

# A Relaxation Approach to Feature Selection for Linear Mixed-Effects Models

Aleksei Sholokhov, James V. Burke, Peng Zheng, Damian Santomauro, and Aleksandr Aravkin

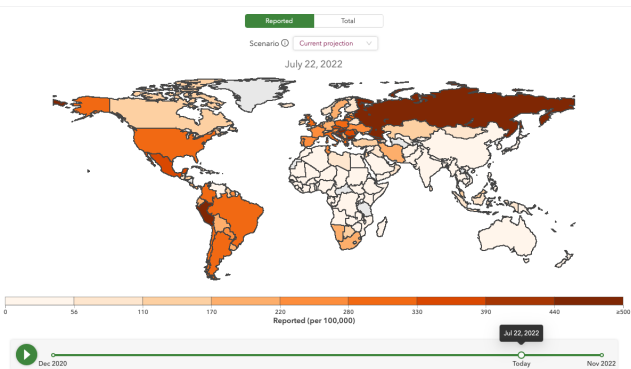
Thursday 8<sup>th</sup> June, 2023



# Feature Selection for Mixed-Effect Models

## Mixed-effect models

- ▶ Used for analyzing **combined data** across a range of **groups**.
- ▶ Use covariates to separate the **population variability** from the **group variability**.
- ▶ **Borrow strength** across groups to estimate key statistics.



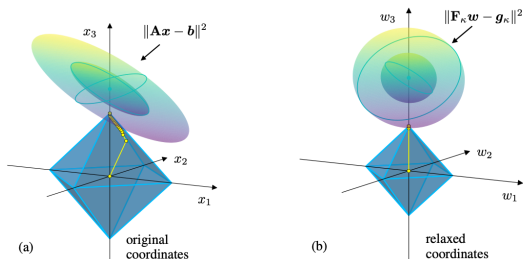
<sup>1</sup> Picture is taken from [covid19.healthdata.org](https://covid19.healthdata.org)

# Feature Selection for Mixed-Effect Models

Practitioners:

- ▶ Often seek **sparse models** that select **most informative** covariates.
- ▶ Want to be **flexible but efficient** in using various sparsity-promoting terms.
- ▶ Want a library to be **universal and compatible** with e.g. sklearn.

Sparse Relaxed Regularized Regression ( $\mathcal{SR3}$ )<sup>2</sup> showed great results for t linear models:



**Goal:** create a feature selection library that uses a relaxation approach for feature-selection in mixed-effect models.

<sup>2</sup>Zheng and Aravkin, "Relax-and-split method for nonconvex inverse problems".

# Linear Mixed-Effect (LME) Models

Dataset:  $m$  groups  $(X_i, Z_i, y_i)$ ,  $i = 1, \dots, m$ , each has  $n_i$  observations

- ▶  $X_i \in \mathbb{R}^{n_i \times p}$  – group  $i$  design matrix for fixed features
- ▶  $Z_i \in \mathbb{R}^{n_i \times q}$  – group  $i$  design matrix for random features
- ▶  $y_i \in \mathbb{R}^{n_i}$  – group  $i$  observations

Model:

$$y_i = X_i \beta + Z_i u_i + \varepsilon_i$$

$$\varepsilon_i \sim \mathcal{N}(0, \Lambda_i)$$

$$u_i \sim \mathcal{N}(0, \Gamma)$$

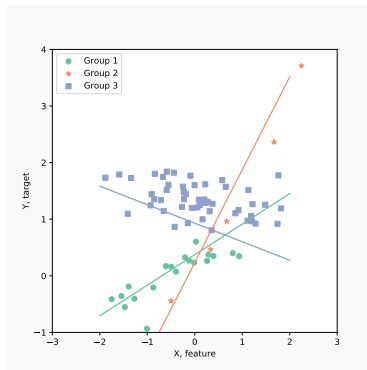
Equivalently:

$$y_i = X_i \beta + \omega_i$$

$$\omega_i \sim \mathcal{N}(0, Z_i \Gamma Z_i^T + \Lambda_i)$$

Simplifying assumption:

$$\Gamma = \text{diag}(\gamma)$$



# Notation

$$\begin{aligned}y_i &= X_i\beta + Z_iu_i + \varepsilon_i \quad i = 1 \dots m \\ \varepsilon_i &\sim \mathcal{N}(0, \Lambda_i) \\ u_i &\sim \mathcal{N}(0, \Gamma)\end{aligned}\tag{1}$$

- ▶  $p$  – number of fixed features,  $q$  – number of random effects.
- ▶  $\beta \in \mathbb{R}^p$  – fixed effects, or mean effects
- ▶  $u_i \in \mathbb{R}^q$  – random effects
- ▶  $\Gamma \in \mathbb{R}^{q \times q}$  – covariance matrix of random effects, often  $\Gamma = \text{diag}((\gamma))$
- ▶  $\varepsilon_i \in \mathbb{R}^{n_i}$  – observation noise
- ▶  $\Lambda_i \in \mathbb{R}^{n_i \times n_i}$  – covariance matrix for noise

Unknowns:  $\beta$ ,  $u_i$ ,  $\gamma$ , sometimes  $\Lambda_i$ .

# Likelihood for Mixed Models

Optimization problem:

$$\mathcal{FS} - \mathcal{LM}\mathcal{E} \quad \min_{\beta \in \mathbb{R}^p, \gamma \in \mathbb{R}_+^q} \mathcal{L}(\beta, \gamma) + R(\beta, \gamma) \quad (2)$$

Where  $\mathcal{L}$ :

$$\begin{aligned} \mathcal{L}(\beta, \gamma) = & \sum_{i=1}^m \frac{1}{2} (y_i - X_i \beta)^T (Z_i \Gamma Z_i^T + \Lambda_i)^{-1} (y_i - X_i \beta) + \\ & + \frac{1}{2} \log \det (Z_i \Gamma Z_i^T + \Lambda_i), \quad \Gamma = \text{diag}((\gamma)) \end{aligned} \quad (3)$$

- ▶  $\mathcal{L}(\beta, \gamma)$  is smooth on its domain, quadratic w.r.t.  $\beta$  and  $\bar{\eta}$ -weakly-convex w.r.t.  $\gamma$ .
- ▶  $R(\beta, \gamma)$  is closed, proper, with easily computed *prox operator*

# Regularization

- ▶  $R(\beta, \gamma)$  is closed, proper, with easily computed *prox operator*

$$\text{prox}_{\alpha R + \delta_{\mathcal{C}}}(\tilde{\beta}, \tilde{\gamma}) := \underset{(\beta, \gamma) \in \mathcal{C}}{\text{argmin}} R(\beta, \gamma) + \frac{1}{2\alpha} \|(\beta, \gamma) - (\tilde{\beta}, \tilde{\gamma})\|_2^2, \quad (4)$$

where  $\mathcal{C} := \mathbb{R}^P \times \mathbb{R}_+^q$

Examples:

- ▶  $R(x) = \lambda \sum_{j=1}^p w_j \|x_j\|_1$  – LASSO and Adaptive LASSO penalties<sup>3</sup>
- ▶  $R(x) = \lambda \|x\|_0 - \ell_0$  penalty<sup>4</sup>
- ▶  $R(x)$  – SCAD penalty<sup>(5)</sup>

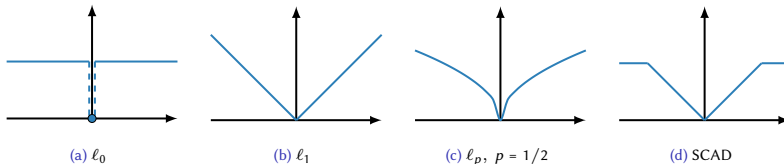


Figure: Four commonly-used regularizers which promote sparsity

<sup>3</sup>Bondell, Krishna, and Ghosh, “Joint Variable Selection for Fixed and Random Effects in Linear Mixed-Effects Models”; Lin, Pang, and Jiang, “Fixed and random effects selection by REML and pathwise coordinate optimization”.

<sup>4</sup>Vaida and Blanchard, “Conditional Akaike information for mixed-effects models”; Jones, “Bayesian information criterion for longitudinal and clustered data”.

<sup>5</sup>J. Fan and Li, “Variable Selection via Nonconcave Penalized Likelihood and its Oracle Properties”; Y. Fan and Li, “Variable selection in linear mixed effects models”.

## SR3-Relaxation for Mixed-Effect Models ( $\mathcal{MSR3}$ )

Original problem  $\mathcal{FS} - \mathcal{LME}$ :

$$\min_{\beta \in \mathbb{R}^p, \gamma \in \mathbb{R}_+^q} \mathcal{L}(\beta, \gamma) + R(\beta, \gamma) \quad (5)$$

Relaxed problem  $\mathcal{MSR3}$ :

$$\min_{\beta, \tilde{\beta} \in \mathbb{R}^p, \gamma, \tilde{\gamma} \in \mathbb{R}_+^q} \mathcal{L}(\beta, \gamma) + \phi_\mu(\gamma) + \kappa_\eta(\beta - \tilde{\beta}, \gamma - \tilde{\gamma}) + R(\tilde{\beta}, \tilde{\gamma}) \quad (6)$$

where the *relaxation*  $\kappa_\eta$  decouples the likelihood and the regularizer

$$\kappa_\eta(\beta - \tilde{\beta}, \gamma - \tilde{\gamma}) := \frac{\eta}{2} \|\beta - \tilde{\beta}\|_2^2 + \frac{\eta}{2} \|\gamma - \tilde{\gamma}\|_2^2, \quad \eta > \bar{\eta} \quad (7)$$

and the *perspective mapping*  $\phi_\mu$  replaces  $\gamma \geq 0$  with a log-barrier

$$\phi_\mu(\gamma) := \begin{cases} -\mu \sum_{i=1}^q \ln(\gamma_i / \mu), & \mu > 0 \\ \delta_{\mathbb{R}_+^q}(\gamma), & \mu = 0 \\ +\infty, & \mu < 0 \end{cases} \quad (8)$$



# Value Function Reformulation

$\mathcal{MSR3}$ -relaxation replaces the original likelihood  $\mathcal{L}$  with a *value function*  $u_{\eta,\mu}$ :

$$\begin{aligned} v_{\eta,\mu}(\tilde{\beta}, \tilde{\gamma}) &:= \min_{(\beta, \gamma)} \mathcal{L}_{\eta,\mu}((\beta, \gamma), (\tilde{\beta}, \tilde{\gamma})) \\ &:= \min_{(\beta, \gamma)} \mathcal{L}(\beta, \gamma) + \phi_{\mu}(\gamma) + \kappa_{\eta}(\beta - \tilde{\beta}, \gamma - \tilde{\gamma}) \end{aligned} \tag{9}$$

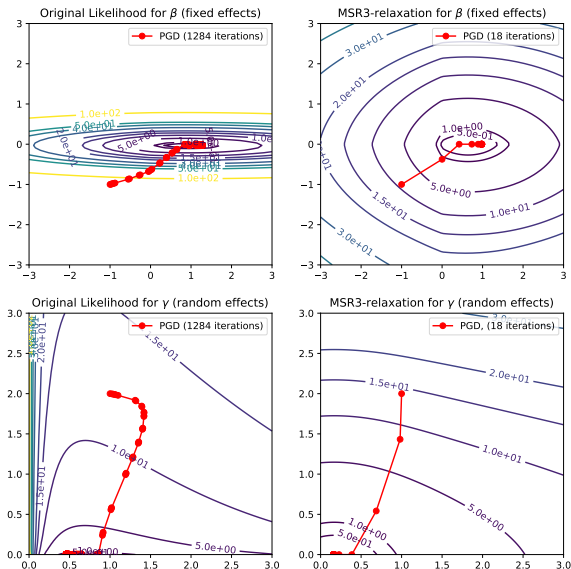
so  $\mathcal{MSR3}$ -formulation (6) becomes

$$\min_{\beta \in \mathbb{R}^p, \gamma \in \mathbb{R}_+^q} v_{\eta,\mu}(\tilde{\beta}, \tilde{\gamma}) + R(\tilde{\beta}, \tilde{\gamma}) \tag{10}$$

When  $\eta$  is larger than the weak-convexity constant

- ▶  $v_{\eta,\mu}$  is well-defined and continuously differentiable.
- ▶ As  $\mu \rightarrow 0$  and  $\eta \rightarrow \infty$ , cluster points of solutions to  $\mathcal{MSR3}$  are first-order stationary points for  $\mathcal{FS} - \mathcal{LM}\mathcal{E}$
- ▶  $v_{\eta,\mu}$  don't need to be evaluated precisely.

# Value Function Reformulation



**Figure:** Comparison of the level-sets for the original likelihood (left) and  $MSR3$ -likelihood (right), for fixed (top) and random (bottom) effects.

# Designing an Algorithm

$G_{\nu,\eta}$  encodes both gradient of a Lagrangian (lines 1-2) and the complementarity condition (line 3):

$$G_{\nu,\eta}((\beta, \gamma, \nu), (\tilde{\beta}, \tilde{\gamma})) := \begin{bmatrix} \nabla_{\beta} \mathcal{L}(\beta, \gamma) + \eta(\beta - \tilde{\beta}) \\ \nabla_{\gamma} \mathcal{L}(\beta, \gamma) + \eta(\gamma - \tilde{\gamma}) - \nu \\ \nu \odot \gamma - \mu \mathbf{1} \end{bmatrix} \quad (11)$$

We apply Newton method to  $G$  while geometrically decreasing  $\mu$ .

**Lemma:** For every  $(\mu, \eta) \in \mathbb{R}_+ \times \mathbb{R}_{++}$ ,

$$\begin{aligned} (\hat{\beta}, \hat{\gamma}) &= \underset{(\beta, \gamma)}{\operatorname{argmin}} \mathcal{L}_{\eta, \mu}((\beta, \gamma), (\tilde{\beta}, \tilde{\gamma})) \\ &\iff \\ \exists \hat{\nu} \in \mathbb{R}_+^q \text{ s.t. } G_{\nu, \eta}((\beta, \gamma, \hat{\nu}), (\tilde{\beta}, \tilde{\gamma})) &= 0 \end{aligned} \quad (12)$$

If  $\mu > 0$ , then  $\hat{\nu} = -\nabla \phi_{\mu}(\hat{\gamma})$ , and if  $\mu = 0$ , then  $\hat{\nu}$  is the unique KKT multiplier associated with the constraint  $0 \leq \gamma$ .

# MSR3-fast Algorithm

---

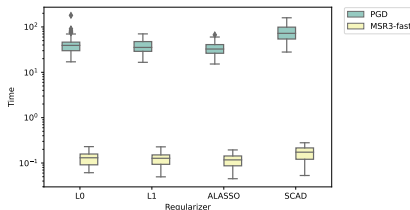
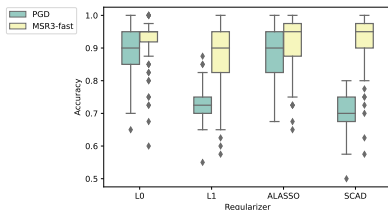
---

```
1 progress  $\leftarrow$  True;  iter = 0;
2  $\beta^+, \tilde{\beta}^+ \leftarrow \beta_0$ ;   $\gamma^+, \tilde{\gamma}^+ \leftarrow \gamma_0$ ;   $v^+ \leftarrow \mathbf{1} \in \mathbb{R}^q$ ;   $\mu \leftarrow \frac{v^{+T} \gamma^+}{10q}$ 
3 while iter < max_iter and  $\|G_\mu(\beta^+, \gamma^+, v^+)\| > \text{tol}$  and progress
4   do
5      $\beta \leftarrow \beta^+$ ;   $\gamma \leftarrow \gamma^+$ ;   $\tilde{\beta} \leftarrow \tilde{\beta}^+$ ;   $\tilde{\gamma} \leftarrow \tilde{\gamma}^+$ 
6      $[dv, d\beta, d\gamma] \leftarrow \nabla G_\mu((\beta, \gamma, v), (\tilde{\beta}, \tilde{\gamma}))^{-1} G_\mu((\beta, \gamma, v), (\tilde{\beta}, \tilde{\gamma}))$ 
7      $\alpha \leftarrow 0.99 \times \min \left( 1, -\frac{\gamma_i}{d\gamma_i}, \forall i : d\gamma_i < 0 \right)$ 
8      $\beta^+ \leftarrow \beta + \alpha d\beta$ ;   $\gamma^+ \leftarrow \gamma + \alpha d\gamma$ ;   $v^+ \leftarrow v + \alpha dv$ 
9     if  $\|\gamma^+ \odot v^+ - q^{-1} \gamma^{+T} v^+ \mathbf{1}\| > 0.5 q^{-1} v^{+T} \gamma^+$  then continue;
10    else
11       $\tilde{\beta}^+ = \text{prox}_{\alpha R}(\beta^+)$ ;   $\tilde{\gamma}^+ = \text{prox}_{\alpha R + \delta_{\mathbb{R}_+}}(\gamma^+)$ ;   $\mu = \frac{1}{10} \frac{v^{+T} \gamma^+}{q}$ 
12    end
13    progress = ( $\|\beta^+ - \beta\| \geq \text{tol}$  or  $\|\gamma^+ - \gamma\| \geq \text{tol}$  or  $\|\tilde{\beta}^+ - \tilde{\beta}\| \geq \text{tol}$  or
14       $\|\tilde{\gamma}^+ - \tilde{\gamma}\| \geq \text{tol}$ )
15    iter += 1
16 end
17 return  $\tilde{\beta}^+, \tilde{\gamma}^+$ 
```

---

# Application to Synthetic Problems

- ▶ The number of fixed effects  $p$  and random effects  $q$  is 20.
- ▶  $\beta = \gamma = \frac{1}{2}[1, 2, 3, \dots, 10, 0 \dots, 0]$
- ▶ 9 groups with sizes  $[10, 15, 4, 8, 3, 5, 18, 9, 6]$
- ▶  $X_i \sim \mathcal{N}(0, I)^p$ ,  $Z_i = X_i$ ,  $\varepsilon_i \sim \mathcal{N}(0, 0.3^2 I)$
- ▶ Each experiment is repeated 100 times.
- ▶ Grid-search for  $\eta \in [10^{-4}, 10^2]$ , golden search for  $\lambda \in [0, 10^5]$
- ▶ Final model is chosen to maximize BIC



- +  $MSR3$ -relaxation improves feature selection performance of the original likelihood.
- +  $MSR3$ -fast optimization accelerates the compute time by  $\sim 10^2$ .
- Initialization of  $\eta$  is problem-specific

# Comparison to Other Libraries

Algorithm	MSR3-Fast ( $\ell_1$ )	glmLasso <sup>67</sup>	lmmLasso <sup>89</sup>	PGD ( $\ell_1$ )
Accuracy, %	<b>88</b>	48	66	73
FE Accuracy, %	<b>86</b>	52	47	56
RE Accuracy, %	<b>91</b>	45	84	<b>91</b>
Time, sec	<b>0.19</b>	1.37	11.51	38.39
Iterations, num	34	50	-	7693

---

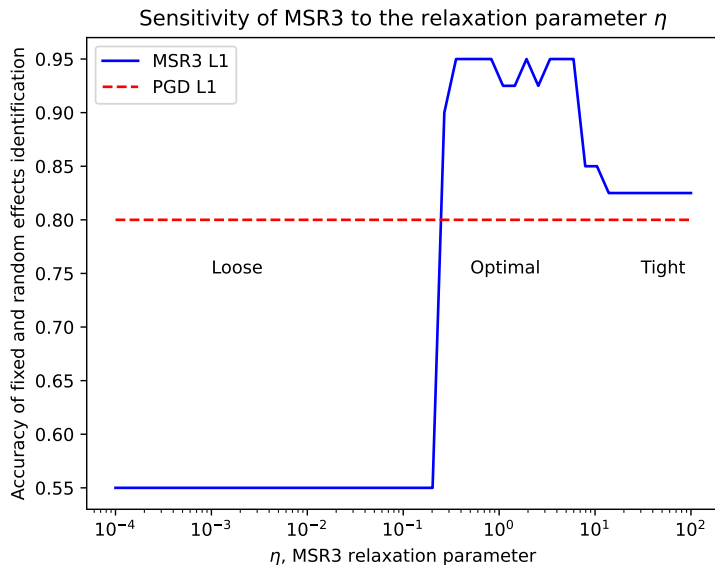
<sup>6</sup><https://rdrr.io/cran/glmLasso/man/glmLasso.html>

<sup>7</sup>Groll and Tutz, “Variable selection for generalized linear mixed models by L 1-penalized estimation”.

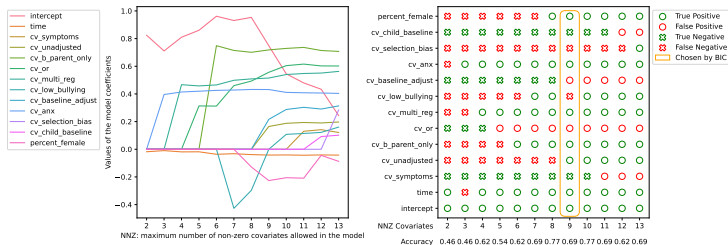
<sup>8</sup><https://rdrr.io/cran/lmmLasso/>

<sup>9</sup>Schellldorfer, Bühlmann, and DE GEER, “Estimation for high-dimensional linear mixed-effects models using l1-penalization”.

## Choice of $\eta$



# $\ell_0$ -based Covariate Selection for Bullying Study from GBD



**Figure:** Fixed and random covariate selection for Bullying dataset from <sup>10</sup>. The model selected 9 covariates, 7 of which were historically significant, and did not select 4 covariates, 1 of which was historically significant.

<sup>10</sup>Institute for Health Metrics and Evaluation (IHME). Bullying Victimization Relative Risk Bundle GBD 2020. Seattle, United States of America (USA), 2021.



The code is available on GitHub: <https://github.com/aksholokhov/pysr3>

- ▶ All estimators are fully compatible to `sklearn` library.
- ▶ Implements SR3 for linear, generalized-linear, and linear mixed-effect models.
- ▶ Has tutorials, tests, and documentation.

# References I



Bondell, Howard D., Arun Krishna, and Sujit K. Ghosh. “Joint Variable Selection for Fixed and Random Effects in Linear Mixed-Effects Models”. In: Biometrics 66.4 (Dec. 2010), pp. 1069–1077. ISSN: 0006341X. DOI: 10.1111/j.1541-0420.2010.01391.x. arXiv: NIHMS150003. URL: <http://doi.wiley.com/10.1111/j.1541-0420.2010.01391.x>.



Fan, Jianqing and Runze Li. “Variable Selection via Nonconcave Penalized Likelihood and its Oracle Properties”. In: Journal of the American Statistical Association 96.456 (Dec. 2001), pp. 1348–1360. ISSN: 0162-1459. DOI: 10.1198/016214501753382273. URL: <http://www.tandfonline.com/doi/abs/10.1198/016214501753382273>.



Fan, Yingying and Runze Li. “Variable selection in linear mixed effects models”. In: The Annals of Statistics 40.4 (Aug. 2012), pp. 2043–2068. ISSN: 0090-5364. DOI: 10.1214/12-AOS1028. URL: <http://projecteuclid.org/euclid.aos/1351602536>.



Groll, Andreas and Gerhard Tutz. “Variable selection for generalized linear mixed models by L1-penalized estimation”. In: Statistics and Computing 24.2 (2014), pp. 137–154.



Jones, Richard H. “Bayesian information criterion for longitudinal and clustered data”. In: Statistics in Medicine 30.25 (Nov. 2011), pp. 3050–3056. ISSN: 02776715. DOI: 10.1002/sim.4323. URL: <http://doi.wiley.com/10.1002/sim.4323>.



Lin, Bingqing, Zhen Pang, and Jiming Jiang. “Fixed and random effects selection by REML and pathwise coordinate optimization”. In: Journal of Computational and Graphical Statistics 22.2 (2013), pp. 341–355. ISSN: 10618600. DOI: 10.1080/10618600.2012.681219.



Schelldorfer, Jürg, Peter Bühlmann, and SARA VAN DE GEER. “Estimation for high-dimensional linear mixed-effects models using l1-penalization”. In: Scandinavian Journal of Statistics 38.2 (2011), pp. 197–214.

# References II



Vaida, Florin and Suzette Blanchard. “Conditional Akaike information for mixed-effects models”. In: Biometrika 92.2 (June 2005), pp. 351–370. ISSN: 1464-3510. DOI: 10.1093/biomet/92.2.351. URL: <http://academic.oup.com/biomet/article/92/2/351/233128/Conditional-Akaike-information-for-mixedeffects>.



Zheng, Peng and Aleksandr Aravkin. “Relax-and-split method for nonconvex inverse problems”. In: Inverse Problems 36.9 (2020). ISSN: 13616420. DOI: 10.1088/1361-6420/aba417.