

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

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

Friday 22<sup>nd</sup> July, 2022



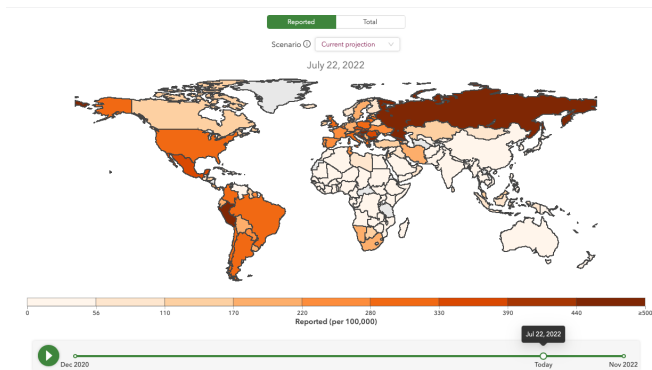
APPLIED MATHEMATICS  
UNIVERSITY *of* WASHINGTON



IHME

# 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.
- ▶ Often seek **sparse models** that only use **most informative** covariates.



1

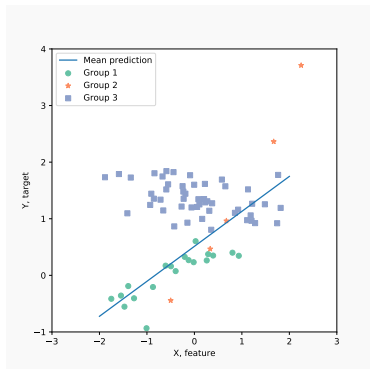
<sup>1</sup>Picture is taken from [covid19.healthdata.org](https://covid19.healthdata.org)

# Linear Mixed-Effect (LME) Models

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

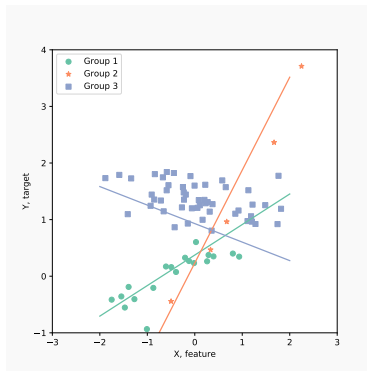
- ▶  $X_i \in \mathbb{R}^{n_i \times p}$  – group  $i$  design matrix for fixed features
- ▶  $Z_i \in \mathbb{R}^{n_i \times q}$  – group  $i$  design matrix for random effects
- ▶  $y_i \in \mathbb{R}^{n_i}$  – group  $i$  observations

Standard Linear Regression:



$$y = X\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \Lambda)$$

Linear Mixed-Effect Model:



$$y_i = X_i\beta + Z_i u_i + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \Lambda_i)$$
$$u_i \sim \mathcal{N}(0, \Gamma)$$

## 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

Negative log-likelihood:

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

Maximum likelihood estimates for  $\beta$  and  $\gamma$  solve the problem:

$$\mathcal{LM}\mathcal{E} \quad \min_{\beta \in \mathbb{R}^p, \gamma \in \mathbb{R}_+^q} \mathcal{L}(\beta, \gamma)\tag{3}$$

To select covariates we add a sparsity-promoting regularizer  $R(\beta, \gamma)$

$$\mathcal{FS} - \mathcal{LM}\mathcal{E} \quad \min_{\beta \in \mathbb{R}^p, \gamma \in \mathbb{R}_+^q} \mathcal{L}(\beta, \gamma) + R(\beta, \gamma)\tag{4}$$

- ▶  $\mathcal{L}(\beta, \gamma)$  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, convex, with easily computed *prox operator*

# Regularization

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

$$\text{prox}_{\alpha R + \delta_{\mathcal{C}}}(\tilde{\beta}, \tilde{\gamma}) := \underset{(\beta, \gamma) \in \mathcal{C}}{\operatorname{argmin}} R(\beta, \gamma) + \frac{1}{2\alpha} \|(\beta, \gamma) - (\tilde{\beta}, \tilde{\gamma})\|_2^2, \quad (5)$$

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 [1, 5]
- ▶  $R(x) = \lambda \|x\|_0 - \ell_0$  penalty [6, 4]
- ▶  $R(x)$  – SCAD penalty ([2, 3])

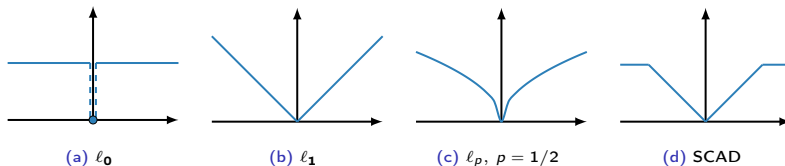


Figure: Four commonly-used regularizers which promote sparsity

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

Original problem  $\mathcal{FS} - \mathcal{LM}\mathcal{E}$ :

$$\min_{\beta \in \mathbb{R}^p, \gamma \in \mathbb{R}_+^q} \mathcal{L}(\beta, \gamma) + R(\beta, \gamma) \quad (6)$$

Relaxed problem  $\mathcal{MSR3}$ :

$$\min_{\beta, \tilde{\beta} \in \mathbb{R}^p, \gamma, \tilde{\gamma} \in \mathbb{R}_+^q} \mathcal{L}(\beta, \gamma) + \phi_\mu(\gamma) + \kappa_\eta(\beta - \tilde{\beta}, \gamma - \tilde{\gamma}) + R(\tilde{\beta}, \tilde{\gamma}) \quad (7)$$

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

$$\kappa_\eta(\beta - \tilde{\beta}, \gamma - \tilde{\gamma}) := \frac{\eta}{2} \|\beta - \tilde{\beta}\|_2^2 + \frac{\eta}{2} \|\gamma - \tilde{\gamma}\|_2^2, \quad \eta > \bar{\eta} \quad (8)$$

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

$$\phi_\mu(\gamma) := \begin{cases} -\mu \sum_{i=1}^q \ln(\gamma_i/\mu), & \mu > 0 \\ \delta_{\mathbb{R}_+^q}(\gamma), & \mu = 0 \\ +\infty, & \mu < 0 \end{cases} \quad (9)$$

# Value Function Reformulation

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

$$\begin{aligned} u_{\eta,\mu}(\tilde{\beta}, \tilde{\gamma}) &:= \min_{(\beta, \gamma)} \mathcal{L}_{\eta,\mu}((\beta, \gamma), (\tilde{\beta}, \tilde{\gamma})) \\ &:= \min_{(\beta, \gamma)} \mathcal{L}(\beta, \gamma) + \phi_{\mu}(\gamma) + \kappa_{\eta}(\beta - \tilde{\beta}, \gamma - \tilde{\gamma}) \end{aligned} \quad (10)$$

so  $\mathcal{MSR3}$ -formulation (7) becomes

$$\min_{(\tilde{\beta}, \tilde{\gamma}) \in \mathcal{C}} u_{\eta,\mu}(\tilde{\beta}, \tilde{\gamma}) + R(\tilde{\beta}, \tilde{\gamma}) \quad (11)$$

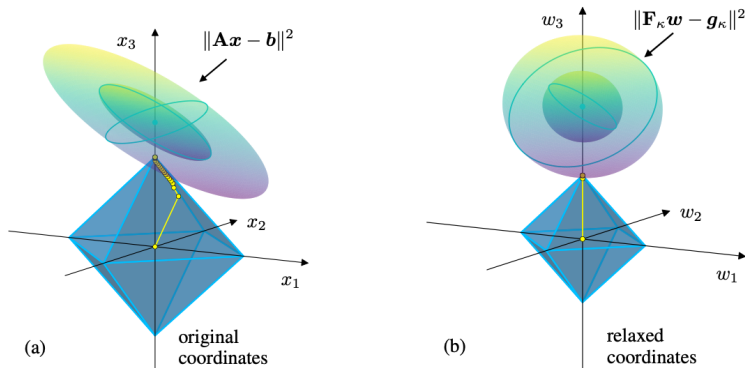
When  $\eta$  is larger than the weak-convexity constant

- ▶  $u_{\eta,\mu}$  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  $\mathcal{FS} - \mathcal{LM}\mathcal{E}$

**Key observation:** in practice, we don't need accurate solutions for (10): a few Newton iterations keep the solution close to the central path.

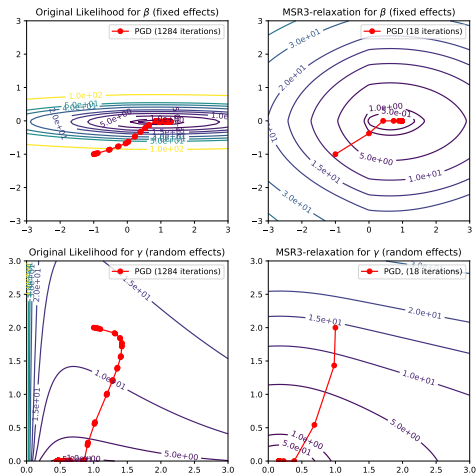


# Value Function Reformulation



**Figure:** Picture from [7]: for a linear problem, value function relaxation “squashes” level-sets simplifying the optimization landscape.

# Value Function Reformulation



**Figure:** Comparison of the level-sets for the original likelihood (left) and  $\mathcal{MSR3}$ -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):

$$G_{\nu,\eta}((\beta, \gamma, \nu), (\tilde{\beta}, \tilde{\gamma})) := \begin{bmatrix} \nabla_{\beta} \mathcal{L}(\beta, \gamma) + \eta(\beta - \tilde{\beta}) \\ \nabla_{\gamma} \mathcal{L}(\beta, \gamma) + \eta(\gamma - \tilde{\gamma}) - \nu \\ \nu \odot \gamma - \mu \mathbf{1} \end{bmatrix} \quad (12)$$

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

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

$$\begin{aligned} (\hat{\beta}, \hat{\gamma}) &= \underset{(\beta, \gamma)}{\operatorname{argmin}} \mathcal{L}_{\eta, \mu}((\beta, \gamma), (\tilde{\beta}, \tilde{\gamma})) \\ &\iff \\ \exists \hat{\nu} \in \mathbb{R}_+^q \text{ s.t. } &G_{\nu, \eta}((\beta, \gamma, \hat{\nu}), (\tilde{\beta}, \tilde{\gamma})) = 0 \end{aligned} \quad (13)$$

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 ← True;  iter = 0;
2  $\beta^+, \tilde{\beta}^+ \leftarrow \beta_0; \quad \gamma^+, \tilde{\gamma}^+ \leftarrow \gamma_0; \quad v^+ \leftarrow 1 \in \mathbb{R}^q; \quad \mu \leftarrow \frac{v^{+T} \gamma^+}{10q}$ 
3 while iter < max_iter and  $\|G_\mu(\beta^+, \gamma^+, v^+)\| > \text{tol}$  and progress
4   do
5      $\beta \leftarrow \beta^+; \quad \gamma \leftarrow \gamma^+; \quad \tilde{\beta} \leftarrow \tilde{\beta}^+; \quad \tilde{\gamma} \leftarrow \tilde{\gamma}^+$ 
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      $\beta^+ \leftarrow \beta + \alpha d\beta; \quad \gamma^+ = \gamma + \alpha d\gamma; \quad v^+ \leftarrow v + \alpha dv$ 
9     if  $\|\gamma^+ \odot v^+ - q^{-1} \gamma^{+T} v^+ 1\| > 0.5 q^{-1} v^{+T} \gamma^+$  then continue;
10    else
11       $\tilde{\beta}^+ = \text{prox}_{\alpha R}(\beta^+); \quad \tilde{\gamma}^+ = \text{prox}_{\alpha R + \delta_{\mathbb{R}^+}}(\gamma^+); \quad \mu = \frac{1}{10} \frac{v^{+T} \gamma^+}{q}$ 
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  $\tilde{\beta}^+, \tilde{\gamma}^+$ 
```

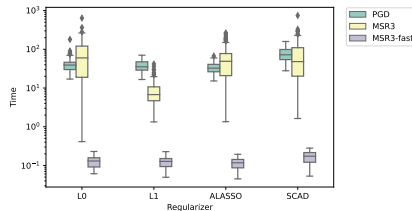
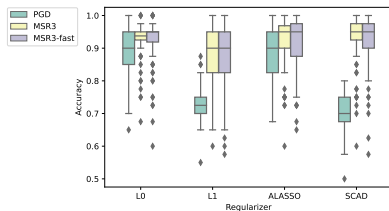
# Application to Synthetic Problems

## The Experiment

- ▶ 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

Regularizer	Model Metric	PGD	MSR3	MSR3-fast
L0	Accuracy	0.89	<b>0.92</b>	<b>0.92</b>
	Time	41.68	88.54	<b>0.13</b>
L1	Accuracy	0.73	<b>0.88</b>	<b>0.88</b>
	Time	38.39	9.13	<b>0.13</b>
ALASSO	Accuracy	0.88	<b>0.92</b>	0.91
	Time	34.55	65.19	<b>0.12</b>
SCAD	Accuracy	0.71	<b>0.93</b>	0.92
	Time	77.62	84.67	<b>0.17</b>

# Application to Synthetic Problems



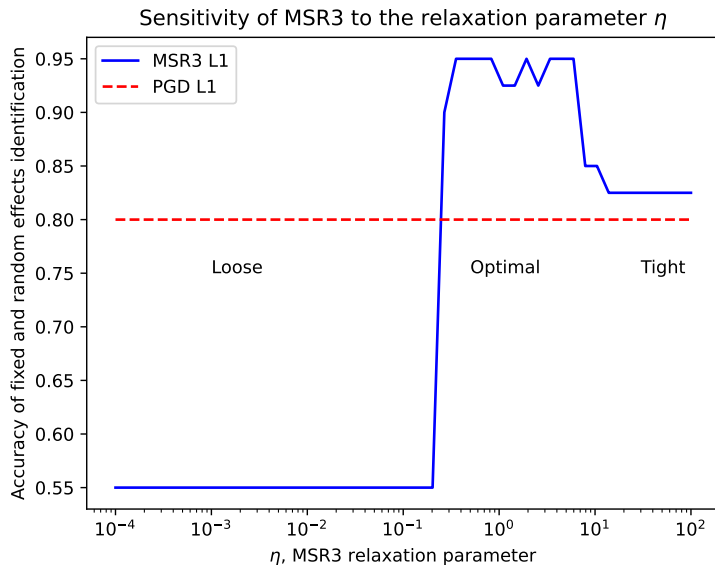
## Benefits:

- ▶ *MSR3*-relaxation improves feature selection performance than the original likelihood.
- ▶ *MSR3*-fast optimization accelerates the compute time by  $\sim 10^2$ .

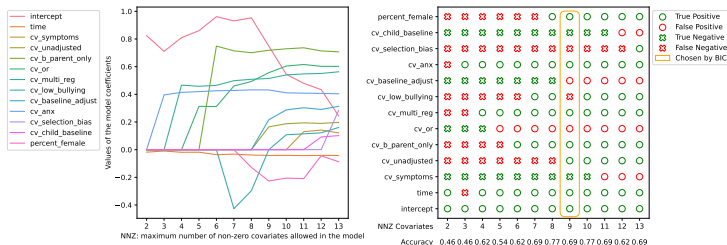
## Challenge:

- ▶ Initialization of  $\eta$  is problem-specific

## Choice of $\eta$



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



**Figure:** Fixed and random covariate selection for Bullying dataset from [?]. The model selected 9 covariates, 7 of which were historically significant, and did not select 4 covariates, 1 of which was historically significant.



# Thank You!

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:

- [1] Howard D. Bondell, Arun Krishna, and Sujit K. Ghosh. Joint Variable Selection for Fixed and Random Effects in Linear Mixed-Effects Models. Biometrics, 66(4):1069–1077, dec 2010.
- [2] Jianqing Fan and Runze Li. Variable Selection via Nonconcave Penalized Likelihood and its Oracle Properties. Journal of the American Statistical Association, 96(456):1348–1360, dec 2001.
- [3] Yingying Fan and Runze Li. Variable selection in linear mixed effects models. The Annals of Statistics, 40(4):2043–2068, aug 2012.
- [4] Richard H. Jones. Bayesian information criterion for longitudinal and clustered data. Statistics in Medicine, 30(25):3050–3056, nov 2011.
- [5] Bingqing Lin, Zhen Pang, and Jiming Jiang. Fixed and random effects selection by REML and pathwise coordinate optimization. Journal of Computational and Graphical Statistics, 22(2):341–355, 2013.
- [6] Florin Vaida and Suzette Blanchard. Conditional Akaike information for mixed-effects models. Biometrika, 92(2):351–370, jun 2005.
- [7] Peng Zheng and Aleksandr Aravkin. Relax-and-split method for nonsmooth nonconvex problems. 1:1–38, feb 2018.