

Exact Bayesian Geostatistics Using Predictive Stacking

BY LU ZHANG

Division of Biostatistics, Department of Population and Public Health Sciences, University of Southern California, USA

lzhang63@usc.edu

WENPIN TANG

Department of Industrial Engineering and Operations Research, Columbia University, USA

wt2319@columbia.edu

AND SUDIPTO BANERJEE

Department of Biostatistics, University of California Los Angeles, Los Angeles, USA

sudipto@ucla.edu

SUMMARY

We develop Bayesian predictive stacking for geostatistical models. Our approach builds an augmented Bayesian linear regression framework that subsumes the realisations of the spatial random field and delivers exact analytically tractable posterior inference conditional upon certain spatial process parameters. We subsequently combine such inference by stacking these individual models across the range of values of the hyper-parameters. We devise stacking of means and posterior densities in a manner that is computationally efficient without the need of iterative algorithms such as Markov chain Monte Carlo (MCMC) and can exploit the benefits of parallel computations. We offer novel theoretical insights into the resulting inference within an infill asymptotic paradigm and through empirical results showing that stacked inference is comparable to full sampling-based Bayesian inference at a significantly lower computational cost.

Some key words: Bayesian inference; conjugate spatial models; Gaussian processes; Geostatistics; stacking.

1. INTRODUCTION

Geostatistics (Cressie, 1993; Chilés & Delfiner, 1999; Zimmerman & Stein, 2010; Banerjee, 2019) refers to the study of a spatially distributed variable of interest, which in theory is defined at every point over a bounded study region of interest. Customary geostatistical modelling proceeds from a latent stochastic process over space that specifies the probability law for the point-referenced measurements on the variable as a partial realisation of the process over a finite set of locations. Inference is sought for the underlying spatial process posited to be generating the data, which is subsequently used for spatial predictions over the domain (“kriging” Stein, 1999) to better understand the scientific phenomenon under study. The spatial process, is often assumed to be stationary and empirically estimated from measurements at sampled locations using the “variogram” and characterised by parameters representing the sill, the nugget, the range and, possibly, the smoothness of the process. We collectively refer to these as process parameters.

Formal likelihood-based inference for this process is, however, thwarted by the absence of classical consistent estimators of the process parameters in a customarily preferred infill asymptotic paradigm (see, e.g., Stein, 1999; Zhang, 2004; Zhang & Zimmerman, 2005; Kaufman & Shaby, 2013; Tang et al., 2021). Bayesian inference for geostatistical processes (Handcock & Stein, 1993; Berger et al., 2001; Banerjee et al., 2014; Li et al., 2023), while not relying upon asymptotic inference, is also not entirely straightforward. Specifically, irrespective of how many spatial locations yield measurements, the likelihood does not mitigate the effect of the prior distributions on the inference. This is undesirable since proper prior elicitation for spatial process parameters is challenging. Objective priors for spatial process models have also been pursued, but interpreting such information in practice and their implications in scientific contexts are not uncontroversial. The related question of how effectively (or poorly) the realised data can identify these process parameters (in an exact sense from finite samples) has also been the subject of commentary (see, e.g., Hodges, 2013; Bose et al., 2018; De Oliveira & Han, 2022).

The aforementioned, rather substantial, literature has focused extensively on inference for geostatistical parameters and their sensitivity to modelling assumptions from diverse perspectives. It is, therefore, not unreasonable to pursue methods that will yield robust inference for the spatial process and for spatial predictions of the outcome at arbitrary points (“kriging”) while circumventing inference on the weakly identified parameters. Instead of seeking families of prior distributions for such parameters, recent efforts at computationally efficient algorithms for geostatistical models have proposed multi-fold cross-validation methods (Finley et al., 2019) to fix the values of weakly identified parameters. However, the metrics for ascertaining optimal values of such parameters are somewhat arbitrary and may not offer robust inference. Instead, our current contribution develops Bayesian predictive stacking of geostatistical models.

Stacking is a model averaging procedure for generating predictions (Wolpert, 1992; Breiman, 1996a; Clyde & Iversen, 2013). Stacking methods and algorithms in diverse data analytic applications are rapidly evolving and a comprehensive review is beyond the scope of this manuscript. Significant developments of stacking methodology in Bayesian analysis have been achieved in recent years (Le & Clarke, 2017; Yao et al., 2018, 2020, 2021), but, to the best of our knowledge, developments in the context of spatial data analysis are lacking. Stacking can be regarded as an alternative to Bayesian model averaging (Madigan et al., 1996; Hoeting et al., 1999). Assume that there are G candidate models $\mathcal{M} = \{M_1, \dots, M_G\}$. Following Bernardo & Smith (1994), Bayesian model comparison customarily considers three settings: (i) \mathcal{M} -closed where a true data generating model exists and is included in \mathcal{M} ; (ii) \mathcal{M} -complete where a true model exists but is not included in \mathcal{M} ; and (iii) \mathcal{M} -open where we do not assume the existence of a true data generating model. While Bayesian model averaging is often the preferred solution in the first setting, predictive stacking has advantages in the \mathcal{M} -complete and \mathcal{M} -open settings. Given the complex nature of spatial dependence, the \mathcal{M} -closed assumption is rarely tenable in practical geostatistics and we prefer stacking to model averaging.

Much of the theoretical properties of conventional stacking rely upon properties of exchangeable models that are not enjoyed by geostatistical models. Inferential behaviour of posterior distributions is, therefore, studied within an infill asymptotic paradigm. We absorb spatial process realisations into a convenient augmented Bayesian linear regression framework (Section 2). We offer some novel theoretical insights into the consistency of the posterior and posterior predictive distributions. Each model $M_g \in \mathcal{M}$ is indexed according to fixed values of certain spatial covariance parameters so that the exact posterior distribution is analytically tractable. Section 3 develops geostatistical predictive stacking by combining the exact inference across the G models in \mathcal{M} . We develop (i) stacking of means, which combines posterior predictive means, $E_g(y(s_0) | y_{obs}; M_g)$, from each of the G models; and (ii) stacking of posterior predictive den-

sities (Yao et al., 2018), which combines posterior predictive densities $p(y(s_0) | y_{obs}; M_g)$ for $g = 1; 2; \dots; G$, where $y(s_0)$ is the random variable denoting model predictions of the outcome at an arbitrary point s_0 and y_{obs} denotes the observed data on the outcome. We obtain exact Bayesian inference by exploiting exact distribution theory and avoiding iterative algorithms such as Markov chain Monte Carlo (MCMC). These methods are evaluated theoretically as well as empirically through simulation experiments (Section 4) demonstrating that stacked inference is comparable to full Bayesian inference using MCMC at significantly less computational expense. An illustrative data analysis is presented in Section 5 to further corroborate results seen in the simulations. We conclude with some discussions and pointers to future work in Section 6.

2. BAYESIAN HIERARCHICAL SPATIAL PROCESS MODELS

2.1. Conjugate Bayesian spatial model

We explore a regression model for a spatially indexed outcome $y(s)$ at a location s in a bounded region $D \subset \mathbb{R}^d$,

$$y(s) = x(s)^T \beta + z(s) + \epsilon(s); \quad (2.1)$$

where $x(s)$ is a $p \times 1$ vector of spatially referenced predictors, β is a $p \times 1$ vector of slopes measuring the trend, $z(s) \sim \text{GP}(0; \Sigma(\cdot; \cdot))$ is a zero-centred spatial Gaussian process on \mathbb{R}^d with spatial correlation function $R(\cdot; \cdot)$ depending on spatial range parameter ρ , and Σ is a scale (spatial variance) parameter. The white noise process $\epsilon(s) \sim N(0; \sigma^2)$ with variance σ^2 captures measurement error.

Let $\mathcal{S} = \{s_1; \dots; s_n\}$ be a set of n spatial locations, each $s_i \in D$, yielding measurements $y = (y(s_1); \dots; y(s_n))^T$ with known values of predictors at these locations collected in the $n \times p$ matrix $X = (x(s_1); \dots; x(s_n))^T$. We also define the finite-dimensional realisation of the spatial process as $z = (z(s_1); \dots; z(s_n))^T$, and let $R(\cdot) = (R(s_i; s_j))_{1 \leq i, j \leq n}$ be the $n \times n$ spatial correlation matrix constructed from the correlation function. A customary Bayesian hierarchical model is constructed as

$$\begin{aligned} y | z; \rho; \sigma^2 &\sim N(X\beta + z; \sigma^2 I_n); & z | \rho &\sim N(0; \Sigma(\rho)) \\ \beta &\sim N(\mu; \Sigma); & \rho &\sim \text{IG}(a; b); \end{aligned} \quad (2.2)$$

where we fix the spatial correlation parameters ρ and the noise-to-spatial variance ratio $\sigma^2 := \frac{\sigma^2}{\Sigma}$, and μ, Σ, a , and b are fixed hyper-parameters specifying the prior distributions for β and ρ . This ensures closed-form conjugate marginal posterior and posterior predictive distributions (Finley et al., 2019; Banerjee, 2020). In order to harness familiar results from conjugate Bayesian linear regression, we cast the spatial model in (2.2) into an augmented linear system:

$$\begin{aligned} \begin{pmatrix} y \\ z \end{pmatrix} &= \begin{pmatrix} X & I_n \\ 0 & R \end{pmatrix} \begin{pmatrix} \beta \\ \rho \end{pmatrix} + \begin{pmatrix} \epsilon \\ z \end{pmatrix} \\ &= \begin{pmatrix} X & I_n \\ 0 & R \end{pmatrix} \begin{pmatrix} \beta \\ \rho \end{pmatrix} + \begin{pmatrix} \epsilon \\ z \end{pmatrix} \end{aligned} \quad (2.3)$$

where $\begin{pmatrix} \epsilon \\ z \end{pmatrix} \sim N(0; \Sigma_y)$, $\Sigma_y = \begin{pmatrix} \Sigma & 0 \\ 0 & R \end{pmatrix}$ and $\rho \sim \text{IG}(ja; b)$.

LEMMA 1. The posterior distribution of $(z; y)$ from (2.2) is

$$p(z; y) = \frac{IG(z; a; b)}{p(z; y)} \frac{N(\hat{z}; M)}{p(\hat{z}; y)}; \quad (2.4)$$

where $a = a + n/2$, $b = b + \frac{1}{2}(y - X\hat{z})^\top V^{-1}(y - X\hat{z})$, $M^{-1} = X^\top V^{-1}X$ and $\hat{z} = M^{-1}X^\top V^{-1}y$. The posterior distribution $p(z; y)$ is a multivariate Student's t with degrees of freedom $2a$, location \hat{z} and scale matrix $(b/a)M$.

Proof. The proof follows from a straightforward adaptation of familiar results from the Normal-Gamma family of distributions (see (Murphy, 2015) and Section 3.1 in (Banerjee, 2020)) to the setting in (2.3).

Furthermore, let $\tilde{s} = \{s_1, \dots, s_m\}$ be a set of m unknown points in D , z and y be the $m \times 1$ vectors with elements $z(s_i)$ and $y(s_i)$ for $i = 1, 2, \dots, m$. Let $X = (x(s_1), \dots, x(s_m))^\top$ be the $m \times p$ matrix that carries the values of predictors at \tilde{s} and let $J(\tilde{s}) = (R(s; s^h))_{s \in \tilde{s}, s^h \in \tilde{s}}$. Then, spatial predictive inference follows from the posterior distribution

$$p(z; y) = \int p(y|z; \tilde{s})p(z|z; \tilde{s})p(\tilde{s}; y) d\tilde{s}; \quad (2.5)$$

which is again a multivariate t distribution with degrees of freedom $2a$, location \tilde{s} and scale matrix $(b/a)M$ where

$$\tilde{s} = W\hat{z}; M = WMW^\top + M_2; M_1 = I_m - J^\top(\tilde{s})R^{-1}(\tilde{s})J(\tilde{s});$$

$$W = \begin{pmatrix} 0 & J^\top(\tilde{s})R^{-1}(\tilde{s}) \\ X & J^\top(\tilde{s})R^{-1}(\tilde{s}) \end{pmatrix}; M_2^{-1} = \begin{pmatrix} \frac{1}{2}I_m + M_1^{-1} & \frac{1}{2}I_m \\ \frac{1}{2}I_m & \frac{1}{2}I_m \end{pmatrix}.$$

Specifically, the predictive distributions $p(z(s_0)|y)$ and $p(y(s_0)|y)$ are also available in analytic form as non-central t distributions for any single point $s_0 \in D$. Bayesian inference can proceed from exact posterior samples obtained from (2.4) as follows. We first draw values of $z \sim IG(a; b)$ followed by a single draw of $\hat{z} \sim N(\hat{z}; M)$ for each drawn value of z . This yields samples $\tilde{s}; y$ from (2.4). Predictive inference for the latent process $z(s_0)$ and the outcome $y(s_0)$ is obtained by sampling from (2.5) by drawing a value of $z \sim N(z; V_{zj})$ with $z(\tilde{s}) := J^\top(\tilde{s})R^{-1}(\tilde{s})z$ and $V_{zj} := R^{-1}(\tilde{s}) - J^\top(\tilde{s})R^{-1}(\tilde{s})J^{-1}(\tilde{s})$ for each value of $\tilde{s}; y$ drawn above (see Section 3.4 in Banerjee (2020)), then drawing a value of $y \sim N(X\hat{z} + z; V_{yj})$ for each drawn value of \tilde{s}, z and y .

This tractability is only possible if the range decay α and the noise-to-spatial variance ratio σ^2 are fixed. While the data can inform about these parameters, they are inconsistently estimable (Zhang, 2004) often resulting in poorer convergence. Alternate approaches using K -fold cross-validation have been explored with limited success (Finley et al., 2019). Therefore, our approach will conduct the exact inference using (2.4) and (2.5) and stack the inference over the different fixed values of $\tilde{s}; y$. Next, we investigate the concentration of the posteriors and the effect of incorrectly specified hyperparameters on the posterior inference.

2.2. Posterior inference for spatial process models

Let L be the Cholesky decomposition of V such that $V = L L^T$, and L a non-singular square matrix such that $R(\cdot)^{-1} = L^T L^{-1}$. The linear model (2.3) can be rewritten as

$$\begin{aligned} \begin{pmatrix} y \\ X \end{pmatrix} &= \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \begin{pmatrix} X \\ 0 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix}^{-1} \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \begin{pmatrix} y \\ X \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix}^{-1} \begin{pmatrix} 1 \\ 0 \end{pmatrix} \\ &= \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \begin{pmatrix} X \\ 0 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix}^{-1} \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \begin{pmatrix} y \\ X \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix}^{-1} \begin{pmatrix} 1 \\ 0 \end{pmatrix} \end{aligned} \quad (2.6)$$

where $y \sim N(0; \sigma^2 I_{2n+p})$. To explore posterior concentrations within an in-fill paradigm, where we assume a fixed size of our spatial domain D and an increasing number of spatial locations inside D , we will use the concept of equivalence of probability measures.

Assumption 1 (Equivalence). Let P_0 be the probability distribution of the process $y(s)$ defined by the model (2.1) with true parameter values $f_0; \sigma_0^2; \rho_0; \sigma_0^2 g$. For each $\epsilon > 0$, there is $\eta > 0$ such that the probability distribution of the process $y(s)$ defined by the model (2.1) with parameter values $f_0; \sigma_0^2; \rho_0; \sigma_0^2 g$ is equivalent to P_0 .

Note that the above Assumption holds when the latent process $z(s)$ follows a Matérn model in dimension $d \geq 1; 2; 3g$, which we provide more details in Section 2.3. In this case, the probability distribution of the process $y(s)$ defined by the model (2.1) with parameters $\sigma_0; \frac{\sigma_0^2}{2}; \rho_0; \sigma_0^2$, with ν a smoothness parameter of the Matérn model, is equivalent to P_0 (see e.g. Tang et al. (2021, Section 2.1)). The following theorem explores the posterior (in)consistency of parameter inference.

THEOREM 1 (POSTERIOR INFERENCE (IN)CONSISTENCY). Let P_0 be the probability measure of the model (2.1) with parameter values $f_0; \sigma_0^2; \rho_0; \sigma_0^2 g$, and let Assumption 1 hold. Let $H = X_y^T (X_y^T X_y)^{-1} X_y^T$ be the $(2n+p) \times (2n+p)$ orthogonal projector onto the column space of X_y and let $H = \begin{pmatrix} H_{11} & H_{12} \\ H_{12}^T & H_{22} \end{pmatrix}$ be a 2×2 partition of H so that H_{22} is the lower right $n \times n$ block formed by rows and columns indexed from $n+p+1$ to $2n+p$. Assume that $\text{Tr}(H_{22}) = n!$ as $n \rightarrow \infty$. Then under P_0 ,

$$\lim_{n \rightarrow \infty} p(\sigma^2 | y(\cdot)) = \delta_{\sigma_0^2}; \quad (2.7)$$

where $y(\cdot) = (y(s_1); y(s_2); \dots; y(s_n))^T$, $\sigma^2 := \frac{\sigma_0^2}{2} + \eta^2(1 - \rho_0)$ and $\delta_{\sigma_0^2}$ denotes the Dirac measure at σ_0^2 .

Proof. See Section A in the Appendix.

Theorem 1 suggests that the posterior distribution of the scale parameter σ^2 in the spatial process does not necessarily concentrate to the true generating value. We also provide conditions for the posterior consistency of σ^2 in the following corollary.

COROLLARY 1. Let P_0 be the probability distribution of the model (2.1) with parameter values $f_0; \sigma_0^2; \rho_0; \sigma_0^2 g$ under Assumption 1. Define $U = (X_y^T X_y)^{-1} = \begin{pmatrix} U_{11} & U_{12} \\ U_{12}^T & U_{22} \end{pmatrix}$, where U_{11} is $p \times p$, $B = \begin{pmatrix} X^T X & X^T \\ X & I_n \end{pmatrix}^{-1} U = \begin{pmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{pmatrix}$, where B_{11} is $p \times p$, $C = U_{11} V^{-1} U_{11}$, and $D = U_{12} R(\cdot)^{-1} U_{12}^T$ and assume the following additional conditions:

$$\text{Tr}(U_{11}); \text{Tr}(B_{11}); \text{Tr}(C); \text{Tr}(D) \rightarrow \infty \text{ as } n \rightarrow \infty; \quad (2.8)$$

Then, $\lim_{n \rightarrow \infty} \rho(\hat{y}(s_0) | y(s_0)) = \rho_0$ under \mathbb{P}_0 , where $\hat{y}(s_0)$ and $y(s_0)$ are defined as in Theorem 1.

Proof. See Section B in the Appendix.

Figure 1 summarises some numerical experiments to empirically explore $\text{Tr}(H_{22})=n^{-1}$ and (2.8) for some typical examples. The study domain D is $[0; 1]^2$, and locations in \mathcal{S} are chosen uniformly on D . We generate data using the Matérn covariogram for $z(s)$ (see (2.12) in Section 2.3). We consider two types of predictors $x(s)$. For the first type, titled ‘‘with intercept’’, $x(s)$ consists of an constant 1 for intercept and a predictor generated by a standard normal. For the second type, labelled ‘‘without intercept’’, $x(s)$ is composed of two predictors sampled from a standard normal. We explore the trends of the target quantities with different hyper-parameter values in the covariogram of $z(s)$ and different types of $x(s)$ as sample size increases. Figure 1 reveals clearly that $\text{Tr}(H_{22})=n^{-1}$ increases as sample size increases for all examples. Since $\text{Tr}(H_{22})=n^{-1}$ is bounded above by 1, the condition $\text{Tr}(H_{22})=n^{-1}$ for some constant is likely to hold for all examples. On the other hand, condition (2.8) is observed to be more sensitive to the choice of $x(s)$. It does not seem to hold for the cases when \mathcal{S} contains an intercept. Examples on the one-dimensional $[0; 1]$ exhibit similar behaviour. Moreover, we prove in Theorem 3 that in the special case where there is no trend term $x(s)^\top$, and the latent process $z(s)$ follows a Matérn model in dimension $d \geq 1; 2; 3g$, the posterior distribution of $\hat{z}(s_0)$ converges to the degenerate distribution at $\hat{z}(s_0) = z(s_0)$. Hence, in general, this result will be inconsistent unless the noise-to-spatial variance ratio σ^2 is fixed to be $\frac{\sigma^2}{\delta^2} = \frac{\sigma^2}{\delta_0^2} = \frac{\sigma^2}{\delta_0^2}$.

Next, we consider Bayesian posterior predictive inference at a new location $s_0 \in D$. Let $Z_n(s_0)$ be a random variable distributed as $p(z(s_0) | y(s_0))$ and $Y_n(s_0)$ be distributed as $p(y(s_0) | y(s_0))$. We study the prediction error $E_0(Z_n(s_0) - z(s_0))^2$ for the latent process, and $E_0(Y_n(s_0) - y(s_0))^2$ for the response variable. Let X_y and y_y be as in (2.6) and U is as in Corollary 1. The following theorem quantifies these posterior prediction errors.

THEOREM 2 (POSTERIOR PREDICTIVE CONSISTENCY). *Consider $s_0 \in D$. For any given $\delta > 0$, denote $\text{COV}(z; z(s_0) | y(s_0))$ and $R^{-1}(s_0)$ by $J_{\cdot; n}^{-2}$ and $R_{\cdot; n}$, respectively. Let $Z_n(s_0)$ have the density $p(z(s_0) | y(s_0))$, and $Y_n(s_0)$ have the density $p(y(s_0) | y(s_0))$, and let*

$$F_n = J_{\cdot; n}^\top R_{\cdot; n}^{-1} (U_n X_y^\top y_y)_{[p+1; p+n]} \quad \text{and}$$

$$G_n = 1 - J_{\cdot; n}^\top R_{\cdot; n}^{-1} (R_{\cdot; n} + U_{n[p+1; p+n; p+1; p+n]}) R_{\cdot; n}^{-1} J_{\cdot; n}$$

where we use the suffix n to highlight the dependency of \cdot_n , and we remove the suffix n of X_y and y_y when there is no confusion. Under the assumptions in Theorem 1, we have

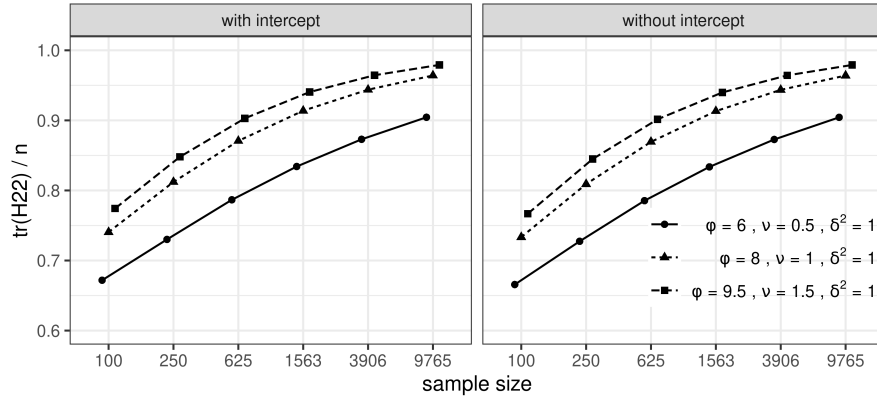
$$E_0(Z_n(s_0) - z(s_0))^2 = E_{1; n} + E_{2; n} + o(1); \quad (2.9)$$

where E_0 on the left side supports $z(s_0)$, $y(s_0)$ and the external randomisation of $p(z(s_0) | y(s_0))$, and $E_{1; n} = E_0(Z_n(s_0) - F_n)^2$ and $E_{2; n} = \sigma^2 G_n$. Moreover, if $E_{1; n}; E_{2; n} \rightarrow 0$ as $n \rightarrow \infty$, then posterior inference for the latent process is consistent in the sense that $E_0(Z_n(s_0) - z(s_0))^2 \rightarrow 0$ as $n \rightarrow \infty$, while posterior predictive inference for $y(s_0)$ satisfies

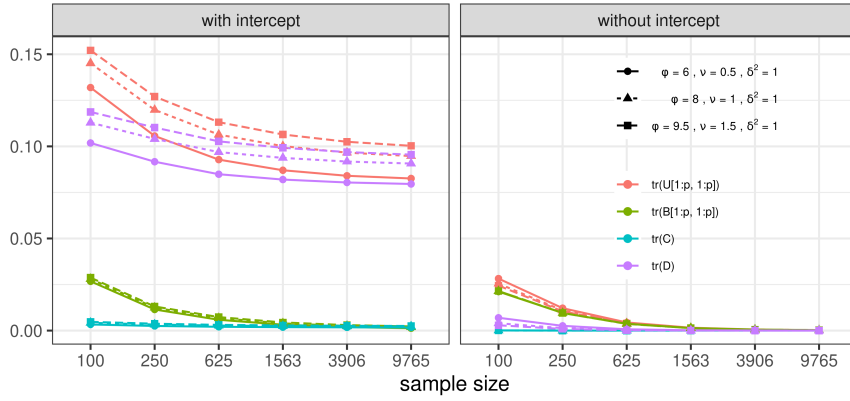
$$E_0(Y_n(s_0) - y(s_0))^2 \rightarrow \frac{\sigma^2}{\delta^2} + \sigma^2 \quad \text{as } n \rightarrow \infty; \quad (2.10)$$

Proof. See Section C in the Appendix.

In the decomposition (2.9), the term $E_{1; n}$ arises in the deviation from the posterior mean, while the term $E_{2; n}$ is from the posterior uncertainty. The conditions $E_{1; n}; E_{2; n} \rightarrow 0$ is analytically intractable in the general case. If, however, we detrend the outcome so that $x(s)^\top = 0$, then $E_{1; n}$ and $E_{2; n}$ are simplified to inherit a rich structure under the Matérn covariance model for the latent process $z(s)$. In this case, we provide evidence to support the posterior predictive consistency for the latent process in the next subsection.



(a)



(b)

Fig. 1: The top row (a) provides plots of $\text{Tr}(H_{22})=n$ when $x(s)$ consists of an intercept only (left), and with additional covariates (right) for different parameter values. The bottom row (b) plots the quantities $\text{Tr}(U_{11})$, $\text{Tr}(B_{11})$, $\text{Tr}(C)$, and $\text{Tr}(D)$ with an intercept only (left) and with covariates (right)

2.3. Conjugate Bayesian model without trend

We now consider a zero-centred spatial process for $y(s)$ at a location $s \in D$ modelled as

$$y(s) = z(s) + \epsilon(s); \tag{2.11}$$

Let $D \subset \mathbb{R}^d$ be a compact domain with $d \geq 1; 2 \leq \nu < 3/2$. We also assume that the correlation function $R(\cdot; \cdot)$ is specified by the isotropic Matérn covariogram

$$R(s; s^j) := \frac{(\nu - 1/2)!}{(\nu - 1/2)!} \frac{K(\nu - 1/2)(\sqrt{2} |s - s^j|)}{(\sqrt{2} |s - s^j|)^{\nu - 1/2}}; \tag{2.12}$$

where $\nu > 0$ is a smoothness parameter, $\Gamma(\cdot)$ is the Gamma function, and $K(\cdot)$ is the modified Bessel function of the second kind of order $\nu - 1/2$ (Abramowitz & Stegun, 1965, Section 10). It is

known that the spectral density of the (isotropic) Matérn model without nugget is

$$f_{\text{Matérn}}(u) = C \frac{u^{2\nu}}{(u^2 + \nu^2)^{\nu + d/2}} \quad \text{for some } C > 0: \quad (2.13)$$

Fix the smoothness parameter $\nu > 0$, call the process (2.11) the Matérn model with parameter values $f^2; \sigma^2; g$. Since there is no x^\top term, the conjugate Bayesian model (2.2) simplifies to

$$y \mid z \sim N(0; \sigma^2(R(\nu) + \nu^2 I_n)); \quad z \sim \text{IG}(a; b); \quad (2.14)$$

and the corresponding augmented linear regression becomes

$$\begin{pmatrix} y \\ 0 \\ \{z\} \end{pmatrix} = \begin{pmatrix} I_n \\ I_n \\ -\{z\} \end{pmatrix} \begin{pmatrix} z \\ x \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \\ -\{z\} \end{pmatrix}; \quad (2.15)$$

where $N(0; \sigma^2 V_y)$ and $V_y = \begin{pmatrix} I_n & 0 \\ 0 & R(\nu) \end{pmatrix}$.

The following theorem shows the posterior inconsistency of the scale σ^2 under the conjugate model (2.14).

THEOREM 3 (POSTERIOR INFERENCE FOR THE MATÉRN MODEL). *Assume that the location set $S = \{s_1; \dots; s_n\}$ satisfies*

$$\max_{s \in D} \min_{i \in [n]} |s - s_i| \geq n^{-\frac{1}{d}}; \quad (2.16)$$

Let P_0 be the probability distribution of the Matérn model (2.11) with the true parameter values $(\frac{\sigma^2}{\sigma_0}; \sigma_0; \frac{\sigma^2}{\sigma_0})$. Under P_0 ,

$$\lim_{n \rightarrow \infty} p(\sigma^2 \mid y(\cdot)) = \frac{\sigma^2}{\sigma_0}; \quad (2.17)$$

Consequently, $\lim_{n \rightarrow \infty} p(\sigma^2 \mid y(\cdot)) = \frac{\sigma^2}{\sigma_0}$.

Proof of Theorem 3. See Section D in the Appendix.

Note that the posterior inference of σ^2 is independent of the range decay ν chosen in the conjugate Bayesian model (2.14). The scale σ^2 is posterior inconsistent unless the noise-to-spatial variance ratio $\frac{\sigma^2}{\sigma_0} = \frac{\sigma_0}{\sigma_0}$, while the nugget σ^2 is posterior consistent.

Posterior prediction for $z(\cdot)$ and $y(\cdot)$: Recall the decomposition (2.9) for the posterior prediction error for the latent process $z(s)$. The following proposition simplifies the two sources of error $E_{1;n}$ and $E_{2;n}$, when the outcome $y(s)$ follows a Matérn model in the presence of a nugget.

THEOREM 4 (POSTERIOR PREDICTIVE CONSISTENCY FOR THE MATÉRN MODEL). *Let $s_0 \in D$. Then, we have the decomposition (2.9), where $E_{1;n}$ is the prediction error of the best linear predictor for a Matérn model with parameters $f^2; \sigma^2; g$ satisfying $\frac{\sigma^2}{\sigma_0} = \frac{\sigma_0}{\sigma_0}$, and*

$$E_{2;n} := \frac{\sigma^2}{2} \mathbf{1} + J_{\cdot;n}^\top (I_n + \sigma^2 R(\nu))^{-1} I_n R(\nu) \mathbf{1}_{\cdot;n}; \quad (2.18)$$

Moreover, if $E_{1;n}; E_{2;n} \rightarrow 0$ as $n \rightarrow \infty$, then the latent process $z(s)$ is posterior predictive consistent in the sense that $E_0(Z_n(s_0) - z(s_0))^2 \rightarrow 0$ as $n \rightarrow \infty$ and, hence, $E_0(Y_n(s_0) - y(s_0))^2 \rightarrow 2 \frac{\sigma^2}{\sigma_0}$ as $n \rightarrow \infty$.

Proof. See Section E in the Appendix.

Theorem 4 reveals that the posterior mean of the conjugate Bayesian model (2.14) is identified as the best linear predictor of any Matérn model with parameters $f^2; \sigma^2; g$ provided $\frac{\sigma^2}{\sigma_0} = \frac{\sigma_0}{\sigma_0}$. This observation connects the Bayesian modelling to a frequentist approach in

that the deviation error $E_{1;n}$ is viewed as the prediction error of the best linear predictor of a Matérn model in the presence of a nugget. Next we provide evidence to support the posterior predictive consistency for the latent process in the sense that $E_{1;n}; E_{2;n} \rightarrow 0$ as $n \rightarrow \infty$.

The deviation error $E_{1;n}$: It is expected that the prediction error of the best linear predictor of a Matérn model in the presence of a nugget tends to 0 as long as the fill distance condition (2.16) holds. However, it is hard to prove this statement in the general case. To provide some ideas, we consider a particular one-dimensional example.

We assume without loss of generality that $D = [0; 1]$, and $s_i = fi/n; 0 \leq i \leq n$ so the fill distance condition (2.16) is satisfied. Also let $s_{i-1} = fi/n; 1 \leq i \leq n$ be the infinite evenly spaced grid. Define

$$\hat{z} := E(z(0) | y(s); s \in \{0, \dots, n\}) \quad (\text{resp. } \hat{z}_1 := E(z(0) | y(s); s \in \{0, \dots, n\}));$$

be the best linear predictor of $z(0)$ based on the observations $y(s)$ on $\{0, \dots, n\}$ (resp. $\{0, \dots, n\}$). Note that \hat{z} (resp. \hat{z}_1) may well be computed using misspecified parameter values $f; \sigma^2$. Let

$$e := E_0(z(0) - \hat{z})^2 \quad (\text{resp. } e_1 := E_0(z(0) - \hat{z}_1)^2); \tag{2.19}$$

be the prediction error of the best linear predictor based on the observations on the finite grid $\{0, \dots, n\}$ (resp. the infinite grid $\{0, \dots, n\}$). We make the following assumption which relates e to e_1 as $n \rightarrow \infty$.

Assumption 2 (Infinite approximation). The difference $e - e_1 \rightarrow 0$ as $n \rightarrow \infty$.

Basically, this assumption suggests that as the grid becomes finer and finer, observations outside any fixed bounded interval have little impact on the prediction of $z(0)$. However, it is not easy to prove such a statement which is conjectured in (Stein, 1999, page 97). The plan is to prove $e_1 \rightarrow 0$, and by Assumption 2 it implies that $e \rightarrow 0$ as $n \rightarrow \infty$. This program is recorded in the following proposition.

PROPOSITION 1 (PREDICTION ERROR FOR THE MATÉRN MODEL). *Let $D = [0; 1]$ and $s_i = fi/n; 0 \leq i \leq n$. Under Assumption 2, we have*

$$e \rightarrow 0 \quad \text{as } n \rightarrow \infty; \tag{2.20}$$

where e is the prediction error of the best linear predictor defined as in (2.19).

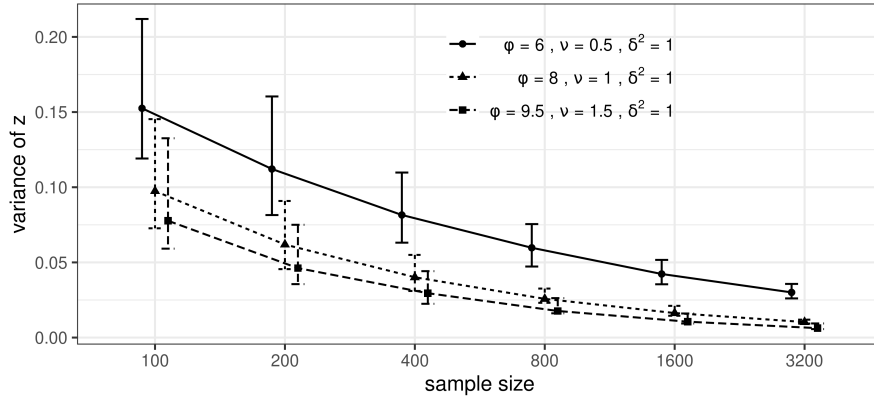
Proof. See Section F in the Appendix.

The posterior variance $E_{2;n}$: The error term $E_{2;n}$ defined by (2.18) is also analytically intractable. Here, we provide a numerical study to investigate the behaviour of $E_{2;n}$ when $n \rightarrow \infty$ in the general case. We first generate the location set \mathcal{S} by uniformly sampling n locations in $[0; 1]$ or $[0; 1]^2$, then we compute $E_{2;n}$ for every location in \mathcal{S} with $\sigma^2 = 1; \nu = 1$ and different values of d and α . We expand the location set \mathcal{S} sequentially to sets with larger sample sizes by adding locations that are uniformly sampled in the study domain. For each expanded set \mathcal{S} , we recompute the $E_{2;n}$ for all locations in \mathcal{S} . Figure 2 plots the median, the 2.5th and 97.5th percentiles of $E_{2;n}$ for different sample sizes. The values of $E_{2;n}$ for points in a fixed domain, shown in Figure 2, decrease rapidly as the sample size increases, although the rate diminishes when d increases from 1 to 2. This suggests that the decreasing rate is related to dimension.

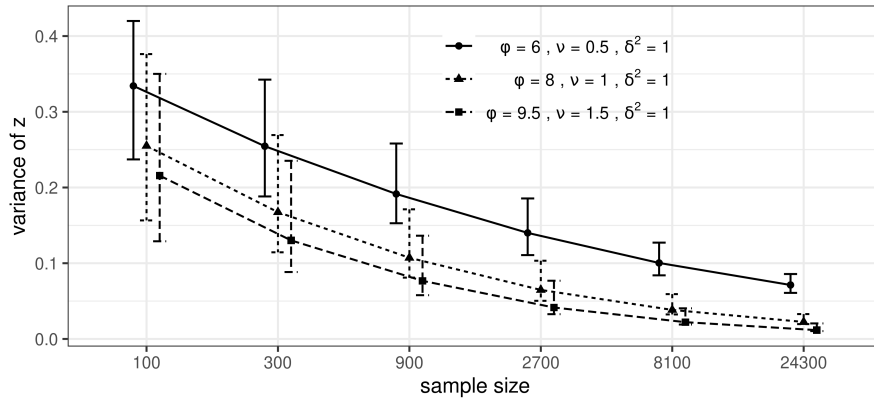
3. STACKING ALGORITHMS FOR SPATIAL MODELS

3.1. Stacking of means

\mathbb{R} The Bayes predictor for $y(s_0)$ under model M_g , for each $g = 1; \dots; G$, is $E_g(y(s_0) | y) = \int y(s_0) p(y(s_0) | y; M_g) dy(s_0)$, where $E_g(y(s_0) | y)$ is the expectation of the posterior predictive



(a)



(b)

Fig. 2: The median of $E_{2;n}$ for locations uniformly sampled on $[0; 1]$ (a) and $[0; 1]^2$ (b). The error bars indicate the 97.5th and 2.5th percentiles. The sample size n ranges from 100 to 3,200 and from 100 to 24,300 for the experiments on $[0; 1]$ and $[0; 1]^2$, respectively.

mean computed from M_g . Stacking methods will combine the G Bayes predictors as a weighted average

$$\sum_{g=1}^G w_g E_g(y(s_0) | y); \tag{3.1}$$

where w_1, \dots, w_G are the weights for combination. Define the leave-one-out (LOO) Bayes predictor for $y(s_i)$ under model M_g as

$$\hat{y}_g(s_i) = E_g(y(s_i) | y_{-i}; M_g) = \int y(s_i) p(y(s_i) | y_{-i}; M_g) p(y_{-i}; M_g) dy_{-i}$$

where y_{-i} is the data without the i th observation. Stacking determines the optimal weights as

$$\arg \min_w \sum_{i=1}^n \sum_{g=1}^G w_g \hat{y}_g(s_i)^2; \tag{3.2}$$

The following proposition studies the posterior prediction error using the stacking predictor, where we only require that the weights be bounded.

PROPOSITION 2. Let $s_0 \in D$, and $w := (w_1, \dots, w_G)$ be the stacking weights defined by (3.2) such that $|w_g|$ is bounded for each $1 \leq g \leq G$. Let the assumptions in Theorem 2 or Theorem 4 hold for each model M_g . We have

$$E_0 @ y(s_0) - \sum_{g=1}^G w_g E_g(y(s_0) | y^A) \leq \frac{1}{n} \sum_{i=1}^n E_{1,g;i} \leq \frac{1}{n} \sum_{i=1}^n \sum_{g=1}^G w_g^2 \text{ as } n \rightarrow \infty \tag{3.3}$$

Proof. See Section G in the Appendix.

The following theorem shows that the stacking predictor asymptotically minimises the posterior prediction error.

THEOREM 5. Let $s_0 \in D$, and $w := (w_1, \dots, w_G)$ be the stacking weights defined by (3.2) such that $|w_g|$ is bounded for each $1 \leq g \leq G$. Let

$$E_{1,g;i} := E_0(z(s_i) - \hat{y}_g(s_i))^2; \quad 1 \leq g \leq G \text{ and } 1 \leq i \leq n;$$

be the deviation error for the latent process $z(s)$ by leaving the i^{th} observation out under the model M_g . Assume that for each $1 \leq g \leq G$,

$$\frac{1}{n} \sum_{i=1}^n E_{1,g;i} \rightarrow 0 \text{ as } n \rightarrow \infty \tag{3.4}$$

Also, if the assumptions in Theorems 2 or 4 hold for each model M_g , then as $n \rightarrow \infty$ we obtain

$$E_0 @ y(s_0) - \sum_{g=1}^G w_g E_g(y(s_0) | y^A) = E_0 @ \frac{1}{n} \sum_{i=1}^n @ y(s_i) - \sum_{g=1}^G w_g \hat{y}_g(s_i) \rightarrow 0 \tag{3.5}$$

Proof. See Section H in the Appendix.

Equation (3.4) implies that for each model M_g the average deviation error for the latent process goes to 0 as the sampling resolution becomes finer. This condition is consistent with the fact that in the one-dimensional grid we typically have $E_{1,g;i} \sim \min(2^{-\nu}, 2^{-\frac{2\nu}{d+1}})$ (see Proposition 1); hence $\frac{1}{n} \sum_{i=1}^n E_{1,g;i} \rightarrow 0$ as $n \rightarrow \infty$.

Theorem 5 holds for candidate models with a misspecified smoothness parameter in the Matérn kernel for $d = 1$. Note that Clyde & Iversen (2013, p.487) proves a similar result, while the established theoretical results about stacking assume exchangeability which is not available in geostatistical models. The innovation in our theoretical results emerge from studying the behaviour of the posterior and predictive distributions within the infill asymptotic paradigm. The proof of Theorem 5 can be extended to the case where the LOO Bayes prediction $\hat{y}_g(s_i)$ is replaced by a much cheaper prediction based on K -fold cross-validation.

3.2. Stacking of predictive densities

Following the generalised Bayesian stacking framework established in Yao et al. (2018), we devise a second stacking algorithm for spatial analysis, which we refer to as *stacking of predictive densities*. This algorithm finds the distribution in the convex hull $C = \{ \sum_{g=1}^G w_g p(\cdot | M_g) : \sum_{g=1}^G w_g = 1; w_g \geq 0 \}$ that is optimal according to some proper scoring functions. Here $p(\cdot | M_g)$ refers to the distribution of interest under model M_g . Let $S_1^G = \{ w \in [0, 1]^G : \sum_{g=1}^G w_g = 1 \}$

and $p_t(\cdot|y)$ be the true posterior predictive distribution. Using the logarithmic score (corresponding to the KL divergence), we seek w so that

$$\max_{w \in \mathcal{S}_1^G} \frac{1}{n} \sum_{i=1}^n \log \sum_{k=1}^G w_k p(y(s_i)|y_{-i}; M_k) \quad (3.6)$$

The optimal distribution $\sum_{g=1}^G w_g p(y(s_i)|y_{-i}; M_g)$ provides a ‘‘likelihood’’ of observing $y(s_i)$ on location s_i given other data. Therefore, $\sum_{i=1}^n \sum_{g=1}^G w_g p(y(s_i)|y_{-i}; M_g)$ serves as a pseudo-likelihood that measures the performance of prediction based on the weighted average of the LOO predictors for all observed locations. Also, it provides an approximation of the expected log point-wise predictive density (ELPD)

$$\sum_{i=1}^n \int p_t(y(s_i)) \log \sum_{g=1}^G w_g p(y(s_i)|y_{-i}; M_g) dy(s_i) \quad (3.7)$$

It is worth pointing out that the weights for stacking of means are not necessarily positive, while those for stacking of predictive densities must be non-negative. Relaxing this restriction imparts greater flexibility in prediction and improves prediction accuracy for stacking of means. On the other hand, we can no longer interpret negative weights as a reflection of our prior beliefs regarding \cdot . In our subsequent experiments, we restrict the stacking weights for both algorithms to be non-negative so that the diagnostic metric MLPD (the average log pointwise predictive density, introduced in Section 4) is valid for both algorithms.

3.3. Predictive stacking in finite sample spatial analysis

A critical step in solving stacking weights is the computation of the Bayes predictor and predictive density. Computing the exact LOO Bayes predictor and predictive densities for all observed locations $\{s_1, \dots, s_n\}$ requires refitting a model n times. For a Gaussian latent variable model with the number of parameters larger than the sample size n , there are limited choices for approximating LOO predictors accurately without the onerous computation (see, e.g., Vehtari et al., 2016). Instead of the LOO cross-validation, we prefer using the much cheaper K -fold cross-validation to generate the predictions. Using K -fold cross-validation instead of LOO in stacking is first implemented in Breiman (1996b), where the simulation evinces that 10-fold cross-validation can provide more efficient predictors than LOO cross-validation. Here, we use $K = 10$ in the implementations in Section 4.

We offer the pseudo-code algorithms for stacking of means and stacking of predictive densities using the conjugate Bayesian spatial regression model in Section I. In the algorithms, we first partition the data into K -folds based on locations. Letting $X = [x(s_1) : \dots : x(s_n)]^T$ be the design matrix, we use $X[k], y[k], s_n[k]$ to denote the predictors, response and observed locations from k -th fold, respectively, and $X[-k], y[-k], s_n[-k]$ to denote the respective data not in k -th fold. The values of the prefixed hyper-parameters $f; \dots; g$ of the conjugate Bayesian spatial regression model are picked from the grid G_{all} , which is expanded over the grids of candidate values as the cartesian product $G \times \dots \times G^2$ for $f; \dots; g$. For stacking of means, we compute the expected latent process $z_{(\cdot; \cdot; 2)}^{(k)}$ on locations in fold k for each $f; \dots; g$ from G_{all} and then generate the corresponding expectation $y_{(\cdot; \cdot; 2)}^{(k)}$ to obtain the stacking weights.

For stacking of predictive densities, we compute the log point-wise predictive density, $(\log p_{(\cdot; \cdot; 2)}(s))$, of $y(s)$ at location s for locations in each fold for all candidate models to find the stacking weights. The closed form of the log point-wise predictive density is derived in Section J

of the Appendix. We also devise a Monte Carlo algorithm for stacking of predictive densities, which is expounded in Section K of the Appendix. Solving the weights for stacking of means is a quadratic programming problem. We use the `solve.QP()` function offered by the `quadprog` package in the R statistical computing environment and the solver `Mosek` to solve the weights in our R and Julia implementations, respectively. The stacking weights for stacking of predictive densities are calculated by `Mosek` in both our R and Julia implementations.

4. SIMULATION

4.1. Simulation settings

We present two simulations to examine the predictive performance of the proposed stacking algorithms. The data sets for experiments are generated by model (2.1) on locations sampled uniformly over a unit square $[0; 1]^2$, where the correlation function R is specified by the Matérn covariogram (2.12). The sample size n of the simulated data sets ranges from 200 to 900, and we randomly pick $n_h = 100$ observations for checking predictive performance. The vector $x(s)$ consists of an intercept and a single predictor generated from a standard normal distribution. The true value of the parameters for generating the data in the first simulation are $\beta = (1; 2)^\top$, $\sigma^2 = 7$, $\nu^2 = 1$, $\rho^2 = 1$ and $\kappa = 1$. In the second simulation, we alter the hyper-parameters in the Matérn covariogram to $\nu = 20$, $\rho^2 = 1$, $\nu^2 = 0.3$ and $\kappa = 0.5$.

We consider customary values of smoothness in spatial analysis, $\nu \in \{0.5; 1; 1.5; 1.75; 2\}g$. The candidate values for ν are selected so that the “effective spatial range”, which refers to the distance where spatial correlation drops below 0.05, covers 0.1 and 0.6 times ν^2 (the maximum inter-site distance within a unit square) for all candidate values of ν . Here we set $G = \{3; 14; 25; 36\}g$. Finally, we specify $G_2 = \{0.1; 0.5; 1; 2\}g$ as the candidate grid for ν^2 . We assign an $\text{IG}(a; b)$ prior with $a = b = 2$ for ν^2 . The prior of β follows a zero centred Gaussian with the covariance matrix equal to a diagonal matrix whose diagonal elements are 4, i.e., $N(\beta; V)$ where $\beta = \mathbf{0}$ and $V = 4 \cdot I$. For each simulated data set, we use stacking of means and stacking of predictive densities to obtain the expected outcome $\hat{y}(s)$ based on the held out observed locations. The predictive accuracy is evaluated by the mean squared prediction error over a set of n_h hold-out locations in set S_h ($\text{MSPE} = \sum_{s \in S_h} ((\hat{y}(s) - y(s))^2) = n_h$). We also compute the posterior expected values of the latent process $\hat{z}(s)$ for $z(s)$ on all of the n sampled locations in S and evaluate the mean squared error for $z(s)$ ($\text{MSEZ} = \sum_{s \in S} (\hat{z}(s) - z(s))^2 = n$). To further evaluate the distribution of predicted values, we compute the mean log point-wise predictive density for the n_h held out locations ($\text{MLPD} = \sum_{s \in S_h} \log(\sum_{g=1}^G w_g p(y(s) | y; M_g)) = n_h$).

Apart from stacking, we also implemented a fully Bayesian model with priors on the hyper-parameters using Markov chain Monte Carlo (MCMC) sampling for comparison. In addition, we carried out exact Bayesian inference using the conjugate model in Section 2.1 with hyper-parameters fixed at the exact value (denoted as M_0). We use the same priors for ν^2 and κ as those in stacking implementations. For the rest of the priors needed in full MCMC sampling, we assign uniform priors $\text{U}(3; 36)$ for β and $\text{U}(0.25; 2)$ for ρ^2 , and an $\text{IG}(2; 2)$ prior for ν^2 . Sampling is fitted through the `spLM` function in the `spBayes` package in R. The diagnostic metrics are computed based on 1,000 posterior samples retained after convergence was diagnosed over a burn-in period of 10,000 initial iterations. The algorithm for recovering the expected z and the log point-wise predictive density based on the output of `spLM` is presented in Section L of the Appendix. We monitor all diagnostic metrics for prediction for all competing algorithms. To measure uncertainty of the diagnostic metrics, we generate 60 data sets for each sample size in each simulation, fit each data set with the four competing methods and record the diagnostic metrics of each model fitting.

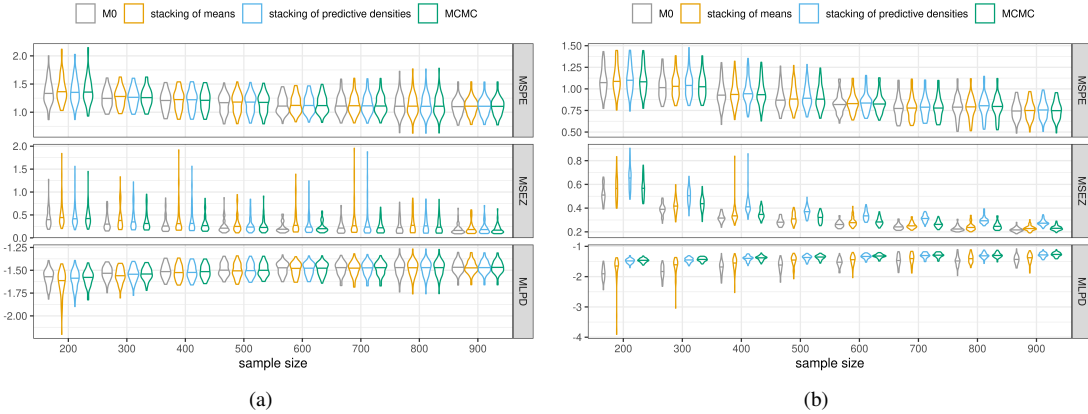


Fig. 3: Distributions of the diagnostic metrics for prediction performance for the first simulation (a) and the second simulation (b). Each distribution is depicted through a violin plot. The horizontal line in each violin plot indicates the median.

4.2. Predictive performances

Interestingly, the candidate algorithms exhibit different behaviours in the two simulation studies. Figure 3 summarises the comparison of the predictive performance. In the first simulation there seem to be no pronounced distinctions between the prediction performance for all competing models. In the second simulation, however, we observe that stacking of means outperforms stacking of predictive densities on having better estimates of the latent process on both the observed and unobserved locations (based on MSEZ), while stacking of predictive densities outperforms stacking of means in terms of the log point-wise predictive density (based on MLPD). These results are expected since we optimise the prediction error in the stacking of means, and we maximise the log predictive densities in the stacking of predictive densities.

Treating the fully Bayesian model with priors on all hyperparameters (fitted using MCMC) as a benchmark, we find that stacking of predictive densities is very competitive in terms of MLPD. The performance of latent process estimation for the full Bayesian model falls between stacking of means and stacking of predictive densities. The prediction accuracy for the response on unobserved locations for all competing algorithms are very close. Based on the medians of the MSPEs for all fittings, stacking of means slightly outperforms the full Bayesian model, and both are slightly better than stacking of predictive densities, The conjugate Bayesian model \mathcal{M}_0 provides the best point estimates for both response and latent processes based on MSPE and MSEZ, while it performs worse in terms of MLPD. These results seem to indicate that, when choosing stacking algorithms, stacking of means is preferred for point estimation and stacking of predictive densities is preferable for interval estimation.

In Figure 4, we compare the counts of the non-zero weights in stacking and find that stacking of means tends to produce a slightly smaller number of non-zero weights than stacking of predictive densities. The number of non-zero weights is small for both stacking algorithms, and this phenomenon is also observed and discussed in Breiman (1996b). On average, there are around 3.7 and 4.5 out of 64 weights that are greater than 0.001 in the simulation studies for stacking of means and stacking of predictive densities, respectively. And the number is relatively consistent when the sample sizes increase. These results suggest that stacking methodology is memory efficient, which is beneficial for large-scale spatial data analysis.

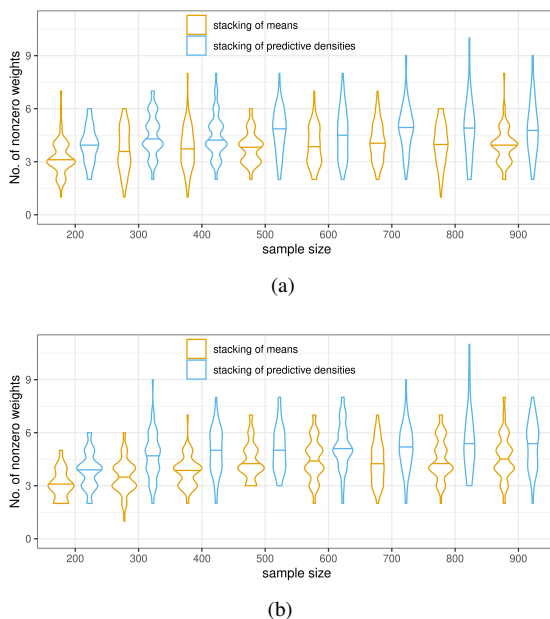


Fig. 4: Distributions of the counts of nonzero weights in the first (a) and the second (b) simulation. The distribution of the counts are described through violin plots whose horizontal lines indicate the medians.

In Section M.1, we pick one run for each simulation study and illustrate the interpolated maps of the predicted response on held out locations and expected latent processes over all locations generated by different fitting algorithms. The predicted response resembles the de-noised response $(X(s)^T + Z(s))$, not $Y(s)$, and the estimated latent process shares a similar pattern with the raw data. The latent process estimated by stacking of predictive densities is observed to be slightly smoother than those estimated by other algorithms, while the predictions of the response on unobserved locations fitted by different fitting algorithms are almost indistinguishable.

4.3. Running time comparisons

We provide both R and Julia code for the two proposed stacking algorithms. The code for the simulation studies are available from the GitHub repository https://github.com/LuZhangstat/spatial_stacking. Comparisons in predictive performances presented above are conducted in R. For the running time comparisons reported here, the stacking algorithms are implemented using Julia-1.6.6 and the MCMC sampling algorithms are run in R-4.2.0. We report the time for obtaining weights for stacking, and we consider the sampling time for $f; ; 2; 2g$ using MCMC (no sampling of $f; zg$ and no predictions). The timing comparisons are based upon experiments on a Windows 10 pro platform, with 64 GB of RAM and a 1-Intel Core i7-7700K CPU @ 4.20GHz processor with 4 cores each and 2 threads per core—totalling 8 possible threads for parallel computing. Figure 5 summarises the running time for the three competing algorithms. On average, the stacking of means is 526 times faster than MCMC in the simulation study. Stacking of predictive densities is slightly slower than stacking of means but is still around 430 times faster than MCMC sampling. These experiments clearly establish that predictive stacking algorithms are efficient alternatives to the full MCMC sampling algorithms for estimating latent spatial processes and predicting geostatistical outcomes.

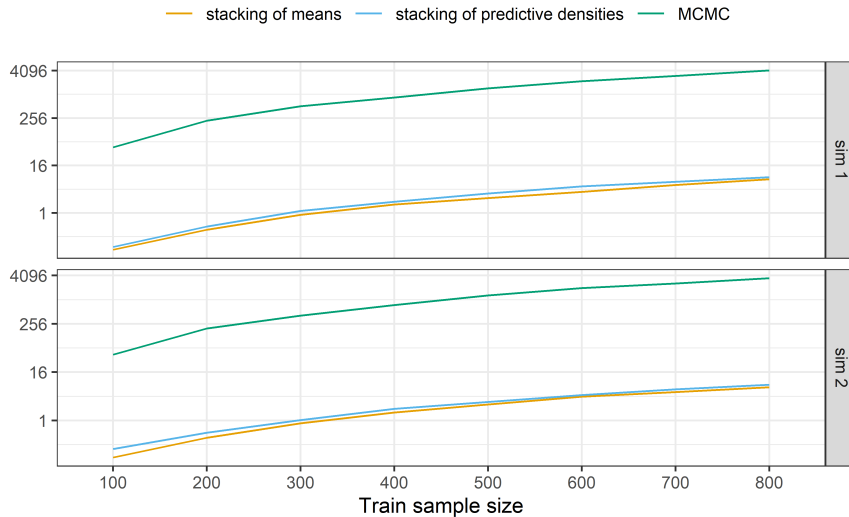


Fig. 5: Running time comparison for stacking and MCMC sampling

4.4. Inference of prefixed hyper-parameters

One limitation of stacking, compared to full Bayesian inference (e.g., using MCMC sampling algorithms), is that it does not provide interval estimates for the prefixed hyper-parameters. Moreover, our experiments show that stacking cannot provide a reliable point estimate of the hyper-parameters. If we treat the grid of the candidate values for the hyper-parameters in our stacking algorithms as a discrete uniform prior, then, intuitively, we should be able to achieve point estimates for those fixed hyper-parameters by a weighted average based on stacking weights. In Figure 6, we compare the point estimates of β based on stacking for the simulation studies. It is clear that stacking of means yields unstable estimates. Stacking of predictive densities has a smaller variance, but the bias can be large. Meanwhile, since β is not identifiable, we observe that the posterior interval estimates for β inferred from MCMC algorithms are wide, showing that the inference for β is relatively unstable for all candidate algorithms in this simulation study. The comparisons of the other two hyper-parameters are provided in the Section M.2.

5. WESTERN EXPERIMENTAL FOREST INVENTORY DATA ANALYSIS

We illustrate our proposed stacking methodology using the Western Experimental Forest (WEF) inventory data from a long-term ecological research site in western Oregon. This data set contains coordinates, species, diameter at breast height (DBH) and other measurements for trees in the experimental forest. The raw data is provided in the R package *spBayes*, and the code is available from the GitHub repository https://github.com/LuZhangstat/spatial_stacking. In this analysis, we focus on the prediction of DBH for living trees using the location and species information. There are 1,954 records in all after cleaning up some of the data. We randomly hold out 500 points for prediction and use the remaining observations to train the model and draw posterior inference.

We first fit a Bayesian linear regression model (BLM) using Hamiltonian Monte Carlo algorithms implemented in the Stan Bayesian modelling environment through the R interface *rstanarm*. We run the model fitting with the default settings, where the default priors are described in

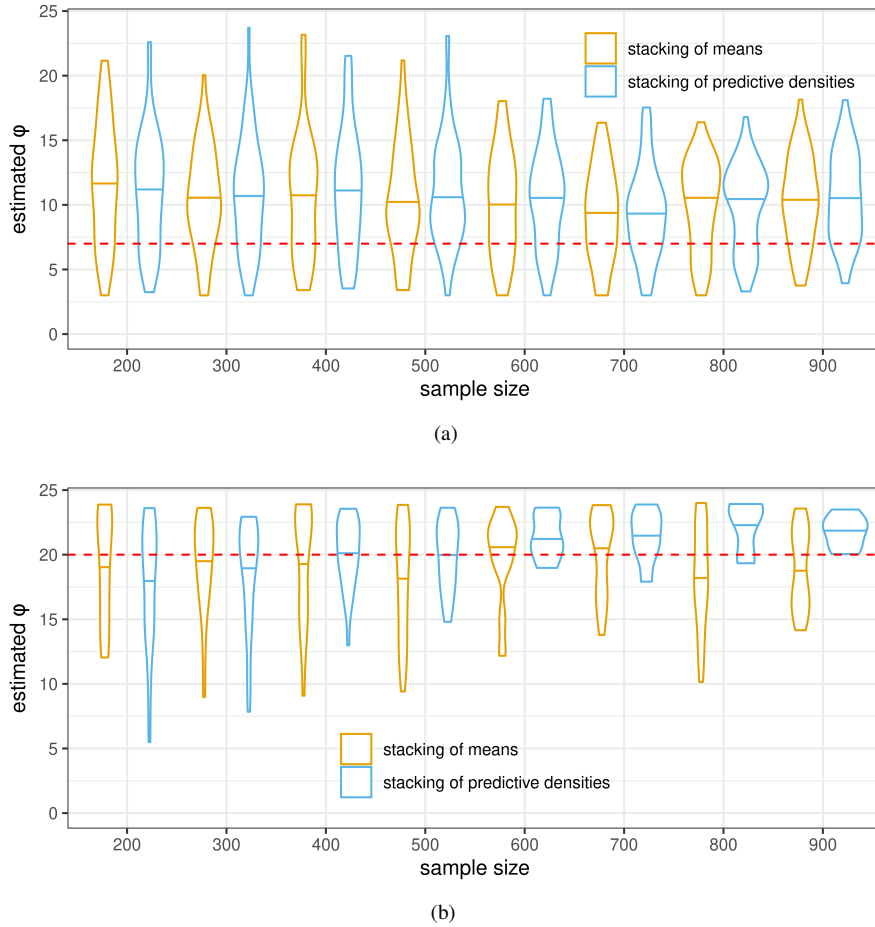


Fig. 6: Distributions of the estimated ϕ in the first (a) and the second (b) simulation. The distribution of the counts are described through violin plots whose horizontal lines indicate the medians. The red dashed horizontal line indicates the actual value of ϕ .

Gabry & Goodrich (2020). Then, we use the same prior for ϕ in our stacking implementations. For the prior of σ^2 in stacking implementations, we set up the hyper-parameters in the inverse-Gamma to have mean at the square of the mean of the default prior for σ^2 in BLM. We use the same candidate values for ρ and σ^2 as our simulation studies in Section 4. The largest and smallest candidate values for ρ are the maximum and minimum values so that the effective spatial range covers 0.1 and 0.6 times the maximum inter-site distance for Matérn covariogram with all candidate smoothness ν values. We pick four candidate values for ν in this data analysis. Finally, we evaluate predictive performance using the root mean squared prediction error (square root of the MSPE defined in Section 4.1) and the average log point-wise predictive density (MLPD defined in Section 4.1) for all held out locations. Table 1 summarises the result. Stacking delivers an RMSPE that is lower by around 9% and yields better MLPD in this data analysis.

Table 1: Summary Statistics of Western experimental forest inventory data analysis

	Bayesian linear regression	stacking of means	stacking of predictive densities
RMSPE	22.86	20.70	20.79
MLPD	-4.55	-4.45	-4.44

6. CONCLUSION AND FUTURE WORK

We develop geostatistical inference using Bayesian stacking. We offer theoretical insights for inferential behaviour of posterior distributions in fixed-domain or infill settings and explore the performance of our methods through simulations and analysis of a forestry data set. The empirical results reveal that our devised stacking methods deliver predictions comparable to full Bayesian inference obtained using MCMC samples, but at significantly lower costs. Our proposed stacking algorithm can be implemented in parallel and made storage efficient and, hence, can become a powerful alternative to full Bayesian inference using MCMC samples. Future directions can build upon our current framework to extend Bayesian stacking for multivariate geostatistics using conjugate matrix-variate normal-Wishart families (Zhang et al., 2021) and conjugate exponential families for non-Gaussian data (Bradley et al., 2020). Another extension is will be in the direction of stacked Bayesian inference for high-dimensional geostatistics (for example, building on the conjugate framework in Banerjee, 2020).

REFERENCES

- ABRAMOWITZ, M. & STEGUN, A. (1965). *Handbook of Mathematical Functions: with Formulas, Graphs, and Mathematical Tables*. Dover.
- BANERJEE, S. (2019). Geostatistics for environmental processes. In *Handbook of Environmental and Ecological Statistics*, A. E. Gelfand, M. Fuentes, J. A. Hoeting & R. L. Smith, eds. CRC press, Boca Raton, FL, pp. 81–96.
- BANERJEE, S. (2020). Modeling massive spatial datasets using a conjugate bayesian linear modeling framework. *Spatial Statistics* **37**, 100417.
- BANERJEE, S., CARLIN, B. P. & GELFAND, A. E. (2014). *Hierarchical Modeling and Analysis for Spatial Data*. CRC Press, Boca Raton, FL.
- BERGER, J. O., OLIVEIRA, V. D. & SANSÓ, B. (2001). Objective bayesian analysis of spatially correlated data. *Journal of the American Statistical Association* **96**, 1361–1374.
- BERNARDO, J. & SMITH, A. (1994). *Bayesian Theory*. Chichester, UK: John Wiley and Sons.
- BOSE, M., HODGES, J. S. & BANERJEE, S. (2018). Toward a diagnostic toolkit for linear models with gaussian-process distributed random effects. *Biometrics* **74**, 863–873.
- BRADLEY, J. R., HOLAN, S. H. & WIKLE, C. K. (2020). Bayesian hierarchical models with conjugate full-conditional distributions for dependent data from the natural exponential family. *Journal of the American Statistical Association* **115**, 2037–2052.
- BREIMAN, L. (1996a). Stacked regressions. *Machine Learning* **24**, 49–64.
- BREIMAN, L. (1996b). Stacked regressions. *Machine learning* **24**, 49–64.
- CHILÉS, J. & DELFINER, P. (1999). *Geostatistics: Modeling Spatial Uncertainty*. John Wiley: New York.
- CLYDE, M. & IVERSEN, E. S. (2013). Bayesian model averaging in the m-open framework. *Bayesian theory and applications* **14**, 483–498.
- CRESSIE, N. (1993). *Statistics for Spatial Data*. Wiley-Interscience, New York, revised ed.
- DE OLIVEIRA, V. & HAN, Z. (2022). On information about covariance parameters in gaussian matern random fields. *Journal of Agricultural, Biological and Environmental Statistics* **27**, 690–712.
- FINLEY, A. O., DATTA, A., COOK, B. C., MORTON, D. C., ANDERSEN, H. E. & BANERJEE, S. (2019). Efficient algorithms for bayesian nearest neighbor gaussian processes. *Journal of Computational and Graphical Statistics* **28**, 401–414.
- GABRY, J. & GOODRICH, B. (2020). Prior distributions for rstanarm models. <http://mc-stan.org/rstanarm/articles/priors.html>.
- HANDCOCK, M. S. & STEIN, M. L. (1993). A bayesian analysis of kriging. *Technometrics* **35**, 403–410.
- HODGES, J. S. (2013). Richly parameterized linear models: Additive, time series, and spatial models using random effects.

HOETING, J. A., MADIGAN, D., RAFTERY, A. E. & VOLINSKY, C. T. (1999). Bayesian model averaging: a tutorial (with comments by M. Clyde, David Draper and E. I. George, and a rejoinder by the authors. *Statistical Science* **14**, 382–417.

KAUFMAN, C. G. & SHABY, B. A. (2013). The role of the range parameter for estimation and prediction in geostatistics. *Biometrika* **100**, 473–484.

LE, T. & CLARKE, B. (2017). A bayes interpretation of stacking for m-complete and m-open settings. *Bayesian Analysis* **12**, 807–829.

LI, C., SUN, S. & ZHU, Y. (2023). Fixed-domain posterior contraction rates for spatial gaussian process model with nugget. *Journal of the American Statistical Association* **0**, 1–21.

MADIGAN, D., RAFTERY, A. E., VOLINSKY, C. & HOETING, J. (1996). Bayesian model averaging. In *Proceedings of the AAAI Workshop on Integrating Multiple Learned Models, Portland, OR*.

MURPHY, K. P. (2015). Conjugate bayesian analysis of the gaussian distribution, 2007. URL <https://www.cs.ubc.ca/~murphyk/Papers/bayesGauss.pdf>.

RASMUSSEN, C. E. & WILLIAMS, C. K. I. (2006). *Gaussian processes for machine learning*. MIT Press, Cambridge, MA.

STEIN, M. L. (1999). *Interpolation of Spatial Data: Some Theory for Kriging*. Springer-Verlag, New York.

TANG, W., ZHANG, L. & BANERJEE, S. (2021). On identifiability and consistency of the nugget in Gaussian spatial process models. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **83**, 1044–1070.

VEHTARI, A., MONONEN, T., TOLVANEN, V., SIVULA, T. & WINTHER, O. (2016). Bayesian leave-one-out cross-validation approximations for gaussian latent variable models. *The Journal of Machine Learning Research* **17**, 3581–3618.

WOLPERT, D. H. (1992). Stacked generalization. *Neural networks* **5**, 241–259.

YAO, Y., PIRRES, G., VEHTARI, A. & GELMAN, A. (2021). Bayesian hierarchical stacking: Some models are (somewhere) useful. *Bayesian Analysis* **1**, 1–29.

YAO, Y., VEHTARI, A. & GELMAN, A. (2020). Stacking for non-mixing bayesian computations: The curse and blessing of multimodal posteriors. *arXiv preprint arXiv:2006.12335*.

YAO, Y., VEHTARI, A., SIMPSON, D. & GELMAN, A. (2018). Using stacking to average bayesian predictive distributions (with discussion). *Bayesian Analysis* **13**, 917–1007.

ZHANG, H. (2004). Inconsistent estimation and asymptotically equal interpolations in model-based geostatistics. *Journal of the American Statistical Association* **99**, 250–261.

ZHANG, H. & ZIMMERMAN, D. L. (2005). Towards reconciling two asymptotic frameworks in spatial statistics. *Biometrika* **92**, 921–936.

ZHANG, L., BANERJEE, S. & FINLEY, A. O. (2021). High-dimensional multivariate geostatistics: A Bayesian matrix-normal approach. *Environmetrics* **32**, e2675.

ZIMMERMAN, D. & STEIN, M. (2010). Classical geostatistical methods. In *Handbook of spatial statistics*, A. E. Gelfand, P. Diggle, P. Guttorp & M. Fuentes, eds. CRC press, Boca Raton, FL, pp. 29–44.

A. PROOF OF THEOREM 1

Proof. For ease of presentation, we denote $y(\cdot)$ by y_n in the proofs. By Assumption 1, there is a probability distribution $P^\theta \in P_0$, which corresponds to the model with parameters $(\mu_0; \Sigma; \tau_0)$ for some $\theta \in \Theta$. Under P^θ , there exists a μ_0 such that

$$y_y \mid X_y \sim N(\mu_0; V^\theta); \tag{A.1}$$

where $V^\theta = \begin{pmatrix} \frac{1}{2} I_n & 0 & 0 \\ 0 & \frac{1}{2} I_p & 0 \\ 0 & 0 & \tau_0 I_n \end{pmatrix}$. Consider the sequence $y_n \sim P^\theta(\cdot \mid y_n)$. Under P^θ , $y_n \mid G(a_{:,n}; b_{:,n})$, where $a_{:,n} = a + n=2$ and

$$\begin{aligned} b_{:,n} &= b + \frac{1}{2} (y_y \mid X_y \mid \mu_0)^\top (I_{2n+p} - H) (y_y \mid X_y \mid \mu_0) \\ &= b + \frac{1}{2} [Q(y_y \mid X_y \mid \mu_0)]^\top \begin{pmatrix} 0 & 0 \\ 0 & I_n \end{pmatrix} [Q(y_y \mid X_y \mid \mu_0)]; \end{aligned} \tag{A.2}$$

The expectation under P^0 for $b_{:,n}$ is

$$\begin{aligned}
E^0(b_{:,n}) &= b + \frac{1}{2} E^0(y_y^T (I_{2n+p} - H) y_y) \\
&= b + \frac{1}{2} \left(X_y^T (I_{2n+p} - H) X_y + \text{Tr}((I_{2n+p} - H) V^0) \right) \\
&= b + \frac{1}{2} \text{Tr}((I_{2n+p} - H) V^0) \quad (\text{since } (I_{2n+p} - H) X_y = 0) \\
&= b + \frac{1}{2} \text{Tr} \begin{pmatrix} Q_{21}^T Q_{21} & \frac{2}{2} I_n & O \\ O & I_p & O \end{pmatrix} + \frac{1}{2} \text{Tr}(Q_{22}^T Q_{22}); \tag{A.3}
\end{aligned}$$

where $Q = \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{pmatrix}$ is an orthogonal matrix such that $H = Q^T \begin{pmatrix} I_{n+p} & 0 \\ 0 & 0 \end{pmatrix} Q$ with Q_{22} being the lower right $n \times n$ block of Q in a 2×2 partition. Using $Q_{21} Q_{21}^T = I_n - Q_{22} Q_{22}^T$ and some further simplification we obtain

$$E^0(b_{:,n}) = b + \frac{n}{2} \left(\frac{2}{2} (1 - \text{Tr}(Q_{22}^T Q_{22})) \right) + \frac{1}{2} \text{Tr}(Q_{22}^T Q_{22}) + \frac{p}{2} \frac{2}{2} \frac{2}{2}; \tag{A.4}$$

Since $E^0(b_{:,n}) = E^0(b_{:,n} - 1)$ and $\text{Tr}(Q_{22}^T Q_{22}) = n - \text{Tr}(H_{22})$, where $\text{Tr}(H_{22}) = n - 1$ as $n \geq 1$, we obtain $\lim_{n \rightarrow \infty} E^0(b_{:,n}) = 1$. The variance of $b_{:,n}$ under P^0 is given by

$$\begin{aligned}
V^0(b_{:,n}) &= E^0[V(b_{:,n})] + V^0[E(b_{:,n})] \\
&= E^0 \left[\frac{b_{:,n}^2}{(a_{:,n} - 1)^2 (a_{:,n} - 2)} \right] + V^0 \left[\frac{b_{:,n}}{a_{:,n} - 1} \right]; \tag{A.5}
\end{aligned}$$

Further note that

$$\begin{aligned}
V^0(b_{:,n}) &= \frac{3}{4} \text{Tr} \left(\frac{2}{2} (I_n - Q_{22}^T Q_{22}) \right) + \frac{1}{2} \text{Tr}(Q_{22}^T Q_{22})^2 + o(n) \\
&= C/n + \text{Tr}(Q_{22}^T Q_{22}) + \text{Tr}(Q_{22}^T Q_{22})^2; \tag{A.6}
\end{aligned}$$

for some $C > 0$ independent of n . Since $Q_{21} Q_{21}^T + Q_{22} Q_{22}^T = I_n$, we have $\text{Tr}(Q_{22}^T Q_{22})$ and $\text{Tr}(Q_{22}^T Q_{22})^2$ are bounded from above by n . Hence,

$$\frac{V^0(b_{:,n})}{n^2} \rightarrow 0 \quad \text{and} \quad \frac{E^0(b_{:,n}^2)}{n^2} \rightarrow 0 \quad \text{for any } \epsilon > 0; \tag{A.7}$$

Combining (A.5) and (A.7) yields $\lim_{n \rightarrow \infty} V^0(b_{:,n}) = 0$. By Chebyshev's inequality, $b_{:,n}$ converges in probability under P^0 , and hence under P_0 , to 1.

B. PROOF OF COROLLARY 1

Proof. Again we denote $y(n)$ by y_n . The conditional posterior distribution $p(j^2; y_n)$ can be derived from Lemma 1 as $N(j^2(M^2)_{[1:p]; 2}(M^2)_{[1:p; 1:p]})$, where $M^2 = V^{-1} + X^T R^{-1}(n) + 2I_n^{-1} X$ and $m = V^{-1} + X^T R^{-1}(n) + I_n^{-1} y_n$. Let $p(j^2; y_n)$. Some straightforward algebra yields

$$E^0(k_n - 0)^2 = \frac{2}{4} \text{Tr}(B_{11}) + E^0(n) \text{Tr}(C) + \frac{1}{2} \text{Tr}(D) + E^0(n) \text{Tr}(U_{11}); \tag{B.1}$$

Since $E^0(n) \rightarrow 2 < 1$ from Theorem 1, the proof follows from (2.8) and (B.1).

Derivation of (B.1)

$$\begin{aligned}
 E^0 k_{n_0}^2 &= E^0 [E^0 f(n_0)^T (n_0) j y_n g] \\
 &= E^0 [E^0 f(n_0) E^0(n_j y_n) + E^0(n_j y_n) E^0(n_0)]^T (n_0) [E^0(n_j y_n) + E^0(n_0)] j y_n g \\
 &= E^0 [E^0 f(n_0) E^0(n_j y_n)]^T (n_0) [E^0(n_j y_n)] j y_n g \\
 &\quad + 2 E^0 E^0(n_j y_n) E^0(n_0)^T (n_0) [E^0(n_j y_n)] j y_n g \\
 &\quad + E^0 [E^0 f(n_j y_n) E^0(n_0)]^T (E^0(n_j y_n) E^0(n_0)) j y_n g \\
 &= E^0 [\text{Tr} f V^0(n_j y_n) g] + E^0 [(E^0(n_j y_n) E^0(n_0))^T (E^0(n_j y_n) E^0(n_0))]
 \end{aligned} \tag{B.2}$$

Let $n_0 = \rho(2j y_n)$. By the definition of n_0 , we have

$$\begin{aligned}
 V^0(n_j y_n) &= V^0[E^0(j^2; y_n) j y_n] + E^0(V^0[j^2; y_n] j y_n) \\
 &= \underbrace{V^0[(M m)_{[1:p]} j y_n]}_{=0} + E^0(V^0(j^2; y_n) j y_n) \\
 &= E^0(n_j y_n) (M)_{[1:p]; [1:p]} = E^0(n_j y_n) U_{11}; \\
 E^0(n_j y_n) &= E^0 f E^0(j^2; y_n) j y_n g = (M m)_{[1:p]}.
 \end{aligned} \tag{B.3}$$

With $M = (X_y^T X_y)^{-1} = U$; $m = X_y^T y_y$, under P^0 ,

$$\begin{aligned}
 (M m)_{[1:p]} j^2 &= N(0; [M X_y^T V^0 X_y M]_{[1:p]; [1:p]}) \\
 &= N(0; U \left(\frac{2}{4} X^T X + \frac{2}{4} V^{-1} \right) U \left(\frac{2}{4} X^T \right)^T \left(\frac{2}{4} I_n + \frac{2}{4} R^{-1}(n) \right) U_{[1:p]; [1:p]})
 \end{aligned} \tag{B.4}$$

Therefore, $E^0[E^0(n_j y_n)] = 0$, and

$$\begin{aligned}
 &E^0 [(E^0(n_j y_n) E^0(n_0))^T (E^0(n_j y_n) E^0(n_0))] \\
 &= \text{Tr} f V^0[E^0(n_j y_n)] g \\
 &= \text{Tr} f V^0[(M m)_{[1:p]}] g \\
 &= \text{Tr} f V^0[E^0(M m)_{[1:p]} j^2] + E^0 V^0[(M m)_{[1:p]} j^2] g
 \end{aligned} \tag{B.5}$$

We separate the variance $V^0[(M m)_{[1:p]} j^2]$ as presented in (B.4) into the following three parts

$$\underbrace{\frac{2}{4} U \begin{pmatrix} X^T X & X^T \\ X & I_n \end{pmatrix} U}_{B_{11}} + \underbrace{\frac{2}{4} U \begin{pmatrix} V^{-1} & 0 \\ 0 & 0 \end{pmatrix} U}_{C} + \underbrace{\frac{2}{4} U \begin{pmatrix} 0 & 0 \\ 0 & R^{-1}(n) \end{pmatrix} U}_{D} :$$

Then

$$E^0 [(E^0(n_j y_n) E^0(n_0))^T (E^0(n_j y_n) E^0(n_0))] = \frac{2}{4} \text{Tr}(B_{11}) + E^0(n) \text{Tr}(C) + \frac{2}{4} \text{Tr}(D); \tag{B.6}$$

Together with (B.2), (B.3), we can obtain (B.1).

C. PROOF OF THEOREM 2

Proof. For ease of presentation, we denote $y(\cdot|n)$ by y_n . Observe that

$$\begin{aligned} E_0(Z_n(s_0) - z(s_0))^2 &= E_0 f Z_n(s_0) - E(z(s_0)|y_n) + E(z(s_0)|y_n) - z(s_0) g^2 \\ &= E_0 f z(s_0) - E(z(s_0)|y_n) g^2 + E_0 f V(z(s_0)|y_n) g; \end{aligned} \quad (\text{C.1})$$

where the second equality follows from the fact that $z(s_0) - E(z(s_0)|y_n)$ is independent of y_n . Note that

$$p(z(s_0)|y_n) = \int p(z(s_0)|y_n; z; \cdot) p(y_n; z) p(z|y_n) d^2 z; \quad (\text{C.2})$$

By standard Gaussian conditioning (see e.g. (Rasmussen & Williams, 2006, Section 2.2)),

$$p(z(s_0)|y_n; z; \cdot) = N(J_{:,n}^T R_{:,n}^{-1} z; \cdot (1 - J_{:,n}^T R_{:,n}^{-1} J_{:,n})): \quad (\text{C.3})$$

By (C.2), (C.3) and Lemma 1, the posterior predictive mean is

$$E(z(s_0)|y_n) = F_n; \quad (\text{C.4})$$

Further by the law of total variance and Theorem 1, we get

$$\begin{aligned} V(z(s_0)|y_n) &= E f V(z(s_0)|y_n; z; \cdot) g + V f E(z(s_0)|y_n; z; \cdot) g \\ &= \int (1 - J_{:,n}^T R_{:,n}^{-1} J_{:,n}) + 2 J_{:,n}^T R_{:,n}^{-1} U_{n[p+1:p+n; p+1:p+n]} R_{:,n}^{-1} J_{:,n} = 2 G_n; \end{aligned} \quad (\text{C.5})$$

Combining (C.1), (C.4) and (C.5) yields the decomposition (2.9), and hence the posterior predictive consistency for $z(s_0)$ holds if $E_{1;n} E_{2;n} \rightarrow 0$ as $n \rightarrow \infty$.

Further, let y_n have the density $p(y_n)$. Under conditions (2.8), by Corollary 1 we have $E_0 |j| \rightarrow 0$ as $n \rightarrow \infty$. As a result,

$$E_0 (Y_n(s_0) - y(s_0))^2 = E_0 (X(s_0)^T (I_n - \alpha))^2 + E_{1;n} + E_{2;n} + \frac{\sigma_0^2}{\alpha} + 2 E_0 (z|y_n);$$

which clearly converges to $\frac{\sigma_0^2}{\alpha} + 2 \sigma_0^2$ as $n \rightarrow \infty$.

D. PROOF OF THEOREM 3

The proof of this theorem breaks into several lemmas. Recall the definition of b from Lemma 1. The lemma below provides a simple expression for $b_{:,n}$, which is specific to the conjugate model (2.14).

LEMMA 2. We have $b_{:,n} = b + \frac{1}{2} y (I_n + R(\cdot))^{-1} y(\cdot)$.

Proof. Note that

$$M^{-1} = X^T V_y^{-1} X = 2 I_n + R(\cdot) \quad \text{and} \quad m = X^T V_y^{-1} y = 2 y; \quad (\text{D.1})$$

By the Woodbury matrix identity,

$$\begin{aligned} m^T M^{-1} m &= 2 y^T (I_n + 2 R^{-1}(\cdot))^{-1} y \\ &= 2 y^T y - y^T (2 I_n + R(\cdot))^{-1} y = y^T V_y^{-1} y - y^T (2 I_n + R(\cdot))^{-1} y; \end{aligned}$$

which yields the desired result.

The next lemma studies the asymptotic behaviour of b in the case where the range decay $\alpha = \alpha_0$, that is α is fixed at the true value.

LEMMA 3. Let $\alpha = \alpha_0$, and assume that $\max_{s \in D} \min_{1 \leq i \leq n} |s_j - s_i| \geq n^{-\frac{1}{d}}$. Then

$$\frac{b_{:,n} - b}{n} \rightarrow \frac{\sigma_0^2}{2}; \quad \mathbb{P}_0\text{-almost surely}; \quad (\text{D.2})$$

Proof. Let Q_n be the orthogonal matrix such that $Q_n R_o(n) Q_n^T = \begin{pmatrix} \lambda_1^{(n)} & & \\ & \ddots & \\ & & \lambda_n^{(n)} \end{pmatrix}$, where $\lambda_i^{(n)}$ is the i -th largest eigenvalue of matrix $R_o(n)$. Thus, under P_0 ,

$$Q_n y(n) \sim N\left(\begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}, \begin{pmatrix} \frac{\sigma^2}{\lambda_1^{(n)} + 2} & & \\ & \ddots & \\ & & \frac{\sigma^2}{\lambda_n^{(n)} + 2} \end{pmatrix}\right)$$

By Lemma 2, we get

$$2(b_{:,n} - b) = \sum_{i=1}^n \frac{\frac{\sigma^2}{\lambda_i^{(n)} + 2}}{\frac{\sigma^2}{\lambda_i^{(n)} + 2} + \frac{\sigma^2}{2}} u_i^2; \tag{D.3}$$

where $u_i \stackrel{i.i.d.}{\sim} N(0, 1)$ for $i = 1, \dots, n$. By Tang et al. (2021, Corollary 2), there exists $C > 0$ independent of n such that $\lambda_i^{(n)} \geq C n i^{-\frac{2}{\sigma}} - 1$ for all $1 \leq i \leq n$. This implies that

$$\sum_{i=1}^n \frac{\frac{\sigma^2}{\lambda_i^{(n)} + 2}}{\frac{\sigma^2}{\lambda_i^{(n)} + 2} + \frac{\sigma^2}{2}} \leq \frac{n \frac{\sigma^2}{2}}{2} \text{ as } n \rightarrow \infty; \tag{D.4}$$

By the law of large numbers, (D.2) follows from (D.3) and (D.4).

Proof of Theorem 3. Let $\theta := \frac{\sigma^2}{2}$, and let P^θ be the probability distribution of the Matérn model with parameters $(\frac{\sigma^2}{2}; \frac{\sigma^2}{2})$. By Tang et al. (2021, Theorem 1), P^θ is equivalent to P . Further by Lemma 3,

$$\frac{b_{:,n} - b}{n} \xrightarrow{P^\theta\text{-almost surely}} \frac{\sigma^2}{2};$$

which also holds P_0 -almost surely. Now by Lemma 1,

$$E_0(\frac{\sigma^2}{2} j y(n)) = \frac{b_{:,n} - b}{a_{:,n}} \frac{\sigma^2}{2} \text{ and } V_0(\frac{\sigma^2}{2} j y(n)) = \frac{b_{:,n} - b}{(a_{:,n} - 1)^2 (a_{:,n} - 2)} \frac{1}{n}; \tag{D.5}$$

which yields (2.17) by Chebyshev's inequality. This implies the posterior consistency of $\hat{\sigma}^2 = \frac{\sigma^2}{2}$.

E. PROOF OF THEOREM 4

Proof. By (D.1), we have

$$p(z|y; \sigma^2) = N((I_n + \sigma^2 R(n))^{-1} y; \sigma^2 (I_n + R(n))^{-1});$$

Combining with (C.4), we get the posterior predictive mean

$$\begin{aligned} E(z(s_0) | y_n) &= J_{:,n}^T R(n)^{-1} E(z | y_n) \\ &= J_{:,n}^T R(n)^{-1} (I_n + \sigma^2 R(n))^{-1} y_n = J_{:,n}^T (\sigma^2 I_n + R(n))^{-1} y_n; \end{aligned}$$

which is the best linear predictor corresponding to a Matérn model with parameter values $f = \frac{\sigma^2}{2}; \frac{\sigma^2}{2} g$ satisfying $\sigma^2 = \frac{\sigma^2}{2} = \frac{\sigma^2}{2}$. Further by Theorem 3, the formula (C.5) reduces to

$$V(z(s_0) | y_n) \leq \frac{\sigma^2}{2} (1 - J_{:,n}^T R(n)^{-1} J_{:,n} + \sigma^2 J_{:,n}^T (\sigma^2 I_n + R(n))^{-1} R(n)^{-1} J_{:,n});$$

which gives (2.18). The rest of the theorem easily follows.

F. PROOF OF PROPOSITION 1

Proof. To simplify the notation, we denote by $\Delta := \frac{1}{n}$ the inter-spacing of the grid. Let f_0 (resp. f) be the spectral density of the Matérn model with true parameter values $(\nu_0, \rho_0, \sigma_0^2, g_0)$ and the smoothness parameter ν_0 (resp. possibly misspecified parameter values (ν, ρ, σ^2, g) the smoothness parameter ν). Define

$$\tilde{f}_0(u) = \prod_{k=1}^{\infty} f_0 \left(\frac{u+2k}{\Delta} \right) \quad \text{and} \quad \tilde{f}(u) = \prod_{k=1}^{\infty} f \left(\frac{u+2k}{\Delta} \right);$$

By (Stein, 1999, Chapter 3, (13)), the prediction error of $y(0)$ based on $y(s)$, $s \in \mathcal{I}_n$ is

$$e_0^2 = \frac{\int_{\mathbb{R}} \tilde{f}_0(u) du}{\int_{\mathbb{R}} \tilde{f}(u) du};$$

and hence the prediction error of $z(0)$ based on $y(s)$, $s \in \mathcal{I}_n$ is

$$e_1 = \frac{\int_{\mathbb{R}} \frac{\tilde{f}_0(u)}{\tilde{f}(u)^2} du}{\int_{\mathbb{R}} \tilde{f}(u) du}.$$

Recall from (2.13) the spectral density of the Matérn model without nugget. We write

$$\tilde{f}_0(u) = \frac{\sigma_0^2}{2} g_0(u) + \frac{\rho_0^2}{2} \quad \text{and} \quad \tilde{f}(u) = \frac{\sigma^2}{2} g(u) + \frac{\rho^2}{2};$$

so that

$$e_1 = \frac{\int_{\mathbb{R}} \frac{\frac{\sigma_0^2}{2} g_0(u) + \frac{\rho_0^2}{2}}{(\frac{\sigma^2}{2} g(u) + \frac{\rho^2}{2})^2} du}{\int_{\mathbb{R}} (\frac{\sigma^2}{2} g(u) + \frac{\rho^2}{2}) du} \quad (F.1)$$

Note that $g(u) \sim cu^{-2-\nu}$ for some $c > 0$. It is known that (see e.g. (Tang et al., 2021, Section 2.3))

$$\int_{\mathbb{R}} \frac{1}{(\frac{\sigma^2}{2} g(u) + \frac{\rho^2}{2})^2} du \sim \frac{4}{2+C} \frac{1}{\sigma^{2+2\nu}} \quad \text{for some } C > 0; \quad (F.2)$$

Furthermore, $g(u) \sim u^{-2}$ for u large and $g(u) \sim 1$ for u small. We prove that $e_1 \rightarrow 0$ as $n \rightarrow \infty$. Further by Assumption 2, we get $e_1 \rightarrow 0$ as $n \rightarrow \infty$. Therefore,

$$\int_{\mathbb{R}} \frac{\frac{\sigma_0^2}{2} g_0(u) + \frac{\rho_0^2}{2}}{(\frac{\sigma^2}{2} g(u) + \frac{\rho^2}{2})^2} du = \frac{4}{4} \frac{\sigma_0^2}{\sigma^2} + O(\min(2-\nu_0, 2-\nu)); \quad (F.3)$$

Combining (F.1), (F.2) and (F.3) yields

$$e_1 = \frac{\sigma_0^2}{\sigma^2} + O(\min(2-\nu_0, 2-\nu)) \left(1 + \frac{C}{2} \frac{\sigma_0^2}{\sigma^2} \right)^{\frac{1}{2}} \sim \frac{\sigma_0^2}{\sigma^2} + O(\min(2-\nu_0, 2-\nu));$$

Thus, we have $e_1 \rightarrow 0$ as $n \rightarrow \infty$, and by Assumption 2 we get (2.20).

G. PROOF OF PROPOSITION 2

Proof. For ease of presentation, we give the proof in the setting of Theorem 4. Note that

$$\begin{aligned}
 E_0 \left[\sum_{g=1}^G w_g E_g(y(s_0) | y) \right]^2 &= E_0 \left[\sum_{g=1}^G w_g (z(s_0) - E_g(z(s_0) | y)) \right]^2 \\
 &= E_0 \left[\sum_{g=1}^G w_g^2 (z(s_0) - E_g(z(s_0) | y))^2 \right] + 2 \sum_{g=1}^{G-1} \sum_{h=g+1}^G w_g w_h E_0 (z(s_0) - E_g(z(s_0) | y)) (z(s_0) - E_h(z(s_0) | y))
 \end{aligned} \tag{G.1}$$

where we apply the Cauchy-Schwarz inequality in (G.1). For each $1 \leq g \leq G$, the term $E_0 (z(s_0) - E_g(z(s_0) | y))^2$ corresponds to the deviation error $E_{1;n}$ for the model M_g , which goes to 0 as $n \rightarrow \infty$. Since $\sum w_g$ is bounded for each g , the bound (3.3) follows readily from (G.1).

H. PROOF OF THEOREM 5

Proof. By Proposition 2, it suffices to prove that

$$E_0 \left[\frac{1}{n} \sum_{i=1}^n y(s_i) - \sum_{g=1}^G w_g \hat{y}_g(s_i) \right]^2 \rightarrow 0 \text{ as } n \rightarrow \infty \tag{H.1}$$

For ease of presentation, we prove (H.1) in the setting of Theorem 4. Note that

$$\begin{aligned}
 E_0 \left[\frac{1}{n} \sum_{i=1}^n y(s_i) - \sum_{g=1}^G w_g \hat{y}_g(s_i) \right]^2 &= E_0 \left[\frac{1}{n} \sum_{i=1}^n \underbrace{(z(s_i) - \hat{y}_g(s_i))}_{B} \right]^2 \\
 &= E_0 \left[\frac{1}{n} \sum_{i=1}^n \sum_{g=1}^G w_g (z(s_i) - \hat{y}_g(s_i)) \right]^2 \\
 &= E_0 \left[\frac{1}{n} \sum_{g=1}^G w_g \sum_{i=1}^n (z(s_i) - \hat{y}_g(s_i)) \right]^2 \\
 &= E_0 \left[\frac{1}{n} \sum_{g=1}^G w_g^2 \sum_{i=1}^n (z(s_i) - \hat{y}_g(s_i))^2 \right] + 2 \sum_{g=1}^{G-1} \sum_{h=g+1}^G w_g w_h E_0 \left[\frac{1}{n} \sum_{i=1}^n (z(s_i) - \hat{y}_g(s_i)) (z(s_i) - \hat{y}_h(s_i)) \right]
 \end{aligned} \tag{H.2}$$

By the boundedness of $\sum w_g$ and the condition (3.4), the limit (H.1) follows easily from (H.2). Taking a closer look at (H.2), we can see that B summarises the average squared prediction errors for $z(s_i); i = 1, \dots, n$, in a LOO cross-validation. Hence, the stacking weights obtained from (3.2) minimises the average squared prediction errors for the latent process over the observed locations.

I. PSEUDO-CODES FOR STACKING ALGORITHMS

Algorithm 1. Stacking of means.

Data: X, y, n : Design matrix, outcome and location set

V, a, b : Prior parameters

$G; G; G_2$: Grids of $s; s^0; s^2$

Result: G_{all} : Grid spanned by $G; G; G_2$

w : Stacking weights

Compute $X_{prod}^{(k)} = X^T[k]X[k]$, $X_y^{(k)} = X^T[k]y[k]$ and record the number of observations n_k in fold k for $k = 1; \dots; K$, where $X[k]y[k]$ denotes the predictors and response for observations not in fold k ;

foreach $f; g$ in grid expanded by grid of f and g **do**

for $k = 1$ to K **do**

Calculate $R_{(k; k)}^{-1} = fR^{-1}(s; s^0; s^2)g_{s; s^0; s^2}^{-1} n[k]$ $O((n - n_k)^3)$;

Store $R_{(k; k)} = fR(s; s^0; s^2)g_{s; s^0; s^2} n[k]$ $O(n - n_k)$;

for z in the grid of z^2 **do**

Compute the Cholesky decomposition L of $M^{-1} = L L^T = X_{prod}^{(k)} + V^{-1} X^T[k] + R_{(k; k)}^{-1} + I_{n - n_k}$ $O(n^3)$;

Compute $u = m = V^{-1} X_y^{(k)} + y[k]$ $O(n - p)$;

Update $u = L^{-1}u$;

Update $u = L^{-T}u$, where $u = (z^T; z^T)^T$ $O(n^2)$;

Compute the expected latent process on locations in fold k

$z_{(z; z^2)}^{(k)} = R_{(k; k)} R_{(k; k)}^{-1} z$;

Compute expected outcome on locations in fold k $y_{(z; z^2)}^{(k)} = X[k]z_{(z; z^2)}^{(k)} + z_{(z; z^2)}^{(k)}$;

Compute Stacking weights w based on $\bar{y}_{(z; z^2)}^{(k)} g_{(z; z^2)}^{k=1; \dots; K} 2G_{all}$;

Algorithm 2. Stacking of predictive densities

Data: X, y, n : Design matrix, outcome and location set
 μ, V, a, b : Prior parameters
 $G_1; G_2; G_3$: Grids of $s_1; s_2; s_3$
 J : number of samples for log point-wise predictive density estimation
Result: G_{all} : Grid spanned by $G_1; G_2; G_3$
 w : Stacking weights

Compute $X_{prod}^{(k)} = X^T[k]X[k]$, $X_y^{(k)} = X^T[k]y[k]$, $ky[k]k^2 = y^T[k]y[k]$ and record the number of observations n_k in fold k for $k = 1; \dots; K$, where $X[k]y[k]$ denotes the predictors and response for observations not in fold k ;

foreach $f; g$ in grid expanded by grid of s_1 and s_2 **do**

for $k = 1$ to K **do**

Calculate $R_{(k; k)}^{-1} = fR^{-1}(s; s^0; \dots; s^J)g_{s; s^0; 2; n[k]}$ $O((n - n_k)^3)$;

Store $R_{(k; k)} = fR(s; s^0; \dots; s^J)g_{s; 2; S[k]; s^0; 2; n[k]}$ $O(n - n_k)$;

for s_3 in the grid of s_3 **do**

Compute the Cholesky decomposition L of $M^{-1} = L L^T = \begin{matrix} X_{prod}^{(k)} + V^{-1} & X^T[k] \\ X[k] & R_{(k; k)}^{-1} + I_{n - n_k} \end{matrix}$ $O(n^3)$;

Compute $u = m = \begin{matrix} V^{-1} + X_y^{(k)} \\ y[k] \end{matrix}$ $O(n - p)$;

Update $u = L^{-1}u$;

Compute $b = b + 0.5(X_y[k]k^2 + u^T V^{-1} u)$ and $a = a + 0.5(N - n_k)$;

Update $u = L^{-T}u$, where $u = (\tau; z^T)^T$ $O(n^2)$;

Generate the posterior expected outcome on locations in fold k
 $\hat{y}_{(s_1; s_2)}^k = X[k]u + R_{(k; k)}^{-1}z$

Compute $lpc = 0.5 \log(2) + \log(a + 1/2) - \log(a) + a \log b$;

foreach $s \in S[k]$ **do**

Generate $h_s = x(s) R_{(k; k)}^{-1}(s) R^{(inv)}$ where $R_{(k; k)}^{-1}(s)$ is the row with elements $fR(s; s^0; \dots; s^J)g_{s^0; 2; n[k]}$ in $R_{(k; k)}$

Compute $V_s = kL^{-1}h_s^T k^2 + 2$.

Compute the log point-wise predictive density of $y(s)$ at location s , by
 $lpc_{(s_1; s_2)}(s) = 0.5 \log(V_s) - (a + 1/2) \log fb + (y(s) - \hat{y}_{(s_1; s_2)}^k)^2 = (2V_s)g(J.1)$;

Compute Stacking weights w based on $f/lpc_{(s_1; s_2)}(s)g_{(s_1; s_2)}^{s_3; n; 2; G_{all}}$;

J. DERIVE THE CLOSED FORM OF POINT-WISE PREDICTIVE DENSITY

In this subsection, we derive the posterior predictive density of the outcome $y(s_0)$ on location s_0 . We follow the notations in Section 2. First, we know that $y(s_0) | z; y_1$ follows a Gaussian with mean $x(s_0) + U(\cdot)R^{-1}(\cdot)z$ and variance 2 . Since the conditional posterior distribution $z | y_1$ follows $N(Mm; 2M)$, the conditional posterior distribution $y(s_0) | z; y_1$ still follows a Gaussian $N(x(s_0) + U(\cdot)R^{-1}(\cdot)z; 2V_{s_0})$ where

$$s_0 = \left| \begin{array}{c} x(s_0) \ U(\cdot)R^{-1}(\cdot) \\ \hline \underbrace{\hspace{10em}}_{h_{s_0}} \end{array} \right\} Mm; V_{s_0} = h_{s_0}M h_{s_0}^T + 2;$$

Next, through equation (2.4)

$$\begin{aligned}
p(y(s_0)|j; y_{(-j)}) &= \int p(y(s_0)|j; y_{(-j)})p(y_{(-j)}|y)dy \\
&= \int N(y(s_0)|j; s_0, \sigma^2 V_{s_0})IG(y_{(-j)}|a; b)dy \\
&= \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left[-\frac{(y(s_0) - s_0)^2}{2\sigma^2 V_{s_0}}\right] \frac{b^a}{\Gamma(a)} \int \exp\left[-\frac{b}{2}(y_{(-j)} - s_0)^2\right] dy \\
&= \frac{b^a}{(2\pi\sigma^2)^{1/2} \Gamma(a)} \int \exp\left[-\frac{1}{2}\left(b + \frac{(y(s_0) - s_0)^2}{2\sigma^2 V_{s_0}}\right) y_{(-j)}^2\right] dy \\
&= \frac{(a+1/2)b^a}{(2\pi\sigma^2)^{1/2} \Gamma(a)} \left[b + \frac{(y(s_0) - s_0)^2}{2\sigma^2 V_{s_0}} \right]^{-a}
\end{aligned}$$

The log point-wise predictive density is

$$\begin{aligned}
\log p(y(s_0)|j; y_{(-j)}) &= -0.5 \log(2\pi\sigma^2) - a \log b - (a+1/2) \log \left[b + \frac{(y(s_0) - s_0)^2}{2\sigma^2 V_{s_0}} \right] \\
&\quad + \log \Gamma(a+1/2) - \log \Gamma(a)
\end{aligned} \tag{J.1}$$

K. STACKING OF PREDICTIVE DENSITIES ALGORITHM (MONTE CARLO VERSION)

We present a Monte Carlo algorithm to estimate the log of point-wise predictive density for outcome in fold k given observations not in fold k . For each k , we generate J posterior samples of θ^2 and $(\tau; Z^T)^T = (\tau, \sigma^2)$, (i.e., $F^{(j)}; G^{(j)}$) for $j = 1, \dots, J$, using data not in fold k . Then we calculate the corresponding expected outcome for location s in fold k , $y_{(\cdot; \cdot; 2)}^{k,j}(s)$ for $j = 1, \dots, J$. Next, we compute the predictive density of $y(s)$ conditional on the prediction $y_{(\cdot; \cdot; 2)}^{k,j}(s)$ and the nugget (variance of the noise process, which equals the product of σ^2 and the j -th posterior sample $\theta^{2(j)}$) for each j . The conditional predictive distribution of $y(s)$ follows

$$p(y(s)|j; y_{(-j)}^{k,j}) = N(y(s)|j; y_{(\cdot; \cdot; 2)}^{k,j}, \sigma^2 \theta^{2(j)}) \tag{K.1}$$

Finally, the log point-wise predictive density (LPD) of $y(s)$ at location s is estimated by

$$\begin{aligned}
\log p_{(\cdot; \cdot; 2)}(s) &= \log \int p(y(s)|j; y_{(-j)}^{k,j})p(y_{(-j)}^{k,j}|k)d y_{(-j)} \\
&= \log \int \frac{1}{J} p_{(\cdot; \cdot; 2)}(y(s)|j; y_{(-j)}^{k,j}) dy_{(-j)}
\end{aligned} \tag{K.2}$$

and we can compute the stacking weights based on the estimated LPDs. The following is the Monte Carlo version of the stacking of predictive densities algorithm.

Algorithm 3. Stacking of predictive densities (Monte Carlo Version)

Data: X, y, n : Design matrix, outcome and location set
 V, a, b : Prior parameters
 $G_1; G_2; G_3$: Grids of $s; s^0; s^2$
 J : number of samples for log point-wise predictive density estimation

Result: G_{all} : Grid spanned by $G_1; G_2; G_3$
 w : Stacking weights

Compute $X_{prod}^{(k)} = X^T[k]X[k]$, $X_y^{(k)} = X^T[k]y[k]$, $ky[k]k^2 = y^T[k]y[k]$ and record the number of observations n_k in fold k for $k = 1; \dots; K$, where $X[k]y[k]$ denotes the predictors and response for observations not in fold k ;

foreach $f; g$ in grid expanded by grid of s and s^2 **do**

for $k = 1$ to K **do**

Calculate $R_{(k; k)}^{-1} = fR^{-1}(s; s^0; s^2; n[k])g_{s; s^0; s^2; n[k]}$ $O((n - n_k)^3)$;

Store $R_{(k; k)} = fR(s; s^0; s^2; n[k])g_{s; s^0; s^2; n[k]}$ $O(n - n_k)$;

for s^2 in the grid of s^2 **do**

Compute the Cholesky decomposition L of $M^{-1} = L L^T = \begin{matrix} X_{prod}^{(k)} + V^{-1} & X^T[k] \\ X[k] & R_{(k; k)}^{-1} + I_{n - n_k} \end{matrix}$ # $O(n^3)$;

Compute $u = m = \begin{matrix} V^{-1} + X_y^{(k)} \\ y[k] \end{matrix}$ $O(n - p)$;

Update $u = L^{-1}u$;

Compute $b = b + 0.5(X_y^{(k)}k^2 + X^T[k]V^{-1}u^T u)$ and $a = a + 0.5(N - n_k)$;

Generate $v^{(1); \dots; v^{(J)}}$ Inverse-Gamma($a; b$);

Generate $v^{(j)} \sim N(M^{-1}m; v^{(j)}M)$ by taking $v^{(j)} \sim N(0; v^{(j)}I_{n - n_k + p})$ and computing $v^{(j)} = L^{-T}(v^{(j)} + u)$ for $j = 1; \dots; J$;

Generate the posterior samples of the expected outcome on locations in fold k
 $y_{(s; s^2)}^{(k; j)} = X[k]v^{(j)} + R_{(k; k)}^{-1}R_{(k; k)}v^{(j)}$, $j = 1; \dots; J$;

foreach $s \in S[k]$ **do**

Compute the posterior samples of the log-density of observation $y(s)$ at location s , $p_{(s; s^2)}(y(s)|j^{(2(j)); v^{(j)})}$ by taking the density of $N(y_{(s; s^2)}^{(k; j)}(s); v^{(j)})$ at $y(s)$ for $j = 1; \dots; J$ (K.1);

Compute the expected log point-wise predictive density of $y(s)$ at location s , by
 $lp_{(s; s^2)}(s) = \log \frac{1}{J} \sum_{j=1}^J p_{(s; s^2)}(y(s)|j^{(2(j)); v^{(j)})}$ (K.2);

Compute Stacking weights w based on $f/p_{(s; s^2)}(s)g_{(s; s^2)}^{s^2; n}G_{all}$;

L. RECOVER EXPECTED Z AND LOG POINT-WISE PREDICTIVE DENSITY FOR MCMC SAMPLING

The package *spBayes* doesn't record the posterior samples of the latent process $Z(s)$. In this subsection, we illustrate how to recover the expected $Z(s)$ for both observed and unobserved location and compute MLPD for the simulation studies based on the outputs returned by *spLM*. To achieve our goal, we need to recover the posterior samples of $Z(s)$ on all locations given the recorded MCMC samples of parameters $s; s^2; s^2$ and s . Denote $Z(s)$ on observed and unobserved locations as Z_o and Z_u , respectively, and

denote $z(s)$ on all locations as z . Based on (2.2),

$$\begin{aligned} p(z | \{Z_o, Z_u\}, y, X) &\propto N(y | Xz + \mu_0, \Sigma_0) N(z | \mu_0, \Sigma_0) N(z | \mu_0, \Sigma_0) \\ &\propto \exp \left\{ -\frac{1}{2} (y - Xz - \mu_0)^T \Sigma_0^{-1} (y - Xz - \mu_0) - \frac{1}{2} z^T \Sigma_0^{-1} z \right\} \\ &\propto \exp \left\{ -\frac{1}{2} z^T \Sigma_0^{-1} z - \frac{1}{2} (y - Xz - \mu_0)^T \Sigma_0^{-1} (y - Xz - \mu_0) \right\} \\ &\propto N(z | M_z, m_z); \end{aligned}$$

where Σ_0 combines the observed and unobserved location sets and

$$M_z = \Sigma_0^{-1} \mu_0 + \Sigma_0^{-1} X^T (y - X\mu_0 - \mu_0); \quad m_z = \Sigma_0^{-1} (y - X\mu_0 - \mu_0);$$

Let $f^{(j)}; z^{(j)}; y^{(j)}; X^{(j)}$ for $j = 1; \dots; J$ denote the recorded MCMC samples. We generate posterior samples for z using the above full conditional posterior distribution for each iteration j and then compute the average as the expected z . We further compute the LPD of $y(s)$ any held out location s by

$$\begin{aligned} \log p(s) &= \log \int p(y(s) | z(s)) p(z(s) | y) dz(s) \\ &= \log \left[\frac{1}{J} \sum_{j=1}^J p(y(s) | z^{(j)}(s)) \right] \\ &= \log \left[\frac{1}{J} \sum_{j=1}^J N(y(s) | X^{(j)T} z^{(j)}(s) + \mu_0, \Sigma_0) \right] \end{aligned}$$

M. FIGURES FOR SIMULATION STUDIES

M.1. Interpolated maps for the simulation studies

M.2. Distributions of the estimated and Σ_0

N. COMPUTE THE STACKING WEIGHTS FOR STACKING OF MEANS (IN R CODE)

Let us format the expected outcome $\hat{y}_{(s)}^{(k)}; g_{(s)}^{(k)}$ computed in Algorithm 1 by an $N \times G$ matrix \hat{Y} . Each column of \hat{Y} stores the expected outcome $\hat{y}_{(s)}^{(k)}; g_{(s)}^{(k)}$ for each candidate model, and it shares the same order of observed locations as the outcome y . Let $w = (w_1; w_2; \dots; w_G)^T$ be the stacking weights, we need to find the weights that satisfy

$$\arg \min_w f(y - \hat{Y}w)^T (y - \hat{Y}w)g;$$

under the constrain $\sum_{g=1}^G w_g = 1$. Now we modify it into a quadratic programming (QP) problem.

$$\begin{aligned} &(y - \hat{Y}w)^T (y - \hat{Y}w) \\ &= f(y - \sum_{g=1}^G w_g \hat{Y}_g) \sum_{g=1}^G w_g (\hat{Y}_g - \hat{Y}_G)^T f(y - \sum_{g=1}^G w_g \hat{Y}_g) \sum_{g=1}^G w_g (\hat{Y}_g - \hat{Y}_G)g \\ &= (y - \hat{Y}w)^T (y - \hat{Y}w) \end{aligned}$$

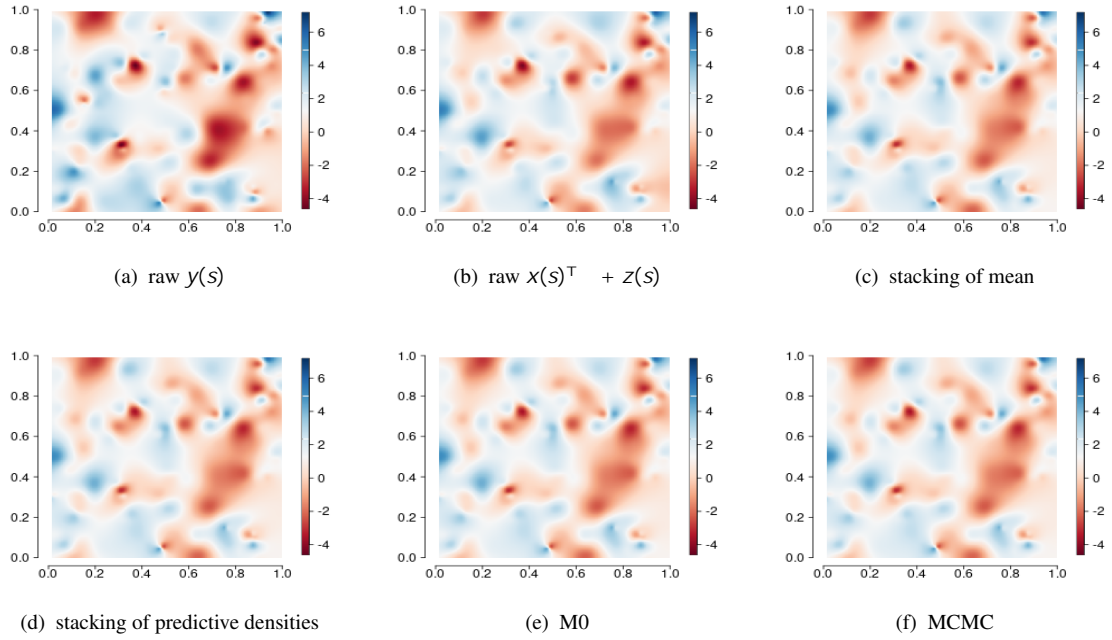


Fig. 7: Interpolated maps of (a) the response $y(s)$, (b) the denoised response $x(s) + z(s)$ and (c-f) the expected $y(s)$ on the $n_h = 100$ held out locations generated by all competing algorithms for one data set in the first simulation.

where $y = (y_1, \dots, y_G)^\top$, $\hat{Y}_G = [(\hat{Y}_1, \hat{Y}_G) : \dots : (\hat{Y}_{G-1}, \hat{Y}_G)]$, and $w = (w_1, \dots, w_{G-1})^\top$. And the QP problem has constraints $\sum_{g=1}^{G-1} w_g = 1$ and $w_g \geq 0$ for $g = 1, \dots, G-1$.

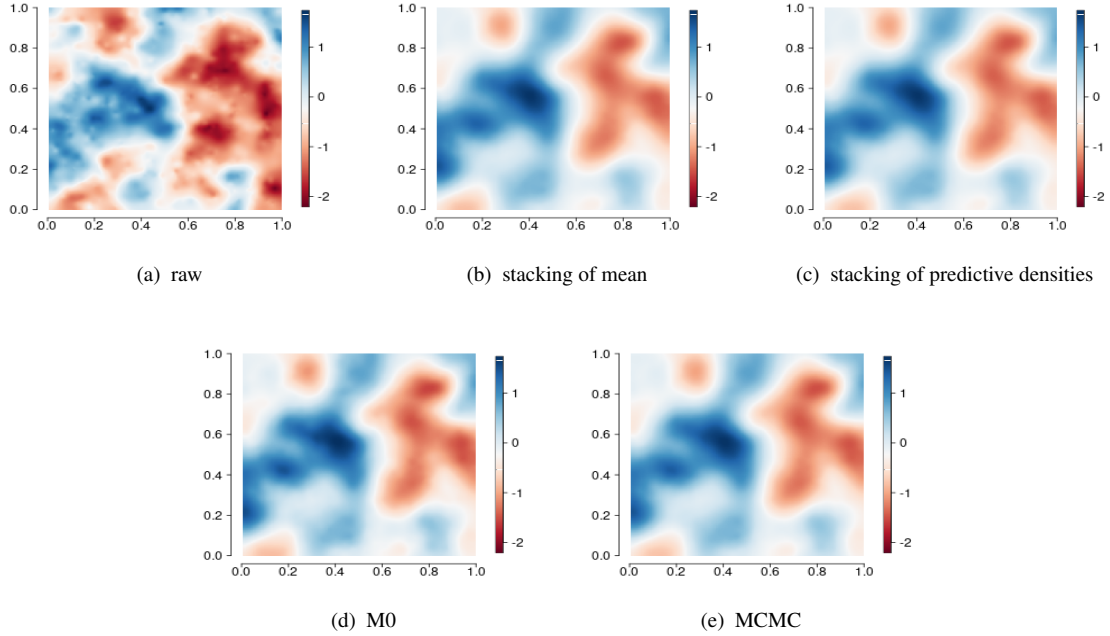


Fig. 8: Interpolated maps of (a) the latent process $Z(s)$ and (b-g) the expected $Z(s)$ on all $n = 900$ sampled locations generated by all competing algorithms for one data set in the first simulation. The $n = 900$ locations include both observed and unobserved locations

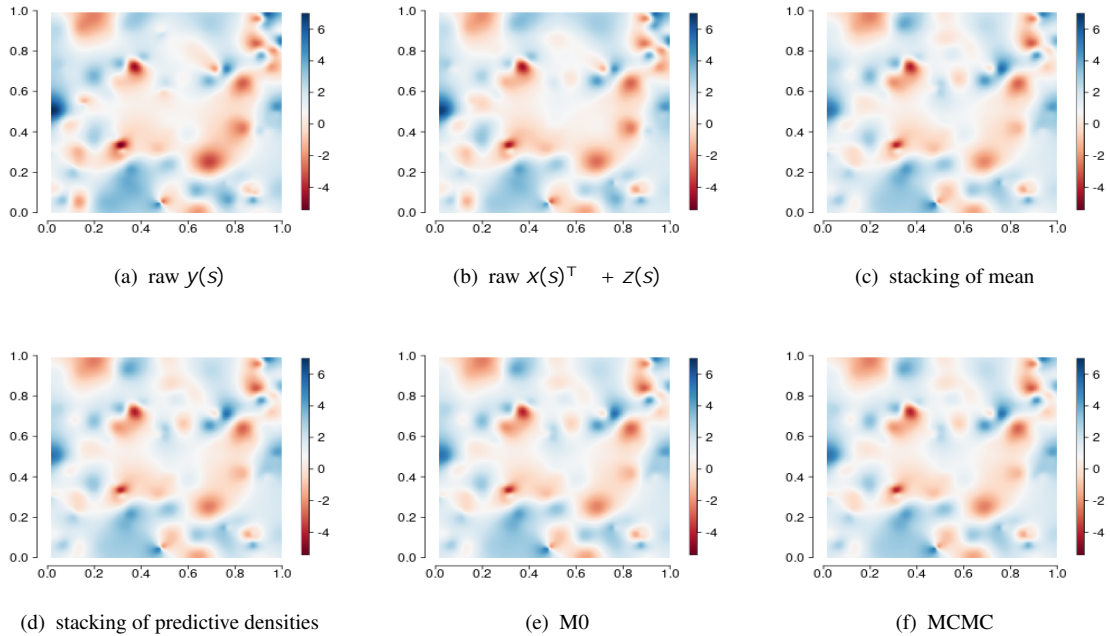


Fig. 9: Interpolated maps of (a) the response $y(s)$, (b) the denoised response $x(s) + Z(s)$ and (c-f) the expected $y(s)$ on the $n_h = 100$ held out locations generated by all competing algorithms for one data set in the second simulation.

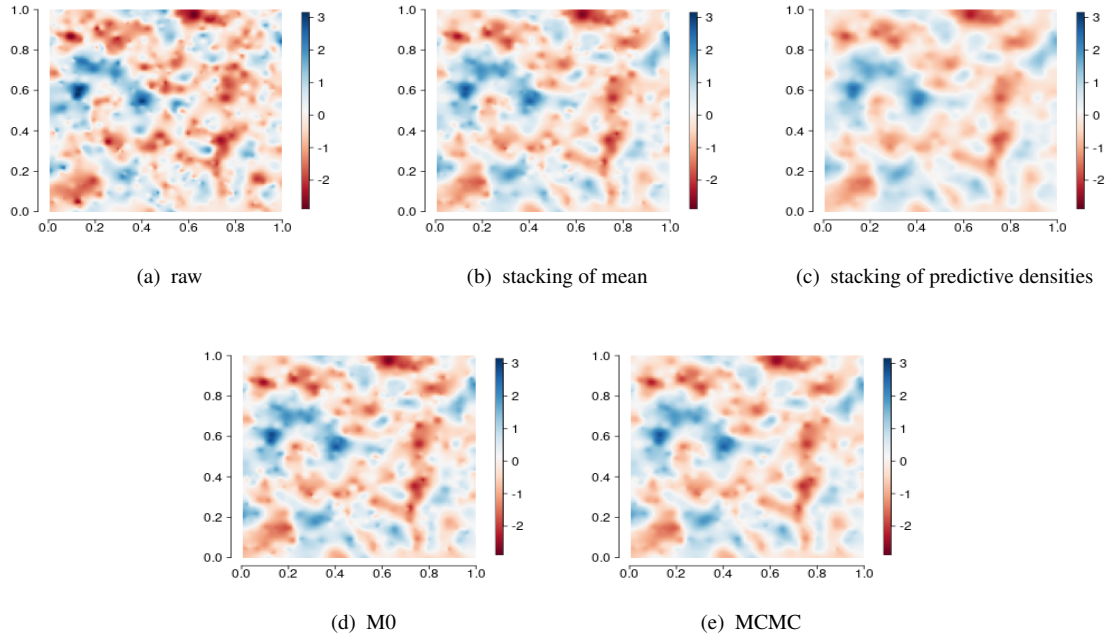
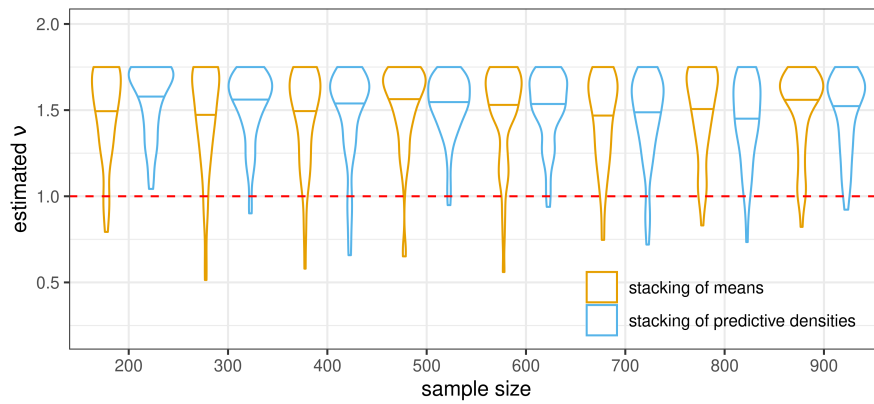
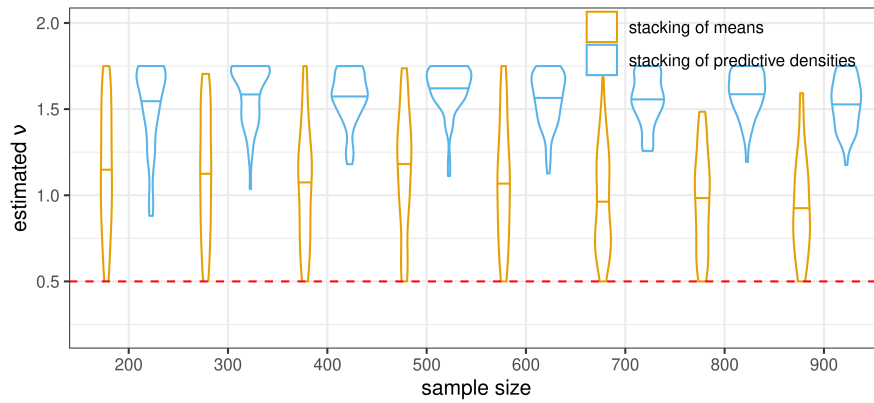


Fig. 10: Interpolated maps of (a) the latent process $Z(s)$ and (b-g) the expected $Z(s)$ on all $n = 900$ sampled locations generated by all competing algorithms for one data set in the second simulation. The $n = 900$ locations include both observed and unobserved locations



(a)



(b)

Fig. 11: Distributions of the estimated v in the first (a) and the second (b) simulation. The distribution of the counts are described through violin plots whose horizontal lines indicate the medians. The red dashed horizontal line indicates the actual value of v .

