



A fast regression via SVD and marginalization

Philip Greengard¹ · Andrew Gelman¹ · Aki Vehtari²

Received: 23 February 2021 / Accepted: 10 July 2021

© The Author(s), under exclusive licence to Springer-Verlag GmbH Germany, part of Springer Nature 2021

Abstract

We describe a numerical scheme for evaluating the posterior moments of Bayesian linear regression models with partial pooling of the coefficients. The principal analytical tool of the evaluation is a change of basis from coefficient space to the space of singular vectors of the matrix of predictors. After this change of basis and an analytical integration, we reduce the problem of finding moments of a density over $k + 2$ dimensions, to finding moments of a 2-dimensional density, where k is the number of coefficients. Moments can then be computed using, for example, MCMC, the trapezoid rule, or adaptive Gaussian quadrature. An evaluation of the SVD of the matrix of predictors is the dominant computational cost and is performed once during the precomputation stage. We demonstrate numerical results of the algorithm.

Keywords Bayesian Regression · Singular Value Decomposition · Marginalization · Fast Algorithms

1 Introduction

Linear regression is a ubiquitous tool for statistical modeling in a range of applications including social sciences, epidemiology, biochemistry, and environmental sciences (Gelman et al. 2013; Gelman and Hill 2007; Greenland 2000; Merlo et al. 2005; Bardini et al. 2017).

A common bottleneck for applied statistical modeling workflow is the computational cost of model evaluation. Since posterior distributions in statistical models are often high dimensional and computationally intractable, various techniques have been used to approximate posterior moments. Standard approaches often involve a variety of techniques including Markov chain Monte Carlo (MCMC) or using a suitable approximation of the posterior.

✉ Philip Greengard
pg2118@columbia.edu

¹ Columbia University, New York, USA

² Aalto University, Espoo, Finland

25 In this paper, we describe an approach for reducing the computational costs for a
 26 particular class of regression models — those that contain parameters $\theta \in \mathbb{R}^k$ such
 27 that θ has a normal prior and normal likelihood. These models represent only a subset
 28 of regression models that appear in applications. We focus our attention in this paper
 29 on normal-normal models because they have well known analytical properties and
 30 are more computationally tractable than the vast majority of multilevel models. A
 31 broader class of models, including logistic regression, contain distributions that are
 32 less amenable to the techniques of this paper and will require other analytical and
 33 computational tools. Mathematically, marginalization of normal-normal parameters
 34 is well-known and has been applied to the posterior by, for example, Lindley and
 35 Smith (1972). Our contribution is to provide a stable, accurate, and fast algorithm for
 36 marginalization.

37 The primary numerical tool used in the algorithm is the singular value decomposi-
 38 tion (SVD) of the data matrix. As a mathematical and statistical tool, SVD has been
 39 known since at least 1936 (see Eckart and Young (1936)). Use of the SVD as a practical
 40 and efficient numerical algorithm only started gaining popularity much later, with the
 41 first widely used scheme introduced in Golub and Kahan (1965). Due in large part to
 42 advances in computing power, use of the SVD as a tool in applied mathematics, statis-
 43 tics, and data science has been gaining significant popularity in recent years, however
 44 efficient evaluation of SVDs and related matrix decompositions is still an active area
 45 of research (see Hastie et al. 2015; Halko et al. 2011; Shamir et al. 2016).

46 Similar schemes to ours are used in the software packages lme4 (Bates et al. 2015)
 47 and INLA (Rue et al. 2017). There are several differences between the problems they
 48 address and their computational techniques, and those that we shall discuss here. While
 49 lme4 finds maximum likelihood and restricted maximum likelihood estimates, our goal
 50 is to find posterior moments. The software package INLA uses Laplace approximation
 51 on the posterior for a general choice of likelihood functions, whereas our algorithm
 52 is focused on fast and accurate solutions for only a particular class of densities: those
 53 with normal-normal parameters.

54 The approach presented in this paper analytically marginalizes the normal-normal
 55 parameters of a model using a change of variables. After marginalization, posterior
 56 moments can be computed using standard techniques on the lower-dimensional den-
 57 sity. In particular, for a model that contains $k + 2$ total variables, k of which are
 58 normal-normal, our scheme converts the problem of evaluating expectations of a den-
 59 sity in $k + 2$ dimensions to finding expectations of a 2-dimensional density. After
 60 marginalization, we evaluate the 2-dimensional posterior density in $O(k)$ operations.

61 We illustrate our scheme on the problem of evaluating the marginal expectations
 62 of the unnormalized density

$$63 \quad q(\sigma_1, \sigma_2, \beta) = \sigma_1^{-(k+1)} \sigma_2^{-n} \exp \left(-\gamma (\log(\sigma_1))^2 - \frac{\sigma_2^2}{2} \right. \\ 64 \quad \left. - \frac{\|X\beta - y\|^2}{2\sigma_2^2} - \frac{\|\beta\|^2}{2\sigma_1^2} \right), \quad (1)$$

65 where $\gamma > 0$ is a constant, $\sigma_1, \sigma_2 > 0$, and $\beta \in \mathbb{R}^k$. We assume that X is a fixed
 66 $n \times k$ matrix, $y \in \mathbb{R}^n$ is fixed, and the normalizing constant of (1) is unknown. For
 67 fixed $n, k \in \mathbb{N}$, the algorithm is nearly identical when X is an $n \times k$ matrix to when
 68 X is a $k \times n$ matrix. In the case where $k \gg n$, Kwon et al. (2011) also use SVD for
 69 marginalization. There are three main distinctions between their method and ours. (i)
 70 Our method applies to $n \times k$ matrices X for $k < n$ and $k > n$. (ii) We use the SVD
 71 to analytically compute conditional second moments with respect to β , not only first
 72 moments. (iii) While they use MCMC for computing posterior moments, we use a
 73 high-order quadrature scheme.

74 Using the standard notation of Bayesian models, density q is the unnormalized
 75 posterior of the model

$$\begin{aligned}
 \sigma_1 &\sim \text{lognormal}(0, \sqrt{\gamma}) \\
 \sigma_2 &\sim \text{normal}^+(0, 1) \\
 \beta &\sim \text{normal}(0, \sigma_1) \\
 y &\sim \text{normal}(X\beta, \sigma_2).
 \end{aligned}
 \tag{2}$$

77 In Appendix A, we include Stan code that can be used to sample from density (1) using
 78 MCMC. We also include Stan code that samples from the marginalized 2-dimensional
 79 posterior obtained via the algorithm of this paper.

80 Statistical model (2) is a standard model of Bayesian statistics and appears when
 81 seeking to model an outcome, y , as a linear combination of related predictors, the
 82 columns of X . In Gelman and Hill (2007), these models are described in detail and
 83 are used in the estimation of the distribution of radon levels in houses in Minnesota.
 84 See (Dias et al. 2013; Rover et al. 2020) for further examples.

85 Density (1) is also closely related to posterior densities that appear in genome-wide
 86 association studies (GWAS; see Zhu and Stephens 2017; Meuwissen, et al. 2001;
 87 Azevedo et al. 2015) which can be used to identify genomic regions containing genes
 88 linked with a specific trait, such as height. Using the notation of (1), each row of matrix
 89 X corresponds to a person, each column of X represents a genomic location, entries
 90 of X indicate genotypes, and y corresponds to the trait. Due to technical advances in
 91 genome sequencing over the last ten years, it is now feasible to collect large amounts
 92 of sequencing data. GWAS models can contain data on up to millions of people and
 93 often between hundreds and thousands of genome locations (see Linner et al. 2019).
 94 As a result, efficient computational tools are required for model evaluation.

95 The number of operations required by the scheme of this paper scales like $O(nk^2)$
 96 with a small constant. The key analytical tool is a change of variables of β such that
 97 the terms,

$$-\frac{1}{2\sigma_2^2} \|X\beta - y\|^2 - \frac{1}{2\sigma_1^2} \|\beta\|^2,
 \tag{3}$$

99 in (1) are converted to a diagonal quadratic form in \mathbb{R}^k . After that change of vari-
 100 ables, expectations over q are analytically converted from integrals over \mathbb{R}^{k+2} to
 101 integrals over \mathbb{R}^2 . The remaining 2-dimensional integrals can be computed to high
 102 accuracy using classical numerical techniques including, for example, adaptive Gaus-
 103 sian quadrature or even the 2-dimensional trapezoid rule.

104 The tools used in this paper to evaluate the expectations of (1) can also be used
 105 in the evaluation of expectations of multilevel and multigroup posterior distributions
 106 including, for example, the two-group posterior of the form,

$$q(\sigma_1, \sigma_2, \sigma_3, \beta) = \exp\left(-\frac{1}{2\sigma_1^2}\|X\beta - y\|^2 - \frac{1}{2\sigma_2^2}\sum_{i=1}^{k_1}\beta_i^2 - \frac{1}{2\sigma_3^2}\sum_{i=k_1+1}^{k_1+k_2}\beta_i^2\right), \quad (4)$$

108 where X is a $n \times k$ matrix, $y \in \mathbb{R}^n$, k_1 and k_2 are non-negative integers satisfying
 109 $k_1 + k_2 = k$, and $\sigma_1, \sigma_2, \sigma_3 > 0$.

110 The structure of this paper is as follows. In the following section we describe the
 111 analytic integration that transforms (1) from a $k + 2$ -dimensional problem to a 2-
 112 dimensional problem. Section 3 includes formulas that will allow for the evaluation
 113 of posterior moments using the 2-dimensional density. In Sects. 4 and 5 we provide
 114 formulas for evaluating covariances of (1). In Sect. 6, we discuss the numerical results
 115 of the implementation of the algorithm. Conclusions and generalizations of the algo-
 116 rithm of this paper are presented in Sect. 7. Appendix A provides Stan code that can
 117 be used to sample from (1), and Appendix B includes proofs of the formulas of this
 118 paper.

119 2 Analytic integration of β

120 In this section, we describe how we analytically marginalize the normal-normal param-
 121 eter β of density (1). We include proofs of all formulas in Appendix B.

122 We start by marginalizing β using a change of variables that converts the quadratic
 123 forms in (1) into diagonal quadratic forms. The resulting integral in the new variable,
 124 z , is Gaussian, and the coefficients of z_i and z_i^2 are available analytically. The change
 125 of variables is given by the right orthogonal matrix of the singular value decomposition
 126 (SVD) of X . That is, we set

$$z = V^t \beta \quad (5)$$

128 where the SVD of X is

$$X = UDV^t. \quad (6)$$

130 We define λ_i to be the i^{th} element of the diagonal of D . The elements of diagonal need
 131 not be sorted. After this change of variables, we obtain the following identity for the
 132 last two terms of (1). A proof can be found in Lemma 5 in Appendix B.

Formula 2.1

$$\begin{aligned} & -\frac{1}{2\sigma_2^2}\|X\beta - y\|^2 - \frac{1}{2\sigma_1^2}\|\beta\|^2 \\ & = a_0 + \sum_{i=1}^k a_{2,i} \left(z_i - \frac{a_{1,i}}{2a_{2,i}}\right)^2 + \frac{a_{1,i}^2}{4a_{2,i}} \end{aligned} \quad (7)$$

134 where

$$135 \quad a_{2,i} = \frac{\lambda_i^2}{2\sigma_2^2} + \frac{1}{2\sigma_1^2}, \quad (8)$$

$$136 \quad a_{1,i} = \frac{w_i}{\sigma_2^2}, \quad (9)$$

137 and

$$138 \quad a_0 = -\frac{y^t y}{2\sigma_2^2} \quad (10)$$

139 where

$$140 \quad w = V^t X^t y. \quad (11)$$

141 After performing the change of variables $z = V^t \beta$ and using (7), we now have an
 142 expression for density (1) in a form that allows us to use the well-known properties of
 143 a Gaussian with diagonal covariance. The following identity uses these properties and
 144 provides a formula for analytically reducing expectations of (1) from integrals over
 145 $k + 2$ dimensions to integrals over 2 dimensions. After the formula is applied, we have
 146 a new density, \tilde{q} , over only 2 dimensions. See Theorem 1 in Appendix B for a proof.

147 **Formula 2.2** For all $\sigma_1, \sigma_2 > 0$ we have

$$148 \quad \int_{\mathbb{R}^k} q(\sigma_1, \sigma_2, \beta) d\beta = \tilde{q}(\sigma_1, \sigma_2) \quad (12)$$

149 where $\tilde{q}(\sigma_1, \sigma_2)$ is defined by the formula

$$150 \quad \tilde{q}(\sigma_1, \sigma_2) = \sigma_1^{-(k+1)} \sigma_2^{-n} \exp\left(-\gamma \log^2(\sigma_1) - \frac{\sigma_2^2}{2}\right. \\ 151 \quad \left.+ a_0 + \sum_{i=1}^k \frac{a_{1,i}^2}{4a_{2,i}}\right) \prod_{i=1}^k \frac{1}{\sqrt{2a_{2,i}}} \quad (13)$$

152 where $a_{2,i}$ is defined in (8), $a_{1,i}$ is defined in (9), a_0 is defined in (10), and γ is a
 153 constant.

154 In (58) we provide a formula for \tilde{q} in the case where both scale parameters have
 155 half-normal priors.

156 **Remark 1** Certain Bayesian models might contain correlated priors on β that will
 157 result in posteriors such as (28) of Sect. 4. For such models, we perform the change
 158 of variables that uses the fact that two diagonal forms over β can be simultaneously
 159 diagonalized.

160 We include in Fig. 1 a plot of the density of q as a function of σ_1 and β_1 for fixed σ_2
 161 and randomly chosen X and y . Figure 2 shows a plot of q as a function of σ_2 and β
 162 for fixed σ_1 . Figure 3 provides an illustration of \tilde{q} , obtained after the change of variables
 163 and marginalization described in this section.

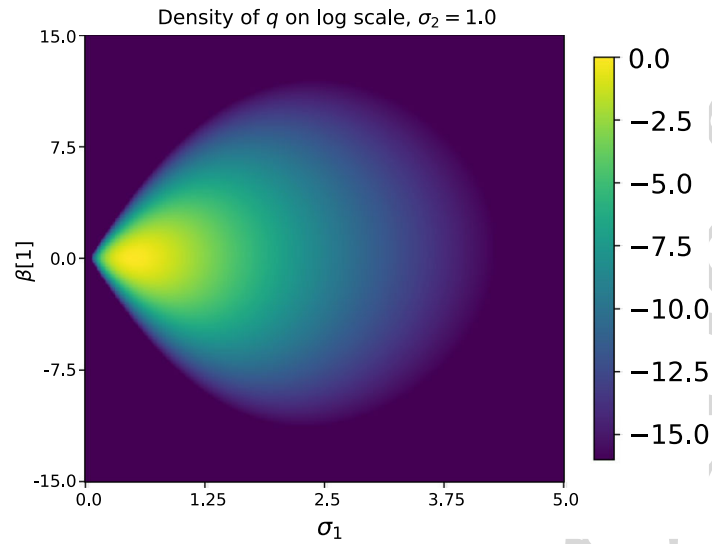


Fig. 1 Density of q (see (1)) with respect to σ_1 and β_1 , where $\gamma = 8$, $n = 100$, $k = 10$, and data were randomly generated

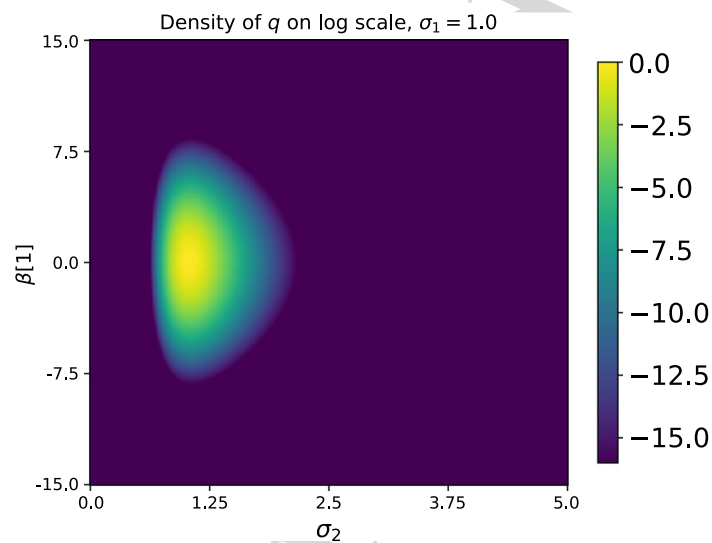


Fig. 2 Density of q (see (1)) with respect to σ_2 and β_1 , where $\gamma = 8$, $n = 100$, $k = 10$, and data were randomly generated

164 **3 Evaluation of posterior means**

165 Now that we have reduced the $k + 2$ -dimensional density q to the 2-dimensional
 166 density \tilde{q} , it remains to recover the posterior moments of q using \tilde{q} . We first observe
 167 that moments of σ_1 and σ_2 with respect to q are equivalent to moments of σ_1 and σ_2
 168 over \tilde{q} . That is,

169
$$\mathbb{E}_q(\sigma_1) = \mathbb{E}_{\tilde{q}}(\sigma_1) \tag{14}$$

170 and

171
$$\mathbb{E}_q(\sigma_2) = \mathbb{E}_{\tilde{q}}(\sigma_2). \tag{15}$$

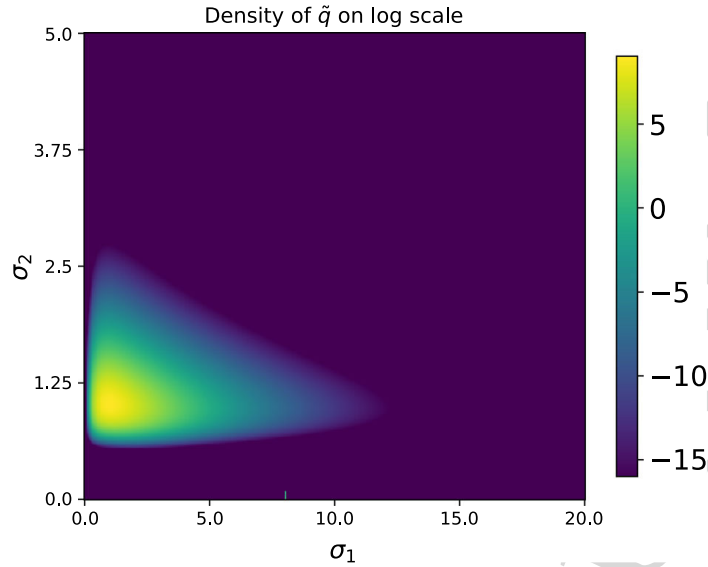


Fig. 3 Log density of \tilde{q} (see (13)) using the same q as Fig. 1, where $n = 100$, $k = 10$, and data were randomly generated

172 As for moments of β , we use (13) and standard properties of Gaussians to obtain
 173 the following formula.

174 **Formula 3.1** For all $\sigma_1, \sigma_2 > 0$,

$$175 \int_{\mathbb{R}^k} z_i q(\sigma_1, \sigma_2, \beta) d\beta = \frac{a_{1,i}}{2a_{2,i}} \tilde{q}(\sigma_1, \sigma_2) \quad (16)$$

176 where q is defined in (1), \tilde{q} is defined in (13), $a_{2,i}$ is defined in (8), and $a_{1,i}$ is defined
 177 in (9).

178 As an immediate consequence of (16), we are able to evaluate the posterior expectation
 179 of z as an expectation of a 2-dimensional density:

$$180 \mathbb{E}_q(z_i) = \mathbb{E}_{\tilde{q}}\left(\frac{a_{1,i}}{2a_{2,i}}\right). \quad (17)$$

181 We then transform those expectations back to expectations over the desired basis, β
 182 using the matrix V computed in (6). Specifically, using linearity of expectation and
 183 (17), we know

$$184 \begin{aligned} \mathbb{E}_q((\beta_1, \dots, \beta_k)^t) &= \mathbb{E}_q(VV^t(\beta_1, \dots, \beta_k)^t) \\ 185 &= V\mathbb{E}_q(V^t(\beta_1, \dots, \beta_k)^t) \\ 186 &= V\mathbb{E}_q((z_1, \dots, z_k)^t) \\ 187 &= V\mathbb{E}_{\tilde{q}}\left(\left(\frac{a_{1,1}}{2a_{2,1}}, \dots, \frac{a_{1,k}}{2a_{2,k}}\right)^t\right). \end{aligned} \quad (18)$$

188 **4 Covariance of β**

189 In addition to facilitating the rapid evaluation of posterior means, the change of vari-
 190 ables described in Sect. 2 is also useful for the evaluation of higher moments.

191 Equation (7) shows that after the change of variables from β to z , the resulting
 192 density is a Gaussian in z with a diagonal covariance matrix. Additionally, for each
 193 z_i , using Eq. (7) and standard properties of Gaussians, we have the following identity.

194 **Formula 4.1** For all $\sigma_1, \sigma_2 > 0$, we have

195
$$\int_{\mathbb{R}^k} (z_i - \mu_{z_i})^2 q(\sigma_1, \sigma_2, \beta) d\beta = (2a_{2,i})^{-1} \tilde{q}(\sigma_1, \sigma_2) \quad (19)$$

196 where μ_{z_i} is the expectation of z_i , \tilde{q} is defined in (13), and $a_{2,i}$ is defined in (8).

197 The second moments of the posterior of β are obtained as a linear transformation of
 198 the posterior variances of z . In particular, denoting the expectation of z by μ_z , we have

199
$$\begin{aligned} \mathbb{E}(\beta\beta^t) &= VV^t \mathbb{E}(zz^t) V^t \\ &= V \mathbb{E}(zz^t) V^t \\ &= V(E((z - \mu_z)(z - \mu_z)^t) + \mu_z \mu_z^t) V^t \end{aligned} \quad (20)$$

202 We observe that due to the independence of all z_i ,

203
$$\mathbb{E}((z - \mu_z)(z - \mu_z)^t) \quad (21)$$

204 is diagonal and we can therefore evaluate the $k \times k$ posterior covariance matrix of β by
 205 evaluating $\text{var}(z_i)$ and μ_{z_i} for $i = 1, \dots, k$ and then applying two orthogonal matrices.
 206 Specifically, combining Formula 4.1, (17), and (20), we obtain

207
$$\begin{aligned} \text{cov}(\beta) &= V \mathbb{E}_{\tilde{q}} \left(\left((2a_{2,1})^{-1}, \dots, (2a_{2,k})^{-1} \right)^t \right) V^t \\ &\quad + V E_{\tilde{q}} \left(\frac{a_{1,i}}{2a_{2,i}} \right) E_{\tilde{q}} \left(\frac{a_{1,i}}{2a_{2,i}} \right)^t V^t - \mu_{\beta} \mu_{\beta}^t. \end{aligned} \quad (22)$$

209 **5 Variance of σ_1 and σ_2**

210 Higher moments of σ_1 and σ_2 with respect to q can be evaluated directly as higher
 211 moments of σ_1 and σ_2 with respect to \tilde{q} . That is, for all $j \in \{2, 3, \dots\}$, we have

212
$$\mathbb{E}_q((\sigma_1 - \mu_{\sigma_1})^j) = \mathbb{E}_{\tilde{q}}((\sigma_1 - \mu_{\sigma_1})^j) \quad (23)$$

213 and

214
$$\mathbb{E}_q((\sigma_2 - \mu_{\sigma_2})^j) = \mathbb{E}_{\tilde{q}}((\sigma_2 - \mu_{\sigma_2})^j). \quad (24)$$

215 In particular, for $j = 2$, we obtain

$$216 \quad \text{var}_q(\sigma_1) = \text{var}_{\tilde{q}}(\sigma_1) \quad (25)$$

217 and

$$218 \quad \text{var}_q(\sigma_2) = \text{var}_{\tilde{q}}(\sigma_2). \quad (26)$$

Algorithm 1: *Evaluation of posterior expectations of normal-normal models*

- 1 Compute SVD of matrix X
 - 2 Compute w (see (11))
 - 3 Compute $V^t \mathbb{1}$ (see (9))
 - 4 Construct evaluator for density \tilde{q} of (13)
 - 5 Evaluate first and second moments with respect to \tilde{q} : $\mathbb{E}_{\tilde{q}}(\sigma_1)$, $\mathbb{E}_{\tilde{q}}(\sigma_2)$, $\mathbb{E}_{\tilde{q}}\left(\frac{a_{1,i}}{2a_{2,i}}\right)$
 - 6 Compute $\mathbb{E}(\beta)$ via formula (18)
-

219 6 Numerical experiments

220 Algorithm 1 was implemented in Fortran. We used the GFortran compiler on a 2.6
 221 GHz 6-Core Intel Core i7 MacBook Pro. All examples were run in double precision
 222 arithmetic. The matrix X and vector y were randomly generated as follows. Each entry
 223 of X was generated with an independent Gaussian with mean 0 and variance 1. The
 224 vector y was created by first randomly generating a vector $\beta \in \mathbb{R}^k$, each entry of
 225 which is an independent Gaussian with mean 0 and variance 1. The vector y was set to
 226 the value of $X\beta + \epsilon$ where $\epsilon \in \mathbb{R}^n$ contains standard normal iid entries. We generated
 227 y this way in order to ensure that the $\mathbb{E}(\beta_i)$ were not all small in magnitude. We set
 228 γ of (1) to 8 for all subsequent experiments and note that in practice the value of γ
 229 would be set according to some problem-specific knowledge.

230 In Table 1 and Fig. 5, we compare the performance of Algorithm 1 to two alterna-
 231 tive schemes for computing posterior expectations — one in which we analytically
 232 marginalize via Eq. (12) and then integrate the 2-dimensional density via MCMC
 233 using Stan. In the other, we use Stan’s MCMC sampling on the full $k + 2$ dimensional
 234 posterior. When using MCMC with Stan, we took 10,000 posterior draws. In Table 1
 235 and Fig. 5 we denote Algorithm 1 by “SVD-Trap”. The algorithm that uses Stan on
 236 the marginal 2-dimensional density is labeled “SVD-MCMC”, and “MCMC” corre-
 237 sponds to the algorithm that uses only MCMC sampling in Stan. We observe that both
 238 the time for evaluation and the accuracy is drastically improved when using Algorithm
 239 1 over full MCMC and MCMC with marginalization. In particular, for large n , the
 240 algorithm of this paper is faster by a factor of thousands compared to full MCMC via
 241 Stan.

242 In the appendix, we include Stan code to sample from the marginal density \tilde{q} of
 243 (13).

Table 1 Accuracy of evaluation of expectations of q (see (1)) using three different algorithms: (i) SVD-Trap: Algorithm 1 of this paper, (ii) SVD-MCMC: marginalization with MCMC integration of \tilde{q} using Stan, and (iii) MCMC: full MCMC integration of q using Stan

n	k	SVD-Trap max error	SVD-MCMC max error	MCMC max error
100	100	0.9×10^{-14}	0.4×10^{-4}	0.1×10^{-1}
200	100	0.9×10^{-14}	0.3×10^{-2}	0.8×10^{-2}
500	100	0.9×10^{-13}	0.2×10^{-2}	0.8×10^{-2}
1000	100	0.2×10^{-13}	0.6×10^{-3}	0.7×10^{-2}
5000	100	0.4×10^{-13}	0.2×10^{-3}	0.3×10^{-2}
10,000	100	0.2×10^{-13}	0.4×10^{-3}	0.2×10^{-2}

Table 2 Scaling of computation times for evaluation of expectations of q (see (1)) using Algorithm 1

n	k	max error	precompute time (s)	integrate time (s)	total (s)
50	5	0.22×10^{-13}	0.01	0.01	0.02
100	10	0.26×10^{-13}	0.02	0.01	0.03
500	20	0.30×10^{-13}	0.04	0.01	0.05
1000	50	0.34×10^{-13}	0.09	0.03	0.12
5000	100	0.37×10^{-13}	0.29	0.05	0.34
10000	500	0.26×10^{-13}	14	0.3	14.2
10,000	1000	0.39×10^{-13}	54	0.6	54.5

244 **Remark 2** In the numerical integration stage of algorithm 1, we use the trapezoid
 245 rule with 200 nodes in each direction. See Sect. C for a brief description of the 2-
 246 dimensional trapezoid rule. Because the integrand is smooth and vanishes near the
 247 boundary, convergence of the integral is super-algebraic when using the trapezoid rule
 248 (see Stoer and Bulirsch 1992). A rectangular grid with 200 points in each direction
 249 is satisfactory for obtaining approximately double precision accuracy. In problems
 250 with large numbers of non-normal-normal parameters, MCMC algorithms such as
 251 Hamiltonian Monte Carlo or other methods can be used.

252 In Tables 1 and 2, n and k represent the size of the $n \times k$ random matrix X .

253 The column labeled “max error” provides the maximum absolute error of the expect-
 254 ations of σ_1 , σ_2 , and β_i for $i \in \{1, 2, \dots, n\}$. The true solution was evaluated using
 255 trapezoid rule with 500 nodes in each direction in extended precision.

256 In Table 2, “Precompute time (s)” denotes the time in seconds of all computations
 257 until numerical integration. These times are dominated by the cost of SVD (36). The
 258 total time of the numerical integration in addition to the matrix-vector product (18) is
 259 given in “integrate time (s).” The final column of Table 2, “total time (s)”, provides
 260 the total time of precomputation and integration.

261 Notably, Table 2 demonstrates that the dominant cost of the algorithm of this paper is
 262 the SVD in the precomputation stage. Additionally, even for large regression problems
 263 with 10,000 observations and 1000 predictors, evaluation time is under a minute.

Fig. 4 Scaling of computation times for evaluation of posterior expectations of q (see (1)) using Algorithm 1 as a function of k with $n = 10,000$

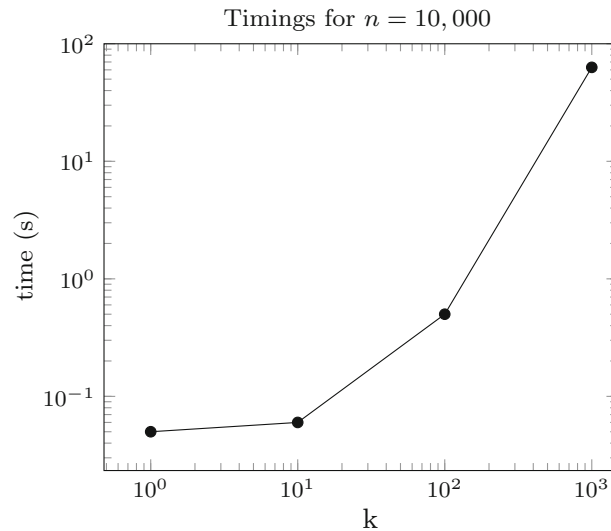
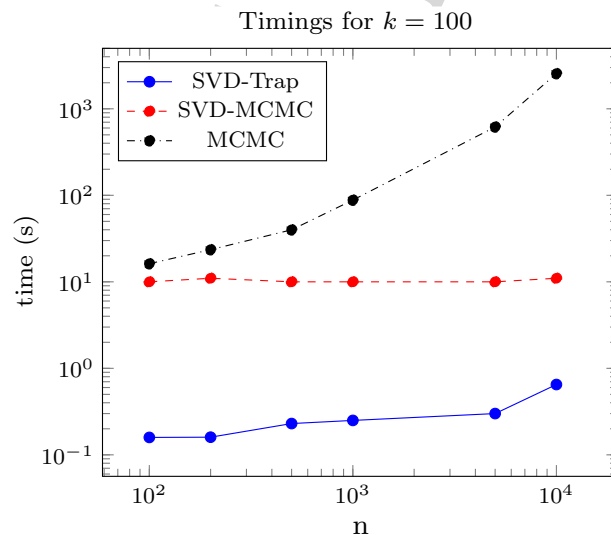


Fig. 5 Scaling of computation times for evaluation of posterior expectations of q (see (1)) as a function of sample size n with $k = 100$. The three algorithms compared are (i) SVD-Trap: Algorithm 1 of this paper, (ii) SVD-MCMC: marginalization with MCMC integration of \tilde{q} using Stan, and (iii) MCMC: full MCMC integration of q using Stan



264 7 Generalizations and conclusions

265 In this paper, we present a numerical scheme for the evaluation of the expectations of a
 266 particular class of distributions that appear in Bayesian statistics; posterior distributions
 267 of linear regression problems with normal-normal parameters.

268 The tools used in the numerical scheme of this paper generalize to several related
 269 classes of distributions that appear frequently in Bayesian statistics. We list several
 270 examples of posterior densities whose expectations can be evaluated using this method.
 271 1. The choice of priors for σ_1 , and σ_2 in this document were log normal and half-
 272 normal. This choice did not substantially impact the algorithm and can be generalized.
 273 Adaptive Gaussian quadrature (see, e.g. Trefethen 2020) can be used for the numerical
 274 integration step of the algorithm for a more general choice of prior on σ_1 and σ_2 .

275 2. Regression problems with multiple groups such as the two-group model with pos-
276 terior

$$277 \exp\left(-\frac{1}{2\sigma_1^2}\|X\beta - y\|^2 - \frac{1}{2\sigma_2^2}\sum_{i=1}^{k_1}\beta_i^2 - \frac{1}{2\sigma_3^2}\sum_{i=k_1+1}^{k_1+k_2}\beta_i^2\right) \quad (27)$$

278 where X is a $n \times k$ matrix, $y \in \mathbb{R}^n$, and k_1 and k_2 are non-negative integers satisfying
279 $k_1 + k_2 = k$.

280 3. Regression problems with correlated priors on β :

$$281 \exp\left(-\frac{1}{2\sigma_2^2}\|X_1\beta - y\|^2 - \frac{1}{2\sigma_1^2}\|X_2\beta\|^2\right) \quad (28)$$

282 For regression problems with large numbers of non-normal-normal parameters,
283 marginal expectations can be computed using, for example, MCMC in Stan. For such
284 problems, the algorithm of this paper would convert an MCMC evaluation from $k + m$
285 dimensions to m dimensions, where k is the number of normal-normal parameters.

286 **Acknowledgements** The authors are grateful to Ben Bales, Bob Carpenter, and Mitzi Morris for many
287 useful discussions.

288 A Code

289 The following Stan code allows for sampling from the distribution corresponding to
290 the probability density function proportional to (1).

```
291 data {
292   int n;
293   int k;
294   vector[n] y;
295   matrix[n,k] X;
296 }
297 parameters {
298   real<lower=0> sigma1;
299   real<lower=0> sigma2;
300   vector<offset=0, multiplier=sigma1>[k] beta;
301 }
302 model {
303   y ~ normal(X*beta, sigma2);
304   beta ~ normal(0, sigma1);
305   sigma1 ~ lognormal(0, 0.25);
306   sigma2 ~ normal(0, 1);
307 }
```

308 The following Stan program samples from the marginal density \tilde{q} (see (13)). The
309 data input `yty` corresponds to $y^t y$ of (10), `lam` is the vector of singular values of X ,

310 and w is the vector w in Eq. (11). We include R code for computing yty , lam , and w
 311 after the following Stan code.

312 functions {
 313 real q_tilde_lpdf(real sig1, real sig2, vector w,
 314 vector lam, real yty, int k,
 315 int n) {
 316 vector[min(n,k)] a2 = lam^2/(sig2^2) + 1/(sig1^2);
 317 real sol = sum(w^2 ./a2)/2/sig2^4 - sum(log(a2))/2
 318 - yty/(2*sig2^2);
 319 sol += -min(n,k)*log(sig1) - n*log(sig2);
 320 return sol;
 321 }
 322 }
 323 data {
 324 int n;
 325 int k;
 326 vector[min(n,k)] w;
 327 vector[min(n,k)] lam;
 328 real yty;
 329 matrix[min(n,k),k] V;
 330 }
 331 parameters {
 332 real<lower=0> sigma1;
 333 real<lower=0> sigma2;
 334 }
 335 model {
 336 sigma1 ~ q_tilde(sigma2, w, lam, yty, k, n);
 337 sigma1 ~ lognormal(0, 0.25);
 338 sigma2 ~ normal(0, 1);
 339 }
 340 generated quantities {
 341 vector[k] beta;
 342 {
 343 vector[min(n,k)] zvar = 1 ./ (2*(lam^2 ./ (2*sigma2^2)
 344 + 1/(2*sigma1^2)));
 345 vector[min(n,k)] zmu = w./sigma2^2 .* zvar;
 346 vector[min(n,k)] z =
 347 to_vector(normal_rng(zmu, sqrt(zvar)));
 348 beta = V * z;
 349 }
 350 }

351 The following is a sample of code from R that can be used for the precomputation
 352 stage of Algorithm 1.

353 `udv <- svd(X)`

```

354 V <- udv$d
355 lam <- as.vector(udv$d)
356 w <- t(V) %*% t(X) %*% Y
357 w <- as.vector(w)
358 yty <- t(Y) %*% Y
359 yty <- yty[1]

```

B Proofs

In this appendix, we include proofs of the formulas provided in this paper. For increased readability, this appendix is self-contained.

B.1 Mathematical preliminaries and notation

In this section, we introduce notation and elementary mathematical identities that will be used throughout the remainder of this section.

We define $C \in \mathbb{R}$ by the Eq.

$$C = \int_{\sigma_1 \in \mathbb{R}^+} \int_{\sigma_2 \in \mathbb{R}^+} \int_{\beta \in \mathbb{R}^k} q(\sigma_1, \sigma_2, \beta) d\beta d\sigma_2 d\sigma_1, \tag{29}$$

and define $\mathbb{E}(\sigma_1)$, $\mathbb{E}(\sigma_2)$, and $\mathbb{E}(\beta_i)$ by the formulas

$$\mathbb{E}(\sigma_1) = \frac{1}{C} \int_{\sigma_1 \in \mathbb{R}^+} \int_{\sigma_2 \in \mathbb{R}^+} \int_{\beta \in \mathbb{R}^k} \sigma_1 q(\sigma_1, \sigma_2, \beta) d\beta d\sigma_2 d\sigma_1, \tag{30}$$

$$\mathbb{E}(\sigma_2) = \frac{1}{C} \int_{\sigma_1 \in \mathbb{R}^+} \int_{\sigma_2 \in \mathbb{R}^+} \int_{\beta \in \mathbb{R}^k} \sigma_2 q(\sigma_1, \sigma_2, \beta) d\beta d\sigma_2 d\sigma_1, \tag{31}$$

and

$$\mathbb{E}(\beta_i) = \frac{1}{C} \int_{\sigma_1 \in \mathbb{R}^+} \int_{\sigma_2 \in \mathbb{R}^+} \int_{\beta \in \mathbb{R}^k} \beta_i q(\sigma_1, \sigma_2, \beta) d\beta d\sigma_2 d\sigma_1 \tag{32}$$

for $i \in \{1, 2, \dots, k\}$.

We provide algorithms for the evaluation of (29), (30), (31), and (32).

We will be denoting by $\mathbb{1}$ the vector of ones

$$\mathbb{1} = (1, 1, \dots, 1)^t. \tag{33}$$

We denote the i^{th} component of a vector v by v_i .

The following two well-known identities give the normalizing constant and expectation of a Gaussian distribution.

Lemma 1 For all $\sigma > 0$ we have

$$\sqrt{2\pi}\sigma = \int_{\mathbb{R}} e^{-\frac{(\beta-\mu)^2}{2\sigma^2}} d\beta \tag{34}$$

382 **Lemma 2** For all μ in \mathbb{R} and $\sigma > 0$, we have

$$383 \quad \mu\sqrt{2\pi}\sigma = \int_{\mathbb{R}} \beta e^{-\frac{(\beta-\mu)^2}{2\sigma^2}} d\beta \quad (35)$$

384 **B.2 Analytic integration of β**

385 We denote the singular value decomposition of X by

$$386 \quad X = UDV^t \quad (36)$$

387 where U is an orthogonal $n \times k$ matrix, V is an orthogonal $k \times k$ matrix, and D is a
388 $k \times k$ diagonal matrix. We define $z \in \mathbb{R}^k$ by the formula

$$389 \quad z = V^t \beta. \quad (37)$$

390 The following lemma, which will be used in the proof of Lemma 5, gives an expression
391 for the second to last term of the exponent in (1) after a change of variables.

392 **Lemma 3** For all $\beta \in \mathbb{R}^k$, and $y \in \mathbb{R}^n$,

$$393 \quad -\frac{1}{2\sigma_2^2} \|X\beta - y\|^2 = -\frac{y^t y}{2\sigma_2^2} + \sum_{i=1}^k -\frac{\lambda_i^2}{2\sigma_2^2} z_i^2 + \frac{w_i}{\sigma_2^2} z_i \quad (38)$$

394 where

$$395 \quad w = V^t X^t y, \quad (39)$$

396 z is defined in (37), and λ_i is the i^{th} entry on the diagonal of D (see (36)).

397 **Proof** Clearly,

$$398 \quad \|X\beta - y\|^2 = \beta^t X^t X \beta - 2y^t X \beta + y^t y. \quad (40)$$

399 Substituting (36) and (37) into (40), we obtain

$$400 \quad \begin{aligned} \|X\beta - y\|^2 &= \beta^t (UDV^t)^t (UDV^t) \beta - 2y^t X V V^t \beta + y^t y \\ &= (\beta^t V) D^2 (V^t \beta) - 2y^t (V^t X^t)^t z + y^t y. \end{aligned} \quad (41)$$

401 where z is defined in (37). Substituting (39) and (37) into (41), we have

$$402 \quad \|X\beta - y\|^2 = z^t D^2 z - 2w^t z + y^t y \quad (42)$$

403 Equation (38) follows immediately from (42). □

404 The following lemma provides an equation for the last term of the exponent in (1).
405 The identity will be used in Lemma 5.

406 **Lemma 4** For all $\sigma_1 > 0$,

$$407 \quad -\frac{\|\beta\|^2}{2\sigma_1^2} = \sum_{i=1}^k -\frac{z_i^2}{2\sigma_1^2} \quad (43)$$

408 where $\beta \in \mathbb{R}^k$, z is defined in (37), and V is defined in (36).

409 **Proof** Clearly,

$$410 \quad \frac{\|\beta\|^2}{2\sigma_1^2} = \frac{1}{2\sigma_1^2} (Vz)^t (Vz) = \frac{z^t z}{2\sigma_1^2} \quad (44)$$

411 where V is defined in (36). Equation (43) follows immediately from (44). \square

412 The following formula combines Lemmas 3 and 4 to convert the final two terms of
413 (1) into a Gaussian in k dimensions.

414 **Lemma 5**

$$415 \quad -\frac{\|X\beta - y\|^2}{2\sigma_2^2} - \frac{\|\beta\|^2}{2\sigma_1^2} = a_0 + \sum_{i=1}^k a_{2,i} \left(z_i - \frac{a_{1,i}}{2a_{2,i}}\right)^2 + \frac{a_{1,i}^2}{4a_{2,i}} \quad (45)$$

416 where

$$417 \quad a_{2,i} = \frac{\lambda_i^2}{2\sigma_2^2} + \frac{1}{2\sigma_1^2}, \quad (46)$$

$$418 \quad a_{1,i} = \frac{w_i}{\sigma_2^2} \quad (47)$$

419 and

$$420 \quad a_0 = -\frac{y^t y}{2\sigma_2^2} \quad (48)$$

421 where z is defined in (37), w is defined in (39) and V is defined in (36).

422 **Proof** By combining Lemmas 3 and 4, we have

$$423 \quad -\frac{1}{2\sigma_2^2} \|X\beta - y\|^2 - \frac{1}{2\sigma_1^2} \|\beta\|^2 = a_0 + \sum_{i=1}^k \left(a_{1,i} z_i - a_{2,i} z_i^2\right). \quad (49)$$

424 We obtain Eq. (45) by completing the square in Eq. (49). \square

425 The following theorem is the principal analytical apparatus of this note. It provides
426 a formula for the k -dimensional integrals that appear in (29), (30), and (31).

427 **Theorem 1** For all $\sigma_1, \sigma_2 > 0$

$$428 \quad \int_{\mathbb{R}^k} q(\sigma_1, \sigma_2, \beta) d\beta = \tilde{q}(\sigma_1, \sigma_2) \quad (50)$$

429 where $\tilde{q}(\sigma_1, \sigma_2)$ is defined by the formula

$$430 \quad \tilde{q}(\sigma_1, \sigma_2) = \sigma_1^{-(k+1)} \sigma_2^{-n} \exp\left(-\log^2(\sigma_1) - \frac{\sigma_2^2}{2}\right) \\ 431 \quad + a_0 + \sum_{i=1}^k \frac{a_{1,i}^2}{4a_{2,i}} \sqrt{2\pi}^k \prod_{i=1}^k \frac{1}{\sqrt{2a_{2,i}}} \quad (51)$$

432 where $a_{2,i}$ is defined in (46), $a_{1,i}$ is defined in (47) and a_0 is defined in (48).

433 **Proof** Using (1), clearly

$$434 \quad \int_{\mathbb{R}^k} q(\sigma_1, \sigma_2, \beta) d\beta = \sigma_1^{-(k+1)} \int_{\mathbb{R}^k} \exp\left(-\log^2(\sigma_1) - \frac{\sigma_2^2}{2}\right) \\ 435 \quad - \frac{1}{2\sigma_2^2} \|X\beta - y\|^2 - \frac{1}{2\sigma_1^2} \|\beta\|^2 \Big) d\beta \quad (52)$$

436 Performing the change of variables (37) and substituting (45) into (52), we have

$$437 \quad \int_{\mathbb{R}^k} q(\sigma_1, \sigma_2, \beta) d\beta = \sigma_1^{-(k+1)} \exp\left(-\log^2(\sigma_1) - \frac{\sigma_2^2}{2} + a_0\right) \\ 438 \quad + \sum_{i=1}^k \frac{a_{1,i}^2}{4a_{2,i}} \int_{\mathbb{R}^k} \exp\left(\sum_{i=1}^k a_{2,i} \left(z_i - \frac{a_{1,i}}{2a_{2,i}}\right)^2\right) dz \quad (53)$$

439 Since the integrand on the right side of (53) is a Gaussian in z_i , Eq. (50) follows from
440 applying Lemma 1 to (53). □

441 **Remark 3** When adjusting the priors on the scale parameter to both become half-
442 normal, we have the model

$$443 \quad \sigma_1 \sim \text{normal}^+(0, 1) \quad (54)$$

$$444 \quad \sigma_2 \sim \text{normal}^+(0, 1) \quad (55)$$

$$445 \quad \beta \sim \text{normal}(0, \sigma_1) \quad (56)$$

$$446 \quad 447 \quad y \sim \text{normal}(X\beta, \sigma_2). \quad (57)$$

448 For the corresponding posterior, we note that \tilde{q} becomes

$$449 \quad \int_{\mathbb{R}^k} q(\sigma_1, \sigma_2, \beta) d\beta = \exp\left(-\frac{\sigma_1^2 + \sigma_2^2}{2} + a_0\right) \\ 450 \quad + \sum_{i=1}^k \frac{a_{1,i}^2}{4a_{2,i}} \int_{\mathbb{R}^k} \exp\left(\sum_{i=1}^k a_{2,i} \left(z_i - \frac{a_{1,i}}{2a_{2,i}}\right)^2\right) dz \quad (58)$$

451 The following theorem provides a formula for the expectation of z (see (37)). We
 452 use this formula, in combination with an orthogonal transformation, to obtain the
 453 expectation of β .

454 **Theorem 2** For all $\sigma_1 > 0$ and $\sigma_2 \in \mathbb{R}$,

$$455 \int_{\mathbb{R}^k} (V^t x)_i q(\sigma_1, \sigma_2, \beta) d\beta = \frac{a_{1,i}}{2a_{2,i}} \tilde{q}(t) \quad (59)$$

456 where q is defined in (1), \tilde{q} is defined in (51), $a_{2,i}$ is defined in (46), $a_{1,i}$ is defined in
 457 (47), a_0 is defined in (48).

458 **Proof** Combining (53) and (37), we have

$$459 \int_{\mathbb{R}^k} (V^t \beta)_i q(\sigma_1, \sigma_2, \beta) d\beta$$

$$460 = \exp\left(-\log^2(\sigma_1) - \frac{\sigma_2^2}{2} + a_0 + \sum_{i=1}^k \frac{a_{1,i}^2}{4a_{2,i}}\right)$$

$$461 \times \int_{\mathbb{R}^k} z_i \exp\left(\sum_{i=1}^k a_{2,i} \left(z_i - \frac{a_{1,i}}{2a_{2,i}}\right)^2\right) dz. \quad (60)$$

462 Applying Lemma 2 to (60), we obtain (59). □

463 C Trapezoid rule

464 The trapezoid rule (see, e.g. Stoer and Bulirsch 1992) is a quadrature scheme that is
 465 used to approximate the integral

$$466 \int_a^b f(x) dx \quad (61)$$

467 with the sum

$$468 \sum_{k=1}^n \frac{f(x_{k-1}) + f(x_k)}{2} \Delta_x \quad (62)$$

469 where $\Delta_x = (b - a)/(n - 1)$ and

$$470 x_k = a + k \frac{b - a}{n} \quad (63)$$

474 for $k = 0, \dots, n$. In the 2-dimensional analogue of the trapezoid rule we approximate
475 the integral

$$476 \int_c^d \int_a^b f(x, y) dx dy \quad (64)$$

478 with the sum

$$479 \sum_{k=1}^n \frac{g(y_{k-1}) + g(y_k)}{2} \Delta_y \quad (65)$$

481 where

$$482 g(y) = \sum_{k=1}^m \frac{f(x_{k-1}, y) + f(x_k, y)}{2} \Delta_x \quad (66)$$

484 and

$$485 \Delta_y = (d - c)/(n - 1), \quad (67)$$

$$486 y_k = c + k \frac{d - c}{n}, \quad (68)$$

$$487 \Delta_x = (b - a)/(m - 1), \quad (69)$$

$$488 x_k = a + k \frac{b - a}{m}. \quad (70)$$

490 References

- 491 Bardini R, Politano G, Benso A, Di Carlo S (2017) Multi-level and hybrid modelling approaches for systems
492 biology. *Comput Struct Biotechnol J* 15:396–402
- 493 Bates D, Martin M, Ben B, Steve W (2015) Fitting linear mixed-effects models using lme4. *J Stat Softw*
494 Carpenter B, Gelman A, Hoffman MD, Lee D, Goodrich B, Betancourt M, Brubaker M, Guo J, Li P, Riddell
495 A (2017) Stan: a probabilistic programming language. *J Stat Softw* 76(1):1–32
- 496 Dias S, Sutton AJ, Welton NJ, Ades AE (2013) Evidence synthesis for decision making 3: heterogeneity-
497 subgroups, meta-regression, bias, and bias-adjustment. *Med Decis Mak* 33(5):618–640
- 498 Eckart C, Young G (1936) The approximation of one matrix with another of lower rank. *Psychometrika* 1:3
- 499 Ferreira AC et al (2015) Ridge, lasso and Bayesian additive-dominance genomic models. *BMC Genet*
500 16:105
- 501 Gelman A, Hill J (2007) *Data analysis using regression and multilevel/hierarchical models*. Cambridge
502 University Press, Cambridge, UK
- 503 Gelman A, Carlin JB, Stern SH, Dunson BD, Vehtari A, Rubin BD (2013) *Bayesian Data Analysis*, 3rd
504 edn. CRC, New York, U.S
- 505 Golub G, Kahan W (1965) Calculating the singular values and psuedo-inverse of a matrix. *J SIAM Numer*
506 *Anal* 2:3
- 507 Greenland S (2000) Principles of multilevel modelling. *Int J Epidemiol* 29(1):158–167
- 508 Halko N, Martinsson PG, Tropp JA (2011) Finding structure with randomness: probabilistic algorithms for
509 constructing approximate matrix decompositions. *SIAM Rev* 53:2
- 510 Hastie T, Rahul M, Lee JD, Zadeh R (2015) Matrix completion and low-rank SVD via fast alternating least
511 squares. *J Mach Learn Res* 16(1):3367–3402

- 512 Kwon S, Yan X, Cui J, Yao J, Yang K, Tsiand D, Li X, Rotter J, Guo X (2011) Application of Bayesian
 513 regression with singular value decomposition method in association studies for sequence data. *BMC*
 514 *Proc* 5:9
- 515 Lindley DV, Smith AFM (1972) Bayes estimates for the linear model. *J Royal Stat Soc B* 34:1–41
- 516 Linner K et al (2019) Genome-wide association analyses of risk tolerance and risky behaviors in over 1
 517 million individuals identify hundreds of loci and shared genetic influences. *Nat Genet* 51:2
- 518 Merlo J, Chaix B, Yang M, Lynch J, Rastam L (2005) A brief conceptual tutorial of multilevel analysis in
 519 social epidemiology: linking the statistical concept of clustering to the idea of contextual phenomenon.
 520 *J Epidemiol Community Health* 59:3367–3402
- 521 Meuwissen TH, Hayes BJ, Goddard ME (2001) Prediction of total genetic value using genome-wide dense
 522 marker maps. *Genet* 157:4
- 523 Rover C., Ralf B., Sofia D., Christopher H.S., Heinz S., Sibylle S., Sebastian W., Tim F. (2020) “On
 524 weakly informative prior distributions for the heterogeneity parameter in Bayesian random-effects
 525 meta-analysis.” [arXiv:2007.08352v3](https://arxiv.org/abs/2007.08352v3)
- 526 Rue H, Riebler A, Sorbye SH, Illian JB, Simpson DP, Lindgren FK (2017) Bayesian computing with INLA:
 527 a review. *Annual Rev Stat Appl* 4:395–421
- 528 Shamir O. (2016) Fast Stochastic Algorithms for SVD and PCA: Convergence Properties and Convexity,
 529 Proceedings of the 33rd ICML, New York, NY,
- 530 Stoer J, Bulirsch R (1992) Introduction to numerical analysis, 2nd edn. Springer-Verlag, Berlin
- 531 Trefethen LN (2020) Approximation Theory and Approximation Practice. SIAM, Extended
- 532 Xiang Z, Matthew S (2017) Bayesian large-scale multiple regression with summary statistics from genome-
 533 wide association studies. *Annals Appl Stat* 11:3

534 **Publisher’s Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps
 535 and institutional affiliations.