must therefore coincide with $\mathbf{H}_{\rho,AA}$ for any symmetric selection $\mathbf{H}_\rho(\theta^\lambda)$. The nonsingularity of $\mathbf{J}_{\lambda,AA}^*$ follows from that of \mathbf{J}_λ^* and the block structure.

For the ℓ_1 -penalty, each coordinate map $\theta_j \mapsto |\theta_j|$ is affine on any interval not containing the origin. Under sign/gap stability, $\theta_j^\lambda \neq 0$ for $j \in A$ and the neighborhood of θ_A^λ can be chosen so that $\text{sgn}(\theta_j)$ is constant for $j \in A$. Hence ρ is locally linear in θ_A and the second derivative vanishes, giving $\mathbf{H}_{\rho,AA} = \mathbf{0}$. \square

The next lemma records a standard penalized Fisher expansion on the active coordinates. Let $\hat{\theta}_n^{\text{pen}}$ be any measurable solution of the penalized estimating equation

$$\mathbf{0} \in \mathbf{U}_n(\hat{\theta}_n^{\text{pen}}) + \lambda_n \partial \rho(\hat{\theta}_n^{\text{pen}}) \quad (23)$$

such that $\hat{\theta}_n^{\text{pen}} \rightarrow_p \theta^\lambda$ and

$$\Pr\left(\text{sgn}(\hat{\theta}_{n,j}^{\text{pen}}) = \text{sgn}(\theta_j^\lambda) \text{ for all } j \in A\right) \rightarrow 1.$$

The latter is the usual sign/gap stability requirement for the active set.

Lemma D.3 (Subgradient Fisher expansion on the active set). *Suppose Assumptions D.1, 2.2 and 2.3 hold, and let $\hat{\theta}_n^{\text{pen}}$ satisfy (23) with the above sign/gap stability. Then*

$$\sqrt{s_n} (P_A \hat{\theta}_n^{\text{pen}} - P_A \theta^\lambda) \rightarrow_d \mathcal{N}\left(\mathbf{0}, (\mathbf{J}_{\lambda,AA}^*)^{-1} \mathbf{K}_{AA}^* (\mathbf{J}_{\lambda,AA}^*)^{-1}\right), \quad n \rightarrow \infty.$$

of Lemma D.3. On the event where $\text{sgn}(\hat{\theta}_{n,j}^{\text{pen}}) = \text{sgn}(\theta_j^\lambda)$ for all $j \in A$ and $\hat{\theta}_n^{\text{pen}}$ lies in the neighborhood \mathfrak{N}_A from Lemma D.2, the restriction of ρ to the active coordinates is C^2 and the subgradient in those coordinates is a singleton: there exists a measurable choice $\nabla_A \rho(\theta)$ such that

$$\partial \rho(\theta) \cap \mathbb{R}^A = \{\nabla_A \rho(\theta)\}, \quad \theta_A \in \mathfrak{N}_A.$$

Projecting (23) onto A and using this selection gives the smooth estimating equation

$$\mathbf{0} = P_A \mathbf{U}_n(\hat{\theta}_n^{\text{pen}}) - \lambda_n \nabla_A \rho(\hat{\theta}_n^{\text{pen}}).$$

A first-order Taylor expansion of both \mathbf{U}_n and $\nabla_A \rho$ at θ^λ yields

$$\mathbf{0} = P_A \mathbf{U}_n(\theta^\lambda) - \lambda_n \nabla_A \rho(\theta^\lambda) + \mathbf{J}_{AA}^*(\hat{\theta}_{n,A}^{\text{pen}} - \theta_A^\lambda) + \lambda_n \mathbf{H}_{\rho,AA}(\hat{\theta}_{n,A}^{\text{pen}} - \theta_A^\lambda) + o_p(\|\hat{\theta}_n^{\text{pen}} - \theta^\lambda\|)$$

on A , uniformly on events of probability tending to one. Using the population inclusion $\mathbf{0} \in \Psi(\theta^\lambda) - \lambda \partial \rho(\theta^\lambda)$ and the definition of \mathbf{J}^* , this can be written as

$$\mathbf{0} = P_A \{\mathbf{U}_n(\theta^\lambda) - \Psi(\theta^\lambda)\} + \mathbf{J}_{\lambda,AA}^*(\hat{\theta}_{n,A}^{\text{pen}} - \theta_A^\lambda) + o_p(\|\hat{\theta}_n^{\text{pen}} - \theta^\lambda\|).$$

By Assumption 2.2 (iv) and the definition of \mathbf{K}^* ,

$$\sqrt{s_n} P_A \{\mathbf{U}_n(\theta^\lambda) - \Psi(\theta^\lambda)\} \rightarrow_d \mathcal{N}(\mathbf{0}, \mathbf{K}_{AA}^*).$$

Since $\hat{\theta}_n^{\text{pen}} \rightarrow_p \theta^\lambda$ and $\mathbf{J}_{\lambda,AA}^*$ is nonsingular by Lemma D.2, the remainder term is $o_p(s_n^{-1/2})$ and the usual M -estimation argument yields

$$\sqrt{s_n}(\hat{\theta}_{n,A}^{\text{pen}} - \theta_A^\lambda) = -(\mathbf{J}_{\lambda,AA}^*)^{-1} \sqrt{s_n} P_A \{\mathbf{U}_n(\theta^\lambda) - \Psi(\theta^\lambda)\} + o_p(1).$$

The desired normal limit follows by Slutsky's theorem. \square

Lemmas D.2 and D.3 show that, under the non-smooth setting, the active coordinates behave as in the smooth case, with curvature matrix $\mathbf{J}_{\lambda,AA}^*$ and variability \mathbf{K}_{AA}^* . The proofs of Proposition 3.1, Theorem 3.2, Lemma 4.1, Lemma 4.2, and Proposition 4.3, use only local quadratic expansions and the central limit theorem. Repeating those arguments with all matrices and vectors restricted to A yields the same Gaussian limits on the active coordinates after replacing \mathbf{J}^* , \mathbf{J}_λ^* , \mathbf{K}^* , $\mathbf{V}_{\text{target}}^*$ by their AA -blocks. Coordinates in A^c may exhibit boundary phenomena, and their asymptotic distribution need not be Gaussian. Accordingly, our asymptotic statements and calibrated inference for non-smooth penalties are reported conditionally on the active set A .

E Detailed settings for experiment

E.1 Sampling algorithm

We describe the augmentation and Gibbs sampler used to approximate the generalized Bayes posterior Π_n^η in Section 5.1. All notation is as in the main text: in particular, the Huber loss $M_n(\beta)$ is defined there in terms of the whitened residuals $\tilde{r}_{ij}(\beta)$ and effective scale $s_n = n$.

For the Huber loss $\rho_c(u)$ in Section 5.1, the following infimal-convolution representation holds:

$$\rho_c(u) = \min_{t \in \mathbb{R}} \left\{ \frac{1}{2}(u - t)^2 + c|t| - \frac{1}{2}c^2 \right\}.$$

Hence, up to a multiplicative constant,

$$\exp\{-\eta M_n(\beta)\} \propto \int \exp\left\{-\eta \sum_{i=1}^G \sum_{j=1}^{n_i} \left[\frac{1}{2}(\tilde{r}_{ij}(\beta) - t_{ij})^2 + c|t_{ij}| \right]\right\} \prod_{i,j} dt_{ij},$$

where $t_{ij} \in \mathbb{R}$ are latent variables.

We then use the standard normal-exponential mixture representation of the Laplace kernel: for $\kappa > 0$,

$$\exp(-\kappa|t|) \propto \int_0^\infty \frac{1}{\sqrt{2\pi\omega}} \exp\left(-\frac{t^2}{2\omega}\right) \frac{\kappa^2}{2} \exp\left(-\frac{\kappa^2}{2}\omega\right) d\omega.$$

Setting $\kappa = \eta c$ introduces latent scales $\omega_{ij} > 0$. Under this augmentation the generalized posterior is conjugate in each block $(\beta, \mathbf{t}, \omega)$.

Let \mathbf{L}_i be the symmetric square root of the working covariance Σ_i from Section 5.1, and define $\tilde{\mathbf{y}}_i = \mathbf{L}_i^{-1} \mathbf{y}_i$ and $\mathbf{t}_i = (t_{i1}, \dots, t_{in_i})^\top$. The prior is the Gaussian prior with n -dependent scale used in the main text, with log-density

$$\log \pi_n(\beta) = -\lambda_n s_n \rho(\beta) + r_n,$$

where $\rho(\beta) = 2^{-1}(\beta - \mu)^\top \mathbf{Q}(\beta - \mu)$, $\lambda_n \rightarrow \lambda \in [0, \infty)$ and r_n is a normalizing constant.

Given (\mathbf{t}, ω) , the full conditional of β is multivariate normal,

$$\beta \mid \mathbf{t}, \omega, \{\mathcal{D}_i\}_{i=1}^n \sim \mathcal{N}(\mathbf{m}_{\text{post}}, \Lambda_{\text{post}}^{-1}),$$

with

$$\Lambda_{\text{post}} = \eta \sum_{i=1}^G \tilde{\mathbf{X}}_i^\top \tilde{\mathbf{X}}_i + \lambda_n s_n \mathbf{Q}, \quad \mathbf{m}_{\text{post}} = \Lambda_{\text{post}}^{-1} \left\{ \eta \sum_{i=1}^G \tilde{\mathbf{X}}_i^\top (\tilde{\mathbf{y}}_i - \mathbf{t}_i) + \lambda_n s_n \mathbf{Q} \mu \right\}.$$

Conditional on (β, ω) , the latent t_{ij} are independent normals. Writing $\tilde{r}_{ij} = \tilde{r}_{ij}(\beta)$,

$$t_{ij} \mid \beta, \omega_{ij}, \{\mathcal{D}_i\}_{i=1}^n \sim \mathcal{N}(\mu_{t,ij}, \sigma_{t,ij}^2), \quad \sigma_{t,ij}^2 = (\eta + \omega_{ij}^{-1})^{-1}, \quad \mu_{t,ij} = \sigma_{t,ij}^2 \eta \tilde{r}_{ij}.$$

Finally, conditional on t_{ij} the latent scales ω_{ij} follow an inverse-Gaussian distribution. With the parameterization $\text{IG}(\mu, \lambda)$ used in our implementation,

$$\omega_{ij} \mid t_{ij} \sim \text{IG}\left(\mu = \frac{|t_{ij}|}{\eta c}, \lambda = 1\right).$$

A single iteration of the Gibbs sampler consists of the three updates $\omega \rightarrow \mathbf{t} \rightarrow \beta$.

E.2 Simulation settings

We fix $G = 100$, $n_i = 5$ so that $n = 500$, and take $p = 1$ to focus on the slope β . Covariates are generated as $x_{ij} \sim \mathcal{N}(0, 1)$, and the data-generating model is

$$y_{ij} = x_{ij}\beta + b_i + \varepsilon_{ij},$$

with independent $b_i \sim \mathcal{N}(0, \tau^2)$ and $\varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2)$, with true values $\beta = 2$, $\tau^2 = 2$, $\sigma^2 = 1$.

To induce model misspecification, we contaminate the errors by replacing ε_{ij} with $\varepsilon_{ij} + \xi_{ij}$ with probability 0.1, where $\xi_{ij} \sim \mathcal{N}(0, 10^2)$ independently. The loss M_n is the Huber objective in Section 5.1 with tuning constant $c = 1$ and effective scale $s_n = n$. For the ridge penalty $\rho(\beta) = (\beta - \mu)^2/2$ we set $\mu = 0$ and $\lambda = 0.5$. The corresponding penalized population equation (1) has solution β^λ , which we treat as the target.

Because a closed-form expression for β^λ is not available, we approximate it numerically by computing the penalized estimating equation estimator on large simulated data sets with $G = 5,000$ groups and averaging the resulting estimates over 1,000 replications. All coverage probabilities and biases below are evaluated with respect to this pseudo-true value.

For each learning rate η we compare three procedures. First, as a frequentist benchmark, we compute the Huber M -estimator with ridge penalty by minimizing $M_n(\beta) + \lambda s_n \rho(\beta)$ and form Wald intervals based on the sandwich variance estimator $\hat{V}_{\text{target}} = \hat{J}_\lambda^{-1} \hat{K} \hat{J}_\lambda^{-1}$ in Section 5.1. Second, we compute the generalized Bayes posterior based on M_n and the n -dependent Gaussian prior corresponding to the ridge penalty, using the Gibbs sampler described in Section E.1, and form equal-tailed 95% credible intervals from the posterior draws of β . Third, we apply the location-scale calibration of Section 3: from the same posterior draws we estimate the working curvature \hat{H}_0 via $s_n \hat{\Sigma}_{\text{post}}$, estimate the target sandwich covariance \hat{V}_{target} by plugging in the MAP estimator and its empirical score covariance, construct the calibration operator $\hat{\Omega}$ as in (4), and transform the posterior draws to obtain calibrated credible intervals for β .

We vary the learning rate on the grid $\eta \in \{10^a : a \in [\log_{10} 0.01, \log_{10} 100]\}$ using 20 equally spaced points on the log-scale. For each value of η , we generate 200 independent data sets and run the Gibbs sampler for 1,000 iterations, discarding the first 500 as burn-in. From these replications we record, for each method and each η , the empirical coverage probability of the nominal 95% intervals for the pseudo-true value, the mean interval width, the mean bias of the point estimator for the pseudo-true value, and the standard deviation of that bias across replications.