

A Flexible Empirical Bayes Approach to Multiple Linear Regression, and Connections with Penalized Regression

Youngseok Kim

*Department of Statistics
University of Chicago
Chicago, IL 60637, USA*

YOUNGSEOK@UCHICAGO.EDU

Wei Wang

*Department of Statistics
University of Chicago
Chicago, IL 60637, USA*

WEIWANG@GALTON.UCHICAGO.EDU

Peter Carbonetto

*Research Computing Center and Department of Human Genetics
University of Chicago
Chicago, IL 60637, USA*

PCARBO@UCHICAGO.EDU

Matthew Stephens

*Department of Statistics and Department of Human Genetics
University of Chicago
Chicago, IL 60637, USA*

MSTEPHENS@UCHICAGO.EDU

Abstract

We introduce a new empirical Bayes approach for large-scale multiple linear regression. Our approach combines two key ideas: (i) the use of flexible “adaptive shrinkage” priors, which approximate the nonparametric family of scale mixture of normal distributions by a finite mixture of normal distributions; and (ii) the use of variational approximations to efficiently estimate prior hyperparameters and compute approximate posteriors. Combining these two ideas results in fast and flexible methods, with computational speed comparable to fast penalized regression methods such as the Lasso, and with competitive prediction accuracy across a wide range of scenarios. Further, we provide new results that establish conceptual connections between our empirical Bayes methods and penalized methods. Specifically, we show that the posterior mean from our method solves a penalized regression problem, with the form of the penalty function being learned from the data by directly solving an optimization problem (rather than being tuned by cross-validation). Our methods are implemented in an R package, `mr.ash.alpha`, available from <https://github.com/stephenslab/mr.ash.alpha>.

Keywords: Empirical Bayes, variational inference, normal means, penalized linear regression, nonconvex optimization

1 Introduction

Multiple linear regression is one of the oldest statistical methods for relating an outcome variable to predictor variables, dating back at least to the eighteenth century (Stigler, 1984). In recent decades, data sets have grown rapidly in size, with the number of predictor variables often exceeding the number of observations. Fitting even simple models such as

multiple linear regression to large data sets raises interesting research questions. These include questions about regularization and prediction (e.g., how to estimate the parameters to optimize out-of-sample prediction accuracy), and questions about variable selection and inference (e.g., how to choose the coefficients in the regression that are non-zero). In this paper, we focus on the former.

Many different approaches for regularization and prediction have been proposed. Most fall into one of two types: penalized linear regression (PLR) methods based on different penalized least-squares criteria (e.g., Hoerl and Kennard, 1970; Tibshirani, 1996; Fan and Li, 2011; Miller, 2002; Zou and Hastie, 2005; Zhang, 2010; Hazimeh and Mazumder, 2020; Amir et al., 2021) and Bayesian approaches based on different priors and computational methods (e.g., Mitchell and Beauchamp, 1988; George and McCulloch, 1993; Meuwissen et al., 2001; Park and Casella, 2008; Hans, 2009; Carvalho et al., 2010; Li and Lin, 2010; Griffin and Brown, 2010; Guan and Stephens, 2011; Habier et al., 2011; Carbonetto and Stephens, 2012; Zhou et al., 2013; Drugowitsch, 2013; Wang et al., 2020; Ročková and George, 2018; Ray and Szabó, 2022; Zabad et al., 2023).

These different approaches have different strengths and weaknesses. For example, ridge regression (an L_2 -penalty on the coefficients; Hoerl and Kennard, 1970; Tikhonov, 1963) is simple, involving a convex optimization problem and a single tuning parameter, and performs well in “dense” settings (many predictors with non-zero effects). However, it does not do well in “sparse” settings where a small number of non-zero coefficients dominate. The Lasso (an L_1 -penalty on the coefficients; Tibshirani, 1996) is similarly computationally convenient, and behaves better than ridge regression in sparse settings. However, prediction accuracy of the Lasso is limited by its tendency to “overshrink” large effects (e.g., Su et al., 2017). The Elastic Net (Zou and Hastie, 2005) combines some of the advantages of ridge regression and the Lasso, and in its most general form includes both as special cases; however, the Elastic Net also introduces an additional tuning parameter that results in a non-trivial additional computation expense.

Nonconvex penalties—examples include the L_0 -penalty (Miller, 2002; Hazimeh and Mazumder, 2020), the smoothly clipped absolute deviation (SCAD) penalty (Fan and Li, 2011) and the minimax concave penalty (MCP) (Zhang, 2010)—can also give better prediction performance in sparse settings, but this comes with the challenge of solving a nonconvex optimization problem.

Bayesian methods, by using flexible priors, have the potential to achieve excellent prediction accuracy in both sparse and dense settings (e.g., Park and Casella, 2008; Hans, 2009; Griffin and Brown, 2010; Li and Lin, 2010; Guan and Stephens, 2011; Zhou et al., 2013; Zeng et al., 2018), but have some practical drawbacks; notably, model fitting typically involves performing Markov chain Monte Carlo (MCMC) with a potentially high computational burden. Further, convergence of the Markov chain can be difficult to diagnose, particularly for non-expert users. In summary, when choosing among existing methods, one must confront tradeoffs between prediction accuracy, flexibility and computational convenience.

In this paper, we develop an approach to multiple linear regression that aims to combine the best features of existing methods: it is fast, comparable in speed to the cross-validated Lasso; it is flexible, capable of adapting to sparse and dense settings; it is self-tuning, with no need for user-specified hyperparameters; and, in our numerical studies, its prediction

accuracy was competitive with the best methods against which we compared, in a wide range of regression settings, including both dense and sparse effects.

This consistent competitive performance across a wide range of dense/sparse settings is particularly valuable in practice, because in practical applications one does not know whether signals are dense or sparse (or perhaps somewhere in between). Thus, in practice, users will be reluctant to trust a method that does not perform consistently well, even if it performs excellently (even optimally) in certain settings. We believe that the consistently strong performance of our method across settings, combined with its computational speed, make it attractive for practitioners looking to apply large-scale multiple linear regression to real problems.

Our method takes an *empirical Bayes* (EB) approach (Robbins, 1964; Efron, 2019; Hartley and Rao, 1967; Carlin and Louis, 2000; Stephens, 2016; Casella, 2001) to multiple regression; that is, it assigns a prior to the coefficients in the regression method, and this prior is learned from the data. The EB approach is, in many ways, a natural approach for attempting to attain the benefits of Bayesian methods while addressing some of their computational challenges. Indeed, EB for the multiple regression problem is far from new; for example, the Bayesian Lasso (Park and Casella, 2008; Joo, 2017) takes an EB approach with a Laplace prior in which the λ parameter in the Laplace prior is estimated by maximizing the (marginal) likelihood. Other EB approaches use a normal prior (Nebebe and Stroud, 1986), a point-normal (“spike-and-slab”) prior (George and Foster, 2000), and a point-double-exponential prior (Yuan and Lin, 2005). (See also van de Wiel et al. 2019 for a review of other EB approaches.) However, previous EB approaches to multiple regression have either focussed on relatively inflexible priors, or have been met with considerable computational challenges.

Here, we propose a different EB approach that is both more flexible than these previous EB approaches and more computationally scalable. This new EB approach has two key components. First, to increase flexibility, we borrow the “adaptive shrinkage” priors used in Stephens (2016); specifically, we use the “scale mixture of normals” priors. This prior family includes most of the popular priors that have been used in Bayesian regression, including normal, Laplace, point-normal, point-Laplace, point- t , normal-inverse-gamma, Dirichlet-Laplace and horseshoe priors (Hoerl and Kennard, 1970; George and McCulloch, 1997; Meuwissen et al., 2001, 2009; Habier et al., 2011; Anirban Bhattacharya and Dunson, 2015). Increasing model flexibility typically means greater computational expense, but in this case the use of the adaptive shrinkage priors actually simplifies many computations, essentially because the scale mixture family is a convex family. Second, to make computations tractable, we adapt the variational approximation methods for multiple regression from Carbonetto and Stephens (2012). The main limitation of the variational approximation approach is that, in sparse settings with very highly correlated predictors, it will give only one of the correlated predictors a non-negligible coefficient (Carbonetto and Stephens, 2012). This limitation, which is shared by several other existing methods, including the Lasso and L_0 -penalized regression, does not greatly affect prediction accuracy. However, it does limit the conclusions that can be drawn about the selected variables. Consequently, other methods (e.g., Wang et al., 2020) may be preferred when the main goal is variable selection for scientific interpretation rather than prediction. Since our approach combines EB ideas with variational approximations, we refer to it as a “variational EB” (VEB) approach.

While variational methods have previously been used to fit Bayesian linear regression models (Girolami, 2001; Logsdon et al., 2010; Carbonetto and Stephens, 2012; Wang et al., 2020; You et al., 2014; Ren et al., 2011), they have not been used to implement an EB method, and not with the flexible class of priors we consider here. (Independently, Zabad et al. 2023 recently used a VEB approach, but with less flexible priors.) Our work arises from the combination of two earlier ideas: (1) the use of variational approximation techniques for fast posterior computation in large-scale linear regression (Carbonetto and Stephens, 2012); and (2) the use of a flexible class of priors (2), which was originally proposed in Stephens (2016) for performing EB inference in a simpler, but related, problem: the “normal means problem” (Efron and Morris, 1973; Johnstone and Silverman, 2004; Sun and Stephens, 2018; Castillo and Van Der Vaart, 2012; Bhadra et al., 2019). The combination of these two ideas results in methods that are simpler, faster, more flexible, and often more accurate than those in Carbonetto and Stephens (2012).

Finally, another key contribution of our paper is to provide a conceptual bridge between PLR methods and Bayesian methods. Specifically, we show that our VEB approach is actually a PLR method, where the penalty function is learned from the data by *directly solving an optimization problem* rather than being tuned by cross-validation (CV). This result is not only conceptually interesting, but opens the door to other potential algorithms (e.g., gradient descent) for the VEB problem. Tuning multiple parameters by solving an optimization problem is also more practical than CV; for example, our VEB approach has a similar computational cost to methods such as the Lasso that tune a single parameter by CV, and is substantially faster than methods such as the Elastic Net that tune two or more parameters by CV.

1.1 Organization of the Paper

The remainder of the paper is organized as follows. Section 2 gives preliminary background on the normal means model and introduces some notation. Section 3 describes our VEB methods and optimization algorithms in detail. Section 4 makes connections between our VEB approach and penalized approaches. Section 5 gives results from numerical studies comparing prediction performance of different methods for multiple linear regression, including our VEB approach. Section 6 summarizes the contributions of this work and discusses future directions.

1.2 Notations and Conventions

We write vectors in bold, lowercase letters (e.g., \mathbf{b}) and matrices are written in bold, uppercase letters (e.g., \mathbf{X}) We use \mathbb{R}^n to denote the set of real-valued vectors of length n , \mathbb{R}_+^n for the set of real non-negative vectors of length n , $\mathbb{R}^{m \times n}$ for the set of real $m \times n$ matrices, and $\mathbb{S}^n = \{\mathbf{x} \in \mathbb{R}_+^n : \sum_{i=1}^n x_i = 1\}$ denotes the n -dimensional simplex. We use \mathbf{x}_j to denote the j th column of matrix \mathbf{X} . We write sets and families in calligraphic font, e.g., \mathcal{G} . We use $N(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ to denote the probability density of the multivariate normal distribution at $\mathbf{x} \in \mathbb{R}^n$ with mean $\boldsymbol{\mu} \in \mathbb{R}^n$ and $n \times n$ covariance matrix $\boldsymbol{\Sigma}$. We use \mathbf{I}_n to denote the $n \times n$ identity matrix. We use $\|\mathbf{x}\| = \sqrt{\mathbf{x}^T \mathbf{x}}$ to denote the L_2 -norm of vector $\mathbf{x} \in \mathbb{R}^n$, and $\|\mathbf{x}\|_\infty = \max\{|x_1|, \dots, |x_n|\}$ denotes the L -infinity norm of \mathbf{x} . We use \triangleq to indicate definitions.

2 Preliminaries: the Empirical Bayes Normal Means Model

The computations for our approach are closely related to those for a simpler model known as the “normal means” (NM) model. This section reviews this model and introduces notation that will be used later.

2.1 The Normal Means Model

The normal means model is a model for a sequence y_1, \dots, y_p of observations in which each observation y_j is normally distributed with unknown mean b_j and known variance σ^2 :

$$y_j \mid b_j, \sigma^2 \sim N(b_j, \sigma^2), \quad j = 1, \dots, p. \quad (1)$$

This can be viewed as a special case of multiple linear regression in which the covariates are orthogonal and the residual variance is known. (Specifically, it is equivalent to (14) below with $\mathbf{X} = \mathbf{I}_n$ and σ^2 known.)

2.2 The Normal Means Model with Adaptive Shrinkage Priors

Stephens (2016) considers an EB version of the NM model that assumes the b_j are *i.i.d.* from some prior distribution g that is to be estimated from the observed data. Specifically, Stephens (2016) considers priors that are scale mixtures of normals; that is, $g \in \mathcal{G}(\sigma_1^2, \dots, \sigma_K^2)$, where

$$\mathcal{G}(\sigma_1^2, \dots, \sigma_K^2) \triangleq \left\{ g = \sum_{k=1}^K \pi_k N(0, \sigma_k^2) : \pi \in \mathbb{S}^K \right\}, \quad (2)$$

and where $0 \leq \sigma_1^2 < \dots < \sigma_K^2 < \infty$ is a pre-specified grid of component variances, and π_1, \dots, π_K are unknown mixture proportions. Typically, $\sigma_1^2 = 0$ so that $\mathcal{G}(\sigma_1^2, \dots, \sigma_K^2)$ includes sparse prior distributions. (We define $N(0, 0)$ to be the Dirac “delta” mass at zero, commonly denoted as δ_0 .) By making the grid of variances sufficiently wide and dense, the prior family $\mathcal{G}(\sigma_1^2, \dots, \sigma_K^2)$ can approximate, with arbitrary accuracy, the nonparametric family of all the scale mixtures of zero-mean normal distributions. This nonparametric family, which we denote by \mathcal{G}_{SMN} , is very flexible, and includes most popular distributions used as priors in Bayesian regression models, including normal (“ridge regression”) (Hoerl and Kennard, 1970), point-normal (“spike and slab”) (Chipman et al., 2001; George and McCulloch, 1993, 1997; Mitchell and Beauchamp, 1988), double-exponential or Laplace (Figueiredo, 2003; Park and Casella, 2008; Hans, 2009; Tibshirani, 1996; Li and Lin, 2010), horseshoe (Carvalho et al., 2010), normal-gamma prior (Griffin and Brown, 2010), normal-inverse-gamma prior (Meuwissen et al., 2009; Habier et al., 2011), mixture of two normals (BSLMM) (Zhou et al., 2013), and the mixture of four zero-centered normals with different variances suggested by Moser et al. (2015).

Stephens (2016) refers to the priors (2) as “adaptive shrinkage” priors. Here, we use these adaptive shrinkage priors, but assume a prior distribution in which the *scaled* coefficients, b_j/σ , are *i.i.d.* from g ,

$$b_j \mid g, \sigma^2 \stackrel{i.i.d.}{\sim} g_\sigma, \quad (3)$$

where $g_\sigma(b_j) \triangleq g(b_j/\sigma)/\sigma$. We use this prior on the scaled coefficients because in the regression setting it may help to reduce issues with multi-modality; see Park and Casella 2008 for example. The scaled prior (3) also provides computational benefits in the “fully Bayesian” regression setting; see, for example, Chipman et al. 2001; George and McCulloch 1997; Liang et al. 2008 for arguments in favour of the scaled prior. All our methods can also be applied, with minor modifications, to work with the unscaled prior $b_j \stackrel{i.i.d.}{\sim} g$.

2.2.1 AUGMENTED-VARIABLE REPRESENTATION

It is helpful to think of $g \in \mathcal{G}(\sigma_1^2, \dots, \sigma_K^2)$ as determining a set of mixture proportions $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)$ and component variances $\sigma_1^2, \dots, \sigma_K^2$ of a normal mixture. When $g \in \mathcal{G}(\sigma_1^2, \dots, \sigma_K^2)$, the prior (3) can also be written as

$$b_j \mid g, \sigma^2 \stackrel{i.i.d.}{\sim} \sum_{k=1}^K \pi_k N(0, \sigma^2 \sigma_k^2). \quad (4)$$

It is also convenient for some of the derivations to introduce the standard augmented-variable representation of this mixture:

$$\begin{aligned} p(\gamma_j = k \mid g) &= \pi_k \\ b_j \mid g, \sigma^2, \gamma_j = k &\sim N(0, \sigma^2 \sigma_k^2), \end{aligned} \quad (5)$$

where the latent variable $\gamma_j \in \{1, \dots, K\}$ indicates which mixture component gave rise to b_j .

2.3 Empirical Bayes for the Normal Means Model

Stephens (2016) provides EB methods to fit the normal means model. These methods proceed in two steps: estimate g (Step 1); compute the posterior distribution for \mathbf{b} given the estimated g (Step 2). Step 1 is simplified by the use of a fixed grid of variances in (4), which means that only the mixture proportions $\boldsymbol{\pi}$ need to be estimated. This is done by maximizing the marginal log-likelihood:

$$\begin{aligned} \hat{\boldsymbol{\pi}} &= \operatorname{argmax}_{\boldsymbol{\pi} \in \mathbb{S}^K} \log p(\mathbf{y} \mid g, \sigma^2) \\ &= \operatorname{argmax}_{\boldsymbol{\pi} \in \mathbb{S}^K} \sum_{j=1}^p \log \sum_{k=1}^K \pi_k L_{jk}, \end{aligned} \quad (6)$$

where

$$\begin{aligned} L_{jk} &\triangleq p(y_j \mid g, \sigma^2, \gamma_j = k) \\ &= N(y_j; 0, \sigma^2 + \sigma^2 \sigma_k^2), \end{aligned} \quad (7)$$

and $\mathbf{y} \triangleq (y_1, \dots, y_p)$. This is a convex optimization problem, and can be solved efficiently using convex optimization techniques (Koenker and Mizera, 2014; Kim et al., 2020), or

simply by iterating the following Expectation Maximization (EM) updates (Dempster et al., 1977):

$$\text{E-step} \quad \phi_{jk} \leftarrow \phi_k(y_j; g, \sigma^2) \triangleq p(\gamma_j = k \mid y_j, g, \sigma^2) = \frac{\pi_k L_{jk}}{\sum_{k'=1}^K \pi_{k'} L_{jk'}}, \quad (8)$$

$$\text{M-step} \quad \pi_k \leftarrow \frac{1}{p} \sum_{j=1}^p \phi_{jk}, \quad k = 1, \dots, K. \quad (9)$$

The posterior mixture assignment probabilities ϕ_{jk} are sometimes referred to as the “responsibilities”.

Step 2, computing the posterior distribution, is also straightforward, again due to the independence of the observations and the conjugacy of the normal (mixture) prior with the normal likelihood:

$$\begin{aligned} p_{\text{post}}^{\text{NM}}(b_j, \gamma_j = k \mid y_j, g, \sigma^2) &= p(b_j \mid y_j, g, \sigma^2, \gamma_j = k) p(\gamma_j = k \mid y_j, g, \sigma^2) \\ &= \phi_{jk} N(b_j; \mu_{jk}, s_{jk}^2), \end{aligned} \quad (10)$$

where

$$\mu_{jk} \triangleq \mu_k(y_j; g, \sigma^2) = \frac{\sigma_k^2}{1 + \sigma_k^2} \times y_j, \quad (11)$$

$$s_{jk}^2 \triangleq s_k^2(y_j; g, \sigma^2) = \frac{\sigma_k^2}{1 + \sigma_k^2} \times \sigma^2. \quad (12)$$

(Although the posterior variances do not depend on y_j , we write them as $s_k^2(y_j; g, \sigma^2)$ for notational consistency.) Summing the component posterior (10) over k then yields an analytic expression for the posterior of b_j ,

$$p_{\text{post}}^{\text{NM}}(b_j \mid y_j, g, \sigma^2) = \sum_{k=1}^K \phi_{jk} N(b_j; \mu_{jk}, s_{jk}^2). \quad (13)$$

3 Variational Empirical Bayes Linear Regression

3.1 Empirical Bayes Linear Regression

We consider the multiple linear regression model,

$$\mathbf{y} \mid \mathbf{X}, \mathbf{b}, \sigma^2 \sim N(\mathbf{X}\mathbf{b}, \sigma^2 \mathbf{I}_n), \quad (14)$$

where $\mathbf{y} \in \mathbb{R}^n$ is a vector of responses, $\mathbf{X} \in \mathbb{R}^{n \times p}$ is a design matrix whose columns contain predictors $\mathbf{x}_1, \dots, \mathbf{x}_p \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^p$ is a vector of regression coefficients, and $\sigma^2 \geq 0$ is the variance of the residual errors. While an intercept is not explicitly included in (14), it is easily accounted for by centering \mathbf{y} and the columns of \mathbf{X} prior to model fitting (Chipman et al., 2001); see also Section 3.5. To simplify presentation, we will assume throughout the main text of the paper that the columns of \mathbf{X} are rescaled so that $\|\mathbf{x}_j\| = 1$, for $j = 1, \dots, p$. However, all our methods and results can be extended to the unscaled case; Appendix E includes these extensions.

Taking an EB approach, as in the NM model above, we assume the scaled regression coefficients, b_j/σ , are *i.i.d* from some prior g , where g is to be estimated from the observed data. Although our methods apply more generally, we focus on the adaptive shrinkage priors (4) because they are flexible and computationally convenient.

A standard EB approach to fitting the regression model (14) with priors (3) would, similar to above, involve the following two steps:

1. Estimate g, σ^2 by maximizing the marginal likelihood:

$$\begin{aligned} (\hat{g}, \hat{\sigma}^2) &= \operatorname{argmax}_{g \in \mathcal{G}, \sigma^2 \in \mathbb{R}_+} p(\mathbf{y} \mid \mathbf{X}, g, \sigma^2) \\ &= \operatorname{argmax}_{g \in \mathcal{G}, \sigma^2 \in \mathbb{R}_+} \log \int p(\mathbf{y} \mid \mathbf{X}, \mathbf{b}, \sigma^2) p(\mathbf{b} \mid g, \sigma^2) d\mathbf{b}. \end{aligned} \quad (15)$$

2. Infer \mathbf{b} based on the posterior distribution,

$$\hat{p}_{\text{post}}(\mathbf{b}) \triangleq p(\mathbf{b} \mid \mathbf{X}, \mathbf{y}, \hat{g}, \hat{\sigma}^2) \propto p(\mathbf{y} \mid \mathbf{X}, \mathbf{b}, \hat{\sigma}^2) p(\mathbf{b} \mid \hat{g}, \hat{\sigma}^2). \quad (16)$$

Unfortunately, in contrast to the NM model, both steps are computationally impractical due to intractable integrals or very large sums, or both, except in special cases.

3.2 Variational Approximation

To circumvent the intractability of the EB approach, we use a mean-field variational approximation (Blei et al., 2017; Jordan et al., 1999; Wainwright and Jordan, 2008; Logsdon et al., 2010; Carbonetto and Stephens, 2012) to derive a “variational empirical Bayes” (VEB) approach. The idea of VEB inference is mentioned explicitly in Blei et al. (2003), although earlier work implemented similar ideas (e.g., Saul and Jordan 1996; Ghahramani and Hinton 2000; see also van de Wiel et al. 2019). To describe the VEB approach, it is convenient to rewrite the two steps of EB as solving a single optimization problem (see also the Supplementary Materials from Wang et al. 2020):

$$(\hat{p}_{\text{post}}, \hat{g}, \hat{\sigma}^2) = \operatorname{argmax}_{q, g \in \mathcal{G}, \sigma^2 \in \mathbb{R}_+} F(q, g, \sigma^2), \quad (17)$$

where the optimization over q is over all possible distributions on $(\mathbf{b}, \boldsymbol{\gamma})$, and

$$F(q, g, \sigma^2) \triangleq \log p(\mathbf{y} \mid \mathbf{X}, g, \sigma^2) - D_{\text{KL}}(q(\mathbf{b}, \boldsymbol{\gamma}) \parallel p(\mathbf{b}, \boldsymbol{\gamma} \mid \mathbf{X}, \mathbf{y}, g, \sigma^2)). \quad (18)$$

Here, $D_{\text{KL}}(q \parallel p)$ denotes the Kullback-Leibler (K-L) divergence from a distribution q to a distribution p (Kullback and Leibler, 1951). To aid in deriving the closed-form updates below, we reuse the augmented-variable representation from the NM model, $\boldsymbol{\gamma} \triangleq (\gamma_1, \dots, \gamma_p)$, in which γ_j was defined in (5). The function F is often called the “evidence lower bound” (ELBO) because it is a lower-bound for the “evidence”, $\log p(\mathbf{y} \mid \mathbf{X}, g, \sigma^2)$.

In (17), the optimization over q is generally intractable. The VEB approach addresses this by restricting the family of distributions to be optimized. Specifically, our mean-field VEB approach solves

$$(\hat{q}, \hat{g}, \hat{\sigma}^2) = \operatorname{argmax}_{q \in \mathcal{Q}, g \in \mathcal{G}, \sigma^2 \in \mathbb{R}_+} F(q, g, \sigma^2) \quad (19)$$

where

$$\mathcal{Q} = \left\{ q : q(\mathbf{b}, \boldsymbol{\gamma}) = \prod_{j=1}^p q_j(b_j, \gamma_j) \right\}, \quad (20)$$

is the family of fully-factorized distributions. The resulting \hat{q} factorizes into a product over individual factors $\hat{q}_j(b_j, \gamma_j)$, $j = 1, \dots, p$, and serves as an approximation to the EB posterior, \hat{p}_{post} . With this constraint, the optimization becomes tractable (see Section 3.3 for details).

3.2.1 CONTRAST WITH CARBONETTO AND STEPHENS (2012)

While Carbonetto and Stephens (2012) also use a mean-field approximation for linear regression, their approach to estimating g, σ^2 is quite different—and substantially more complex—than the VEB approach we describe here. In brief, they treat $F(\hat{q}(g, \sigma^2), g, \sigma^2)$ as a direct approximation to the evidence,

$$p(\mathbf{y} \mid \mathbf{X}, g, \sigma^2) \approx \exp\{F(\hat{q}(g, \sigma^2), g, \sigma^2)\},$$

and combine this with a prior distribution on g, σ^2 to arrive at an approximate posterior distribution for g, σ^2 . This approach is computationally burdensome because it requires finding a separate approximation \hat{q} for each g, σ^2 . It also requires specifying a prior on (g, σ^2) , which introduces an additional layer of decision-making for the user. Our VEB approach greatly simplifies this by simultaneously fitting q, g, σ^2 in a single optimization (19). Our approach also uses more flexible prior families than Carbonetto and Stephens (2012). And, as we show in numerical experiments below, our approach generally improves predictive performance.

3.3 Coordinate Ascent Algorithm

We solve (19) for the multiple linear regression model (14) with adaptive shrinkage priors (4) using a simple coordinate-wise approach, outlined in Algorithm 1. While theoretical analysis suggests that coordinate ascent can suffer poor rates of convergence (Beck and Tetrushvili, 2013; Wright, 2015; Hazimeh and Mazumder, 2020), it has the advantage of being guaranteed to converge to a stationary point of the objective under mild conditions (see Proposition 11 in Appendix G). In practice, coordinate ascent has emerged as a simple, fast and reliable approach for optimizing large-scale multiple linear regression models, with both convex and nonconvex penalties (Friedman et al., 2007; Wu and Lange, 2008; Friedman et al., 2010; Mazumder et al., 2011; Breheny and Huang, 2011; Hazimeh and Mazumder, 2020).

In the following, we show that the steps in Algorithm 1 are easy to implement:

- (i) The update for each q_j involves computing a posterior distribution under the NM model.
- (ii) The update for g involves running a single M-step update for the NM model, in which the exact posterior probabilities are replaced with approximate posterior probabilities.
- (iii) The update for the residual variance, σ^2 , has a simple, closed-form solution.

Algorithm 1 Coordinate ascent for fitting VEB model (outline only).

Require: Data $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{y} \in \mathbb{R}^n$,
 and initial estimates $q_1, \dots, q_p, g, \sigma^2$.
repeat
 for $j \leftarrow 1$ to p **do**
 $q_j \leftarrow \operatorname{argmax}_{q_j} F(q, g, \sigma^2)$
 end for
 $g \leftarrow \operatorname{argmax}_{g \in \mathcal{G}} F(q, g, \sigma^2)$
 $\sigma^2 \leftarrow \operatorname{argmax}_{\sigma^2 \in \mathbb{R}_+} F(q, g, \sigma^2)$
until termination criterion is met
return $q_1, \dots, q_p, g, \sigma^2$.

These results are formally stated in the following proposition. (We assume $\mathbf{x}_j^T \mathbf{x}_j = 1$ here to simplify the expressions; more general results and algorithms that do not make this assumption are given in Appendix E.)

Proposition 1 (Coordinate ascent updates for VEB) Assume $\mathbf{x}_j^T \mathbf{x}_j = 1$, for $j = 1, \dots, p$, let $\bar{b}_j = \mathbb{E}(b_j)$, $\bar{\mathbf{b}} = \mathbb{E}(\mathbf{b})$ be the expected values of b_j, \mathbf{b} with respect to q , $\bar{\mathbf{r}} = \mathbf{y} - \mathbf{X}\bar{\mathbf{b}} \in \mathbb{R}^n$ is the vector of expected residuals with respect to q , \mathbf{X}_{-j} denotes the design matrix \mathbf{X} excluding the j th column, q_{-j} is shorthand for all factors $q_{j'}, j' \neq j$, and $\bar{\mathbf{r}}_j \in \mathbb{R}^n$ is the vector of expected residuals accounting for linear effects of all variables other than j ,

$$\bar{\mathbf{r}}_j \triangleq \mathbf{y} - \mathbf{X}_{-j} \bar{\mathbf{b}}_{-j} = \mathbf{y} - \sum_{j' \neq j} \mathbf{x}_{j'} \bar{b}_{j'}. \quad (21)$$

Additionally, we use $\tilde{b}_j \triangleq \mathbf{x}_j^T \bar{\mathbf{r}}_j = \bar{b}_j + \mathbf{x}_j^T \bar{\mathbf{r}}$ to denote the ordinary least squares (OLS) estimate of the coefficient b_j when the residuals $\bar{\mathbf{r}}_j$ are regressed against \mathbf{x}_j . Then we have the following results:

- (i) The coordinate ascent update $q_j^* \triangleq \operatorname{argmax}_{q_j} F(q, g, \sigma^2)$ is obtained by the posterior distribution for b_j, γ_j under the normal means model (1, 3) in which the observation y_j is replaced by the OLS estimate of b_j ; that is,

$$q_j^*(b_j, \gamma_j = k) = p_{\text{post}}^{\text{NM}}(b_j, \gamma_j = k; \tilde{b}_j, g, \sigma^2).$$

In particular, the posterior distribution at the maximum is

$$q_j^*(b_j, \gamma_j = k) = \phi_{jk}^* N(b_j; \mu_{jk}^*, (s_{jk}^2)^*), \quad (22)$$

in which

$$\mu_{jk}^* = \mu_k(\tilde{b}_j; g, \sigma^2) \quad (23)$$

$$(s_{jk}^2)^* = s_k^2(\tilde{b}_j; g, \sigma^2) \quad (24)$$

$$\phi_{jk}^* = \phi_k(\tilde{b}_j; g, \sigma^2). \quad (25)$$

See (8, 10, 11, 13) for the definitions of μ_k, s_k^2, ϕ_k and $p_{\text{post}}^{\text{NM}}$.

(ii) The coordinate ascent update

$$g^* \triangleq \operatorname{argmax}_{g \in \mathcal{G}(\sigma_1^2, \dots, \sigma_K^2)} F(q, g, \sigma^2)$$

is achieved by setting

$$\begin{aligned} g^* &= \sum_{k=1}^K \pi_k^* N(0, \sigma_k^2) \\ \pi_k^* &= \frac{1}{p} \sum_{j=1}^p q_j(\gamma_j = k), \quad k = 1, \dots, K. \end{aligned} \quad (26)$$

(iii) Assuming $\sigma_1^2 = 0$, the coordinate ascent update

$$(\sigma^2)^* \triangleq \operatorname{argmax}_{\sigma^2 \in \mathbb{R}_+} F(q, g, \sigma^2)$$

is achieved with

$$(\sigma^2)^* = \frac{\|\bar{\mathbf{r}}\|^2 + \sum_{j=1}^p \operatorname{Var}_q(b_j) + \sum_{j=1}^p \sum_{k=2}^K \phi_{jk} \mathbb{E}[b_j \mid \gamma_j = k] / \sigma_k^2}{n + p - \sum_{j=1}^p \phi_{j1}}. \quad (27)$$

Or, assuming q_1, \dots, q_p and g have just been updated as in (i) and (ii) above, we have the simpler update formula

$$(\sigma^2)^* = \frac{\|\bar{\mathbf{r}}\|^2 + \bar{\mathbf{b}}^T(\tilde{\mathbf{b}} - \bar{\mathbf{b}}) + \sigma^2 p(1 - \pi_1^*)}{n + p(1 - \pi_1^*)}. \quad (28)$$

Proof See Appendix E. ■

Inserting these expressions into Algorithm 1 (and organizing computations to limit redundant operations and memory requirements) yields Algorithm 2. The computational complexity of this algorithm is $O((n+K)p)$ per outer-loop iteration, with memory requirements $O(n+p+K)$ (in addition to storing the data matrix \mathbf{X}), making it tractable for large data sets. The algorithm also exploits the fact that the approximate posterior $q(\mathbf{b}, \boldsymbol{\gamma})$ can be recovered from $\boldsymbol{\pi}, \sigma^2, \bar{\mathbf{b}}$ by running a single round of the coordinate ascent updates for q_1, \dots, q_p . Because of this, the algorithm is initialized simply by providing an initial estimate of $\bar{\mathbf{b}}$; the full q is not needed. See Section 3.5.3 for more on initialization.

3.4 Accuracy of VEB and Exactness for Orthogonal Predictors

Carbonetto and Stephens (2012) note that their variational approximation approach provides the exact posterior distribution when the columns of \mathbf{X} are orthogonal. Here we extend this result, showing that in this special case the VEB method recovers the standard EB method.

Proposition 2 When \mathbf{X} has orthogonal columns, the VEB approach (19) is mathematically equivalent to the exact EB approach (15, 16).

Algorithm 2 Coordinate ascent for fitting VEB model (in detail).

Require: Data $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{y} \in \mathbb{R}^n$; number of mixture components, K ;
 prior variances, $\sigma_1^2 < \dots < \sigma_K^2$, with $\sigma_1^2 = 0$; initial estimates $\bar{\mathbf{b}}, \boldsymbol{\pi}, \sigma^2$.
 $\bar{\mathbf{r}} = \mathbf{y} - \mathbf{X}\bar{\mathbf{b}}$ (compute mean residuals)
 $t \leftarrow 0$
repeat
 for $j \leftarrow 1$ to p **do**
 $\bar{\mathbf{r}}_j = \bar{\mathbf{r}} + \mathbf{x}_j \bar{b}_j$ (disregard j th effect in residuals)
 $\tilde{b}_j \leftarrow \mathbf{x}_j^T \bar{\mathbf{r}}_j$. (compute OLS estimate)
 for $k \leftarrow 1$ to K **do**
 $\phi_{jk} \leftarrow \phi_k(\tilde{b}_j; g, \sigma^2)$
 $\mu_{jk} \leftarrow \mu_k(\tilde{b}_j; g, \sigma^2)$
 end for (update q_j ; eqs. 23, 25)
 $\bar{b}_j \leftarrow \sum_{k=1}^K \phi_{jk} \mu_{jk}$. (update posterior mean of b_j)
 $\bar{\mathbf{r}} \leftarrow \bar{\mathbf{r}}_j - \mathbf{x}_j \bar{b}_j$. (update mean residuals)
 end for
 for $k \leftarrow 1$ to K **do**
 $\pi_k \leftarrow \sum_{j=1}^p \phi_{jk} / p$. (update g ; eq. 26)
 end for
 $\sigma^2 \leftarrow \frac{\|\bar{\mathbf{r}}\|^2 + \bar{\mathbf{b}}^T (\tilde{\mathbf{b}} - \bar{\mathbf{b}}) + \sigma^2 p (1 - \pi_1)}{n + p(1 - \pi_1)}$ (update σ^2 ; eq. 28)
 $t \leftarrow t + 1$
until termination criterion is met
return $\bar{\mathbf{b}}, \boldsymbol{\pi}, \sigma^2$

Proof See Appendix G. ■

In brief, this result follows from the fact that, when \mathbf{X} has orthogonal columns, the (exact) posterior distribution for \mathbf{b} factorizes as (20), and therefore the mean-field assumption is not an approximation. By contrast, the “conditional maximum likelihood” (CML) approach to approximating EB inference (George and Foster, 2000; Yuan and Lin, 2005) is not exact even in the case of orthogonal columns.

Proposition 2 suggests that our VEB method should be accurate when the columns of \mathbf{X} are close to orthogonal. It also suggests that the approximation may be less accurate for very highly correlated columns. However, Carbonetto and Stephens (2012) note that even in this setting the estimated hyperparameters (here, g) can be accurate. They also note that, when two predictors are highly correlated and predictive of \mathbf{y} , the fully-factorized variational approximation tends to give just one of them an appreciable estimated coefficient. This is similar to the behavior of many PLR methods including the Lasso (but different from a well-mixed MCMC-based Bayesian method). While this behavior is undesirable when the main aim of the analysis is to select variables for scientific interpretation, it does not necessarily harm prediction accuracy. Thus, the VEB approach may perform well for prediction even in settings where the assumptions of the mean-field variational approximation are clearly violated. Our numerical studies (Section 5) confirm this.

3.5 Practical Issues and Extensions

In this subsection, we discuss some practical implementation issues and potential extensions of this method.

3.5.1 INTERCEPT

In multiple regression applications, it is common to include an intercept term that is not regularized in the same way as other variables. A common approach, and the approach we take here, is to center \mathbf{y} and the columns of \mathbf{X} (Chipman et al., 2001; Friedman et al., 2010); that is, the observed responses y_i are replaced with the centered responses $y_i - \sum_{i'=1}^n y_{i'}/n$, and the data x_{ij} are replaced with column-centered values $x_{ij} - \sum_{i'=1}^n x_{i'j}/n$.

3.5.2 SELECTION OF GRID FOR PRIOR VARIANCES

Following Stephens (2016), we choose a grid $\{\sigma_1^2, \dots, \sigma_K^2\}$ that is sufficiently broad and dense so that results do not change much if the grid is made broader and denser; the aim is to choose a $\mathcal{G}(\sigma_1^2, \dots, \sigma_K^2)$ that closely approximates the non-parametric family \mathcal{G}_{SMN} . Specifically, we set the lower end of the grid to be $\sigma_1^2 = 0$, which is a point mass at zero, and we set the largest prior variance to be $\sigma_K^2 \approx n$ so that the prior variance of $\mathbf{x}_j b_j$ is close to σ^2 (recall, we assumed $\mathbf{x}_j^T \mathbf{x}_j = 1$, so $\text{Var}(\mathbf{x}_j) \approx 1/n$ when \mathbf{x}_j is centered). We have found that 20 grid points spanning this range to be good enough to achieve reliable prediction performance across many settings (see Section 5). Based on this, our default grid is the sequence $\sigma_k^2 = n(2^{(k-1)/K} - 1)^2$, $k = 1, \dots, 20$. In rare cases, $\text{Var}(\mathbf{x}_j b_j)$ may be larger than σ^2 for some j . In that case, we may need a larger σ_K^2 to avoid underestimating, or “overshrinking”, the effect b_j . Therefore, we suggest checking that the final estimate of π_K is negligible and, if not, the grid can be made wider.

3.5.3 INITIALIZATION AND UPDATE ORDER

Except in special cases, maximizing F is a nonconvex optimization problem, and so although Algorithm 2 is guaranteed to converge, the solution obtained may depend on initialization of $\bar{\mathbf{b}}$, $\boldsymbol{\pi}$, σ^2 , as well as the order in which the coordinate ascent updates cycle through the coordinates $j \in \{1, \dots, p\}$ (e.g., Carbonetto and Stephens, 2012; Ray and Szabó, 2022). Therefore, we experimented with different initialization and update orderings.

The simplest initialization for $\bar{\mathbf{b}}$ is the “null initialization” $\bar{\mathbf{b}} = (0, \dots, 0)^T$. We also consider initializing $\bar{\mathbf{b}}$ to the Lasso solution $\hat{\mathbf{b}}^{\text{lasso}}$ in which the Lasso penalty is chosen via cross-validation. Given $\bar{\mathbf{b}}$, we initialize the residual variance as $\sigma^2 = \|\mathbf{y} - \mathbf{X}\bar{\mathbf{b}}\|^2/n$, and we initialize the mixture weights to $\boldsymbol{\pi} = (1/K, \dots, 1/K)$. In numerical experiments (Appendix C.2) we found that the Lasso initialization often improved predictive performance when columns of \mathbf{X} were highly correlated. In other cases, the null and the Lasso initializations performed similarly. With the Lasso initialization, we did not find any systematic benefit to different update orderings. Therefore, our default approach is to use the Lasso initialization with updates performed in the natural order, $1, 2, \dots, p$.

3.5.4 TERMINATION CRITERION

We stop iterating when estimates of the prior distribution g stabilize; specifically, we stop at iteration t if $\|\boldsymbol{\pi}^{(t)} - \boldsymbol{\pi}^{(t-1)}\|_\infty < K \times 10^{-8}$.

3.5.5 COMPUTING THE ELBO

Although Algorithm 2 optimizes the ELBO, F , it does not require computation of the ELBO. Still, it can be useful to compute the ELBO, for example to monitor progress of the coordinate ascent updates, or to compare fits obtained from different runs of the algorithm. Appendix E includes analytic expressions for the ELBO that may be useful in such cases.

3.5.6 EXTENSION TO OTHER MIXTURE PRIOR FAMILIES

We have focussed on the normal mixture prior family $\mathcal{G} = \mathcal{G}(\sigma_1^2, \dots, \sigma_K^2)$ because it makes computations simple, fast, and numerically stable. Furthermore, this prior family includes most prior distributions previously used for multiple regression, and so we expect it to suffice for many practical applications. However, Algorithm 2 could be adapted to accommodate other prior families of the form $\mathcal{G} = \{g = \sum_{k=1}^K \pi_k g_k : \boldsymbol{\pi} \in \mathbb{S}^K\}$, with fixed mixture components g_1, \dots, g_K . When choosing \mathcal{G} , there are two important practical points to consider: (i) the convolution of g_k with a normal likelihood should be numerically tractable, ideally with an analytic expression; (ii) the posterior mean in the normal means model with prior $b_j \sim g_k$ should be easy to compute. Examples of fixed mixture components g_k satisfying (i) and (ii) include point masses, uniform distributions and Laplace distributions.

3.5.7 INFERENCE AND VARIABLE SELECTION

We have focussed on developing flexible multiple regression methods for accurate *prediction*, which requires only a point estimate for \mathbf{b} (e.g., the posterior mean, $\bar{\mathbf{b}}$). However, the approximate posterior distributions from our method could also be used for inference, that is, to assess uncertainty in the estimated \mathbf{b} . For example, to assess significance of each variable in the regression, it is easy to compute the local false sign rate, *lfsr* (Stephens, 2016), which quantifies confidence in the sign of the effect (and which we generally prefer to the closely related *local false discovery rate*; see Stephens 2016). However, caution is warranted in settings with highly correlated variables: in such settings, the approximate posterior distribution from the fully-factored variational approximation will often be inaccurate. While predictive performance is quite robust to this issue, inference is more sensitive (Carbonetto and Stephens, 2012). Thus, other methods may be preferred for inference with highly correlated variables; see Wang et al. (2020) for further discussion.

4 Connecting VEB and Penalized Linear Regression

This section shows that the approximate posterior mean computed by our VEB approach also solves a PLR with a nonconvex penalty, where the form of the penalty is flexible and is automatically learned from the data without need for cross-validation.

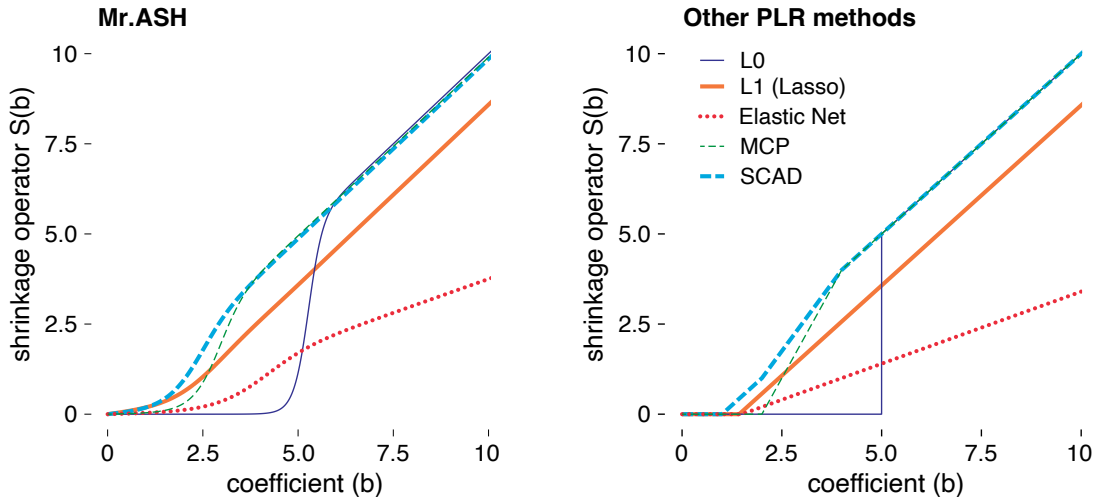


Figure 1: Examples of posterior mean shrinkage operators for different $g \in \mathcal{G}(\sigma_1^2, \dots, \sigma_K^2)$ (left-hand panel) and σ^2 that were chosen to mimic the shrinkage operators from some commonly used penalties (right-hand panel).

4.1 Penalties and Shrinkage Operators

Penalized linear regression (PLR) methods estimate the regression coefficients by minimizing a penalized squared-loss function:

$$\underset{\mathbf{b} \in \mathbb{R}^p}{\text{minimize}} h_\rho(\mathbf{b}) \triangleq \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 + \sum_{j=1}^p \rho(b_j), \quad (29)$$

for some penalty function $\rho : \mathbb{R} \rightarrow \mathbb{R}$. As mentioned above, the PLR problem (29) is often tackled using coordinate descent algorithms; that is, by iterating over the coordinates of b_1, \dots, b_p sequentially, at each iteration solving (29) for one coordinate b_j while keeping the remaining coordinates fixed. Assuming the columns of \mathbf{X} are scaled such that $\mathbf{x}_j^T \mathbf{x}_j = 1$, the solution for the j th coordinate is obtained by

$$b_j \leftarrow S_\rho(b_j + \mathbf{x}_j^T (\mathbf{y} - \mathbf{X}\mathbf{b})), \quad (30)$$

where

$$S_\rho(t) \triangleq \underset{\theta \in \mathbb{R}}{\text{argmin}} \frac{1}{2} (t - \theta)^2 + \rho(\theta) \quad (31)$$

is the “shrinkage operator” (or “univariate proximal operator”; Parikh and Boyd 2014) for penalty ρ . Studying the shrinkage operator S_ρ is often helpful for understanding the behaviour of its corresponding penalty, ρ . For example, the shrinkage operators for the L_0 and L_1 penalties both have a “thresholding” property in which coefficient estimates t less than some value are driven to zero. This thresholding property tends to produce sparse solutions to (29); see Figure 1 for an illustration. Table 3 in Appendix A gives some commonly used penalty functions and their corresponding shrinkage operators.

In the next section, we show that our VEB approach can be interpreted as solving a PLR problem with a penalty function $\rho_{g_\sigma, \sigma}$ that depends on the prior g_σ and the residual variance σ^2 . This penalty function has a corresponding shrinkage operator, denoted by $S_{g_\sigma, \sigma}$, that has a particularly simple form and interpretation: *it is the posterior mean under a normal means model*. It is formally defined as follows.

Definition 3 (Normal Means Posterior Mean Operator) Define the *normal means posterior mean operator*, $S_{f, \sigma} : \mathbb{R} \rightarrow \mathbb{R}$, as the mapping

$$S_{f, \sigma}(y) \triangleq \mathbb{E}_{\text{NM}}(b \mid y, f, \sigma^2), \quad (32)$$

where \mathbb{E}_{NM} denotes the posterior expectation under the following normal means model with prior f and variance σ^2 ,

$$\begin{aligned} y \mid b, \sigma^2 &\sim N(b, \sigma^2) \\ b &\sim f. \end{aligned} \quad (33)$$

From (13), $S_{f, \sigma}$ has a simple analytic form when $f = g_\sigma$, $g \in \mathcal{G}(\sigma_1^2, \dots, \sigma_K^2)$:

$$S_{g_\sigma, \sigma}(y) = \sum_{k=1}^K \phi_k(y; g, \sigma^2) \mu_k(y; g, \sigma^2). \quad (34)$$

It is easy to show that $S_{g_\sigma, \sigma}$ is an odd function and is monotonic in y . Also, $S_{g_\sigma, \sigma}$ is a shrinkage operator, in that $|S_{g_\sigma, \sigma}(y)| \leq |y|$; see Lemma 12 in Appendix G. Indeed, given a suitable choice of prior family, the shrinkage operator $S_{g_\sigma, \sigma}$ can qualitatively mimic the behavior of many commonly used shrinkage operators (Figure 1).

The behavior of $S_{g_\sigma, \sigma}$ naturally depends on the prior. For example, the more mass the prior places near zero, the stronger the shrinkage toward zero. Our VEB method estimates the prior from within a flexible family capable of capturing a wide range of scenarios; consequently, the corresponding shrinkage operator is also estimated from a flexible family of shrinkage operators. This process is analogous to estimating the tuning parameters in regular PLR methods which is usually done by cross-validation (CV). However, our VEB approach dispenses with CV and makes it possible to efficiently tune across a much wider range of shrinkage operators.

4.2 VEB as Penalized Linear Regression

The first step to connecting our VEB approach to a PLR is to write it as an optimization over the regression coefficients, \mathbf{b} , rather than an optimization over the (approximate) posterior distributions, q . To do this, we define an objective function:

$$h(\bar{\mathbf{b}}, g, \sigma^2) \triangleq - \left\{ \max_{q \in \mathcal{Q}, \mathbb{E}_q[\mathbf{b}] = \bar{\mathbf{b}}} F(q, g, \sigma^2) \right\}. \quad (35)$$

The constraint $\mathbb{E}_q[\mathbf{b}] = \bar{\mathbf{b}}$ means that the expected value of \mathbf{b} with respect to q is $\bar{\mathbf{b}}$. The negative sign is introduced to align with the convention that PLRs are usually minimization problems, as in (29).

Any algorithm for optimizing F over q (and possibly g, σ^2) also provides a way to optimize h over $\bar{\mathbf{b}}$ (and possibly g, σ^2), as formalized in the following proposition.

Proposition 4 (Computing Posterior Mean as an Optimization Problem) Let \hat{q} , \hat{g} , $\hat{\sigma}^2$ be a solution to

$$\hat{q}, \hat{g}, \hat{\sigma}^2 = \operatorname{argmax}_{q \in \mathcal{Q}, g \in \mathcal{G}, \sigma^2 \in \mathcal{T}} F(q, g, \sigma^2),$$

where \mathcal{Q} is the variational mean-field family of approximate posterior distributions (20), \mathcal{G} is any family of prior distributions on $b \in \mathbb{R}$, and \mathcal{T} is any subset of \mathbb{R}_+ . (This general formulation allows, as a special case, g, σ^2 to be fixed by taking \mathcal{G} and \mathcal{T} to be singleton sets.) Let $\hat{\mathbf{b}}$ denote the expected value of \mathbf{b} with respect to \hat{q} . Then $\hat{\mathbf{b}}, \hat{g}, \hat{\sigma}^2$ also solves the following optimization problem:

$$\hat{\mathbf{b}}, \hat{g}, \hat{\sigma}^2 = \operatorname{argmin}_{\bar{\mathbf{b}} \in \mathbb{R}^p, g \in \mathcal{G}, \sigma^2 \in \mathcal{T}} h(\bar{\mathbf{b}}, g, \sigma^2).$$

Proof See Appendix G. ■

The final step to connecting VEB with PLR is to show that h has the form of a penalized squared-loss function.

Theorem 5 (VEB as a Penalized Log-likelihood) The objective function h defined in (35) has the form of a PLR,

$$h(\bar{\mathbf{b}}, g, \sigma^2) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\bar{\mathbf{b}}\|^2 + \frac{1}{\sigma^2} \sum_{j=1}^p \rho_{g\sigma, \sigma}(\bar{b}_j) + \frac{n-p}{2} \log(2\pi\sigma^2), \quad (36)$$

in which the penalty function $\rho_{f, \sigma}$ satisfies

$$\rho_{f, \sigma}(S_{f, \sigma}(y)) = -\sigma^2 \ell_{\text{NM}}(y; f, \sigma^2) - \frac{1}{2}(y - S_{f, \sigma}(y))^2, \quad (37)$$

and

$$\rho'_{f, \sigma}(S_{f, \sigma}(y)) = (y - S_{f, \sigma}(y)). \quad (38)$$

Here, $\ell_{\text{NM}}(y; f, \sigma^2) \triangleq \log p(y | f, \sigma^2)$ denotes the marginal log-likelihood under the NM model (33), and $S_{f, \sigma}$ denotes the shrinkage operator (32).

Proof See Appendix G. ■

From this theorem, it follows that the NM posterior mean shrinkage operator $S_{f, \sigma}$ (32) can be also written in the form of (31), a shrinkage operator for the penalty $\rho_{f, \sigma}$. Explicit computation of $\rho_{f, \sigma}(\bar{b})$ for a given \bar{b} in (37) would require computing the inverse shrinkage operator, $y = S_{f, \sigma}^{-1}(\bar{b})$. This inverse exists because the shrinkage operator (34) is strictly increasing; however, we do not have an analytic expression for the inverse, so we do not have an analytic expression for $\rho_{f, \sigma}(\bar{b})$.

4.2.1 SPECIAL CASE WHEN g AND σ^2 ARE FIXED

The special case when g and σ^2 are fixed is particularly simple and helpful for intuition. In this case, the VEB approach is solving a PLR problem with fixed penalty $\rho_{g\sigma, \sigma}$ and

Algorithm 3 Coordinate Ascent Iterative Shrinkage Algorithm for Variational Posterior Mean (with fixed g, σ^2)

Require: $\mathbf{X} \in \mathbb{R}^{n \times p}, \mathbf{y} \in \mathbb{R}^n, \sigma^2 > 0$, prior g , and initial estimates $\bar{\mathbf{b}}$.
repeat
 for $j \leftarrow 1$ to p **do**
 $\bar{b}_j \leftarrow S_{g,\sigma}(\bar{b}_j + \mathbf{x}_j^T(\mathbf{y} - \mathbf{X}\bar{\mathbf{b}}))$
 end for
until convergence criteria is met
return $\bar{\mathbf{b}}$

shrinkage operator $S_{g,\sigma}$. This leads to a simple coordinate ascent algorithm (Algorithm 3). Compare this with the inner loop of Algorithm 2 which maximizes the ELBO, F , over each q_j in turn. The key computation in the inner loop is the computation of the posterior mean, \bar{b}_j . (When g and σ^2 are fixed, computing ϕ_{jk} and μ_{jk} is needed only to compute \bar{b}_j .) Further, by Proposition 1, this value is computed as the posterior mean under a simple NM model, which is given by the shrinkage operator $S_{g,\sigma}$.

In summary, for fixed g, σ^2 , Algorithm 2 can be reframed as a coordinate ascent algorithm for PLR, which is Algorithm 3.

4.2.2 SPECIAL CASE OF A NORMAL PRIOR (RIDGE REGRESSION)

When the prior, g , is a fixed normal distribution with zero mean, the NM posterior mean shrinkage operator $S_{g,\sigma}$ is the same as the ridge regression (or L_2) shrinkage operator (Table 3) and the penalty function $\rho_{g,\sigma}$ is the L_2 -penalty. Thus, in this special case, Algorithm 3 is solving ridge regression (*i.e.*, PLR with L_2 -penalty), which is a convex optimization problem. Furthermore, in this special case Algorithm 3 converges to the *exact* posterior mean of \mathbf{b} because the posterior is multivariate normal, and therefore the posterior mean is equal to the posterior mode. Thus, in this special case, even though the variational posterior approximation q does not exactly match the true posterior—the true posterior does not factorize as in (20)—the variational posterior mean recovers the true posterior mean.

4.2.3 POSTERIOR MEAN VS. POSTERIOR MODE

Traditional PLR approaches are sometimes motivated from a Bayesian perspective as computing a posterior mode estimate for \mathbf{b} —*i.e.*, a *maximum a posteriori* (MAP) estimate—in which the penalty term corresponds to some prior on \mathbf{b} (Fu, 1998). For example, the Lasso is the MAP with a Laplace (“double-exponential”) prior (Figueiredo, 2003; Park and Casella, 2008; Tibshirani, 1996). By contrast, the variational approach seeks the *posterior mean*, not the posterior mode, and likewise the VEB shrinkage (32) is based on an averaging (mean) operation instead of the usual maximization (mode). Our formulation of the (approximate)

posterior mean as solving a PLR is new, at least as far as we are aware, and this formulation may be useful in other settings.

From a Bayesian decision-theoretic perspective (e.g., Chapter 4 of Berger 1985), the posterior mean for \mathbf{b} has better theoretical support than the posterior mode; not only does it minimize the expected mean-squared error in \mathbf{b} , it also minimizes the expected mean-squared error in the predicted, or “fitted,” responses $\hat{y}_i = (x_{i1}, \dots, x_{ip})^T \mathbf{b}$. Although the posterior mode can have attractive properties—for example, the posterior mode with a Laplace prior is sparse whereas the posterior mean with a Laplace prior is not—the posterior mode has very little support as an estimator for \mathbf{b} , particularly when predicting \mathbf{y} is the main goal. To illustrate this, consider a spike-and-slab prior (Mitchell and Beauchamp, 1988) with non-zero mass at zero: the posterior mode is always $\mathbf{b} = \mathbf{0}$, which will generally provide poor prediction performance. Also, consider that if $\hat{\mathbf{b}}^{\text{mode}}$ is the posterior mode of \mathbf{b} , then $(x_{i1}, \dots, x_{ip})^T \hat{\mathbf{b}}^{\text{mode}}$ is not generally the posterior mode of the fitted value \hat{y}_i .

5 Numerical Experiments

We used simulations to empirically assess the predictive performance of our proposed method and compare it with other methods. We call the proposed method *multiple regression with adaptive shrinkage priors*, or “Mr.ASH”. Mr.ASH is implemented in the R package `mr.ash.alpha`, available at <https://github.com/stephenslab/mr.ash.alpha>. In the experiments, we used `mr.ash.alpha` version 0.1-33 (git commit id 0845778). The source code and analysis steps used to generate the results of our numerical experiments are included in a separate repository on GitHub, <https://github.com/stephenslab/mr-ash-workflow>.

5.1 Methods Compared

We compared Mr.ASH with 12 different methods (Table 1). These include a range of well established and recently developed methods based on both PLR and Bayesian ideas. These methods vary in (1) the choice of penalty or prior (and possibly other modeling assumptions); and (2) the algorithm used to fit the model (e.g., point estimation of \mathbf{b} vs. approximate posterior inference of \mathbf{b} via MCMC or variational inference). When selecting methods to compare, we tried to choose methods that were publicly available as open source software and that were well maintained and documented. Since Mr.ASH and many other popular large-scale linear regression methods are implemented in R, we sought out methods implemented in R (R Core Team, 2019).

Here we give a brief overview of the different methods and point out some of their expected strengths and weaknesses:

- Ridge regression (Hoerl and Kennard 1970) and the Bayesian Lasso (Park and Casella 2008; Perez and de los Campos 2014) are well adapted to dense signals, so should be competitive in such settings. On the other hand, they may perform poorly for sparse signals.
- L0Learn (Hazimeh and Mazumder, 2020) and SuSiE (Wang et al., 2020) are better adapted to sparse signals, so they should perform well in such settings. On the other hand, their assumptions are poorly suited to more dense signals.

- The Lasso (Tibshirani, 1996) is one of the most widely used PLR methods. Computing the Lasso estimator is a convex optimization problem. One of the well-studied issues is that Lasso estimates can suffer from bias by overshrinking the strongest signals (e.g., Su et al., 2017; Javanmard and Montanari, 2018).
- The Elastic Net (Zou and Hastie 2005) is another widely used convex PLR method. It has two tuning parameters. The Elastic Net penalty is more flexible than the Lasso and ridge regression and, indeed, it includes both as special cases. The Elastic Net may therefore perform well across a wider range of settings, at the cost of increased computation in the parameter tuning.
- SCAD, MCP, the Spike-and-Slab Lasso (SSLasso) and the Trimmed Lasso are methods based on nonconvex or adaptive penalties that were designed, in part, to address limitations of the Lasso penalty (Bai et al., 2021; Bertsimas et al., 2017; Breheny and Huang, 2011; Ročková and George, 2018; Amir et al., 2021; Yun et al., 2019). They might therefore outperform the Lasso (and the Elastic Net), possibly at the cost of some additional computation. Since these methods were primarily developed with sparse regression in mind, they may not always perform as well for dense signals.
- BayesB is a Bayesian regression method with a “spike-and-slab” prior, in which the “slab” is a t distribution (Meuwissen et al., 2001; Perez and de los Campos, 2014). It has the potential to perform well for both sparse and dense signals. One concern is that the Markov chain may or may not be simulated long enough so as to adequately explore the posterior distribution.
- `varbvs` (Carbonetto and Stephens, 2012; Carbonetto et al., 2017) and `Mr.ASH` both compute approximate posteriors using the same mean field variational approximation. Compared to `varbvs`, `Mr.ASH` features a more flexible prior, and uses a simpler and more efficient empirical Bayes approach to estimate the prior. `Mr.ASH` also uses an initialization based on the Lasso, which, as we show below, can improve the model fit, particularly when the predictors are strongly correlated, or correlated in complex ways. We expect `Mr.ASH` to outperform `varbvs` some settings, particularly in “dense” settings when many of the predictors affect the outcome, or when the predictors are strongly correlated.

As an additional point of comparison in “sparse” data simulations, we also show results for the “Oracle OLS” method, which is the ordinary least-squares (OLS) estimate of \mathbf{b} conditional on knowledge of which coefficients b_j are non-zero. This can be considered a lower bound on the achievable prediction error when the number of non-zero coefficients is small. (When the number of non-zero coefficients is large, the Oracle OLS will perform poorly, so we do not include the Oracle OLS result in settings with many non-zero coefficients.)

A factor that inevitably complicates comparisons is that most methods have many options and tuning parameters. Even a relatively straightforward method such as the Lasso has multiple tuning parameters that can affect performance, some of which involve tradeoffs in computing effort versus prediction accuracy: number of folds to use in the K-fold cross-validation step; what criterion to use for selecting the optimal penalty strength parameter; whether to “relax” the fit; etc. For each method, we tried to follow the recommendations

method	R package	brief description
PLR methods		
ridge regression	glmnet	PLR with convex L_2 penalty
Lasso	glmnet	PLR with convex L_1 penalty
Elastic Net	glmnet	PLR with linear combination of L_1 and L_2 penalties
SCAD	ncvreg	PLR with nonconvex SCAD penalty
MCP	ncvreg	PLR with minimax concave penalty
L0Learn	L0Learn	PLR with nonconvex L_0 , L_0L_1 or L_0L_2 penalty
SSLasso	SSLASSO	PLR with adaptive penalty based on a Laplace mixture prior
Trimmed Lasso	(none)	PLR with nonseparable “trimmed lasso” penalty
Bayesian and empirical Bayes methods		
BayesB	BGLR	MCMC with spike-and-slab prior (the “slab” is t)
Bayesian Lasso	BGLR	MCMC with scaled Laplace prior
varbvs	varbvs	variational inference with spike-and-slab prior
SuSiE	susieR	variational inference for “SuSiE” model

Table 1: Summary of methods compared in the simulations. All methods are implemented in R packages except the Trimmed Lasso which is implemented in MATLAB.

given in the software documentation or in published papers, and in some cases we changed the default settings to improve performance. However, with so many methods to compare, it was infeasible to find the best settings for each method. Appendix B includes additional information on how the methods were run in these experiments.

5.2 Design of Simulations

To test the methods in a wide variety of settings, we designed five sets of simulations, which we refer to as “Experiment 1” through “Experiment 5.” In each set of simulations, we varied one aspect of the simulation while keeping the other aspects fixed.

- In Experiment 1, we varied the “sparsity level”; that is, the proportion of variables with non-zero coefficients. We denote the sparsity level by s , so that $s = 1$ is the sparsest model (with only a single non-zero coefficient) and $s = p$ is the densest model (all variables affect \mathbf{y}).
- In Experiment 2, we varied the “total signal strength”; specifically, the proportion of variance in the response \mathbf{y} that is explained by \mathbf{X} . We refer to this parameter as “PVE”, short for “proportion of variance explained.”
- In Experiment 3, we considered different distributions for the non-zero coefficients. We use h to denote the distribution that was used to simulate the non-zero coefficients.
- In Experiment 4, we varied the number of predictors, p .

- In Experiment 5, we simulated residual errors (noise) in different ways. This was done to assess how departures from the assumption of normally-distributed noise—an assumption made by most methods—might affect performance.

These experiments focus on the case $p > n$. However, we note that Mr.ASH also performs well in the easier case where $p \ll n$; see Appendix C.1.

In all the experiments, we took the following steps to simulate each data set:

- First, we generated the $n \times p$ design matrix, \mathbf{X} . We considered three types of design matrices: (1) *independent variables*, in which the individual observations x_{ij} were simulated *i.i.d.* from the standard normal distribution; (2) *correlated variables*, in which each row of \mathbf{X} was an independent draw from the multivariate normal distribution with mean zero and a covariance matrix diagonal entries set to 1 and off-diagonal entries set to $\rho \in [0, 1]$; and (3) *real genotype data*, in which \mathbf{X} was a genotype data matrix from the Genotype-Tissue Expression (GTEx) project (GTEx Consortium, 2017). (Specifically, we used the processed genotype data sets generated in Wang et al. 2020.) In these data sets, the variables were genetic variants—specifically, single nucleotide polymorphisms, or “SNPs”—and the SNPs exhibited complex correlation patterns. Some SNP pairs had very strong correlations, approaching 1 or -1 . Each genotype matrix \mathbf{X} contained the genotypes of all SNPs within 1 Megabase (Mb) of a gene’s transcription start site after filtering out SNPs with minor allele frequencies less than 5% (see Wang et al. 2020 for details). Among the thousands of data sets used in Wang et al. 2020, we randomly selected 20 data sets for our simulations. Unless otherwise stated, we simulated independent variables with $n = 500$ and $p = 1,000$. For the genotype data sets, $n = 287$, and p ranged from 4,012 to 8,760. Each genotype matrix \mathbf{X} was centered and scaled so that the mean of each column was zero and its standard deviation was 1. The other data matrices were not centered or scaled.
- We chose s , the number of non-zero coefficients. Unless stated otherwise we set $s = 20$. We selected the indices $j \in \{1, \dots, p\}$ of the s non-zero coefficients uniformly at random among the p variables.
- We simulated the s non-zero coefficients b_j *i.i.d.* from some distribution, h . We used the following distributions: standard normal; uniform on $[-1, 1]$; double-exponential (Laplace) distribution centered at zero with variance $2\lambda^2$, $\lambda = 1$ (Gelman et al., 2013); t -distribution with 1, 2, 4 and 8 degrees of freedom; and a point mass (so that all coefficients were the same). Unless stated otherwise, we simulated the coefficients from the standard normal distribution.
- Finally, we simulated the responses $y_i = \sum_{j=1}^p x_{ij}b_j + e_i$, where e_i was drawn from some noise distribution. We used the following noise distributions: normal with mean zero; uniform distribution with mean zero; double-exponential (Laplace) distribution centered at zero; and t -distribution with 1, 2, 4 and 8 degrees of freedom. In all cases, the variance of the noise distribution was adjusted to attain the target PVE; specifically, denoting the variance of the noise distribution by σ^2 , we set it to $\sigma^2 = \text{Var}(\mathbf{X}\mathbf{b}) \times \frac{1-\text{PVE}}{\text{PVE}}$. Unless stated otherwise, we set $\text{PVE} = 0.5$, and simulated $e_i \sim N(0, \sigma^2)$, with σ^2 adjusted to attain the target PVE of 0.5.

We repeated the simulations 20 times for each simulation setting in Experiments 1–5. The test sets used to evaluate the model fits were the same size as the training sets and generated using the same coefficients \mathbf{b} .

5.3 Evaluation

Each method returns $\hat{\mathbf{b}}$, an estimate of the regression coefficients. We evaluated this estimate using the (scaled) root mean squared prediction error in the test data:

$$\text{RMSE-scaled}(\mathbf{y}_{\text{test}}, \hat{\mathbf{b}}) \triangleq \frac{\text{RMSE}(\mathbf{y}_{\text{test}}, \hat{\mathbf{b}})}{\text{RMSE}(\hat{\mathbf{b}} = \mathbf{0})}, \quad (39)$$

where

$$\text{RMSE}(\mathbf{y}_{\text{test}}, \hat{\mathbf{b}}) \triangleq \frac{1}{\sqrt{n}} \|\mathbf{y}_{\text{test}} - \mathbf{X}_{\text{test}} \hat{\mathbf{b}}\|. \quad (40)$$

$\text{RMSE}(\hat{\mathbf{b}} = \mathbf{0}) \triangleq \sigma/\sqrt{1 - \text{PVE}}$ denotes the expected RMSE for the “null predictor”, $\hat{\mathbf{b}} = \mathbf{0}$. The value of “RMSE-scaled” will vary from $\sqrt{1 - \text{PVE}}$ (for the “oracle” predictor) to approximately 1 (for the null predictor); however, values greater than 1 are possible if $\hat{\mathbf{b}}$ performs worse than the null predictor.

5.4 Results

We present the results of Experiments 1–5 in Sections 5.4.1–5.4.5, and we give a high-level summary in Section 5.4.6.

5.4.1 EXPERIMENT 1—VARYING THE SPARSITY LEVEL

In the first set of simulations, we varied s , the number of non-zero predictors, which controls the model sparsity when p is fixed. The results of these simulations, summarized in Figures 2 and 3, highlight differences in performance between the methods that were better suited to a particular level of sparsity versus the methods that better adapted to different levels of sparsity. For example, SuSiE, L0Learn with the L_0 penalty and the Trimmed Lasso better adapted to sparse settings and often performed poorly in dense settings; by contrast, ridge regression and the Bayesian Lasso are better adapted to dense settings, and as expected performed worse in sparse settings.

Other methods performed more consistently across different sparsity levels. In particular, the Lasso and Elastic Net performed somewhat similarly, with the Elastic Net usually performing slightly better, but in many settings there was a noticeable gap in performance between these two methods compared with the best performing method. The nonconvex, penalty-based methods MCP and SCAD performed similarly to one another, and were competitive in many settings, except in denser-signal data sets in some scenarios (e.g., in the low-dimension settings, with $p = 200$ non-zero coefficients). The Spike-and-Slab Lasso (SSLasso) with the adaptive penalty performed competitively across a range of sparse to dense settings when predictor variables were independent, but it performed less well (and sometimes very poorly) in settings with correlated predictors.

Mr.ASH was competitive at all sparsity levels, consistently achieving the best performance or close to the best in almost all simulation settings. Mr.ASH tended to outperform

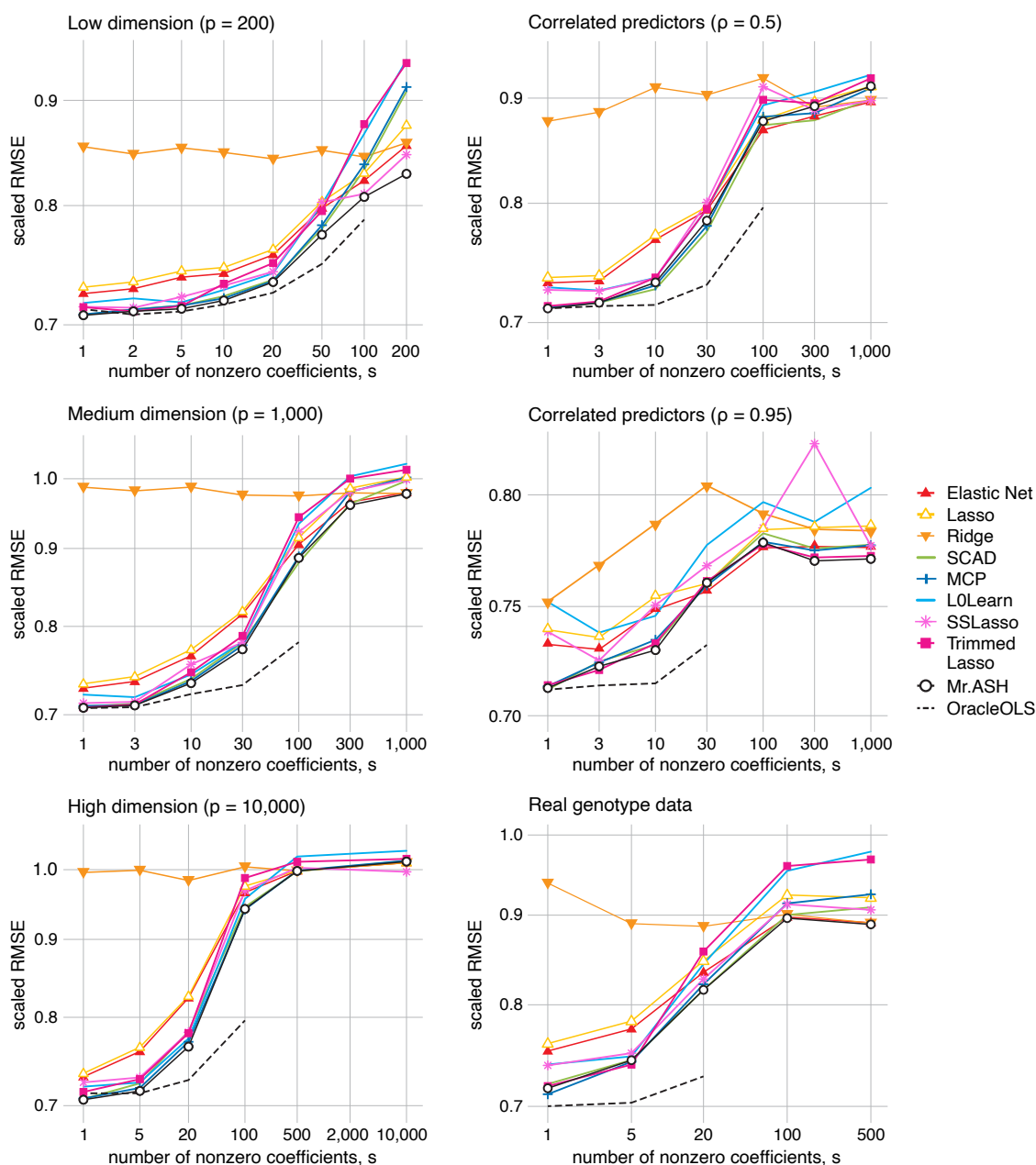


Figure 2: Results from Experiment 1 in which the sparsity level, s , was varied. Each point shows the prediction error (scaled RMSE) averaged over 20 simulations.

other methods in data sets with correlated predictors, with a couple of exceptions: the Bayesian Lasso had better prediction accuracy in settings with the densest signals, and SuSiE was better in some of the sparse simulation settings with genotype data. One possible explanation for this is that the mean-field variational approximation used by Mr.ASH is less appropriate in settings with strongly correlated predictors, whereas the Bayesian Lasso

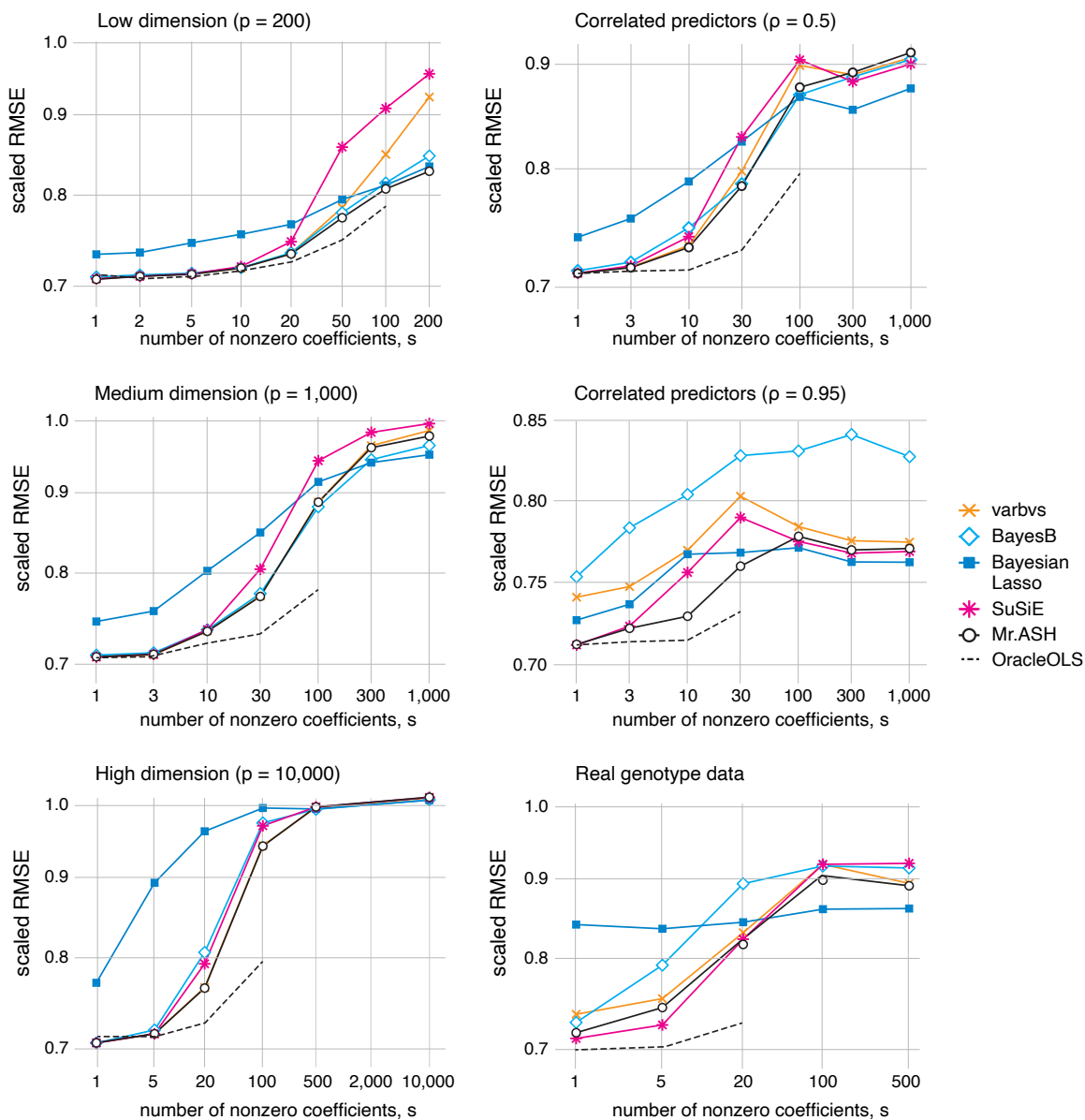


Figure 3: Results from Experiment 1 in which the sparsity level, s , was varied. Each point shows prediction error (scaled RMSE) averaged over 20 simulations. The Mr.ASH results are the same as in Figure 2 and provide a common point of reference.

and SuSiE do not have this issue (Bayesian Lasso uses MCMC to compute the posterior estimate; and SuSiE uses a different variational approximation that can deal with strong correlations).

Overall, these results illustrate the versatility of Mr.ASH and the effectiveness of the VEB approach to adapt to different sparsity levels by adapting the prior—and therefore penalty—to the data.

Two other results stand out. First, BayesB performed poorly in the simulations with correlated predictors. In principle, correlated predictors should not cause problems for Bayesian methods such as BayesB. We suspect that this reflects failure of the MCMC to converge and that the performance of BayesB (and the Bayesian Lasso) would improve with longer MCMC runs.

Second, although varbvs is based on the same variational approximation as Mr.ASH, its prediction accuracy was generally worse than Mr.ASH. This is particularly evident in some settings with dense signals, probably because the default settings in varbvs are designed to favor sparse priors. Also, varbvs performed worse than Mr.ASH in some data sets with correlated predictors, probably because varbvs does not put as much effort into initialization. (See Appendix C.2 for investigations of the impact of initialization on the performance of Mr.ASH.)

5.4.2 EXPERIMENT 2—VARYING TOTAL SIGNAL STRENGTH

We did not expect strong systematic differences in performance as the total signal strength (PVE) varied. Results for sparse simulations (top half of Figure 4) generally matched these expectations; the performance curves of the different methods generally did not cross as the PVE varies. And, among the methods compared, Mr.ASH consistently performed among the best. However, results in dense simulations (bottom half of Figure 4), showed a different pattern: some methods that performed competitively at moderate PVE were no longer competitive at high PVE. This included Mr.ASH, which was unexpected because Mr.ASH should, in principle, adapt well to both sparse and dense signals. Further, while Mr.ASH includes ridge regression as a special case, the ridge regression estimates yielded better predictions than Mr.ASH in these dense data sets, suggesting a failure of the variational EB approach to appropriately adapt the prior. We believe that this failure occurred because the fully factorized (mean-field) variational method has a tendency in some settings to favor sparse priors over dense priors; under sparse priors, the true posterior distribution comes close to factorizing so the gap between the ELBO and the true evidence is small (*i.e.*, the ELBO is a tight lower bound), whereas under dense priors with large p there are stronger dependencies in the posterior, especially when the PVE is large. Thus, the factorization assumption is more strongly violated and the ELBO underestimates the evidence more. Since our EB approach seeks to optimize the ELBO in place of the true evidence, this creates a bias towards estimating sparse priors rather than dense priors. In this sense, the VEB approach could be characterized as leaning towards the “bet on sparsity” principle (Hastie et al., 2009, p. 610), which argues that one should “use a procedure that does well in sparse problems, since no procedure does well in dense problems.” These simulations therefore suggest a situation (dense signal, high PVE, $p = 2n$) where this principle fails.

The good performance of BayesB in these simulations with dense signals and high PVE suggests an advantage of MCMC over the variational approximation in dense scenarios in which many predictors have a small effect on the regression outcome. Consistent with

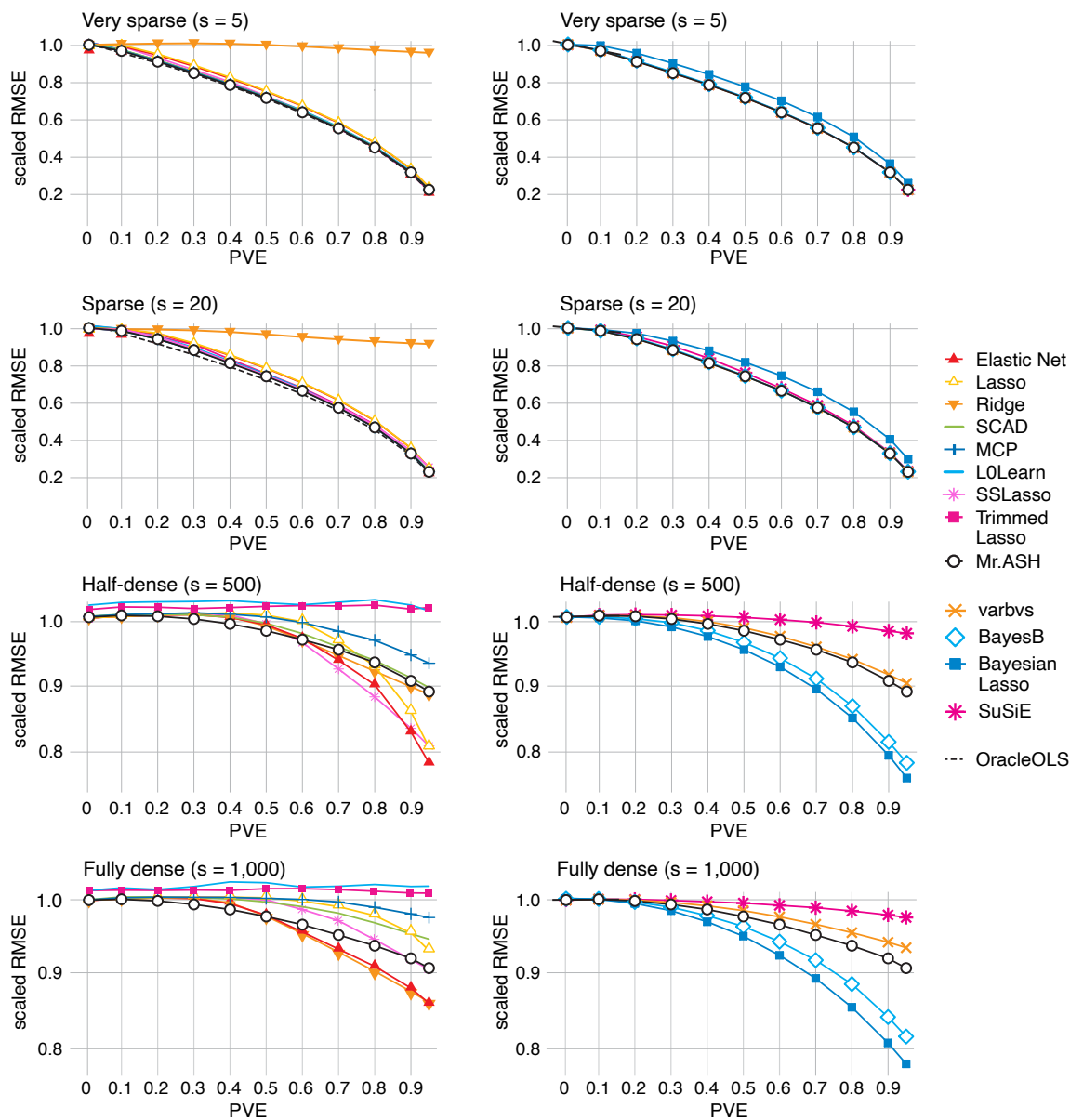


Figure 4: Results from Experiment 2 in which total signal strength (PVE) was varied. Each point shows prediction error (scaled RMSE) averaged over 20 simulations. Left panels show PLR methods; right panels show Bayes-based methods. The Mr.ASH results are included in all plots to provide a common reference point.

Experiment 1, ridge regression produced comparatively poor predictions in sparse settings, whereas L0Learn, SuSiE and the Trimmed Lasso had worse accuracy in dense settings.

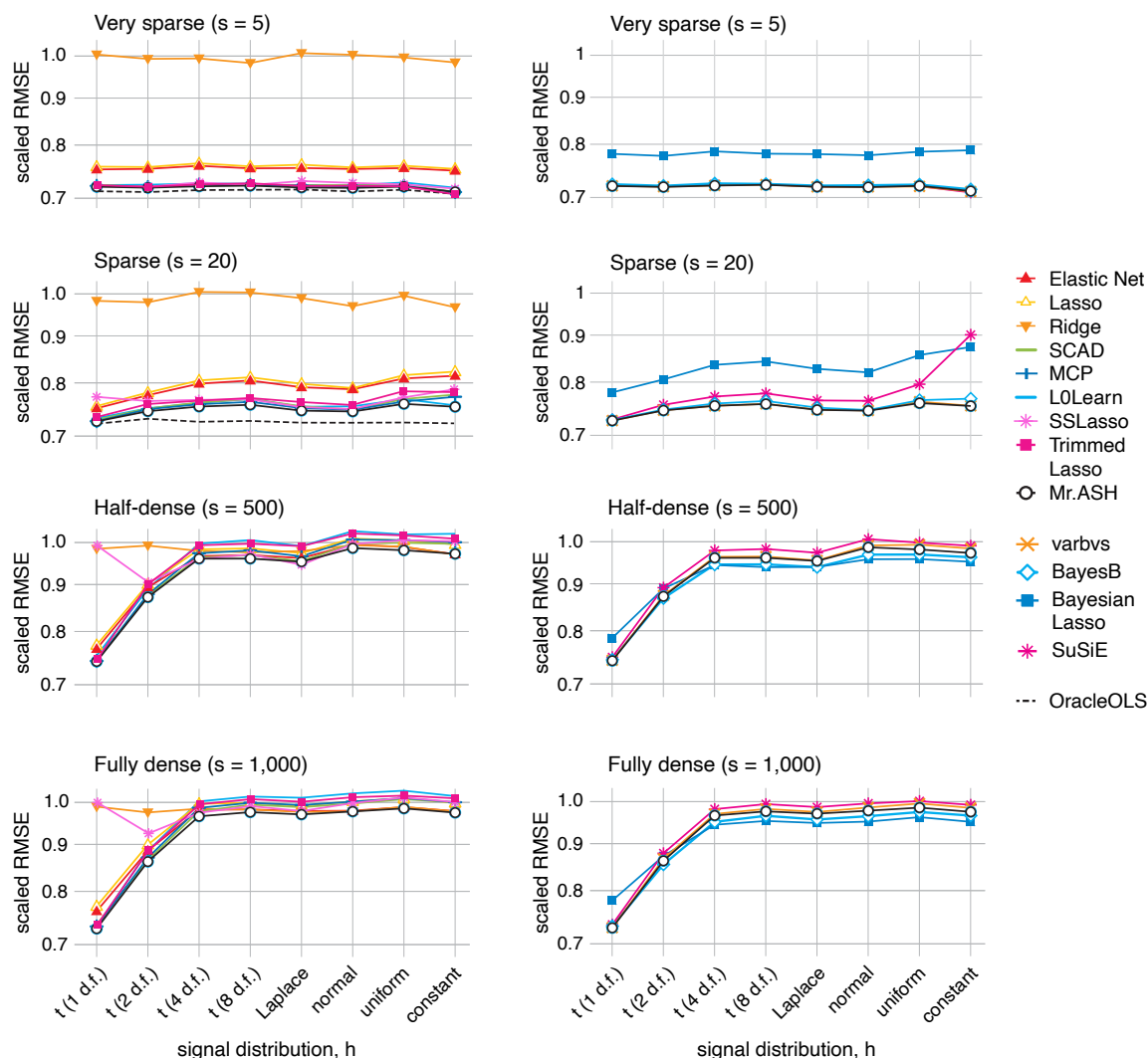


Figure 5: Results from Experiment 3 in which the signal distribution (h) was varied. Each point shows the prediction error (scaled RMSE) averaged over the 20 simulations at that setting. Left panels show PLR methods; right panels show Bayes-based methods. The Mr.ASH results are included in all the plots to provide a common reference point.

5.4.3 EXPERIMENT 3—VARYING SIGNAL DISTRIBUTION

For most methods, the distribution used to simulate the coefficients, h , had only a small impact on performance (Figure 5). There were a few exceptions to this. For example, in dense settings the ridge regression method struggled with long-tailed effect-size distributions such as the t -distribution with 1 degree of freedom. Consider that, when h is long-tailed, often a small number of coefficients will dominate, so the normal prior in ridge regression

is poorly suited for this case. The SSLasso also performed poorly in this setting for reasons that remain unclear to us.

5.4.4 EXPERIMENT 4—VARYING THE NUMBER OF PREDICTORS

This experiment assessed how prediction accuracy and computational effort change with the number of predictors, p (Figure 6). In general, running times for different methods increased similarly with p . For large p , SuSiE, L0Learn and the Lasso were the fastest methods. The Elastic Net was considerably slower than other methods because we tuned both of the Elastic Net parameters by CV whereas other methods tuned by CV involved tuning only one parameter. (The Elastic Net could be run faster by tuning only one parameter but at the risk of losing accuracy in some settings.) Fitting the Mr.ASH prior involved tuning a large number of parameters, but because it tuned these parameters via VEB rather than by CV, Mr.ASH ended up being roughly as fast as methods that tune a single parameter by CV. This is an important benefit of the VEB approach.

The long running times for the Trimmed Lasso were due to running the algorithm at several different settings of its “target sparsity level” parameter (k), and the Trimmed Lasso was slow when k was large (even with the settings that were suggested in the software documentation for larger data sets). When k was small, say, less than 10, the Trimmed Lasso running times were comparable to the other methods.

An important detail is that a large fraction of the effort in running Mr.ASH was due to running the Lasso (which was used to initialize Mr.ASH). The Lasso initialization often greatly reduced the number of iterations required for the Mr.ASH coordinate ascent updates to converge, so the total running time of Mr.ASH with Lasso initialization was often not much greater than Mr.ASH with null initialization.

Although MCMC methods have a reputation for being slow, here the MCMC-based methods—the Bayesian Lasso and BayesB—had similar running times to other methods. The running time of these methods also increased approximately linearly with p because the defaults in the software implementations set the number of iterations proportional to p . (The per-iteration cost is independent of p when the model configurations are sparse.) However, as p increased, the relative prediction performance of these methods got worse. For the Bayesian Lasso, this was probably due to the fact that we fixed s (the number of non-zero coefficients) to simulate the data sets, so simulations with larger p were based on sparser models, and the Bayesian Lasso tends to be less competitive in sparser settings. For BayesB, the reduction in prediction accuracy may instead reflect a failure of the Markov chain to adequately explore the posterior distribution, and perhaps better performance could have been achieved by running the MCMC longer (at increased computational cost). Nonetheless, it is interesting that BayesB obtained better prediction accuracy than the Lasso at roughly the same computational effort.

To summarize, Mr.ASH consistently achieved the best prediction accuracy in these simulations, with computational effort comparable with the fastest methods such as L0Learn and Lasso.

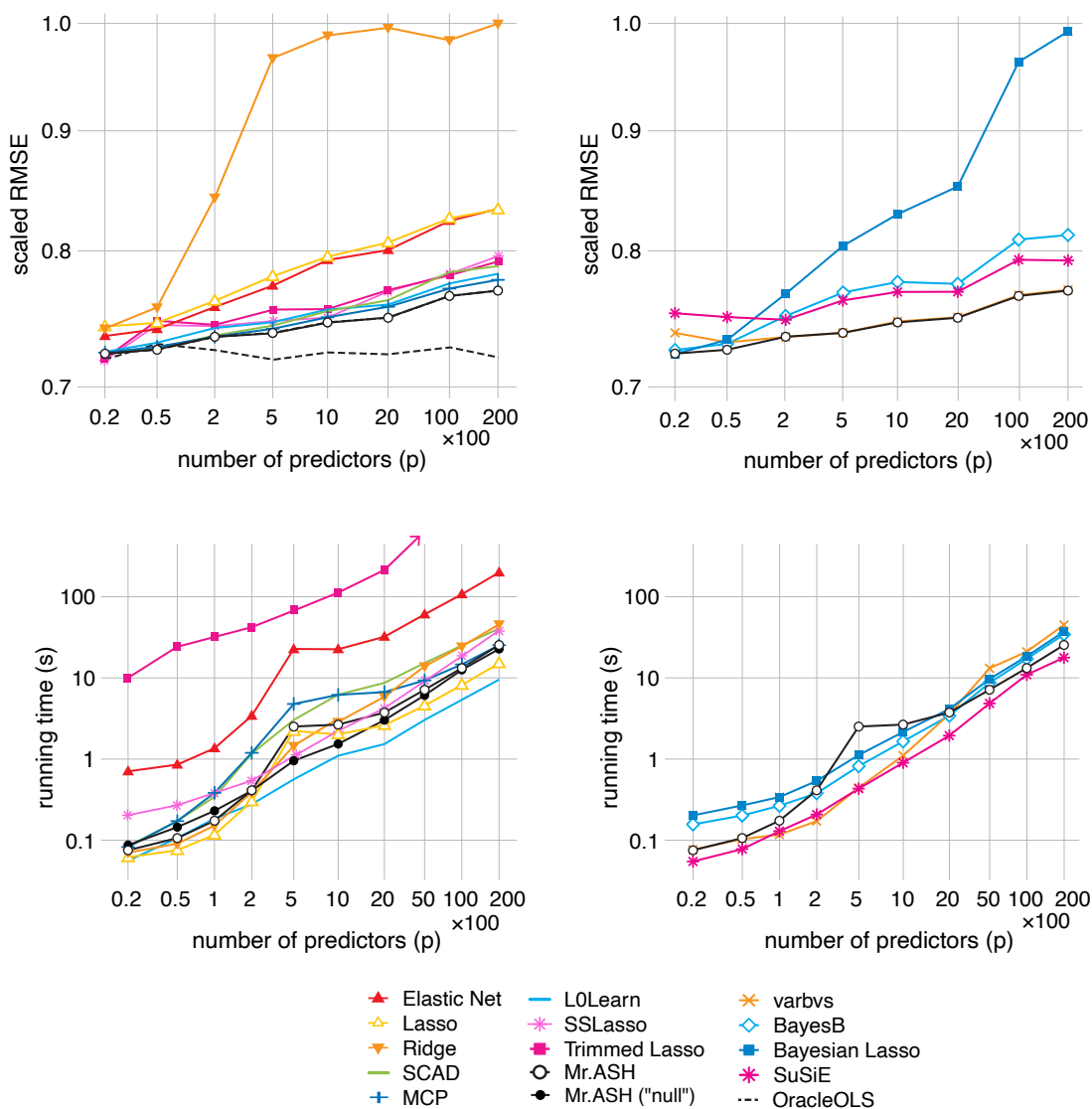


Figure 6: Results from Experiment 4 in which the number of predictors, p , was varied. Each point shows the scaled RMSE (top row) and the running time (bottom row) averaged over 20 simulations. Left panels show PLR methods; right panels show Bayes-based methods. The Mr.ASH results are included in all plots to provide a common reference point. Running times for all methods except Trimmed Lasso are from running the methods in R 3.5.1 (R Core Team, 2019) on a machine with a quad-core 2.6 GHZ Intel Core i7 processor and 16 GB of memory. R was installed from macOS binaries and we used the BLAS libraries that were distributed with R. The Trimmed Lasso was run using MATLAB 9.13 on machines with 4 Intel Xeon E5-2680v4 (“Broadwell”) processors and 24 GB of memory. The Mr.ASH running times include running the Lasso to initialize the estimates; the running times for Mr.ASH with a “null” initialization ($\mathbf{b} = \mathbf{0}$) are also shown for comparison.

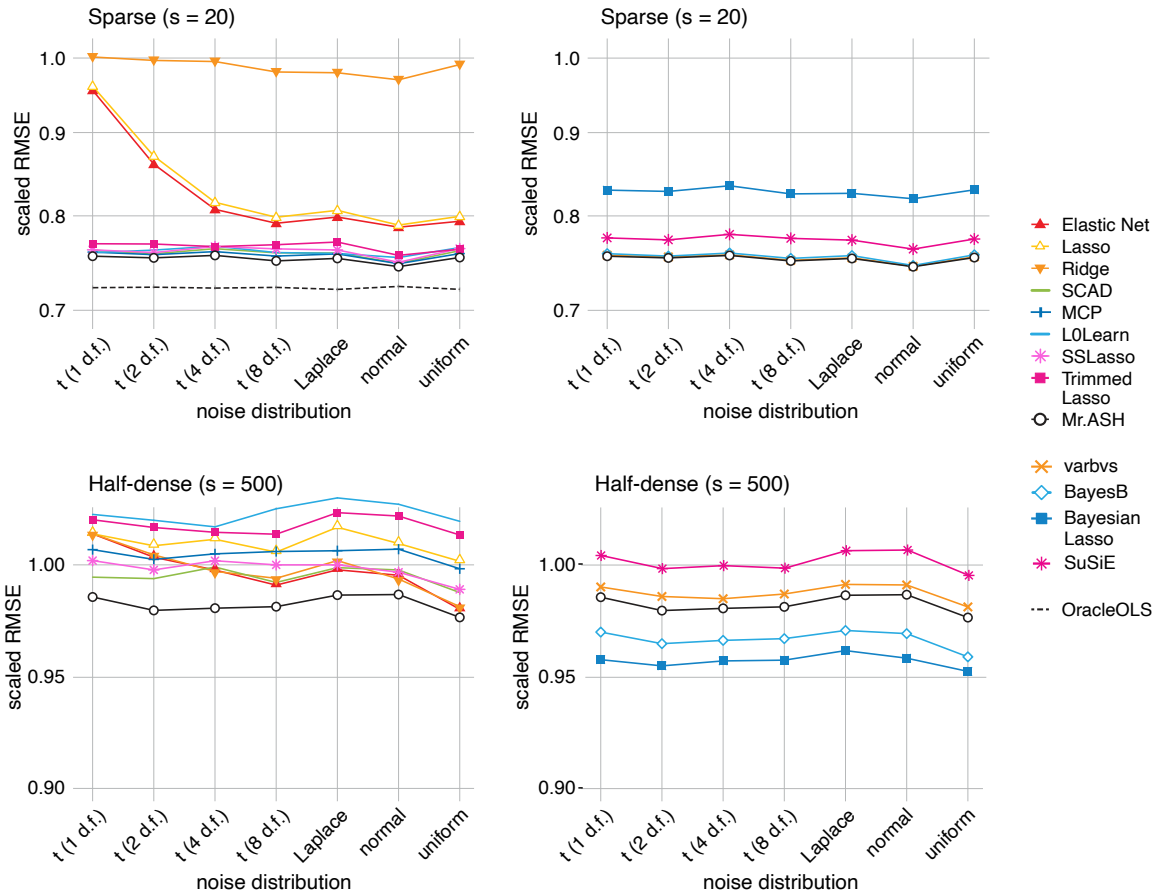


Figure 7: Results from Experiment 5 in which the noise distribution was varied. Each point shows the prediction error (scaled RMSE) averaged over 20 simulations. Left panels show PLR methods; right panels show Bayes-based methods. The Mr.ASH results are included in all plots to provide a common reference point.

5.4.5 EXPERIMENT 5—VARYING NOISE DISTRIBUTION

In our final set of experiments, we simulated data sets with different noise distributions. Reassuringly, most methods were largely insensitive to the noise distribution (Figure 7). However, the Lasso and Elastic Net both performed poorly in sparse settings when the noise was very heavy-tailed (t distribution with small degrees of freedom). We do not have an explanation for this result, which we have not seen previously.

5.4.6 SUMMARY OF THE RESULTS

The five experiments highlighted some differences in performance and behaviour among the multiple regression models. In Figure 8, we give a higher level summary of the results across all five of the experiments. To produce this summary, for each simulation t we calculated the relative prediction accuracy as the ratio of the RMSE to the best RMSE achieved in

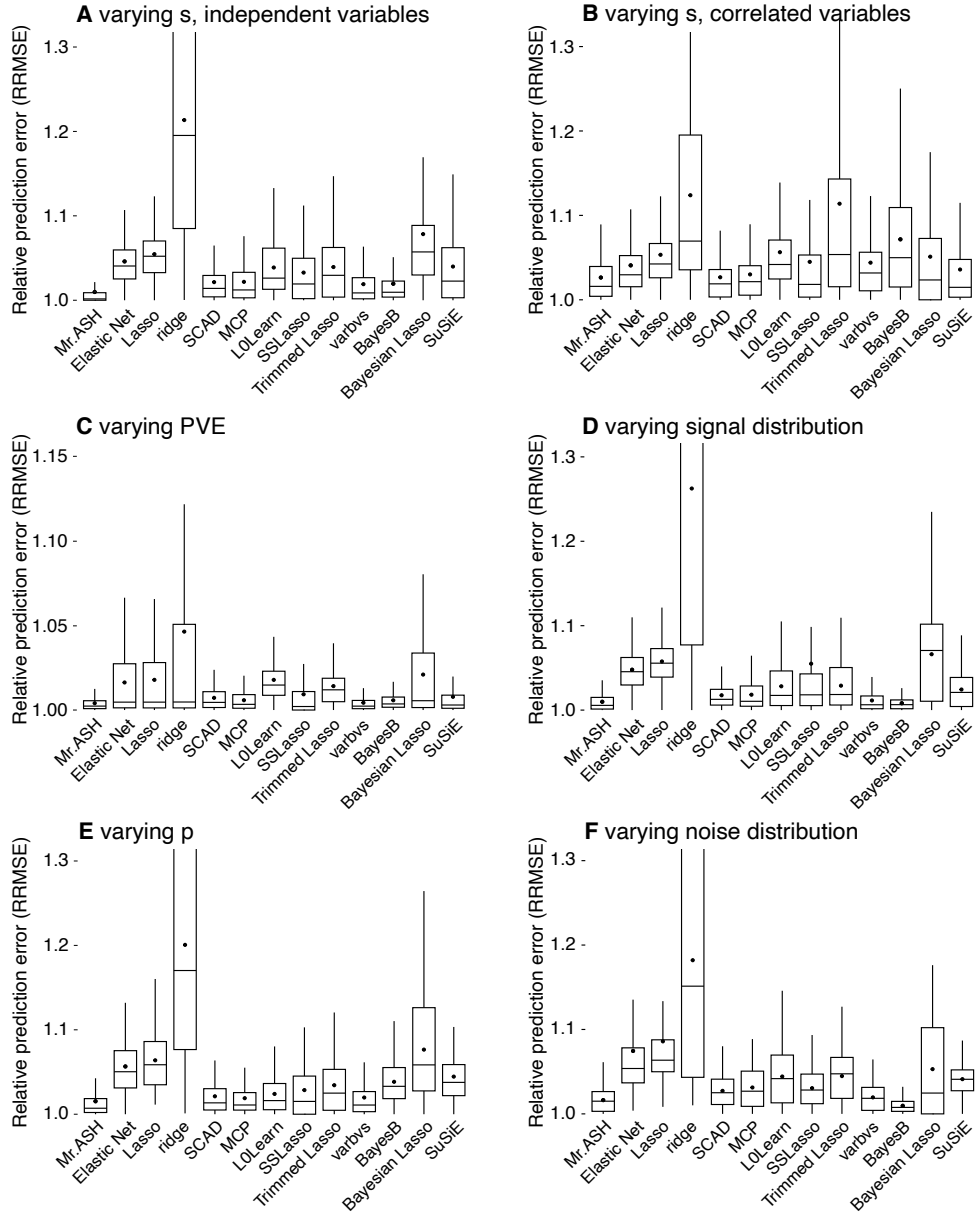


Figure 8: Summary of results from Experiments 1–5. The simulation results are summarized slightly differently in this figure to emphasize common trends. The boxplots show the distribution of RMSEs relative to the best performing method in each simulation (see the text for details). The horizontal line inside each box depicts the median; the dot depicts the mean; and the upper and lower lines depict the interquartile range.

that simulation among all the methods compared, $\text{RRMSE}_{tm} \triangleq \text{RMSE}_{tm} / \min_{m'} \text{RMSE}_{tm'}$, where RMSE_{tm} is the root mean squared error (40) generated by model m for the test set

method	running time (s)
L0Learn	8.12
Mr.ASH (“null”)	14.96
SSLasso	16.19
Lasso	17.98
Bayesian Lasso	18.24
Mr.ASH	18.56
MCP	21.19
SuSiE	22.54
BayesB	23.77
ridge regression	32.88
SCAD	33.59
varbvs	52.26
Elastic Net	223.66
Trimmed Lasso	609.54

Table 2: Average running times in Experiments 1–3.

in simulation t . Defined in this way, the RRMSE can never be smaller than 1 and the most accurate method in a given simulation has an RRMSE of 1.

Figure 8 highlights the consistently good prediction accuracy of Mr.ASH compared with the other methods across a range of settings; in all the panels (A–F), Mr.ASH’s performance was the best, or close to the best, among all methods compared, and was rarely much worse than the best method. The benefits of Mr.ASH were more mixed in the simulations with correlated predictors (Panel B). Yet, even in these simulations, Mr.ASH remained competitive, achieving an average prediction accuracy that was among the best. This good accuracy was also obtained efficiently with a computational effort that was not much higher than the fastest methods such as L0Learn and the Lasso (Table 2).

6 Discussion

We have presented a new VEB method for multiple linear regression with a focus on fast and accurate prediction. This VEB method combines flexible shrinkage priors with variational methods for efficient posterior computations. Variational methods and EB methods are sometimes criticized because of their tendency to understate uncertainty compared with “fully Bayesian” methods; see Morris (1983), Wang and Titterton (2005) and references therein for discussion. However, for some applications uncertainty is of secondary importance compared with speed and accuracy of point estimates. For example, speed and accuracy is often important when multiple regression is used simply to build an accurate predictor for downstream use (see Gamazon et al. 2015 for one such application). Our VEB approach seems particularly attractive for such uses.

A natural next step would be to produce similar VEB methods for non-Gaussian (*i.e.*, generalized) linear models (McCullagh and Nelder, 1989). Extension of our methods to logistic regression should be possible via additional approximations that allow for efficient analytic computations (Carbonetto and Stephens, 2012; Carbonetto et al., 2017; Jaakkola

and Jordan, 2000; Marlin et al., 2011; Bishop, 2006; Wang and Blei, 2013). Extension to other types of outcome distributions and link functions should also be possible but may require more work.

Our work also illustrates the benefits of an EB approach in an important and well studied statistical problem. While there is much theoretical (Johnstone and Silverman, 2004) and empirical (Efron, 2008) support for the benefits of EB approaches, EB approaches have not been widely adopted outside of specific research topics such as wavelet shrinkage (Johnstone and Silverman, 2005) and moderated estimation in gene expression studies (Smyth, 2004; Lu and Stephens, 2016; Zhu et al., 2019). Recent work has highlighted the potential for EB methods in other applications, including smoothing non-Gaussian data (Xing et al., 2021), multiple testing (Stephens, 2016; Sun and Stephens, 2018; Urbut et al., 2019; Gerard and Stephens, 2020), matrix factorization (Wang and Stephens, 2021) and additive models (Wang et al., 2020). We hope that these examples, including our work here, will inspire readers to apply EB approaches to new problems.

Acknowledgments and Disclosure of Funding

We thank Gao Wang for help with processing the the GTEx data, and for helpful discussions. This work was supported by NIH grant R01HG002585 to Matthew Stephens. We also thank the anonymous reviewers for their constructive suggestions for improvement and the staff at the Research Computing Center for providing the high-performance computing resources used to implement the numerical experiments.

Appendix Appendix A. Shrinkage Operators of Commonly Used Penalties

method	penalty function $\rho(t)$	shrinkage operator $S_\rho(t)$
normal shrinkage (ridge regression, L_2 penalty)	$\lambda t^2/2$	$\frac{t^2}{1+\lambda}$
hard thresholding (best subset, L_0 penalty)	$\lambda \times \mathbb{I}\{ t > 0\}$	$\begin{cases} t & \text{if } t < -\lambda, \\ t & \text{if } t > \lambda, \\ 0 & \text{otherwise} \end{cases}$
soft thresholding (Lasso, L_1 penalty)	$\lambda t $	$S_{\text{soft},\lambda} \triangleq \begin{cases} t + \lambda & \text{if } t < -\lambda, \\ t - \lambda & \text{if } t > \lambda, \\ 0 & \text{otherwise} \end{cases}$
Elastic Net	$(1 - \eta)\lambda t^2/2 + \eta\lambda t $	$S_{\text{soft},\eta\lambda/a}(t/a),$ $a = 1 + (1 - \eta)\lambda$
Minimax Concave Penalty	$\begin{cases} \lambda t - t^2/(2\eta) & \text{if } t \leq \eta\lambda, \\ \eta\lambda^2/2 & \text{otherwise} \end{cases}$	$\begin{cases} \frac{S_{\text{soft},\lambda}(t)}{1-1/(\eta-1)} & \text{if } t \leq \eta\lambda, \\ t & \text{otherwise} \end{cases}$
Smoothly Clipped Absolute Deviation	$\begin{cases} \lambda t & \text{if } t \leq 2\lambda, \\ \lambda^2(\eta + 1)/2 & \text{if } t > \eta\lambda, \\ \frac{\eta\lambda t - (t^2 + \lambda^2)/2}{\eta - 1} & \text{otherwise} \end{cases}$	$\begin{cases} S_{\text{soft},\lambda}(t) & \text{if } t \leq 2\lambda, \\ t & \text{if } t > \eta\lambda, \\ \frac{S_{\text{soft},\eta\lambda/(\eta-1)}(t)}{1-1/(\eta-1)} & \text{otherwise} \end{cases}$

Table 3: Some commonly used penalty functions and their corresponding shrinkage operators.

Appendix Appendix B. Additional Notes on the Methods Compared

Here we give additional details about how the methods were applied to the simulated data sets.

Ridge regression. We used function `cv.glmnet` from the `glmnet` R package to choose the penalty strength by CV. Specifically, we called `cv.glmnet` with `alpha = 0`, `intercept = TRUE` and `standardize = FALSE`; all other settings were kept at their defaults. The setting of the penalty strength parameter λ minimizing the mean CV error (`lambda.min`) was used to make predictions.

Lasso. We fit Lasso models in the the same way that we fit ridge regression models using `glmnet`, except that we set the Elastic Net mixing parameter to `alpha = 1`.

Elastic Net. Elastic Net models were fit similarly to ridge regression and Lasso models, again using the `cv.glmnet` interface from the `glmnet` package. The difference is that the

Elastic Net involves two tuning parameters: the penalty strength parameter λ and the “mixing” parameter α . `glmnet` does not provide an automated way to chose α by CV, so we ran `cv.glmnet` for 11 settings of α ranging from 0 to 1 then we chose (α, λ) minimizing the 10-fold CV error.

SCAD and MCP. We called the `cv.ncvreg` function from the R package `ncvreg` which performs k -fold CV to select the regularization parameter. We called `cv.ncvreg` with `nFolds = 10` and `penalty = "SCAD"` or `penalty = "MCP"`. We kept other settings at their defaults. By default, `cv.ncvreg` standardized \mathbf{X} and added an intercept to the model.

L0Learn. We called the `L0Learn.cvfit` function from the R package `L0Learn` which performs k -fold CV to select the penalty strength parameter λ_0 . We called `L0Learn.cvfit` with `penalty = "L0"` and `nFolds = 10`. We chose the setting of λ with the smallest (mean) CV error. We kept other settings at their defaults. By default, `L0Learn.cvfit` included an intercept in the model.

SSLasso. We called the `SSLASSO` function from the R package `SSLASSO` which fits coefficients paths for spike-and-slab linear regression models over a grid of values for the λ_0 regularization parameter. We used the “adaptive” variant of the Spike-and-Slab Lasso which was recommended over the “separable” variant. This function automatically standardizes \mathbf{X} and the model includes an intercept. `SSLASSO` did not provide an automated way to select λ_0 so we performed 5-fold CV and chose the setting of λ_0 minimizing the CV error. We did not perform CV to choose the “spike” penalty parameter λ_1 ; Ročková and George (2018) showed the performance is less sensitive to the choice of λ_1 .

Trimmed Lasso. We downloaded the MATLAB code for the Trimmed Lasso from <https://github.com/tal-amir/sparse-approximation-gsm> and we compiled the MEX file using gcc 10.2.0. We called function `sparse_approx_gsm_v1.22` with the following settings: `sparse_approx_gsm_v1.22(X,y,k,'profile','fast')` in which the target sparsity level k was one of $\{1, 5, 20, 100, 500, 2000, 10000\}$. The sparsity level was chosen by 5-fold CV, taking the setting that minimized the average CV error.

BayesB and Bayesian Lasso We called the `BGLR` function from the `BGLR` package, which simulates the posterior using a Gibbs sampler. We called `BGLR` with the following settings: `standardize = FALSE`, `nIter = 1500`, `burnIn = 500` and `model = "BayesB"` or `model = "ML"`. We kept other settings at their defaults. By default, an intercept was included in the model.

varbvs. We called the `varbvs` function from the `varbvs` package which fits an approximate posterior distribution for a Bayesian variable selection model using variational inference methods. All settings were kept at their defaults. The model included an intercept and \mathbf{X} was not standardized.

SuSiE. We called the `susie` function from the `susie` package which fits a SuSiE model using the iterative Bayesian stepwise selection (IBSS) algorithm. We called `susie` with `standardize = FALSE` and we set the upper bound on the number of “single effects” to 20.

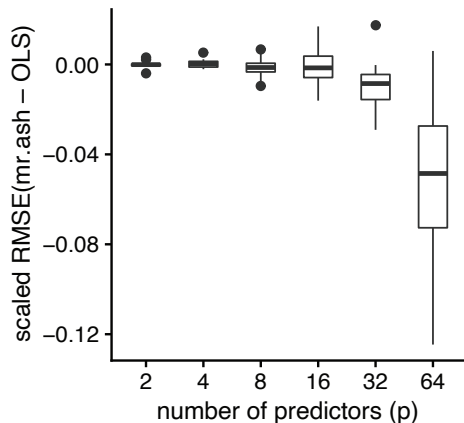


Figure 9: Comparison of Mr.ASH vs. OLS in data sets with $p \geq 64$, $n = 200$. Each box in the boxplot summarizes the difference in the test set prediction error (scaled RMSE) between the Mr.ASH predictions and the ordinary least squares (OLS) estimates across 20 simulations. The horizontal line inside each box depicts the median; the dot depicts the mean; the upper and lower lines depict the interquartile range.

Appendix Appendix C. Additional Experiments

C.1 Simulations with $p < n$

Although the paper focusses on large-scale multiple linear regression with many predictor variables, Mr.ASH can also be applied in settings with $p \ll n$ where one would expect ordinary least squares (OLS) to work well. To illustrate this, we simulated data sets with $\text{PVE} = 0.5$, $n = 200$, $p \leq 64$ and $s = p$ (all predictors had non-zero coefficients). The results show that Mr.ASH performs similarly to OLS when p is very small and outperforms the OLS estimate as p increases (Figure 9).

C.2 Impact of Initialization and Update Order on Prediction Accuracy

Since Mr.ASH is solving a nonconvex optimization problem, and therefore is only guaranteed to converge to a local optimum (except in special cases), the quality of the solution—and hence the accuracy of the predictions—can be sensitive to initialization. This situation is similar to other methods such as SCAD that solve nonconvex optimization problems, but different from methods such as the Lasso that solve a convex optimization problem and are therefore guaranteed to end up with the same final estimates irrespective of initialization (provided of course that the algorithm is given enough time to converge to the solution). Additionally, the order in which the coordinatewise updates are performed can also affect which local solution the Mr.ASH algorithm converges to (Ray and Szabó, 2022). For these nonconvex optimization problems, sometimes a “smart” initialization or update order can lead to a better local solution. From experience, we have found that initializing Mr.ASH to the cross-validated Lasso estimate of \mathbf{b} seems to work well and does not greatly

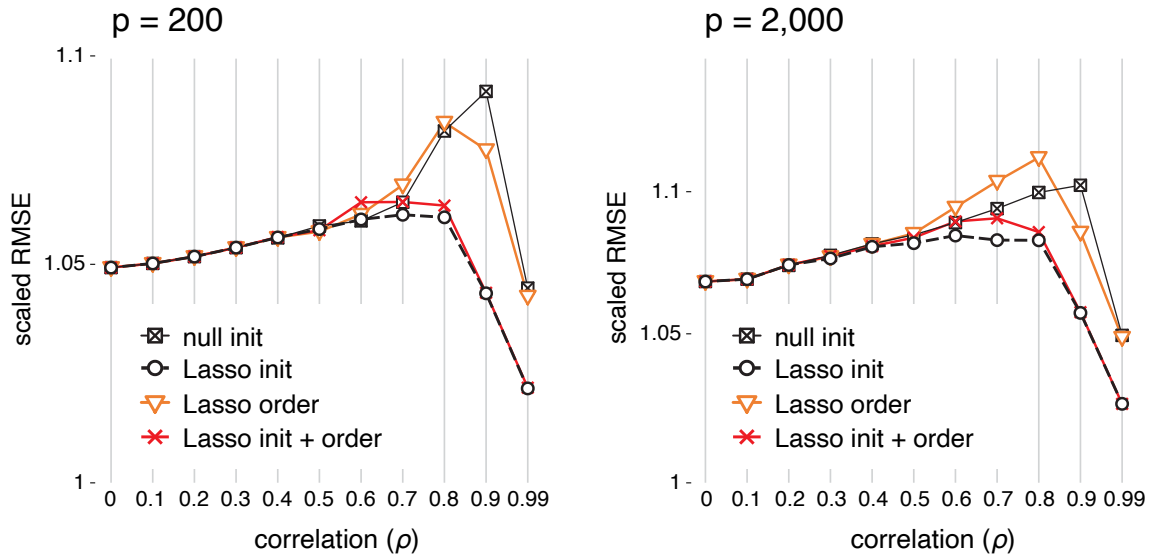


Figure 10: Comparison of Mr.ASH with different initializations and update orders. Each point shows the prediction error (scaled RMSE) averaged over the 20 simulations at that setting.

increase computational effort. In this experiment, we investigated the benefits of a smart initialization.

Specifically, we compared the following four Mr.ASH variants:

- **“Null” initialization.** The posterior mean coefficients $\bar{\mathbf{b}}$ are initialized to zero and the coordinates $j = 1, \dots, p$ are updated in a random order. By “random order”, we mean a random permutation of the indices $1, 2, \dots, p$. A new random permutation is generated for each iteration of the algorithm (that is, for each iteration of the repeat-until loop in Algorithm 2).
- **Lasso initialization.** The posterior mean coefficients are set to $\mathbf{b} = \hat{\mathbf{b}}^{\text{lasso}}$ (see Section 3.5.3), and the coordinates $j = 1, \dots, p$ are updated in a random order.
- **Lasso update order.** The coordinates $j = 1, \dots, p$ are updated in the order that they are estimated to have non-zero coefficients as the strength of the Lasso penalty is decreased. We call this the “Lasso update order,” and it can be understood as the order in which the coefficients “enter the Lasso path” (Su et al., 2017). (Note that determining the Lasso update order typically takes less effort than computing $\hat{\mathbf{b}}^{\text{lasso}}$ because the cross-validation step is avoided.)
- **Lasso initialization and Lasso update order.** Both the Lasso initialization $\bar{\mathbf{b}} = \hat{\mathbf{b}}^{\text{lasso}}$ and Lasso update order are used.

Initialization of σ^2 and $\boldsymbol{\pi}$ is described in Section 3.5.3.

To assess the benefits of these four initialization and update order strategies, we simulated data sets with varying correlation strengths among the predictors, then we compared the performance of the four Mr.ASH variants. The results of these simulations are summarized in Figure 10. The smart initialization and update ordering provided little benefit when the variables were not correlated or only weakly correlated, but produced considerable gains in prediction accuracy when the variables were strongly correlated. Interestingly, once the coefficients were initialized to the Lasso estimates, there was no additional benefit to updating the coordinates using the Lasso update order.

In summary, initializing the coefficients to the cross-validated Lasso estimates is a simple way to improve the performance of Mr.ASH when predictors are strongly correlated.

Appendix Appendix D. More General Formulation of the Normal Means model

In Section 2 we defined the normal means (NM) model for the special case when all observations j have the same variance, σ^2 . This special case was sufficient to develop the VEB methods with the assumption that $\mathbf{x}_j^T \mathbf{x}_j = 1$, $j = 1, \dots, p$. Here, we extend the NM model to allow for observation-specific variances, σ_j^2 , which is needed to generalize the VEB method to cases in which the $\mathbf{x}_j^T \mathbf{x}_j = 1$ assumption no longer holds.

D.1 The Normal Means Model

Let $\text{NM}_p(f, \mathbf{s}^2)$ denote the normal means model with prior f and observation-specific variances $\mathbf{s}^2 = (s_1^2, \dots, s_p^2) \in \mathbb{R}_+^p$:

$$\begin{aligned} y_j \mid b_j, s_j^2 &\sim N(b_j, s_j^2), \\ b_j &\overset{i.i.d.}{\sim} f, \quad j = 1, \dots, p, \end{aligned} \tag{41}$$

such that $y_j, b_j \in \mathbb{R}$, $j = 1, \dots, p$. We assume priors that are mixtures of zero-mean normals, $f \in \mathcal{G}(u_1^2, \dots, u_K^2)$, $u_k^2 \geq 0$, $k = 1, \dots, K$, so that any prior can be written as

$$b_j \sim \sum_{k=1}^K \pi_k N(0, u_k^2),$$

such that $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K) \in \mathbb{S}^K$.

As in (5), it is helpful in the derivations to make use of the latent variable representation:

$$\begin{aligned} p(\gamma_j = k \mid f) &= \pi_k \\ b_j \mid f, \gamma_j = k &\sim N(0, u_k^2), \end{aligned} \tag{42}$$

with $\gamma_j \in \{1, \dots, K\}$, $j = 1, \dots, p$. We write the joint prior for b_j, γ_j as

$$\begin{aligned} p_{\text{prior}}(b_j, \gamma_j = k) &\triangleq p(b_j, \gamma_j = k \mid f) \\ &= \pi_k N(b_j; 0, u_k^2). \end{aligned} \tag{43}$$

In the expressions below we sometimes write the joint prior as $p_{\text{prior}}(f)$ to make its dependence on f explicit.

Note that the definition of the NM model given in the main text (Section 2), with prior (4), is a special case of these definitions and can be obtained with the substitutions $s_j^2 \leftarrow \sigma^2, j = 1, \dots, p$ and $u_k \leftarrow \sigma^2 \sigma_k^2, k = 1, \dots, K$.

D.2 Posterior Distribution under Normal Means Model with One Observation

Let $q^{\text{NM}}(b, \gamma \mid y, s^2, f)$ denote the posterior distribution of b, γ under the normal means model $\text{NM}_1(f, s^2)$ with a single observation ($p = 1$):

$$\begin{aligned} y \mid b, s^2 &\sim N(b, s^2) \\ b &\sim f. \end{aligned} \tag{44}$$

For a mixture of normals prior, $f \in \mathcal{G}(u_1^2, \dots, u_K^2)$, the posterior distribution can be written as

$$p^{\text{NM}}(b \mid y, s^2, f) = \sum_{k=1}^K \phi_{1k} N(b; \mu_{1k}, s_{1k}^2), \tag{45}$$

in which the posterior component means μ_{1k} , variances s_{1k}^2 , responsibilities ϕ_{1k} , and component (marginal) likelihoods L_k are

$$\mu_{1k} \triangleq \mu_{1k}(y; f, s^2) = \frac{u_k^2}{s^2 + u_k^2} \times y \tag{46}$$

$$s_{1k}^2 \triangleq s_{1k}^2(y; f, s^2) = \frac{s^2 u_k^2}{s^2 + u_k^2} \tag{47}$$

$$\phi_{1k} \triangleq \phi_{1k}(y; f, s^2) = \frac{\pi_k L_k}{\sum_{k'=1}^K \pi_{k'} L_{k'}} \tag{48}$$

$$L_k \triangleq L_k(y; f, s^2) = p(y \mid s^2, f, \gamma = k) = \int p(y \mid b, s^2) p(b \mid f, \gamma = k) db \tag{49}$$

The posterior expressions for the NM model given in the main text (10–13) can be recovered from these more general expressions with the substitutions $s_j^2 \leftarrow \sigma^2, j = 1, \dots, p$ and $u_k \leftarrow \sigma^2 \sigma_k^2, k = 1, \dots, K$.

D.3 Evidence Lower Bound for Normal Means Model with One Observation

Given some probability density on $b \in \mathbb{R}$, denoted by q , the ELBO for the normal means model $\text{NM}_1(f, s^2)$ with observation y is

$$\begin{aligned} F_1^{\text{NM}}(q, f, s^2; y) &= \log p(y \mid f, s^2) - D_{\text{KL}}(q \parallel p^{\text{NM}}) \\ &= \mathbb{E}_q[\log p(y \mid b, s^2)] - D_{\text{KL}}(q \parallel p_{\text{prior}}(f)) \\ &= -\frac{1}{2} \log(2\pi s^2) - \frac{1}{2s^2} \mathbb{E}_q[(y - b)^2] - D_{\text{KL}}(q \parallel p_{\text{prior}}(f)). \end{aligned} \tag{50}$$

With this expression, we state the following result.

Lemma 6 (Normal means posterior as maximum of ELBO) The posterior distribution (45) under the NM model $\text{NM}_1(f, s^2)$ with observation y maximizes the ELBO (50); that is,

$$p^{\text{NM}} = \operatorname{argmax}_q -\frac{1}{2s^2} \mathbb{E}_q[(y - b)^2] - D_{\text{KL}}(q \parallel p_{\text{prior}}(f)).$$

From this lemma it follows that any q maximizing the ELBO (50) must have the following form:

$$q(b) = \sum_{k=1}^K \phi_{1k} N(b; \mu_{1k}, s_{1k}^2), \quad (51)$$

with $\phi_{1k} \geq 0$, $\mu_{1k} \in \mathbb{R}$, $s_{1k}^2 > 0$, $k = 1, \dots, K$.

For any q of the form (51), the ELBO (50) has an analytic expression, which we derive in part by making use of the formula for the K-L divergence between two normal distributions (Hastie et al., 2009):

$$F_1^{\text{NM}}(q, f, s^2; y) = \mathbb{E}_q[\log p(y | b, s^2)] - D_{\text{KL}}(q \| p_{\text{prior}}(f)), \quad (52)$$

in which

$$\mathbb{E}_q[\log p(y | b, s^2)] = -\frac{1}{2} \log(2\pi s^2) - \frac{(y - \bar{b})^2}{2s^2} - \frac{1}{2s^2} \sum_{k=1}^K [\phi_{1k}(\mu_{1k}^2 + s_{1k}^2) - \bar{b}^2]$$

and

$$D_{\text{KL}}(q \| p_{\text{prior}}(f)) = \sum_{k=1}^K \phi_{1k} \log \left(\frac{\phi_{1k}}{\pi_k} \right) - \frac{1}{2} \sum_{k=2}^K \phi_{1k} \left[1 + \log \left(\frac{s_{1k}^2}{u_k^2} \right) - \frac{\mu_{1k}^2 + s_{1k}^2}{u_k^2} \right],$$

and where \bar{b} is the posterior mean of b with respect to q , $\bar{b} = \sum_{k=1}^K \phi_{1k} \mu_{1k}$. Here we have assumed that the first component in the prior mixture is a point mass at zero, $\sigma_1^2 = 0$.

D.4 ELBO for Normal Means Model with Multiple Observations

Now we extend the above results for the single-observation NM model to the NM model with multiple observations, $\text{NM}_p(f, \mathbf{s}^2)$. Since the b_j 's are independent under the posterior, the ELBO is simply the sum of the ELBOs for the single-observation NM models:

$$F^{\text{NM}}(q, f, \mathbf{s}^2; y) = \sum_{j=1}^p F_1^{\text{NM}}(q_j, f, s_j^2; y_j). \quad (53)$$

From Lemma 6, the q that maximizes the ELBO is

$$q(\mathbf{b}) = \prod_{j=1}^p q_j(b_j)$$

$$q_j(b_j) = p^{\text{NM}}(b_j | y_j, f, s_j^2),$$

It also follows that any q maximizing the ELBO (53) must have the following form:

$$q(\mathbf{b}) = \prod_{j=1}^p q_j(b_j) \quad (54)$$

$$q_j(b_j) = \sum_{k=1}^K \phi_{1jk} N(b_j; \mu_{1jk}, s_{1jk}^2),$$

in which $\phi_{1jk} \geq 0$, $\mu_{1jk} \in \mathbb{R}$, $s_{1jk}^2 > 0$, $j = 1, \dots, p$, $k = 1, \dots, K$. For any q of the form (54), the analytic expression for the ELBO (53) is easily obtained by applying the analytic expression for the single-observation NM model (eq. 52).

Appendix Appendix E. Derivation of Algorithm 2, Proof of Proposition 1

We prove Proposition 1 by proving a slightly more general proposition that does not require that $\mathbf{x}_j^T \mathbf{x}_j = 1$, $j = 1, \dots, p$.

Proposition 7 Let $d_j = \mathbf{x}_j^T \mathbf{x}_j$, $j = 1, \dots, p$, and let

$$\tilde{b}_j \triangleq \frac{\mathbf{x}_j^T \bar{\mathbf{r}}_j}{\mathbf{x}_j^T \mathbf{x}_j}$$

denote the ordinary least squares (OLS) estimate of the coefficient b_j when the residuals $\bar{\mathbf{r}}_j$ are regressed against \mathbf{x}_j . See Proposition 1 for more definitions. Then we have the following results:

- (i) The coordinate ascent update $q_j^* \triangleq \operatorname{argmax}_{q_j} F(q, g, \sigma^2)$ is obtained by

$$q_j^*(b_j) = p^{\text{NM}}(b_j; \tilde{b}_j, \sigma^2/d_j, g_\sigma),$$

in which p^{NM} , defined in (45), is the posterior distribution of b under the following NM model:

$$\begin{aligned} \tilde{b} \mid b, \sigma^2 &\sim N(b, \sigma^2/d_j) \\ b \mid g, \sigma^2 &\sim g_\sigma. \end{aligned} \tag{55}$$

- (ii) The coordinate ascent update

$$g^* \triangleq \operatorname{argmax}_{g \in \mathcal{G}(\sigma_1^2, \dots, \sigma_K^2)} F(q, g, \sigma^2)$$

is achieved by setting

$$\begin{aligned} g^* &= \sum_{k=1}^K \pi_k^* N(0, \sigma_k^2) \\ \pi_k^* &= \frac{1}{p} \sum_{j=1}^p q_j(\gamma_j = k), \quad k = 1, \dots, K. \end{aligned}$$

- (iii) Using the parameterization of q in (54), and assuming that g is updated as in (ii) above, and $\sigma_1^2 = 0$, the coordinate ascent update

$$(\sigma^2)^* \triangleq \operatorname{argmax}_{\sigma^2 \in \mathbb{R}_+} F(q, g, \sigma^2)$$

is achieved by setting

$$(\sigma^2)^* = \frac{\|\bar{\mathbf{r}}\|^2 + \sum_{j=1}^p \sum_{k=2}^K \phi_{jk}(d_j + 1/\sigma_k^2)(\mu_{1jk}^2 + s_{1jk}^2) - \sum_{j=1}^p d_j \bar{b}_j^2}{n + p(1 - \pi_1^*)}.$$

Additionally, if q_1, \dots, q_p are updated as in (i) above, we obtain the simpler expression

$$(\sigma^2)^* = \frac{\|\bar{\mathbf{r}}\|^2 + \bar{\mathbf{b}}^T \mathbf{D}(\tilde{\mathbf{b}} - \bar{\mathbf{b}}) + \sigma^2 p(1 - \pi_1^*)}{n + p(1 - \pi_1^*)}, \tag{56}$$

where \mathbf{D} is the $p \times p$ diagonal matrix with diagonal entries d_1, \dots, d_p .

Note that Proposition 1 is a special case of this proposition when $d_j = 1$, for $j = 1, \dots, p$.

In the next sections, we prove parts (i), (ii) and (iii) of Proposition 7. These proofs start from the ELBO (18). From Bayes' rule,

$$p_{\text{post}}(\mathbf{b}) = \frac{p(\mathbf{y} \mid \mathbf{X}, \mathbf{b}, \sigma^2) p(\mathbf{b} \mid g, \sigma^2)}{\int p(\mathbf{y} \mid \mathbf{X}, \mathbf{b}, \sigma^2) p(\mathbf{b} \mid g, \sigma^2) d\mathbf{b}},$$

we can write the ELBO as

$$F(q, g, \sigma^2) = \mathbb{E}_q[\log p(\mathbf{y} \mid \mathbf{X}, \mathbf{b}, \sigma^2)] - \sum_{j=1}^p D_{\text{KL}}(q_j \parallel p_{\text{prior}}). \quad (57)$$

Next, using the property that q factorizes over the individual coordinates $j = 1, \dots, p$, we have

$$F(q, g, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \mathbb{E}_q[\|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2] - \sum_{j=1}^p D_{\text{KL}}(q_j \parallel p_{\text{prior}}).$$

E.1 Update for q_j

The coordinate ascent update for q_j involves solving the following optimization problem:

$$q_j^* = \underset{q_j}{\operatorname{argmax}} F(q, g, \sigma^2).$$

From (57), this is equivalent to solving

$$q_j^* = \underset{q_j}{\operatorname{argmax}} \mathbb{E}_q[\log p(\mathbf{y} \mid \mathbf{X}, \mathbf{b}, \sigma^2)] - D_{\text{KL}}(q_j \parallel p_{\text{prior}})$$

By rearranging terms, it can be shown that this is equivalent to solving

$$q_j^* = \underset{q_j}{\operatorname{argmax}} -\frac{d_j}{2\sigma^2} \mathbb{E}_{q_j}[(\tilde{b}_j - b_j)^2] - D_{\text{KL}}(q_j \parallel p_{\text{prior}}). \quad (58)$$

The right-hand side of (58) is the ELBO for the NM model (55); that is, if we ignore constant terms, the ELBO in (50) recovers (58) by making the substitutions $y \leftarrow \tilde{b}_j$, $s^2 \leftarrow \sigma^2/d_j$, $f \leftarrow g_\sigma$. And, therefore, from Lemma 6—and specifically from (51)—we have

$$q_j^*(b_j) = p^{\text{NM}}(b_j \mid \tilde{b}_j, \sigma^2/d_j, g_\sigma) = \sum_{k=1}^K \phi_{1jk} N(b_j; \mu_{1jk}, s_{1jk}^2),$$

in which

$$\begin{aligned} \phi_{1jk} &= \phi_{1k}(\tilde{b}_j, g_\sigma, \sigma^2/d_j) \\ \mu_{1jk} &= \mu_{1k}(\tilde{b}_j, g_\sigma, \sigma^2/d_j) \\ s_{1jk}^2 &= s_{1k}^2(\tilde{b}_j, g_\sigma, \sigma^2/d_j). \end{aligned}$$

This proves part (i) of Proposition 7.

E.2 Update for g

The coordinate ascent update for g involves solving the following optimization problem:

$$g^* \leftarrow \operatorname{argmax}_{g \in \mathcal{G}} F(q, g, \sigma^2). \quad (59)$$

Recall, for the mixture prior with fixed mixture components, fitting g reduces to fitting the mixture weights, $\boldsymbol{\pi}$. Since $\boldsymbol{\pi}$ only appears in the ELBO in the K-L divergence term with respect to the prior, solving (59) is equivalent to solving

$$\boldsymbol{\pi}^* = \operatorname{argmin}_{\boldsymbol{\pi} \in \mathbb{S}^K} \sum_{j=1}^p D_{\text{KL}}(q_j \parallel p_{\text{prior}}),$$

which simplifies further:

$$\boldsymbol{\pi}^* = \operatorname{argmax}_{\boldsymbol{\pi} \in \mathbb{S}^K} \sum_{j=1}^p \sum_{k=1}^K \phi_{jk} \log \pi_k.$$

This has the following analytic solution:

$$\pi_k^* = \frac{1}{p} \sum_{j=1}^p \phi_{jk}, \quad k = 1, \dots, K.$$

This proves part (ii) of Proposition 7.

This update can be thought of as an approximate M-step update for the mixture weights in which the posterior probabilities (the “responsibilities”) are computed approximately using q .

E.3 Update for σ^2

The coordinate ascent update for the residual variance σ^2 is the solution to

$$(\sigma^2)^* = \operatorname{argmax}_{\sigma^2 \in \mathbb{R}_+} F(q, g, \sigma^2).$$

From earlier results, the ELBO for any q parameterized as (54) works out to

$$F(q, g, \sigma^2) = \mathbb{E}_q[\log p(\mathbf{y} \mid \mathbf{X}, \mathbf{b}, \sigma^2)] - \sum_{j=1}^p D_{\text{KL}}(q_j \parallel p_{\text{prior}}) \quad (60)$$

in which

$$\begin{aligned} \mathbb{E}_q[\log p(\mathbf{y} \mid \mathbf{X}, \mathbf{b}, \sigma^2)] &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\bar{\mathbf{b}}\|^2 - \frac{1}{2\sigma^2} \sum_{j=1}^p d_j \left[\sum_{k=1}^K \phi_{1jk} (\mu_{1jk}^2 + s_{1jk}^2) - \bar{b}_j^2 \right] \\ D_{\text{KL}}(q_j \parallel p_{\text{prior}}) &= \sum_{k=1}^K \phi_{1jk} \log \frac{\phi_{1jk}}{\pi_k} - \frac{1}{2} \sum_{k=2}^K \phi_{1jk} \left[1 + \log \frac{s_{1jk}^2}{\sigma^2 \sigma_k^2} - \frac{s_{1jk}^2 + \mu_{1jk}^2}{\sigma^2 \sigma_k^2} \right], \end{aligned}$$

and where $\bar{b}_j = \sum_{k=1}^K \phi_{1jk} \mu_{1jk}$. Taking the partial derivative of F with respect to σ^2 then solving for σ^2 yields the following update:

$$(\sigma^2)^* = \frac{\|\bar{\mathbf{r}}\|^2 + \sum_{j=1}^p \sum_{k=2}^K \phi_{1jk} (d_j + 1/\sigma_k^2) (\mu_{1jk}^2 + s_{1jk}^2) - \sum_{j=1}^p d_j \bar{b}_j^2}{n + p(1 - \pi_1)}.$$

When σ^2 is updated following updates to q , we can simplify this expression by noting the specific form of the posterior means and variances,

$$\begin{aligned} s_{1jk}^2 &= \frac{\sigma^2}{d_j + 1/\sigma_k^2} \\ \mu_{1jk} &= \frac{d_j}{d_j + 1/\sigma_k^2} \times \tilde{b}_j, \end{aligned}$$

for $k = 2, \dots, K$, which gives

$$(\sigma^2)^* = \frac{\|\bar{\mathbf{r}}\|^2 + \sum_{j=1}^p d_j \bar{b}_j (\tilde{b}_j - \bar{b}_j) + \sigma^2 p(1 - \pi_1)}{n + p(1 - \pi_1)}.$$

This proves part (iii) of Proposition 7.

E.4 More detailed VEB algorithm

Further details about the implementation of the VEB algorithm are given in Algorithm 4. Algorithm 4 does not require $\mathbf{x}_j^T \mathbf{x}_j = 1$, $j = 1, \dots, p$.

Appendix Appendix F. VEB as a PLR

F.1 Normal Means as a Penalized Estimation Problem

In this section, we formulate the normal means problem with one observation, $\text{NM}_1(f, s^2)$, as penalized estimation problem. Specifically, we express the posterior mean of b under the NM model as a solution to a penalized least squares problem. The results for the single-observation NM model are used later to derive results for the multiple linear regression model.

Recall, $F_1^{\text{NM}}(q, f, s^2; y)$ denotes the ELBO for the single-observation NM model (50). From this, we define

$$h_1^{\text{NM}}(\bar{b}, f, s^2; y) \triangleq - \max_{q: \mathbb{E}_q(b) = \bar{b}} F_1^{\text{NM}}(q, f, s^2; y). \quad (61)$$

As a reminder, $F_1^{\text{NM}}(q, f, s^2; y)$ attains its maximum over q the exact posterior, $q = p^{\text{NM}}$; analogously, $h_1^{\text{NM}}(\bar{b}, f, s^2; y)$ attains its minimum over \bar{b} at $\bar{b} = S_{f,s}(y)$, the posterior mean of b (see Definition 3). Further, at their respective optima these two functions recover the

Algorithm 4 Coordinate ascent for fitting VEB model (more detailed).

Require: Data $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{y} \in \mathbb{R}^n$; number of mixture components, K ;
 prior variances, $\sigma_1^2 < \dots < \sigma_K^2$, with $\sigma_1^2 = 0$; initial estimates $\bar{\mathbf{b}}, \boldsymbol{\pi}, \sigma^2$.
 $\bar{\mathbf{r}} \leftarrow \mathbf{y} - \mathbf{X}\bar{\mathbf{b}}$ (compute mean residuals)
 $t \leftarrow 0$
for $j \leftarrow 1$ to p **do**
 $d_j = \mathbf{x}_j^T \mathbf{x}_j$
end for
repeat
 for $j \leftarrow 1$ to p **do**
 $\bar{\mathbf{r}}_j = \bar{\mathbf{r}} + \mathbf{x}_j \bar{b}_j$ (disregard j th effect in residuals)
 $\tilde{b}_j \leftarrow \mathbf{x}_j^T \bar{\mathbf{r}}_j / d_j$ (compute OLS estimate)
 for $k \leftarrow 1$ to K **do**
 $\mu_{jk} \leftarrow \frac{d_j}{d_j + 1/\sigma_k^2} \times \tilde{b}_j$ (update q_j)
 $\phi_{jk} \leftarrow \pi_k \times \frac{1}{1 + d_j \sigma_k^2} \times \exp\left(\frac{d_j \tilde{b}_j \mu_{jk}}{2\sigma^2}\right)$
 end for
 for $k \leftarrow 1$ to K **do**
 $\phi_{jk} \leftarrow \phi_{jk} / \sum_{k'=1}^K \phi_{jk'}$
 end for
 $\bar{b}_j \leftarrow \sum_{k=1}^K \phi_{jk} \mu_{jk}$ (update posterior mean of b_j)
 $\bar{\mathbf{r}} \leftarrow \bar{\mathbf{r}}_j - \mathbf{x}_j \bar{b}_j$ (update mean residuals)
 end for
 for $k \leftarrow 1$ to K **do**
 $\pi_k \leftarrow \sum_{j=1}^p \phi_{jk} / p$. (update g ; eq. 26)
 end for
 $(\sigma^2)^* \leftarrow \frac{\|\bar{\mathbf{r}}\|^2 + \bar{\mathbf{b}}^T \mathbf{D}(\tilde{\mathbf{b}} - \bar{\mathbf{b}}) + \sigma^2 p(1 - \pi_1^*)}{n + p(1 - \pi_1)}$ (update σ^2 ; eq. 56)
 $t \leftarrow t + 1$.
until convergence criterion is met
return $\bar{\mathbf{b}}, \boldsymbol{\pi}, \sigma^2$

marginal likelihood:

$$\begin{aligned}
 \log p(y \mid f, s^2) &= F_1^{\text{NM}}(p^{\text{NM}}, f, s^2; y) \\
 &= \max_q F_1^{\text{NM}}(q, f, s^2; y) \\
 &= \max_{\bar{b} \in \mathbb{R}} \max_{q: \mathbb{E}_q(b) = \bar{b}} F_1^{\text{NM}}(q, f, s^2; y) \\
 &= \max_{\bar{b} \in \mathbb{R}} -h_1^{\text{NM}}(\bar{b}, f, s^2; y) \\
 &= -h_1^{\text{NM}}(S_{f,s}(y), f, s^2; y).
 \end{aligned}$$

With these definitions, we can express the posterior mean for b as the solution to a real-valued optimization problem:

$$S_{f,s}(y) = \operatorname{argmin}_{\bar{b}} h_1^{\text{NM}}(\bar{b}, f, s^2; y).$$

The following lemma states that this can be understood as optimizing a penalized loss function and gives an explicit form for the penalty.

Lemma 8 $h_1^{\text{NM}}(\bar{b}, f, s^2; y)$ can be written as a penalized loss function,

$$h_1^{\text{NM}}(\bar{b}, f, s^2; y) = \frac{1}{2s^2}(y - \bar{b})^2 + \frac{1}{s^2}\rho_{f,s}(\bar{b}), \quad (62)$$

in which the penalty is

$$\rho_{f,s}(\bar{b}) \triangleq \min_{q: \mathbb{E}_q(b) = \bar{b}} \frac{s^2}{2} \log(2\pi s^2) + \frac{1}{2} \text{Var}_q(b) + s^2 D_{\text{KL}}(q \| p_{\text{prior}}(f)). \quad (63)$$

For any $y \in \mathbb{R}$, this penalty term satisfies

$$\rho_{f,s}(S_{f,s}(y)) = -s^2 \ell_{\text{NM}}(y; f, s^2) - \frac{1}{2}(y - S_{f,s}(y))^2, \quad (64)$$

and

$$\rho'_{f,s}(S_{f,s}(y)) = (y - S_{f,s}(y)) \quad (65)$$

$$= -s^2 \ell'_{\text{NM}}(y; f, s^2), \quad (66)$$

in which $\ell_{\text{NM}}(y; f, s^2)$ is the marginal log-likelihood $\ell_{\text{NM}}(y; f, s^2) \triangleq \log p(y | f, s^2)$ for the single-observation normal means model, $\text{NM}_1(f, s^2)$.

Proof From (50), we have

$$F_1^{\text{NM}}(q, f, s^2; y) = -\frac{1}{2s^2}(y - \mathbb{E}_q(b))^2 - \left[\frac{1}{2} \log(2\pi s^2) + \frac{1}{2s^2} \text{Var}_q(b) + D_{\text{KL}}(q \| p_{\text{prior}}(f)) \right].$$

Expressions (62) and (63) follow from (61).

Expression (64) is obtained by substituting $\bar{b} = S_{f,s}(y)$ into (62) and rearranging, noting that $h(\bar{b}, f, s^2; y)$ attains its minimum at this \bar{b} and therefore recovers the marginal log-likelihood,

$$h^{\text{NM}}(S_{f,s}(y), f, s^2; y) = -\ell_{\text{NM}}(y; f, s^2).$$

Expression (65) is a consequence of the fact that $h^{\text{NM}}(\bar{b}, f, s^2; y)$ attains its minimum at $\bar{b} = S_{f,s}(y)$ for any $y \in \mathbb{R}$, that is,

$$S_{f,s}(y) = \underset{\bar{b} \in \mathbb{R}}{\text{argmin}} h^{\text{NM}}(\bar{b}, f, s^2; y).$$

Therefore, the derivative of (62) with respect to \bar{b} at $\bar{b} = S_{f,s}(y)$ must be zero. Finally, (66) is obtained by applying Tweedie's formula (Efron, 2011). \blacksquare

F.2 VEB as a Penalized Regression Problem

Here we consider the ELBO for the multiple linear regression model (18). We begin with the following lemma.

Lemma 9 If the distribution $q(\mathbf{b})$ factorizes as $q(\mathbf{b}) = \prod_{j=1}^p q_j(b_j)$, then

$$\mathbb{E}_q[\|\mathbf{r}\|^2] = \|\bar{\mathbf{r}}\|^2 + \sum_{j=1}^p d_j \text{Var}_{q_j}(b_j),$$

where $\bar{\mathbf{b}} \triangleq \mathbb{E}_q[\mathbf{b}]$, $\mathbf{r} \triangleq \mathbf{y} - \mathbf{X}\mathbf{b}$ and $\bar{\mathbf{r}} \triangleq \mathbb{E}_q[\mathbf{r}] = \mathbf{y} - \mathbf{X}\bar{\mathbf{b}}$.

Proof

$$\begin{aligned} \mathbb{E}_q[\|\mathbf{r}\|^2] &= \mathbb{E}_q[\|\bar{\mathbf{r}} + \mathbf{X}(\bar{\mathbf{b}} - \mathbf{b})\|^2] \\ &= \|\bar{\mathbf{r}}\|^2 + \mathbb{E}_q[\|\mathbf{X}(\bar{\mathbf{b}} - \mathbf{b})\|^2] \\ &= \|\bar{\mathbf{r}}\|^2 + \mathbb{E}_q[(\bar{\mathbf{b}} - \mathbf{b})^T \mathbf{X}^T \mathbf{X} (\bar{\mathbf{b}} - \mathbf{b})] \\ &= \|\bar{\mathbf{r}}\|^2 + \text{tr}(\mathbf{X}^T \mathbf{X} \text{Cov}_q(\mathbf{b})) \\ &= \|\bar{\mathbf{r}}\|^2 + \sum_{j=1}^p d_j \text{Var}_{q_j}(b_j). \end{aligned}$$

■

In the following proposition, we express the ELBO for the multiple linear regression model as a penalized loss function.

Proposition 10 The objective function h (35) can be written as a penalized loss function,

$$h(\bar{\mathbf{b}}, g, \sigma^2) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\bar{\mathbf{b}}\|^2 + \sum_{j=1}^p \rho_{g\sigma, s_j}(\bar{b}_j)/s_j^2 + \frac{1}{2} \sum_{j=1}^p \log(d_j) + \frac{n-p}{2} \log(2\pi\sigma^2), \quad (67)$$

using the penalty function $\rho_{f,s}$ defined in (63), and defining $s_j^2 \triangleq \sigma^2/d_j$, $j = 1, \dots, p$. Note that when $d_j = 1$, $k = 1, \dots, p$, (67) simplifies to (36).

Proof From Lemma 9, we have

$$\begin{aligned} h(\bar{\mathbf{b}}, g, \sigma^2) &= \frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \mathbb{E}_q[\|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2] + \sum_{j=1}^p D_{\text{KL}}(q_j \| p_{\text{prior}}(g\sigma)) \\ &= \frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\bar{\mathbf{r}}\|^2 + \frac{1}{2} \sum_{j=1}^p \text{Var}_{q_j}(b_j)/s_j^2 + \sum_{j=1}^p D_{\text{KL}}(q_j \| p_{\text{prior}}(g\sigma)). \end{aligned}$$

Therefore,

$$\begin{aligned}
 h(\bar{\mathbf{b}}, g, \sigma^2) &= \max_{q: \mathbb{E}_q(\mathbf{b}) = \bar{\mathbf{b}}} F(q, g, \sigma^2) \\
 &= \frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\bar{\mathbf{r}}\|^2 \\
 &\quad + \sum_{j=1}^p \frac{1}{s_j^2} \times \left\{ \min_{q_j: \mathbb{E}_{q_j}(b_j) = \bar{b}_j} \frac{1}{2} \text{Var}_{q_j}(b_j) + s_j^2 D_{\text{KL}}(q_j \parallel p_{\text{prior}}(g_\sigma)) \right\} \\
 &= \frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\bar{\mathbf{r}}\|^2 + \sum_{j=1}^p \frac{1}{s_j^2} \left[\rho_{g_\sigma, s_j}(\bar{b}_j) - \frac{s_j^2}{2} \log(2\pi s_j^2) \right] \\
 &= \frac{1}{2\sigma^2} \|\bar{\mathbf{r}}\|^2 + \sum_{j=1}^p \rho_{g_\sigma, s_j}(\bar{b}_j) / s_j^2 + \frac{1}{2} \sum_{j=1}^p \log(d_j) + \frac{n-p}{2} \log(2\pi\sigma^2).
 \end{aligned}$$

■

Appendix Appendix G. Additional Results and Proofs

Proposition 11 (Convergence of cyclic coordinate ascent for VEB) The sequence of iterates

$$\{q^{(t)}, g^{(t)}, (\sigma^2)^{(t)}\}, \quad t = 0, 1, 2, \dots,$$

generated by Algorithm 2 converge monotonically to a stationary point of the ELBO, F (18).

Proof By Proposition 2.7.1 of Bertsekas (1999), the sequence of iterates $\{q^{(t)}, g^{(t)}, (\sigma^2)^{(t)}\}$, $t = 0, 1, 2, \dots$, generated by Algorithm 2 converges monotonically to a stationary point of F provided that F is continuously differentiable and each coordinate update,

$$\begin{aligned}
 q_j^{(t+1)} &= \operatorname{argmax}_{q_j} F(q_1^{(t+1)}, \dots, q_{j-1}^{(t+1)}, q_j, q_{j+1}^{(t)}, \dots, q_p^{(t)}, g^{(t)}, (\sigma^2)^{(t)}), \quad j = 1, \dots, p \\
 g^{(t+1)} &= \operatorname{argmax}_{g \in \mathcal{G}} F(q^{(t+1)}, g, (\sigma^2)^{(t)}) \\
 (\sigma^2)^{(t+1)} &= \operatorname{argmax}_{\sigma^2 \in \mathbb{R}_+} F(q^{(t+1)}, g^{(t+1)}, \sigma^2),
 \end{aligned}$$

is finite and uniquely determined. (See also Luo and Tseng 1992; Tseng 2001 for another treatment of convergence of coordinate ascent under general conditions.) A sufficient condition for F to be continuously differentiable and for the coordinate ascent updates (Proposition 1 or Proposition 7) to have a unique solution is that $0 < \sigma^2 < \infty$, $\pi_k > 0$ for all $k = 1, \dots, K$, and $0 \leq \sigma_1^2 < \dots < \sigma_K^2 < \infty$. ■

G.1 Proof of Proposition 2

The ELBO

$$F(q, g, \sigma^2) = \iint q(\mathbf{b}, \gamma) \log \left\{ \frac{p(\mathbf{y} \mid \mathbf{X}, \mathbf{b}, \sigma^2) p(\mathbf{b}, \gamma \mid g, \sigma^2)}{q(\mathbf{b}, \gamma)} \right\} d\mathbf{b} d\gamma$$

is maximized with respect to q when $q(\mathbf{b}, \gamma) \propto p(\mathbf{y} \mid \mathbf{X}, \mathbf{b}, \sigma^2) p(\mathbf{b}, \gamma \mid g, \sigma^2)$. This follows from the equality condition of Jensen's inequality (Jordan et al., 1999). When the columns of \mathbf{X} are orthogonal, the posterior factorizes over the individual coordinates j ,

$$\begin{aligned} p(\mathbf{b}, \gamma \mid \mathbf{X}, \mathbf{y}, g, \sigma^2) &\propto p(\mathbf{y} \mid \mathbf{X}, \mathbf{b}, \sigma^2) p(\mathbf{b}, \gamma \mid g, \sigma^2) \\ &\propto \prod_{j=1}^p \exp \left\{ -\frac{(b_j - \mathbf{x}_j^T \mathbf{y})^2}{2\sigma^2} \right\} \times p(b_j, \gamma_j \mid g, \sigma^2). \end{aligned}$$

Therefore, when \mathbf{X} has orthogonal columns, the best q , even with the restriction of being fully factorized (20), is able to recover the exact posterior since the exact posterior also factorizes over the coordinates $j = 1, \dots, p$.

G.2 Proof of Proposition 4

First, we note that

$$F(\hat{q}, \hat{g}, \hat{\sigma}^2) = \max_{q \in \mathcal{Q}} F(q, \hat{g}, \hat{\sigma}^2) = -h(\hat{\mathbf{b}}, \hat{g}, \hat{\sigma}^2).$$

Hence, for any $\bar{\mathbf{b}} \in \mathbb{R}^p$, we have

$$\begin{aligned} -h(\bar{\mathbf{b}}, g, \sigma^2) &= \max_{q \in \mathcal{Q}, \mathbb{E}_q[\mathbf{b}] = \bar{\mathbf{b}}} F(q, g, \sigma^2) \\ &\leq \max_{q \in \mathcal{Q}} F(q, g, \sigma^2) \\ &\leq \max_{q \in \mathcal{Q}, g \in \mathcal{G}, \sigma^2 \in \mathcal{T}} F(q, g, \sigma^2) \\ &= F(\hat{q}, \hat{g}, \hat{\sigma}^2) \\ &= -h(\hat{\mathbf{b}}, \hat{g}, \hat{\sigma}^2). \end{aligned}$$

This proves that

$$\hat{\mathbf{b}}, \hat{g}, \hat{\sigma}^2 = \operatorname{argmin}_{\bar{\mathbf{b}} \in \mathbb{R}^p, g \in \mathcal{G}, \sigma^2 \in \mathcal{T}} h(\bar{\mathbf{b}}, g, \sigma^2).$$

G.3 Proof of Theorem 5

The proof of the first part of Theorem 5 follows immediately from Proposition 10 and Lemma 8 by letting $d_j = 1$ for $j = 1, \dots, p$.

The missing piece of the proof is to show that $S_{\rho_{f,\sigma}}(y) = S_{f,\sigma}(y)$. We start with the definition of S_ρ in (31),

$$S_{\rho_{f,\sigma}}(y) = \operatorname{argmin}_{b \in \mathbb{R}} \frac{1}{2}(y - b)^2 + \rho_{f,\sigma}(b).$$

To solve for the argmin on the right-hand side, we differentiate with respect to b and set the derivative to zero,

$$\rho'_{f,\sigma}(b) = y - b.$$

From (38), this derivative vanishes when $b = S_{f,\sigma}(y)$. Therefore, $S_{\rho_{f,\sigma}}(y) = S_{f,\sigma}(y)$.

G.4 Mathematical Properties of the Posterior Mean Shrinkage Operator

Lemma 12 Let f be a symmetric unimodal distribution on \mathbb{R} with a mode at zero, and assume $\sigma^2 > 0$. Then the NM posterior mean operator $S_{f,\sigma}(y)$ defined in (32) is symmetric, non-negative, and non-decreasing on $y \in \mathbb{R}$, and $S_{f,\sigma}(y) \leq y$ on $y \in (0, \infty)$. That is, $S_{f,\sigma}$ is a “shrinkage operator” that shrinks towards zero.

Proof The marginal likelihood for the NM model (33) is

$$\begin{aligned} \ell_{\text{NM}}(y; f, \sigma^2) &\triangleq \log p(y | f, \sigma^2) \\ &= \int p(y | b, \sigma^2) p(b) db \\ &= \int N(y; b, \sigma^2) f(b) db. \end{aligned}$$

By Khintchine’s representation theorem (Dharmadhikari and Joag-Dev, 1988), f can be represented as mixture of uniform distributions,

$$f(b) = \int_0^\infty \frac{\mathbb{I}\{|b| < t\}}{2t} p(t) dt$$

for some (possibly improper) univariate mixing density $p(t)$. Let $p(b | t)$ be the density function of the uniform distribution on $[-t, t]$. Then we have

$$\begin{aligned} p(y | f, \sigma^2) &= \int p(y | b, \sigma^2) p(b | f, \sigma^2) db \\ &= \int_0^\infty \left[\int p(y | b, \sigma^2) p(b | t) db \right] \times p(t) dt \\ &= \int_0^\infty p(y | \sigma^2, t) p(t) dt, \end{aligned} \tag{68}$$

where

$$\begin{aligned} p(y | \sigma^2, t) &= \int p(y | b, \sigma^2) p(b | t) db \\ &= \frac{1}{2t} \left[\Phi\left(\frac{t-y}{\sigma}\right) + \Phi\left(\frac{t+y}{\sigma}\right) - 1 \right], \end{aligned} \tag{69}$$

and where $\Phi(x)$ denotes the normal cumulative distribution function. Note that, from (33), $p(y | b, \sigma^2) = N(y; b, \sigma^2)$. Since $\Phi(t+y) + \Phi(t-y)$ is non-increasing in $y \in (0, \infty)$ for any $t \geq 0$, $p(y | \sigma^2, t)$ is also non-increasing in $y \in (0, \infty)$ for any $t \geq 0$. This implies that $p(y | \sigma^2, t)$ is unimodal with a mode at zero, and therefore $p(y | f, \sigma^2)$ must also be unimodal with a mode at zero since it is a mixture of unimodal distributions that all have modes at zero.

From (66), which was obtained by applying Tweedie’s formula, we have that

$$S_{f,\sigma}(y) = y + \sigma^2 \ell'_{\text{NM}}(y; f, \sigma) \leq y,$$

in which the inequality is obtained by noting that $\ell_{\text{NM}}(y; f, \sigma)$ is non-increasing in $y \in (0, \infty)$. And since $\ell_{\text{NM}}(y; f, \sigma^2)$ is symmetric about zero, the shrinkage operator must be an odd function; *i.e.*, $S_{f,\sigma}(y) = -S_{f,\sigma}(-y)$.

It remains to show that $S_{f,\sigma}(y)$ is non-decreasing on $y \in \mathbb{R}_+$. Since $p(y | b, \sigma^2) = N(y; b, \sigma^2)$ and $p(b | t)$ is the uniform distribution on interval $[-t, t]$, the posterior density is truncated normal:

$$\begin{aligned} p(b | y, \sigma^2, t) &\propto p(y | b, \sigma^2) p(b | t) \\ &= N_{[-t,t]}(b; y, s^2) \end{aligned}$$

where $N_{[-t,t]}(x; \mu, s^2) = N(x; \mu, s^2) \mathbb{I}\{|x| < t\}$ denotes the probability density of the normal distribution with mean μ and variance s^2 truncated to the interval $x \in [-t, t]$. The expected value of the truncated normal,

$$\mathbb{E}[X] = \mu + s \times \frac{N(-t; \mu, s^2) - N(t; \mu, s^2)}{\Phi((t - \mu)/s) - \Phi(-(t + \mu)/s)},$$

is non-decreasing with respect to μ , for all $\mu \in \mathbb{R}$, $t > 0$. To show this, the derivative of the expected value with respect to μ is always positive: $\frac{\partial}{\partial \mu} \mathbb{E}[X] = \text{Var}(X)/s^2 > 0$. Therefore, $S_{f,\sigma}$ is a mixture of non-decreasing functions on \mathbb{R}_+ . ■

References

- T. Amir, R. Basri, and B. Nadler. The Trimmed Lasso: sparse recovery guarantees and practical optimization by the generalized soft-min penalty. *SIAM Journal on Mathematics of Data Science*, 3(3):900–929, 2021.
- N. S. P. Anirban Bhattacharya, Debdeep Pati and D. B. Dunson. Dirichlet–Laplace priors for optimal shrinkage. *Journal of the American Statistical Association*, 110(512):1479–1490, 2015.
- R. Bai, V. Ročková, and E. I. George. Spike-and-slab meets LASSO: A review of the spike-and-slab LASSO. In M. G. Tadesse and M. Vannucci, editors, *Handbook of Bayesian Variable Selection*, pages 81–108. Chapman and Hall/CRC, Boca Raton, FL, 2021.
- A. Beck and L. Tetruashvili. On the convergence of block coordinate descent type methods. *SIAM Journal on Optimization*, 23(4):2037–2060, 2013.
- J. O. Berger. *Statistical Decision Theory and Bayesian Analysis*. Springer, New York, NY, 1985.
- D. P. Bertsekas. *Nonlinear Programming*. Athena Scientific, Belmont, MA, second edition, 1999.
- D. Bertsimas, M. S. Copenhaver, and R. Mazumder. The Trimmed Lasso: sparsity and robustness. *arXiv*, 1708.04527, 2017.
- A. Bhadra, J. Datta, N. G. Polson, and B. Willard. Lasso meets horseshoe: a survey. *Statistical Science*, 34(3):405–427, 2019.
- C. Bishop. *Pattern Recognition and Machine Learning*. Springer, New York, NY, 2006.

- D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
- D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: a review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- P. Breheny and J. Huang. Coordinate descent algorithms for nonconvex penalized regression, with applications to biological feature selection. *Annals of Applied Statistics*, 5(1):232–253, 2011.
- P. Carbonetto and M. Stephens. Scalable variational inference for Bayesian variable selection in regression, and its accuracy in genetic association studies. *Bayesian Analysis*, 7(1):73–108, 2012.
- P. Carbonetto, X. Zhou, and M. Stephens. varbvs: fast variable selection for large-scale regression. *arXiv*, 1709.06597, 2017.
- B. P. Carlin and T. A. Louis. Empirical Bayes: past, present and future. *Journal of the American Statistical Association*, 95(452):1286–1289, 2000.
- C. M. Carvalho, N. G. Polson, and J. G. Scott. The horseshoe estimator for sparse signals. *Biometrika*, 97(2):465–480, 2010.
- G. Casella. Empirical Bayes Gibbs sampling. *Biostatistics*, 2(4):485–500, 2001.
- I. Castillo and A. Van Der Vaart. Needles and straw in a haystack: posterior concentration for possibly sparse sequences. *Annals of Statistics*, 40(4):2069–2101, 2012.
- H. Chipman, E. I. George, and R. E. McCulloch. The practical implementation of Bayesian model selection. In *Model Selection*, volume 38 of *IMS Lecture Notes*, pages 65–116. 2001.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B*, 39(1):1–22, 1977.
- S. Dharmadhikari and K. Joag-Dev. *Unimodality, convexity, and applications*. Academic Press, Boston, MA, 1988.
- J. Drugowitsch. Variational Bayesian inference for linear and logistic regression. *arXiv*, 1310.5438, 2013.
- B. Efron. Microarrays, empirical Bayes and the two-groups model. *Statistical Science*, 32(1):1–22, 2008.
- B. Efron. Tweedie’s formula and selection bias. *Journal of the American Statistical Association*, 106(496):1602–1614, 2011.
- B. Efron. Bayes, oracle Bayes and empirical Bayes. *Statistical Science*, 34(2):177–201, 2019.
- B. Efron and C. Morris. Stein’s estimation rule and its competitors—an empirical Bayes approach. *Journal of the American Statistical Association*, 68(341):117–130, 1973.

- J. Fan and R. Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96(456):1348–1360, 2011.
- M. A. T. Figueiredo. Adaptive sparseness for supervised learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(9):1150–1159, 2003.
- J. Friedman, T. Hastie, H. Höfling, and R. Tibshirani. Pathwise coordinate optimization. *Annals of Applied Statistics*, 1(2):302–332, 2007.
- J. Friedman, T. Hastie, and R. Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33(1):1–22, 2010.
- W. J. Fu. Penalized regressions: The Bridge versus the Lasso. *Journal of Computational and Graphical Statistics*, 7(3):397–416, 1998.
- E. R. Gamazon, H. E. Wheeler, K. P. Shah, S. V. Mozaafari, K. Aquino-Michaels, R. J. Carroll, A. E. Eyler, J. C. Denny, D. L. Nicolae, N. J. Cox, and H. K. Im. A gene-based association method for mapping traits using reference transcriptome data. *Nature Genetics*, 47(9):1091–1098, 2015.
- A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin. *Bayesian Data Analysis*. CRC Press, Boca Raton, FL, third edition, 2013.
- E. I. George and D. P. Foster. Calibration and empirical Bayes variable selection. *Biometrika*, 87(4):731–747, 2000.
- E. I. George and R. E. McCulloch. Variable selection via Gibbs sampling. *Journal of the American Statistical Association*, 88(423):881–889, 1993.
- E. I. George and R. E. McCulloch. Approaches for Bayesian variable selection. *Statistica Sinica*, 9(2):339–373, 1997.
- D. Gerard and M. Stephens. Empirical Bayes shrinkage and false discovery rate estimation, allowing for unwanted variation. *Biostatistics*, 21(1):15–32, 2020.
- Z. Ghahramani and G. E. Hinton. Variational learning for switching state-space models. *Neural Computation*, 12(4):831–864, 2000.
- M. Girolami. A variational method for learning sparse and overcomplete representations. *Neural Computation*, 13(11):2517–2532, 2001.
- J. E. Griffin and P. J. Brown. Inference with normal-gamma prior distributions in regression problems. *Bayesian Analysis*, 5(1):171–188, 2010.
- GTEX Consortium. Genetic effects on gene expression across human tissues. *Nature*, 550(7675):204–213, 2017.
- Y. Guan and M. Stephens. Bayesian variable selection regression for genome-wide association studies and other large-scale problems. *Annals of Applied Statistics*, 5(3):1780–1815, 2011.

- D. Habier, R. L. Fernando, K. Kizilkaya, and D. J. Garrick. Extension of the Bayesian alphabet for genomic selection. *BMC Bioinformatics*, 12:186, 2011.
- C. Hans. Bayesian lasso regression. *Biometrika*, 96(4):835–845, 2009.
- H. O. Hartley and J. N. K. Rao. Maximum-likelihood estimation for the mixed analysis of variance model. *Biometrika*, 54(1–2):93–108, 1967.
- T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer, New York, NY, second edition, 2009.
- H. Hazimeh and R. Mazumder. Fast best subset selection: coordinate descent and local combinatorial optimization algorithms. *Operations Research*, 68(5):1517–1537, 2020.
- A. E. Hoerl and R. W. Kennard. Ridge regression: biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- T. S. Jaakkola and M. I. Jordan. Bayesian parameter estimation via variational methods. *Statistics and Computing*, 10:25–37, 2000.
- A. Javanmard and A. Montanari. Debiasing the lasso: optimal sample size for Gaussian designs. *Annals of Statistics*, 46(6A):2593–2622, 2018.
- I. M. Johnstone and B. W. Silverman. Needles and straw in haystacks: empirical Bayes estimates of possibly sparse sequences. *Annals of Statistics*, 32(4):1594–1649, 2004.
- I. M. Johnstone and B. W. Silverman. Empirical Bayes selection of wavelet thresholds. *Annals of Statistics*, pages 1700–1752, 2005.
- L. Joo. *Bayesian lasso: an extension for genome-wide association study*. PhD thesis, New York University, 2017.
- M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. An introduction to variational methods for graphical models. *Machine Learning*, 37(2):183–233, 1999.
- Y. Kim, P. Carbonetto, M. Stephens, and M. Anitescu. A fast algorithm for maximum likelihood estimation of mixture proportions using sequential quadratic programming. *Journal of Computational and Graphical Statistics, forthcoming*, 29(2):261–273, 2020.
- R. Koenker and I. Mizera. Convex optimization, shape constraints, compound decisions, and empirical Bayes rules. *Journal of the American Statistical Association*, 109(506):674–685, 2014.
- S. Kullback and R. A. Leibler. On information and sufficiency. *Annals of Mathematical Statistics*, 22(1):79–86, 1951.
- Q. Li and N. Lin. The Bayesian elastic net. *Bayesian Analysis*, 5(1):151–170, 2010.
- F. Liang, R. Paulo, G. Molina, M. a. Clyde, and J. O. Berger. Mixtures of g priors for Bayesian variable selection. *Journal of the American Statistical Association*, 103(481):410–423, 2008.

- B. A. Logsdon, G. E. Hoffman, and J. G. Mezey. A variational Bayes algorithm for fast and accurate multiple locus genome-wide association analysis. *BMC Bioinformatics*, 11: 58, 2010.
- M. Lu and M. Stephens. Variance adaptive shrinkage (vash): flexible empirical Bayes estimation of variances. *Bioinformatics*, 32(22):3428–3434, 2016.
- Z. Q. Luo and P. Tseng. On the convergence of the coordinate descent method for convex differentiable minimization. *Journal of Optimization Theory and Applications*, 72:7–35, 1992.
- B. M. Marlin, M. E. Khan, and K. P. Murphy. Piecewise bounds for estimating Bernoulli-logistic latent Gaussian models. In *Proceedings of the 28th International Conference on Machine Learning*, pages 633–640, 2011.
- R. Mazumder, J. H. Friedman, and T. Hastie. Sparsenet: coordinate descent with nonconvex penalties. *Journal of the American Statistical Association*, 106(495):1125–1138, 2011.
- P. McCullagh and J. A. Nelder. *Generalized Linear Models*, volume 37 of *Monographs on Statistics and Applied Probability*. Chapman and Hall, New York, NY, second edition, 1989.
- T. H. Meuwissen, B. J. Hayes, and M. E. Goddard. Prediction of total genetic value using genome-wide dense marker maps. *Genetics*, 157(4):1819–1829, 2001.
- T. H. E. Meuwissen, T. R. Solberg, R. Shepherd, and J. A. Woolliams. A fast algorithm for BayesB type of prediction of genome-wide estimates of genetic value. *Genetics Selection Evolution*, 41:2, 2009.
- A. J. Miller. *Subset Selection in Regression*. Chapman and Hall/CRC, Boca Raton, FL, 2002.
- T. J. Mitchell and J. J. Beauchamp. Bayesian variable selection in linear regression. *Journal of the American Statistical Association*, 83(404):1023–1032, 1988.
- C. N. Morris. Parametric empirical Bayes inference: theory and applications. *Journal of the American Statistical Association*, 78(381):47–55, 1983.
- G. Moser, S. H. Lee, B. J. Hayes, M. E. Goddard, N. R. Wray, and P. M. Visscher. Simultaneous discovery, estimation and prediction analysis of complex traits using a Bayesian mixture model. *PLoS Genetics*, 11(4):e1004969, 2015.
- F. Nebebe and T. Stroud. Bayes and empirical Bayes shrinkage estimation of regression coefficients. *Canadian Journal of Statistics*, 14(4):267–280, 1986.
- N. Parikh and S. Boyd. Proximal algorithms. *Foundations and Trends in Optimization*, 1(3):127–239, 2014.
- T. Park and G. Casella. The Bayesian Lasso. *Journal of the American Statistical Association*, 103(482):681–686, 2008.

- P. Perez and G. de los Campos. Genome-wide regression and prediction with the BGLR statistical package. *Genetics*, 198(2):483–495, 2014.
- R Core Team. *R: a language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria, 2019. URL <https://www.r-project.org>.
- K. Ray and B. Szabó. Variational Bayes for high-dimensional linear regression with sparse priors. *Journal of the American Statistical Association*, 117(539):1270–1281, 2022.
- Q. Ren, S. Banerjee, A. O. Finley, and J. S. Hodges. Variational Bayesian methods for spatial data analysis. *Computational Statistics and Data Analysis*, 55(12):3197–3217, 2011.
- H. Robbins. The empirical Bayes approach to statistical decision problems. *Annals of Mathematical Statistics*, 35(1):1–20, 1964.
- V. Ročková and E. I. George. The Spike-and-Slab LASSO. *Journal of the American Statistical Association*, 113(521):431–444, 2018.
- L. K. Saul and M. I. Jordan. Exploiting tractable substructures in intractable networks. In D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo, editors, *Advances in Neural Information Processing Systems 8*, pages 486–492. MIT Press, 1996.
- G. K. Smyth. Linear models and empirical Bayes methods for assessing differential expression in microarray experiments. *Statistical Applications in Genetics and Molecular Biology*, 3(1), 2004.
- M. Stephens. False discovery rates: a new deal. *Biostatistics*, 18(2):275–294, 2016.
- S. M. Stigler. Studies in the history of probability and statistics XL: Boscovich, Simpson and a 1760 manuscript note on fitting a linear relation. *Biometrika*, 71(3):615–620, 1984.
- W. Su, M. Bogdan, and E. Candès. False discoveries occur early on the lasso path. *Annals of Statistics*, 45(5):2133–2150, 2017.
- L. Sun and M. Stephens. Solving the empirical Bayes normal means problem with correlated noise. *arXiv*, 1812.07488, 2018.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, pages 267–288, 1996.
- A. N. Tikhonov. Solution of incorrectly formulated problems and the regularization method. *Soviet Mathematics—Doklady*, 4:501–504, 1963.
- P. Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of Optimization Theory and Applications*, 109(3):475–494, 2001.
- S. M. Urbut, G. Wang, P. Carbonetto, and M. Stephens. Flexible statistical methods for estimating and testing effects in genomic studies with multiple conditions. *Nature Genetics*, 51(1):187–195, 2019.

- M. A. van de Wiel, D. E. Te Beest, and M. M. Münch. Learning from a lot: empirical Bayes for high-dimensional model-based prediction. *Scandinavian Journal of Statistics*, 46(1): 2–25, 2019.
- M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1(1–2):1–305, 2008.
- B. Wang and D. M. Titterton. Inadequacy of interval estimates corresponding to variational bayesian approximations. In *Proceedings of the 10th International Workshop on Artificial Intelligence and Statistics*, pages 373–380, 2005.
- C. Wang and D. M. Blei. Variational inference in nonconjugate models. *Journal of Machine Learning Research*, 14:1005–1031, 2013.
- G. Wang, A. Sarkar, P. Carbonetto, and M. Stephens. A simple new approach to variable selection in regression, with application to genetic fine mapping. *Journal of the Royal Statistical Society, Series B*, 82(5):1273–1300, 2020.
- W. Wang and M. Stephens. Empirical Bayes matrix factorization. *Journal of Machine Learning Research*, 22(120):1–40, 2021.
- S. J. Wright. Coordinate descent algorithms. *Mathematical Programming*, 151:3–34, 2015.
- T. T. Wu and K. Lange. Coordinate descent algorithms for lasso penalized regression. *Annals of Applied Statistics*, 2(1):224–244, 2008.
- Z. Xing, P. Carbonetto, and M. Stephens. Flexible signal denoising via flexible empirical Bayes shrinkage. *Journal of Machine Learning Research*, 22(93):1–28, 2021.
- C. You, J. T. Ormerod, and S. Müller. On variational Bayes estimation and variational information criteria for linear regression models. *Australian and New Zealand Journal of Statistics*, 56(1):73–87, 2014.
- M. Yuan and Y. Lin. Efficient empirical Bayes variable selection and estimation in linear models. *Journal of the American Statistical Association*, 100(472):1215–1225, 2005.
- J. Yun, P. Zheng, E. Yang, A. Lozano, and A. Aravkin. Trimming the ℓ_1 regularizer: statistical analysis, optimization, and applications to deep learning. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, pages 7242–7251, 2019.
- S. Zabad, S. Gravel, and Y. Li. Fast and accurate Bayesian polygenic risk modeling with variational inference. *American Journal of Human Genetics*, 110(5):741–761, 2023.
- J. Zeng, R. de Vlaming, Y. Wu, M. R. Robinson, L. R. Lloyd-Jones, L. Yengo, C. X. Yap, A. Xue, J. Sidorenko, A. F. McRae, J. E. Powell, G. W. Montgomery, A. Metspalu, T. Esko, G. Gibson, N. R. Wray, P. M. Visscher, and J. Yang. Signatures of negative selection in the genetic architecture of human complex traits. *Nature Genetics*, 50(5): 746–753, 2018.

- C.-H. Zhang. Nearly unbiased variable selection under minimax concave penalty. *Annals of Statistics*, 38(2):894–942, 2010.
- X. Zhou, P. Carbonetto, and M. Stephens. Polygenic modeling with Bayesian sparse linear mixed models. *PLoS Genetics*, 9(2), 2013.
- A. Zhu, J. G. Ibrahim, and M. I. Love. Heavy-tailed prior distributions for sequence count data: removing the noise and preserving large differences. *Bioinformatics*, 35(12):2084–2092, 2019.
- H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society, Series B*, 67(2):301–320, 2005.