

MixLasso: Generalized Mixed Regression via Convex Atomic-Norm Regularization

Ian E.H. Yen¹², Wei-Cheng Lee³, Sung-En Chang³, Kai Zhong⁴, Shou-De Lin³ and Pradeep Ravikumar¹

¹Carnegie Mellon University. ²Snap Inc. ³National Taiwan University. ⁴Amazon Inc.

Abstract

- ▶ In this work, we propose a novel **convex estimator (MixedLasso)** for **Generalized Mixed Regression** models.
- ▶ To best of our knowledge, this is the first method with low-order **polynomial runtime** and **sample complexity** **without restrictive assumptions** on the data distribution for GMR.
- ▶ In experiments, the MixLasso **significantly outperforms** other methods when there is a **larger number** of latent regression functions.

Generalized Mixed Regression Models

- ▶ **Generalized Mixed Regression Model (GMR)** is a generalization of **Mixed Regression**, where each response is an **additive combination of latent regression functions**.

- ▶ In **Generalized Mixed Regression**, each response

$$\mathbf{y}_i = \sum_{k=1}^K z_{i,k} f_k(\mathbf{x}_i) + \omega$$

where $\mathbf{y}_i \in \mathbb{R}$: response, $\mathbf{x}_i \in \mathbb{R}^D$: **explanatory variable**, $z_{i,k} \in \{0, 1\}, k = 1, \dots, K$ **binary latent indicators**, $f_k(\mathbf{x}_i) : \mathbb{R}^D \rightarrow \mathbb{R}$: **is the regression function of k -th component**, and $\omega \in \mathbb{R}$: **noise**.

- ▶ **Standard Mixed Regression** is a special case with $\|\mathbf{z}_i\|_0 = 1$.

- ▶ **Goal** is to find $\mathcal{F} := \{f_k(\mathbf{x})\}_{k=1}^K$ minimizing the risk

$$r(\mathcal{F}) := \mathbb{E} \left[\min_{\mathbf{z} \in \{0,1\}^K} \frac{1}{2} \left(y - \sum_{k=1}^K z_k f_k(\mathbf{x}) \right)^2 \right] \quad (1)$$

which yields a trade-off between risk and size of K .

- ▶ **We focus on** two family of functions:

- ▶ **linear case**: $f_k(\mathbf{x}) := \langle \mathbf{w}_k, \mathbf{x} \rangle$.
- ▶ **Non-linear extension**: $f_k(\mathbf{x})$ lying in a **Reproducing Kernel Hilbert Space (RKHS) \mathcal{H}** with respect to some Mercer kernel $\mathcal{K}(\cdot, \cdot)$.

Related Works & Results

- ▶ **Existing Approaches**:

- ▶ **MCMC, Variational Bayes**:
No finite-time theoretical guarantee.
- ▶ **convex relaxation based on nuclear norm**:
Restricted to two components, with Gaussian assumptions on the input matrix.
- ▶ **Tensor (Spectral) Methods**:
have high sample complexity w.r.t. D or K , and isotropic Gaussian assumptions on the inputs.

- ▶ **This Paper**:

- ▶ A convex estimator — **MixLasso**.
- ▶ Low-order **polynomial runtime** and **sample complexity**.
- ▶ **No restrictive assumption** on the inputs.

Convex Estimation via Atomic Norm (linear function)

- ▶ Regularized Empirical Risk Minimization:

$$\min_{W \in \mathbb{R}^{K \times D}, \mathbf{z}_i \in \{0,1\}^K} \frac{1}{2N} \sum_{i=1}^N (y_i - \mathbf{z}_i^T W \mathbf{x}_i)^2 + \frac{\tau}{2} \|W\|_F^2.$$

- ▶ Given $Z := (\mathbf{z}_i)_{i=1}^N$, the **dual problem w.r.t. W** is:

$$\min_{M=ZZ^T \in \{0,1\}^{N \times N}} \left\{ \max_{\alpha \in \mathbb{R}^N} \underbrace{\frac{-1}{\tau} \text{tr}(\mathcal{D}(\alpha) X X^T \mathcal{D}(\alpha) Z Z^T)}_{g(M)} - \sum_{i=1}^N L^*(y_i, -\alpha_i) \right\}.$$

- ▶ **Key insight**: the function is **convex w.r.t. M** .
- ▶ Enforce structure $M = ZZ^T$ via an atomic norm.

- ▶ Let $\mathcal{S} := \{k \mid \mathbf{z}_k \in \{0,1\}^N\}$. We define **Atomic Norm**:

$$\|M\|_{\mathcal{S}} := \min_{c \geq 0} \sum_{k \in \mathcal{S}} c_k \quad \text{s.t.} \quad M = \sum_{k \in \mathcal{S}} c_k \mathbf{z}_k \mathbf{z}_k^T.$$

- ▶ The **MixLasso** estimator:

$$\min_{M \in \mathbb{R}^{N \times N}} g(M) + \lambda \|M\|_{\mathcal{S}}.$$

- ▶ Equivalently, one can solve the estimator by

$$\min_{\mathbf{c} \in \mathbb{R}_+^{|\mathcal{S}|}} g\left(\sum_{k \in \mathcal{S}} c_k \mathbf{z}_k \mathbf{z}_k^T\right) + \lambda \|\mathbf{c}\|_1$$

Question: How to optimize with $|\mathcal{S}| = 2^N$ variables?

Greedy Coordinate Descent via MAX-CUT

- ▶ At each iteration, we find the **coordinate of steepest descent**:

$$j^* = \underset{j}{\operatorname{argmax}} -\nabla_j f(\mathbf{c}) = \underset{\mathbf{z} \in \{0,1\}^N}{\operatorname{argmax}} \langle -\nabla g(M), \mathbf{z} \mathbf{z}^T \rangle \quad (2)$$

which is a **Boolean Quadratic problem** similar to **MAX-CUT**:

$$\max_{\mathbf{z} \in \{0,1\}^N} \mathbf{z}^T C \mathbf{z}$$

- ▶ Can be solved to a **3/5-approximation** by rounding from a special type of **SDP** with $O(ND)$ iterative solver.

Active-Set Algorithm

0. $\mathcal{A} = \emptyset, \mathbf{c} = 0$.

for $t = 1 \dots T$ **do**

1. Find an **approximate greedy** atom $\mathbf{z} \mathbf{z}^T$ by MAX-CUT-like problem:

$$\max_{\mathbf{z} \in \{0,1\}^N} \langle -\nabla g(M), \mathbf{z} \mathbf{z}^T \rangle.$$

.

2. Add $\mathbf{z} \mathbf{z}^T$ to an **active set \mathcal{A}** .

3. **Refine $\mathbf{c}_{\mathcal{A}}$** via Proximal Gradient Method on:

$$\min_{\mathbf{c} \geq 0} g\left(\sum_{k \in \mathcal{A}} c_k \mathbf{z}_k \mathbf{z}_k^T\right) + \lambda \|\mathbf{c}\|_1$$

4. Eliminate $\{\mathbf{z}_k \mathbf{z}_k^T \mid c_k = 0\}$ from \mathcal{A} .
end for.

- ▶ Finding **approximate greedy** coordinate costs $O(ND)$ (via SDP).
- ▶ Evaluating $\nabla g(M)$: a **least-square problem** of cost $O(D^3 |\mathcal{A}|^3)$.
- ▶ Each iteration costs $\underbrace{O(ND)}_{\text{MAX-CUT}} + \underbrace{O(D^3 |\mathcal{A}|^3)}_{\text{Least-Square}}$

Non-linear Extension

- ▶ For $f_k(\mathbf{x})$ lies in RKHS \mathcal{H} generated by $\mathcal{K}(\cdot, \cdot)$, given $\{\mathbf{z}_i\}_{i=1}^N$:

$$\min_{f_k \in \mathcal{H}} \frac{1}{2N} \sum_{i=1}^N \left(y_i - \sum_{k=1}^K z_{ik} f_k(\mathbf{x}_i) \right)^2 + \frac{\tau}{2} \sum_{k=1}^K \|f_k\|_{\mathcal{H}}^2 \quad (3)$$

- ▶ **Representer Theorem** ensures an expression of the form

$$f_k^*(\mathbf{x}) = \sum_{i=1}^N \alpha_i z_{ik} \mathcal{K}(\mathbf{x}_i, \mathbf{x}), \quad k \in [K],$$

for the minimizers:

- ▶ Similar **MixLasso** estimator with

$$g(M) := \max_{\alpha \in \mathbb{R}^N} \frac{1}{2\tau} \text{tr}(\mathcal{D}(\alpha) Q \mathcal{D}(\alpha) M) - \sum_{i=1}^N L^*(y_i, -\alpha_i) \quad (4)$$

where $Q : N \times N$ is the kernel matrix with $Q_{ij} = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$.

Theoretical Results: Risk Bound

Let W^* be the minimizer of risk (1) with K components and $\|W^*\|_F \leq R$ and $\hat{W} = D(\sqrt{c_{\mathcal{A}}}) W$, then \hat{W} with probability $1 - \rho$ satisfies

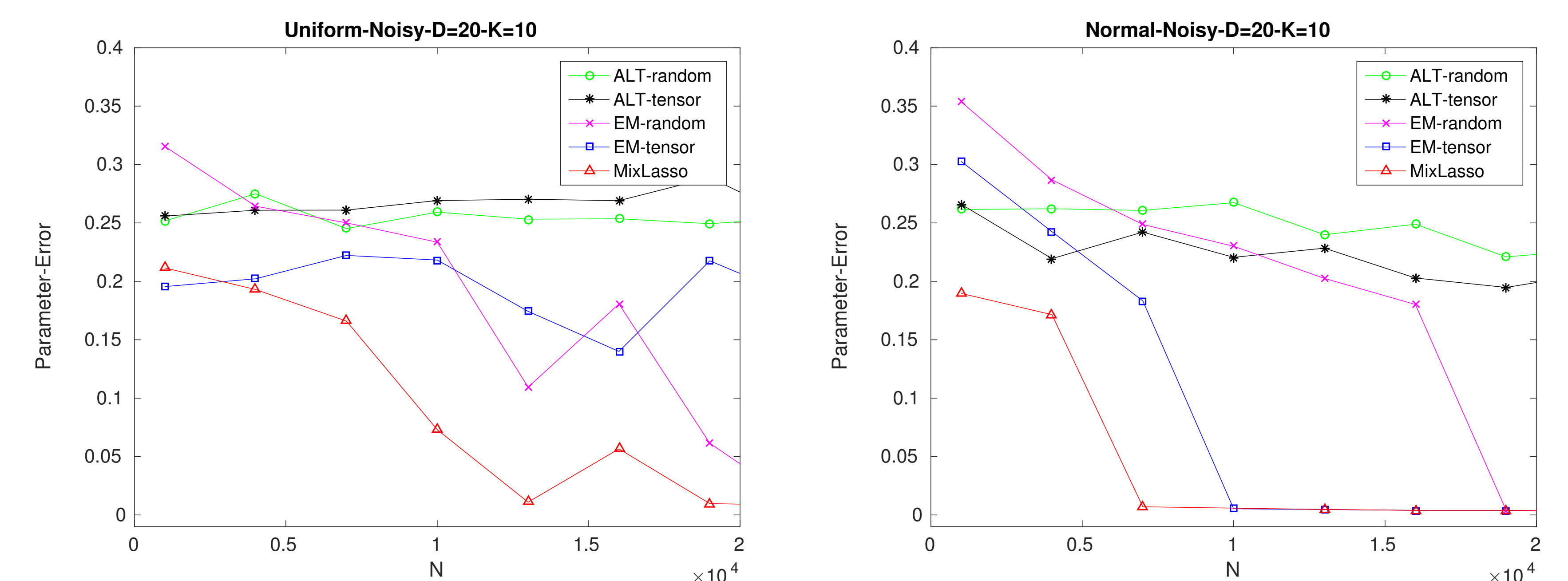
$$r(\hat{W}) \leq r(W^*) + \epsilon$$

as long as

$$t = \Omega\left(\frac{K}{\epsilon}\right) \quad \text{and} \quad N = \Omega\left(\frac{DK}{\epsilon^3} \log\left(\frac{RK}{\epsilon\rho}\right)\right).$$

- ▶ The number of output components $\hat{K} = O(K/\epsilon)$
- ▶ The result trades between **risk** and **sparsity**.
- ▶ **No assumption** on W except that of boundedness.
- ▶ The **sample complexity** is (quasi) linear to D and K .

Experiments on Synthetic data



Experiments on Stock Data

