

Sparse Linear Mixed Model Selection via Streamlined Variational Bayes

Luca Maestrini

The Australian National University
luca.maestrini@anu.edu.au

joint work with

Emanuele Degani, Dorota Toczyłowska and Matt Wand

2022 ISBA World Meeting
Montreal, 26 June – 1 July 2022

Outlook

Introduction

- ▶ linear mixed models in general form
- ▶ two- and three-level models
- ▶ variational approximations

Efficient sparse linear mixed model selection

- ▶ streamlined variational inference
- ▶ global-local shrinkage priors

Simulated and real data experiments

- ▶ simulated data involving three-level models
- ▶ Perinatal data (two-level model)

Motivation

- ▶ A variety of applications (e.g. longitudinal or grouped data) can be studied through **linear mixed-effects models** having fixed and random effects in the linear predictor
- ▶ In some cases, datasets include a large number of predictors, but only a few are effectively relevant
- ▶ These predictors are often introduced as fixed effects in the model to guarantee parsimony and validity of the inferential conclusions, especially in **sparse** model settings
- ▶ A proper **fixed effects selection** procedure is recommended to identify relevant predictors
- ▶ Other approaches are based on
 - random effect selection induced by their covariance matrix decomposition (Chen and Dunson, 2003; Yang, 2013),
 - joint fixed and random effects selection (Kinney and Dunson, 2007; Yang *et al.*, 2020)

General Model Formulation

We study linear mixed models having the general form

$$\begin{aligned} \mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma^2 &\sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma^2\mathbf{I}), \quad \mathbf{u}|\mathbf{G} \sim N(\mathbf{0}, \mathbf{G}), \quad \boldsymbol{\beta} \sim N(\boldsymbol{\mu}_{\boldsymbol{\beta}}, \boldsymbol{\Sigma}_{\boldsymbol{\beta}}), \\ \sigma^2|a_{\sigma^2} &\sim \text{Inverse-}\chi^2(\nu_{\sigma^2}, 1/a_{\sigma^2}), \quad a_{\sigma^2} \sim \text{Inverse-}\chi^2(1, 1/(\nu_{\sigma^2}s_{\sigma^2}^2)), \\ \mathbf{G}|\mathbf{A}_G &\sim p(\mathbf{G}|\mathbf{A}_G), \quad \mathbf{A}_G \sim p(\mathbf{A}_G), \end{aligned}$$

where the auxiliary variable a_{σ^2} and matrix \mathbf{A}_G facilitate the implementation of variational inference (Maestrini and Wand, 2021)

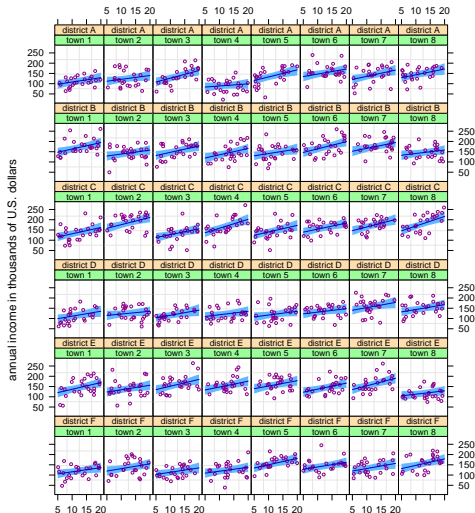
These are Gaussian response models where

- ▶ $\boldsymbol{\beta}$ is a vector of fixed effects
- ▶ \mathbf{u} is a vector of random effects
- ▶ \mathbf{X} is a fixed effects design matrix
- ▶ \mathbf{Z} is a random effects design matrix
- ▶ \mathbf{G} is a random effects covariance matrix

The focus of our work is on [two- and three-level linear mixed models](#)

Multilevel Data: Example

Three-level data with 6 districts, each having 8 towns where 25 residents are randomly selected (Nolan *et al.*, 2020)



Two- and Three-Level Models

A **Two-level** multilevel model for units clustered in m groups:

$$y_i | \beta, u_i, \sigma^2 \stackrel{\text{ind}}{\sim} N(\mathbf{X}_i \beta + \mathbf{Z}_i u_i, \sigma^2 \mathbf{I}), \quad u_i | \Sigma \stackrel{\text{ind}}{\sim} N(\mathbf{0}, \Sigma), \quad 1 \leq i \leq m, \quad + \text{priors}$$

A **three-level** linear mixed model can be defined in terms of observations from the i th group and its j th subgroup as follows:

$$y_{ij} | \beta, u_i^{L1}, u_{ij}^{L2}, \sigma^2 \stackrel{\text{ind}}{\sim} N(\mathbf{X}_{ij} \beta + \mathbf{Z}_{ij}^{L1} u_i^{L1} + \mathbf{Z}_{ij}^{L2} u_{ij}^{L2}, \sigma^2 \mathbf{I}),$$
$$\begin{bmatrix} u_i^{L1} \\ u_{ij}^{L2} \end{bmatrix} | \Sigma^{L1}, \Sigma^{L2} \stackrel{\text{ind}}{\sim} N \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma^{L1} & \mathbf{0} \\ \mathbf{0} & \Sigma^{L2} \end{bmatrix} \right), \quad + \text{priors}$$

In this model

- ▶ Σ^{L1} is the covariance matrix for the group-specific random effect vectors u_i^{L1} , $1 \leq i \leq m$
- ▶ Σ^{L2} is the covariance matrix for the subgroup-specific random effect vectors u_{ij}^{L2} , $1 \leq i \leq m$, $1 \leq j \leq n_i$

Three-level Random Effects Models

The three-level model is a special case of the general specification

$$\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma^2 \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma^2 \mathbf{I}), \quad \mathbf{u}|\mathbf{G} \sim N(\mathbf{0}, \mathbf{G}), \quad + \text{priors}$$

with

$$\mathbf{y} = \underset{1 \leq i \leq m}{\text{stack}} \left(\underset{1 \leq j \leq n_i}{\text{stack}}(\mathbf{y}_{ij}) \right), \quad \mathbf{X} = \underset{1 \leq i \leq m}{\text{stack}} \left(\underset{1 \leq j \leq n_i}{\text{stack}}(\mathbf{X}_{ij}) \right),$$

$$\mathbf{Z} = \underset{1 \leq i \leq m}{\text{blockdiag}} \left(\left[\underset{1 \leq j \leq n_i}{\text{stack}}(\mathbf{Z}_{ij}^{\text{L1}}) \mid \underset{1 \leq j \leq n_i}{\text{blockdiag}}(\mathbf{Z}_{ij}^{\text{L2}}) \right] \right),$$

$$\mathbf{u} = \underset{1 \leq i \leq m}{\text{stack}} \left(\left[(\mathbf{u}_i^{\text{L1}})^T \mid \left(\underset{1 \leq j \leq n_i}{\text{stack}}(\mathbf{u}_{ij}^{\text{L2}}) \right)^T \right]^T \right),$$

$$\mathbf{G} = \underset{1 \leq i \leq m}{\text{blockdiag}} \left(\left[\begin{array}{cc} \boldsymbol{\Sigma}^{\text{L1}} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n_i} \otimes \boldsymbol{\Sigma}^{\text{L2}} \end{array} \right] \right)$$

Three-level Random Effects Models (3)

If, for example, the data is divided into $m = 3$ groups each having $n_1 = 2$, $n_2 = 3$, $n_3 = 2$ subgroups, then

$$\mathbf{Z} = \begin{bmatrix} \mathbf{Z}_{11}^{\text{L1}} & \mathbf{Z}_{11}^{\text{L2}} & & & & & \\ & \mathbf{Z}_{12}^{\text{L1}} & & \mathbf{Z}_{12}^{\text{L2}} & & & \\ & & & & & & \\ & & & \mathbf{Z}_{21}^{\text{L1}} & \mathbf{Z}_{21}^{\text{L2}} & & \\ & & & \mathbf{Z}_{22}^{\text{L1}} & & \mathbf{Z}_{22}^{\text{L2}} & \\ & & & \mathbf{Z}_{23}^{\text{L1}} & & & \mathbf{Z}_{23}^{\text{L2}} \\ & & & & & & & \\ & & & & & & & \mathbf{Z}_{31}^{\text{L1}} & \mathbf{Z}_{31}^{\text{L2}} \\ & & & & & & & \mathbf{Z}_{32}^{\text{L1}} & & \mathbf{Z}_{32}^{\text{L2}} \end{bmatrix}$$

Variational Approximations

Variational approximations are a class of techniques for making approximate model fitting and inference

Idea: Approximate a complex density function p with a simpler q

How to choose q ?

Suppose p is a posterior density function $p(\boldsymbol{\theta}|\mathbf{y})$. Then q is chosen minimize the Kullback–Leibler divergence (or other divergence measures):

- ▶ $KL(q(\boldsymbol{\theta})||p(\boldsymbol{\theta}|\mathbf{y})) = \int q(\boldsymbol{\theta})\log \{q(\boldsymbol{\theta})/p(\boldsymbol{\theta}|\mathbf{y})\} d\boldsymbol{\theta}$
- ▶ $KL(p(\boldsymbol{\theta}|\mathbf{y})||q(\boldsymbol{\theta})) = \int p(\boldsymbol{\theta}|\mathbf{y})\log \{p(\boldsymbol{\theta}|\mathbf{y})/q(\boldsymbol{\theta})\} d\boldsymbol{\theta}$
- ▶ etc...

Variational Bayes

The objective of variational Bayes approximations is to approximate the posterior density function $p(\boldsymbol{\theta}|\mathbf{y})$ with an approximating density q which is easier to obtain

- ▶ A **mean field variational approximation** $q^*(\boldsymbol{\theta})$ to $p(\boldsymbol{\theta}|\mathbf{y})$ is the minimizer of the Kullback–Leibler divergence

$$\int q(\boldsymbol{\theta}) \log \left\{ \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathbf{y})} \right\} d\boldsymbol{\theta} \quad \text{subject to} \quad q(\boldsymbol{\theta}) = \prod_{i=1}^M q(\boldsymbol{\theta}_i),$$

according to a partition $\{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m\}$ of $\boldsymbol{\theta}$

- ▶ If q belongs to an exponential family of distributions,

$$q^*(\boldsymbol{\theta}_i) \propto \exp \left\{ E_{q(\boldsymbol{\theta} \setminus \boldsymbol{\theta}_i)} \log p(\boldsymbol{\theta}_i | \mathbf{y}, \boldsymbol{\theta} \setminus \boldsymbol{\theta}_i) \right\}, \quad i = 1, \dots, M,$$

where $\boldsymbol{\theta} \setminus \boldsymbol{\theta}_i$ are the entries of $\boldsymbol{\theta}$ with $\boldsymbol{\theta}_i$ omitted

MFVB Approximation

For the model in general form, a tractable variational approximation arises from

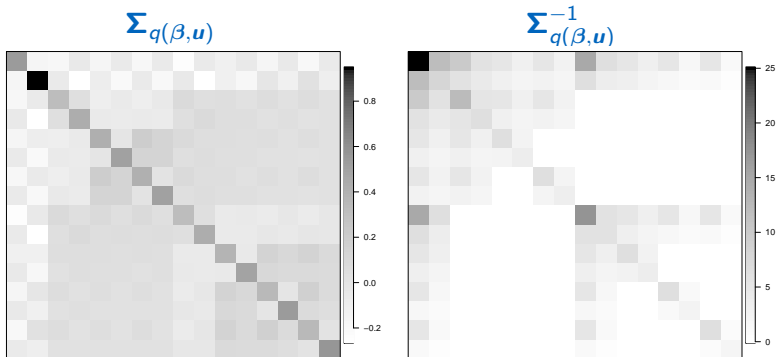
$$p(\boldsymbol{\beta}, \mathbf{u}, \sigma^2, a_{\sigma^2}, \mathbf{G}, \mathbf{A}_G | \mathbf{y}) \approx q(\boldsymbol{\beta}, \mathbf{u}) q(\sigma^2) q(a_{\sigma^2}) q(\mathbf{G}) q(\mathbf{A}_G)$$

Under this restriction, the optimal MFVB approximating densities of the parameters of interest are the following:

$q^*(\boldsymbol{\beta}, \mathbf{u})$ is a $N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})$ density function,
 $q^*(\sigma^2)$ is an Inverse- $\chi^2(\xi_{q(\sigma^2)}, \lambda_{q(\sigma^2)})$ density function,
 $q^*(a_{\sigma^2})$ is an Inverse- $\chi^2(\xi_{q(a_{\sigma^2})}, \lambda_{q(a_{\sigma^2})})$ density function,
 $q^*(\mathbf{G})$ and $q^*(\mathbf{A}_G)$ have form depending on the model specification

Streamlined Variational Inference

- ▶ Note that $q^*(\beta, \mathbf{u})$ is $N(\mu_{q(\beta, \mathbf{u})}, \Sigma_{q(\beta, \mathbf{u})})$
- ▶ At each iteration of the variational algorithm, $\Sigma_{q(\beta, \mathbf{u})}$ has to be computed from a given $\Sigma_{q(\beta, \mathbf{u})}^{-1}$



Linear mixed model with three-level random effects, $m = 2$, $n_1 = 2$, $n_2 = 3$, $\sigma_{ij} = 5$, ($i = 1, 2$; $1 \leq j \leq n_1, n_2$) and $p = q_1 = q_2 = 2$

Naïve MFVB

Let $\mathbf{C} = [\mathbf{X} \ \mathbf{Z}]$. The naïve MFVB algorithm updates for the parameters of $q(\boldsymbol{\beta}, \mathbf{u})$ and $q(\sigma^2)$ are:

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \leftarrow \left(\mu_{q(1/\sigma^2)} \mathbf{C}^T \mathbf{C} + \begin{bmatrix} \boldsymbol{\Sigma}_{\boldsymbol{\beta}}^{-1} & \mathbf{O} \\ \mathbf{O} & E_q(\mathbf{G}^{-1}) \end{bmatrix} \right)^{-1},$$

$$\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \leftarrow \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \left(\mu_{q(1/\sigma^2)} \mathbf{C}^T \mathbf{y} + \begin{bmatrix} \boldsymbol{\Sigma}_{\boldsymbol{\beta}}^{-1} \boldsymbol{\mu}_{\boldsymbol{\beta}} \\ \mathbf{0} \end{bmatrix} \right),$$

$$\xi_{q(\sigma^2)} \leftarrow \nu_{\sigma^2} + n,$$

$$\lambda_{q(\sigma^2)} \leftarrow \mu_{q(1/a_{\sigma^2})} + \|\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}\|^2 + \text{tr} \left\{ \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^T \mathbf{C} \right\}$$

The issues in large dimensions are:

- ▶ \mathbf{Z} has a lot of zero elements that do not need to be stored
- ▶ $E_q(\mathbf{G}^{-1})$ becomes extremely sparse and inversion of $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}^{-1}$ is unfeasible
- ▶ only few sub-blocks of $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$ are effectively useful for computing the other MFVB updates
- ▶ updating the expression for $\lambda_{q(\sigma^2)}$ requires operations involving huge sparse matrices

Streamlined MFVB

Streamlined MFVB:

- ▶ The mean field variational Bayes updates of $\mu_{q(\beta, \mathbf{u})}$ and each of the sub-blocks of $\Sigma_{q(\beta, \mathbf{u})}$ that are relevant for variational inference are expressible as a sparse matrix problem of the form:

$$\mathbf{A}\mu_{q(\beta, \mathbf{u})} = \mathbf{a}$$

- ▶ The non-zero sub-blocks of \mathbf{A} and \mathbf{a} have close (but involved!) forms; see Nolan *et al.* (2020)
- ▶ Efficient matrix inversion and operations arise exploiting the **sparse block-arrowhead** matrix structure of $\Sigma_{q(\beta, \mathbf{u})}^{-1}$ that occurs when two- and three-level models are considered

Global-Local Shrinkage Priors

- ▶ We perform selection of fixed effects through priors on β .
- ▶ Let $\mathbf{X}\beta = \mathbf{X}^R\beta^R + \mathbf{X}^S\beta^S + \mathbf{X}^A\beta^A$, where
 - β^R are fixed effects sharing a predictor with a random effect
 - β^S are fixed effects subject to selection
 - β^A are other fixed effects

- ▶ We use priors

$$\beta^R \sim \mathcal{N}(\mu_{\beta^R}, \Sigma_{\beta^R}), \quad \beta^S \sim p(\beta^S) = \prod_{h=1}^{ps} p(\beta_h^S), \quad \beta^A \sim \mathcal{N}(\mu_{\beta^A}, \Sigma_{\beta^A})$$

where $p(\beta_h^S)$ is the pdf of a **global-local shrinkage prior**

- ▶ These priors admit the scale mixture representation (Polson and Scott, 2011)

$$\beta_h^S | \tau, \zeta_h \stackrel{\text{ind}}{\sim} \mathcal{N}(0, \tau^2 / \zeta_h), \quad \tau \sim p(\tau), \quad \zeta_h \sim p(\zeta_h),$$

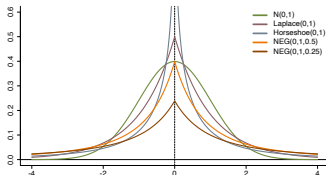
where τ^2 is a **global** variance parameter, and ζ_h is a **local** which may induce a shrinkage effect towards zero to the h th fixed effect

Global-Local Shrinkage Priors and SAVS

- For the fixed effects subject to selection we consider the priors

$$\beta_h^S | \tau \stackrel{\text{iid}}{\sim} \text{Laplace}(0, \tau), \quad \beta_h^S | \tau \stackrel{\text{iid}}{\sim} \text{Horseshoe}(0, \tau), \quad \beta_h^S | \tau \stackrel{\text{iid}}{\sim} \text{NEG}(0, \tau; \lambda)$$

(NEG=Negative-Exponential-Gamma) and use the diffuse prior $\tau \sim \text{Half-Cauchy}(10^5)$



Global-local prior	$p(\beta_h^S \zeta_h, \tau)$	$p(\zeta_h a_{\zeta_h})$	$p(a_{\zeta_h})$
Laplace(0, τ)	$N(0, \tau^2 / \zeta_h)$	Inverse- $\chi^2(2, 1)$	—
Horseshoe(0, τ)	$N(0, \tau^2 / \zeta_h)$	Gamma(1/2, a_{ζ_j})	Gamma(1/2, 1)
NEG(0, $\tau; \lambda$)	$N(0, \tau^2 / \zeta_h)$	Inverse- $\chi^2(2, 2a_{\zeta_j})$	Gamma(λ , 1)

- The posterior densities of irrelevant β_h^S 's are shrunk towards zero but not sparsified
- We perform **fixed-effects selection** via the **signal adaptive variable selector (SAVS)** of Ray and Bhattacharya (2018), which is **hyperparameters-tuning free**:

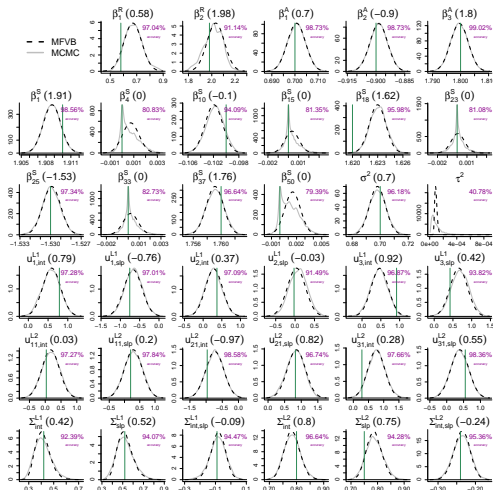
$$\text{if } \|\mathbf{x}_h\|^2 \leq \left| \mu_{q(\beta_h^S)}^* \right|^{-3}, \text{ then } \beta_h^S \text{ is discarded}$$

Simulated Data: Settings

- ▶ 50 simulated datasets from a **three-level linear mixed model** with $m = 100$ groups, each with $n_i = 15$ sub-groups of $o_{ij} = 20$ units each
- ▶ Random intercept and slope for both the group and sub-group levels ($q^{L1} = q^{L2} = p_R = 2$), and $p_A = 3$ additional fixed effects not subject to selection
- ▶ $p_S = 50$ fixed effects subject to selection, where 40 of them are irrelevant (true value equal to zero); $\beta_h^S | \tau \stackrel{\text{iid}}{\sim} \text{Horseshoe}(0, \tau)$, for $1 \leq h \leq 50$
- ▶ Hence, \mathbf{X} has size $30,000 \times 55$, and \mathbf{Z} has size $30,000 \times 3,200$ (96×10^6 cells, of which **99.875% are empty** and do not contain any data)
- ▶ Diffuse priors for all the model parameters
- ▶ Both the naïve and streamlined MFVB approximations are obtained (200 iterations)
- ▶ Comparison with MCMC: Gibbs sampler (5,000 burn-in, 25,000 kept, thinning of size 5)
- ▶ All algorithms are implemented in C++

Simulated Data: Approximation Assessment

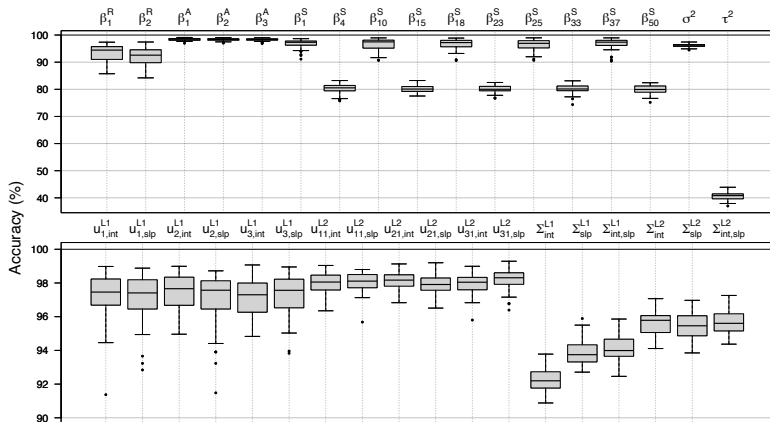
Optimal $q^*(\theta)$ densities obtained via MFVB (black dashed curves) and the corresponding MCMC-based $p(\theta|\mathbf{y})$ densities (grey curves). Vertical lines indicate the parameters true values.



Simulated Data: Accuracy Assessment

Boxplots of the accuracy scores for a selection of model parameters and random effects

$$\text{Accuracy}\{q^*(\theta)\} \equiv (1 - \frac{1}{2} \int |q^*(\theta) - p(\theta|y)| d\theta) \times 100\%$$



Simulated Data: Selection Performance

- ▶ Get MFVB approximations for each of the 50 datasets and priors

$$\beta_h^S \stackrel{\text{iid}}{\sim} \text{N}(0, 10^{10}), \quad \beta_h^S | \tau \stackrel{\text{iid}}{\sim} \text{Laplace}(0, \tau),$$

$$\beta_h^S | \tau \stackrel{\text{iid}}{\sim} \text{Horseshoe}(0, \tau), \quad \beta_h^S | \tau \stackrel{\text{iid}}{\sim} \text{NEG}(0, \tau; 0.25)$$

- ▶ After sparsifying the $q^*(\beta_h^S)$ via the SAVS procedure, the **selection performance** is measured using the F_1 -score:

$$F_1 = \left(\frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \right) \times 100\%, \quad \text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{recall} = \frac{\text{TP}}{\text{TP} + \text{FN}},$$

TP = # of fixed effects with true $\beta_h^S \neq 0$ and selected by SAVS,

TN = # of fixed effects with true $\beta_h^S = 0$ and *not* selected by SAVS,

FP = # of fixed effects with true $\beta_h^S = 0$ and selected by SAVS,

FN = # of fixed effects with true $\beta_h^S \neq 0$ and *not* selected by SAVS

- ▶ Selection performance:
 - **Gaussian**: median F_1 -score 63.5% (q1: 43.5% – q3: 94.2%)
 - **Laplace**: 95.24% (80%–100%)
 - **Horseshoe** and **NEG**: 100% (100%–100%, i.e. perfect fixed effects selection for each of the 50 datasets)

Simulated Data: Speed and Memory Saving

Speed and memory saving assessment for:

- ▶ three-level model
- ▶ p_S fixed effects subject to selection in $\{25, 100, 200\}$
- ▶ number of groups m in $\{10, 50, 100, 200\}$
- ▶ n_i and o_{ij} uniformly from the discrete sets $\{10, \dots, 20\}$ and $\{20, \dots, 30\}$

The table displays average and standard deviation (in brackets) of

- ▶ running time (seconds)
- ▶ total size of required data input (megabytes)

m	p_S	Total runtime of the algorithm (seconds)				Total size of the required data inputs (megabytes)		
		Streamlined MFVB	Naïve MFVB	Naïve MFVB Streamlined MFVB	MCMC	Streamlined MFVB	Naïve MFVB	Naïve MFVB Streamlined MFVB
10	25	0.70 _(0.06)	3.86 _(0.50)	5.54	13.03 _(0.96)	1.39 _(0.07)	11.06 _(1.00)	7.97
	100	3.90 _(0.23)	6.51 _(0.82)	1.67	46.99 _(1.47)	3.68 _(0.24)	12.73 _(1.34)	3.46
	200	9.71 _(1.22)	12.09 _(2.04)	1.24	176.22 _(9.37)	6.60 _(0.52)	15.09 _(1.89)	2.29
50	25	3.29 _(0.11)	365.07 _(37.10)	111.01	63.28 _(3.76)	6.69 _(0.20)	240.40 _(13.92)	35.91
	100	19.66 _(0.87)	415.29 _(51.66)	21.12	138.42 _(5.72)	18.37 _(0.81)	254.04 _(20.84)	13.83
	200	49.25 _(1.47)	501.33 _(39.76)	10.18	305.25 _(6.51)	34.00 _(1.04)	269.34 _(14.90)	7.92
100	25	6.77 _(0.32)	2877.33 _(177.77)	424.72	151.05 _(10.48)	13.56 _(0.34)	967.34 _(42.15)	71.33
	100	40.42 _(1.66)	3172.66 _(116.83)	78.50	266.71 _(10.06)	36.77 _(1.29)	1010.25 _(26.94)	27.47
	200	97.81 _(2.07)	3403.55 _(204.38)	34.80	494.15 _(9.34)	67.53 _(1.31)	1013.90 _(48.67)	15.01
200	25	13.30 _(0.30)	> 5 hours	> 1355	328.87 _(20.04)	26.90 _(0.56)	3817.06 _(113.85)	141.88
	100	79.97 _(1.18)	> 5 hours	> 225	578.88 _(13.45)	73.74 _(1.01)	3941.15 _(97.90)	53.45
	200	197.20 _(4.13)	> 5 hours	> 95	904.64 _(15.77)	135.84 _(2.73)	3991.66 _(102.48)	29.38

Perinatal Data (1)

National Collaborative Perinatal Project data (Klebanoff, 2009)

- ▶ A multisite prospective cohort study which took place in the United States of America between 1959 and 1974.
- ▶ Designed to identify the effects of complications during pregnancy or the perinatal period on birth and child outcomes
- ▶ Focus on predicting the **height-for-age z-score** (standardized measure of the WHO for the height of children after accounting for age) for **37,257 infants** followed longitudinally over their first year of life

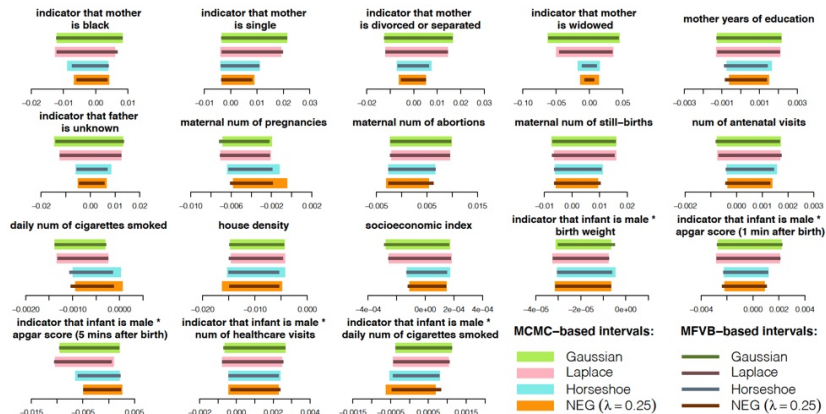
Perinatal Data (1)

National Collaborative Perinatal Project data (Klebanoff, 2009)

- ▶ A multisite prospective cohort study which took place in the United States of America between 1959 and 1974.
- ▶ Designed to identify the effects of complications during pregnancy or the perinatal period on birth and child outcomes
- ▶ Focus on predicting the **height-for-age z-score** (standardized measure of the WHO for the height of children after accounting for age) for **37,257 infants** followed longitudinally over their first year of life
- ▶ We fit a **two-level linear mixed model** with
 - Intercept, age of the infant and its square, plus 38 candidate predictors subject to selection as **fixed effects**
 - Intercept, age of the infant and its square as **random effects**
 - Gaussian, Laplace, Horseshoe and Negative-Exponential-Gamma **priors** for fixed effects subject to selection
- ▶ **MCMC** (5,000 burn-in iterations followed by 25,000 iterations with thinning factor of 5) took **more than 35 minutes** to run
- ▶ **Streamlined MFVB** (200 iterations) took around **2 minutes** to run

Perinatal Data (2)

The 90% high posterior density credible intervals for some of the fixed effects subject to selection



Extensions

More accurate variational approximations can be obtained at higher computational cost when using

- ▶ global-local prior specifications that are not based on auxiliary variables (see Neville *et al.*, 2014)
- ▶ less restrictive mean-field factorizations on the fixed effects, e.g.
 $q(\beta, \mathbf{u}) = q(\beta^R, \mathbf{u})q(\beta^A, \beta^S)$

With relatively little changes in the streamlined MFVB algorithms, it is possible to treat models with

- ▶ other response types (binary data with logit/probit link, Poisson, t, Skew Normal, Skew t, etc...)
- ▶ other types of priors for fixed effects selection (e.g. Bayesian lasso and spike-and-slab priors)

Other extensions include the development of streamlined variational inference for models with

- ▶ unit-specific errors, heteroskedastic covariance structures for groups and sub-groups, higher levels of nesting or crossed random effects
- ▶ priors for selecting random effects (however, this would entail formulating new streamlining strategies for the variational algorithms)

References

- Chen, Z. and Dunson, D. B. (2003). Random effects selection in linear mixed models. *Biometrics*, **59**, 762–769
- Kinney, S. K. and Dunson, D. B. (2007). Fixed and random effects selection in linear and logistic models. *Biometrics*, **63**, 690–698.
- Degani, E., Maestrini, L., Toczyłowska, D. and Wand, M. P. (2022+). Sparse linear mixed model selection via streamlined variational Bayes. [arXiv:2110.07048](https://arxiv.org/abs/2110.07048)
- Klebanoff M. A. (2009). The collaborative perinatal project: a 50-year retrospective. *Paediatric and Perinatal Epidemiology*, **23**, 2–8.
- Maestrini, L. and Wand, M.P. (2021). The Inverse G-Wishart Distribution and Variational Message Passing. *Australian and New Zealand Journal of Statistics*, **63**, 517–541
- Neville, S. E., Ormerod, J. T and Wand, M. P. (2014). Mean field variational Bayes for continuous sparse signal shrinkage: pitfalls and remedies. *Electronic Journal of Statistics*, **8**, 1113–1151
- Nolan, T.H., Menictas, M. and Wand, M.P. (2020). Streamlined computing for variational inference with higher level random effects. *Journal of Machine Learning Research*, **21**, 1–62
- Polson, N. G. and Scott, J. G. (2011). Shrink globally, act locally: sparse Bayesian regularization and prediction. *Bayesian Statistics*, **9**, 501–538
- Yang, M. (2013). Bayesian nonparametric centered random effects models with variable selection. *Biometrical Journal*, **55**, 217–230
- Ray, P. and Bhattacharya, A. (2018). Signal adaptive variable selector for the horseshoe prior. [arXiv:1810.09004](https://arxiv.org/abs/1810.09004)
- Yang, M., Wang, M., and Dong, G. (2020). Bayesian variable selection for mixed effects model with shrinkage prior. *Computational Statistics*, **35**, 227–243

Contacts

E-mail: luca.maestrini@anu.edu.au

Personal website: <https://sites.google.com/view/lucamaestrini>

Twitter: [LucaMaeStats](#)

Article:

Degani, E., Maestrini, L., Toczyłowska, D. and Wand, M. P. (2022+). Sparse linear mixed model selection via streamlined variational Bayes. *arXiv:2110.07048*