

# Pathwise Coordinate Optimization for Nonconvex Sparse Learning

Tuo Zhao

<http://www.princeton.edu/~tuoz>

Department of Computer Science  
Johns Hopkins University

Mar. 25. 2015

# Collaborators

This is joint work with

- Prof. Han Liu at Princeton University,
- Prof. Tong Zhang at Rutgers University and Baidu,
- Xingguo Li at University of Minnesota.

Manuscript:

<http://arxiv.org/abs/1412.7477>

Software Package:

<http://cran.r-project.org/web/packages/picasso/>

# Outline

- Background
- Pathwise Coordinate Optimization
- Computational and Statistical Theories
- Numerical Simulations
- Conclusions

# Background

# Regularized M-Estimation

- Let  $\beta^*$  denote the parameter to be estimated. We solve the following regularized M-estimation problem

$$\min_{\beta \in \mathbb{R}^d} \underbrace{\mathcal{L}(\beta) + \mathcal{R}_\lambda(\beta)}_{\mathcal{F}_\lambda(\beta)},$$

where  $\mathcal{L}(\beta)$  is a smooth loss function, and  $\mathcal{R}_\lambda(\beta)$  is a regularization function with a tuning parameter  $\lambda$ .

- Examples: Lasso, Logistic Lasso (Tibshirani, 1996), Group Lasso (Yuan and Lin, 2006), Graphical Lasso (Yuan and Lin, 2007; Banerjee et al., 2008; Friedman et al. 2008), ...

# Regularization Functions

$\mathcal{R}_\lambda(\boldsymbol{\beta})$  is coordinate separable,

$$\mathcal{R}_\lambda(\boldsymbol{\beta}) = \sum_{j=1}^d r_\lambda(\beta_j).$$

$\mathcal{R}_\lambda(\boldsymbol{\beta})$  is decomposable,

$$\mathcal{R}_\lambda(\boldsymbol{\beta}) = \lambda \|\boldsymbol{\beta}\|_1 + \mathcal{H}_\lambda(\boldsymbol{\beta}) = \lambda \sum_{j=1}^d [|\beta_j| + h_\lambda(\beta_j)].$$

Examples: Smooth Clipped Absolute Deviation (SCAD, Fan and Li, 2001) and Minimax Concavity Penalty (MCP, Zhang, 2010)

# Regularization Functions

For any  $\gamma > 2$ , SCAD is defined as

$$\blacksquare r_\lambda(\beta_j) = \begin{cases} \lambda|\beta_j| & \text{if } |\beta_j| \leq \lambda, \\ \frac{|\beta_j|^2 - 2\lambda\gamma|\beta_j| + \lambda^2}{2(\gamma - 1)} & \text{if } \lambda < |\beta_j| \leq \lambda\gamma, \\ \frac{(\gamma + 1)\lambda^2}{2} & \text{if } |\beta_j| > \lambda\gamma. \end{cases}$$

$$\blacksquare h_\lambda(\beta_j) = \begin{cases} 0 & \text{if } |\beta_j| \leq \lambda, \\ \frac{2\lambda|\beta_j| - |\beta_j|^2 - \lambda^2}{2(\gamma - 1)} & \text{if } \lambda < |\beta_j| \leq \lambda\gamma, \\ \frac{(\gamma + 1)\lambda^2 - 2\lambda|\beta_j|}{2} & \text{if } |\beta_j| > \lambda\gamma. \end{cases}$$

# Regularization Functions

For any  $\gamma > 1$ , MCP is defined as

$$\begin{aligned} \blacksquare r_\lambda(\beta_j) &= \begin{cases} \lambda\left(|\beta_j| - \frac{|\beta_j|^2}{2\lambda\gamma}\right) & \text{if } |\beta_j| \leq \lambda\gamma, \\ \frac{\lambda^2\gamma}{2} & \text{if } |\beta_j| > \lambda\gamma. \end{cases} \\ \blacksquare h_\lambda(\beta_j) &= \begin{cases} -\frac{|\beta_j|^2}{2\gamma} & \text{if } |\beta_j| \leq \lambda\gamma, \\ \frac{\lambda^2\gamma - 2\lambda|\beta_j|}{2} & \text{if } |\beta_j| > \lambda\gamma. \end{cases} \end{aligned}$$



# Regularization Functions

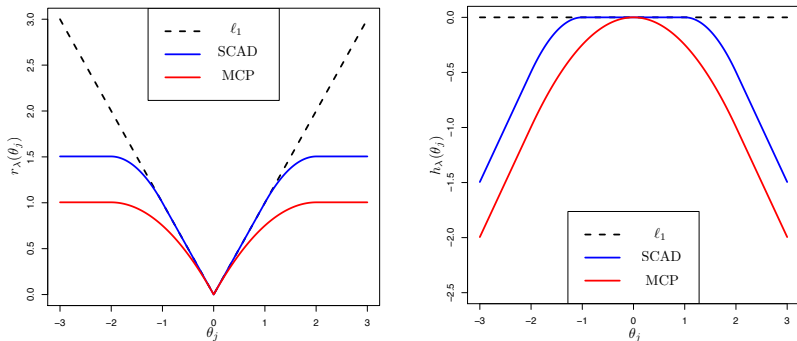


Figure: 1.  $\lambda = 1$  and  $\gamma = 2.01$ .

# Loss Functions

$\mathbf{X} \in \mathbb{R}^{n \times d}$  – design matrix,  $\mathbf{y} \in \mathbb{R}^n$  – response vector.

- Least Square Loss:

$$\mathcal{L}(\boldsymbol{\beta}) = \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2.$$

- Logistic Loss:

$$\mathcal{L}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \left( \log \left[ 1 + \exp(\mathbf{X}_{i*}^T \boldsymbol{\beta}) \right] - y_i \mathbf{X}_{i*}^T \boldsymbol{\beta} \right).$$

- Others: Huber Loss, Multi-category Logistic Loss,...

# Reformulation

We rewrite the regularized M-estimation problem as

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^d} \underbrace{\tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1}_{\mathcal{F}_\lambda(\boldsymbol{\beta})}.$$

- $\tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta})$  is smooth but nonconvex,

$$\tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) = \mathcal{L}(\boldsymbol{\beta}) + \mathcal{H}_\lambda(\boldsymbol{\beta}).$$

- $\lambda \|\boldsymbol{\beta}\|_1$  is nonsmooth but convex.

**Remark:** Amenable to theoretical analysis.

# Randomized Coordinate Descent Algorithm

At the  $t$ -th iteration, we randomly select a coordinate  $j$  from  $d$  coordinates. We then take  $\beta_{\setminus j}^{(t+1)} \leftarrow \beta_{\setminus j}^{(t)}$ , and

- Exact Coordinate Minimization (Fu, 1998)

$$\beta_j^{(t+1)} \leftarrow \arg \min_{\beta_j} \tilde{\mathcal{L}}_\lambda(\beta_j; \beta_{\setminus j}^{(t)}) + \lambda|\beta_j|.$$

- Inexact Coordinate Minimization (Shalev-Shwartz, 2011)

$$\beta_j^{(t+1)} \leftarrow \arg \min_{\beta_j} (\beta_j - \beta^{(t)}) \nabla_j \tilde{\mathcal{L}}_\lambda(\beta^{(t)}) + \frac{L}{2} (\beta_j - \beta^{(t)})^2 + \lambda|\beta_j|,$$

where  $L$  is the step size parameter.

# Examples

- Sparse Linear Regression + MCP:

$$\mathcal{T}_{j,\lambda}(\boldsymbol{\beta}^{(t)}) = \begin{cases} \tilde{\beta}_j^{(t+1)} & \text{if } |\tilde{\beta}_j^{(t+1)}| \geq \gamma\lambda, \\ \frac{\mathcal{S}_\lambda(\tilde{\beta}_j^{(t+1)})}{1 - 1/\gamma} & \text{if } |\tilde{\beta}_j^{(t+1)}| < \gamma\lambda. \end{cases}$$

where  $\tilde{\beta}_j^{(t+1)} = \mathbf{X}_{*j}^T(\mathbf{y} - \mathbf{X}_{*\setminus j}\boldsymbol{\beta}_{\setminus j}^{(t)})/n$ .

- Sparse Logistic Regression + MCP:

$$\mathcal{T}_{j,\lambda}(\boldsymbol{\beta}^{(t)}) = \mathcal{S}_\lambda(\boldsymbol{\beta}^{(t)} - \nabla_j \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}^{(t)})/L)$$

**Remark:** Sublinear Convergence to Local Optima without Statistical Guarantees (Shalev-Shwartz, 2011).

# Pathwise Coordinate Optimization

# Pathwise Coordinate Optimization

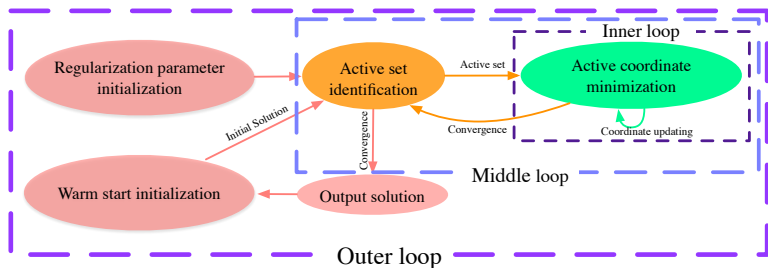
- Much faster than other competing algorithms.
- Very simple implementation.
- Easily scale to large problems.
- **NO** computational analysis in existing literature
- **NO** statistical guarantee on the obtained estimator.

## Our Contribution:

- The **FIRST** pathwise coordinate optimization algorithm with both computational and statistical guarantees.
- The **FIRST** two-step estimator with both computational and statistical guarantees.

# Pathwise Coordinate Optimization

Friedman et al. 2007, Mazumder et al. 2011



**Figure:** 2. The pathwise coordinate optimization framework contains 3 nested loops : (I) Warm start initialization; (II) Active set identification; (III) Active coordinate minimization.



# Restricted Strong Convexity and Smoothness

**Motivation:** For any  $\beta, \beta' \in \mathbb{R}^d$  such that  $|\{j \mid \beta_j \neq 0 \text{ or } \beta'_j \neq 0\}| \leq s$ , we have

- $\tilde{\mathcal{L}}_\lambda(\beta') - \tilde{\mathcal{L}}_\lambda(\beta) - (\beta' - \beta)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta) \geq \frac{C_-(s)}{2} \|\beta' - \beta\|_2^2,$
- $\tilde{\mathcal{L}}_\lambda(\beta') - \tilde{\mathcal{L}}_\lambda(\beta) - (\beta' - \beta)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta) \leq \frac{C_+(s)}{2} \|\beta' - \beta\|_2^2,$

where  $C_-(s), C_+(s) > 0$  are two constants depending on  $s$ .

**Remark:** An algorithm, which can maintain **SPARSE** solutions throughout all iterations, behaves like minimizing a **STRONGLY CONVEX** function. Therefore a linear convergence can be expected.

# Warm Start Initialization (Outer Loop)

- We choose a sequence of **DECREASING** regularization parameters  $\{\lambda_K\}_{K=1}^N$ :

$$\lambda_0 \geq \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{N-1} \geq \lambda_N.$$

- The algorithm yields a sequence of output solutions  $\{\hat{\boldsymbol{\beta}}^{[K]}\}_{K=0}^N$  from sparse to dense,

$$\hat{\boldsymbol{\beta}}^{[K]} \leftarrow \min_{\boldsymbol{\beta}} \tilde{\mathcal{L}}_{\lambda_K}(\boldsymbol{\beta}) + \lambda_K \|\boldsymbol{\beta}\|_1.$$

## Warm Start Initialization (Outer Loop)

- We choose  $\lambda_0 = \|\nabla \mathcal{L}(\mathbf{0})\|_\infty$ , then have

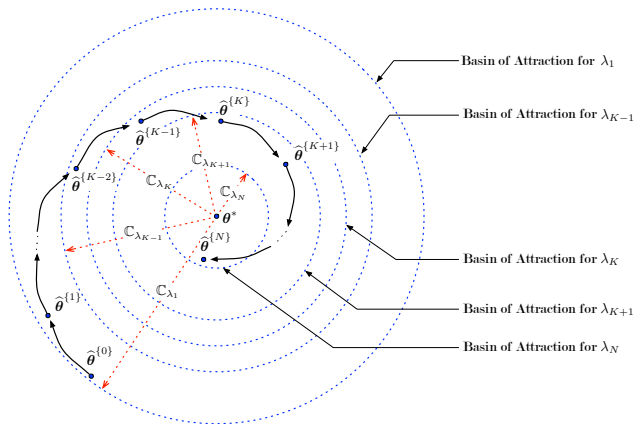
$$\min_{\xi \in \partial \|\mathbf{0}\|_1} \|\nabla \mathcal{L}(\mathbf{0}) + \nabla \mathcal{H}_\lambda(\mathbf{0}) + \lambda_0 \xi\|_\infty = 0 \quad \text{and} \quad \hat{\beta}^{\{0\}} = \mathbf{0}.$$

- The regularization sequence  $\{\lambda_K\}_{K=0}^N$  is geometrically decreasing

$$\lambda_K = \eta \lambda_{K-1} \quad \text{with} \quad \eta \in (0, 1).$$

- When solving the optimization problem with  $\lambda_K$ , we use  $\hat{\beta}^{\{K-1\}}$  as **INITIALIZATION**.

# Geometric Interpretation



**Figure:** 3. Large regularization parameters suppress the overselection of irrelevant variables  $\{j \mid \beta_j^* = 0\}$  and yields highly sparse solutions.

## Active Set Strategy (Friedman et al. 2007)

Define  $\mathcal{A} = \{j \mid \beta_j \neq 0\}$  be the set of indices of nonzero coordinates, and  $\overline{\mathcal{A}} = \{j \mid \beta_j = 0\}$  be the set of indices of zero coordinates. A naive updating scheme is:

- (1) **Active coordinate minimization:** Cyclically update  $\beta_j$ 's in  $\mathcal{A}$  until convergence.
- (2) **Sweeping coordinates:** Check all  $\beta_j$  in  $\mathcal{A}$ , and if any coordinate becomes zero, move it to  $\overline{\mathcal{A}}$ .
- (3) **Adding coordinates:** Update  $\beta_j$ 's over  $\overline{\mathcal{A}}$  for only once, and if any coordinate becomes nonzero, move it to  $\mathcal{A}$ . Then we go back to (1).

**Remark:** Heuristic tricks without theoretical guarantees.

## Active Set Identification (Middle Loop)

For notational simplicity, the outer loop index  $K$  is omitted.

### Greedy Selection:

At the  $m$ -th iteration, we have  $\beta^{[m]}$  and define

$$\mathcal{A}_m = \{j \mid \beta_j^{[m]} \neq 0\} \text{ and } \bar{\mathcal{A}}_m = \{j \mid \beta_j^{[m]} = 0\}.$$

- $\beta^{[m+0.5]} \leftarrow$  Active Coordinate Minimization over  $\mathcal{A}_m$ .
- $k_m \leftarrow \arg \max_{k \in \bar{\mathcal{A}}_m} |\nabla_k \tilde{\mathcal{L}}_\lambda(\beta^{[m+0.5]})|$ .
- $\beta_{k_m}^{[m+1]} \leftarrow \mathcal{T}_{k_m, \lambda}(\beta^{[m+0.5]})$  and  $\beta_{\setminus k_m}^{[m+1]} = \beta_{\setminus k_m}^{[m+0.5]}$ .

**Remark:** Conservative coordinate selection.

## Active Set Identification (Middle Loop)

For notational simplicity, the outer loop index  $K$  is omitted.

### Randomized Selection:

At the  $m$ -th iteration, we have  $\beta^{[m]}$  and define

$$\mathcal{A}_m = \{j \mid \beta_j^{[m]} \neq 0\} \text{ and } \bar{\mathcal{A}}_m = \{j \mid \beta_j^{[m]} = 0\}.$$

- $\beta^{[m+0.5]} \leftarrow$  Active Coordinate Minimization over  $\mathcal{A}_m$ .
- Randomly select  $k_m \in \bar{\mathcal{A}}_m$  such that  $|\nabla_k \tilde{\mathcal{L}}_\lambda(\beta^{[m+0.5]})| \geq \delta\lambda$ .
- $\beta_{k_m}^{[m+1]} \leftarrow \mathcal{T}_{k_m, \lambda}(\beta^{[m+0.5]})$  and  $\beta_{\setminus k_m}^{[m+1]} \leftarrow \beta_{\setminus k_m}^{[m+0.5]}$ .

**Remark:** Conservative coordinate selection.

## Active Set Identification (Middle Loop)

For notational simplicity, the outer loop index  $K$  is omitted.

### Truncated Cyclic Selection:

At the  $m$ -th iteration, we have  $\beta^{[m]}$  and define

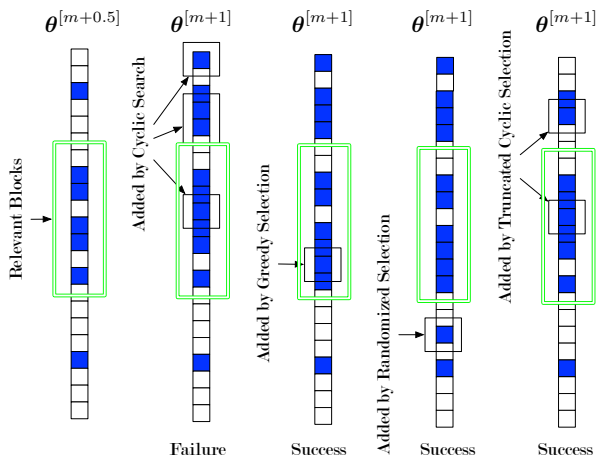
$$\mathcal{A}_m = \{j \mid \beta_j^{[m]} \neq 0\} \text{ and } \bar{\mathcal{A}}_m = \{j \mid \beta_j^{[m]} = 0\}.$$

- $\beta^{[m+0.5]} \leftarrow$  Active Coordinate Minimization over  $\mathcal{A}_m$ .
- For all  $k \in \bar{\mathcal{A}}_m$ , take
 
$$\beta_k^{[m+0.5]} \leftarrow \begin{cases} \mathcal{T}_{k,\lambda}(\beta^{[m+0.5]}) & \text{if } |\nabla_k \tilde{\mathcal{L}}_\lambda(\beta^{[m+0.5]})| \geq \delta\lambda, \\ \beta_k^{[m+0.5]} & \text{if } |\nabla_k \tilde{\mathcal{L}}_\lambda(\beta^{[m+0.5]})| < \delta\lambda. \end{cases}$$
- $\beta^{[m+1]} \leftarrow \beta^{[m+0.5]}$ .

**Remark:** Prevent from the overselection of irrelevant variables.

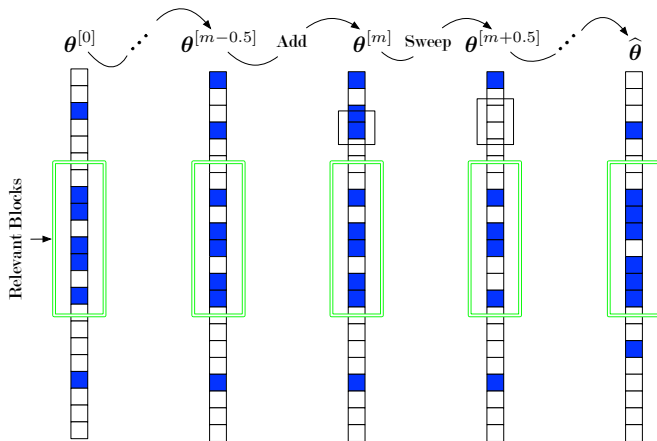


# Active Set Identification (Middle Loop)



**Figure:** 4. The cyclic search in Friedman et al. 2007, Mazumder et al. 2011 may overselect irrelevant variables.

# Active Coordinate Minimization (Inner Loop)



**Figure 5.** The active coordinate minimization not only decreases the objective value, but also sweeps some variables from the active set.

# Computational and Statistical Theories

# Preliminaries

**Definition [Sparse Eigenvalues]:** Given an integer  $s \geq 1$ ,

$$\rho_+(s) = \sup_{\|\mathbf{v}\|_0 \leq s} \frac{\mathbf{v}^T \nabla^2 \mathcal{L}(\boldsymbol{\beta}) \mathbf{v}}{\|\mathbf{v}\|_2^2}, \quad \rho_-(s) = \inf_{\|\mathbf{v}\|_0 \leq s} \frac{\mathbf{v}^T \nabla^2 \mathcal{L}(\boldsymbol{\beta}) \mathbf{v}}{\|\mathbf{v}\|_2^2},$$

and  $\tilde{\rho}_-(s) = \rho_-(s) - \alpha$ .

**Lemma [Restricted Curvature]:** Given  $\rho_+(s) > \rho_-(s) > \alpha$ , for any  $\boldsymbol{\beta}, \boldsymbol{\beta}' \in \mathbb{R}^d$  such that  $|\{j \mid \beta_j \neq 0 \text{ or } \beta'_j \neq 0\}| \leq s$ ,

$$\tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}') - \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) - (\boldsymbol{\beta}' - \boldsymbol{\beta})^T \nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) \leq \frac{\rho_+(s)}{2} \|\boldsymbol{\beta}' - \boldsymbol{\beta}\|_2^2,$$

$$\tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}') - \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) - (\boldsymbol{\beta}' - \boldsymbol{\beta})^T \nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) \geq \frac{\tilde{\rho}_-(s)}{2} \|\boldsymbol{\beta}' - \boldsymbol{\beta}\|_2^2,$$

where  $\mathcal{H}_\lambda(\boldsymbol{\beta}') - \mathcal{H}_\lambda(\boldsymbol{\beta}') - (\boldsymbol{\beta}' - \boldsymbol{\beta})^T \nabla \mathcal{H}_\lambda(\boldsymbol{\beta}) \geq -\frac{\alpha}{2} \|\boldsymbol{\beta}' - \boldsymbol{\beta}\|_2^2$ .

# Preliminaries

**Assumption A:**  $\lambda_N \geq 4 \|\nabla \mathcal{L}(\beta^*)\|_\infty$  and  $\eta \in [23/24, 1)$ .

The regularization parameters are **LARGE** enough to eliminate irrelevant variables (Negahban et al. 2012).

**Assumption B:** Given  $\|\beta^*\|_0 \leq s^*$ , there exists an  $\tilde{s}$  such that

- (1)  $\tilde{s} \geq (484\kappa^2 + 100\kappa)s^*$ ,
- (2)  $\rho_+(s^* + 2\tilde{s} + 2) < +\infty$ ,
- (3)  $\tilde{\rho}_-(s^* + 2\tilde{s} + 2) > 0$ ,

where  $\kappa = \rho_+(s^* + 2\tilde{s} + 2)/\tilde{\rho}_-(s^* + 2\tilde{s} + 2)$ .

The algorithm can tolerate **AT MOST**  $\tilde{s} + 1$  nonzero irrelevant variables throughout all iterations (Bickel, 2009; Zhang, 2009).

# Globally Convergence (Greedy-PICASSO)

Suppose that Assumptions A and B hold. We have the following results:

- **(Solution Sparsity)** Through all iterations of PICASSO, any solution  $\beta$  satisfies  $\|\beta_{\tilde{s}}\|_0 \leq \tilde{s} + 1$ .
- **(Sparse Optimum)** At the  $K$ -th iteration of the outer loop, PICASSO converges to a unique sparse local optimum  $\tilde{\beta}^{\lambda_K}$  satisfying
 
$$\|\tilde{\beta}_{\tilde{s}}^{\lambda_K}\|_0 \leq \tilde{s} \text{ and } \min_{\xi \in \partial \|\tilde{\beta}^{\lambda_K}\|_1} \|\nabla \tilde{\mathcal{L}}_{\lambda_K}(\tilde{\beta}^{\lambda_K}) + \lambda \xi\|_{\infty} = 0.$$
- **(Logarithm Iteration Complexity)** To attain  $\mathcal{F}_{\lambda_N}(\hat{\beta}^{\{N\}}) - \mathcal{F}_{\lambda_N}(\tilde{\beta}^{\lambda_N}) \leq \epsilon$ , the number of active set identification iterations is at most  $\mathcal{O}(N \cdot \log(1/\epsilon))$ .

## Two-step Method

- **Step 1 – Convex Relaxation:** Obtain  $\beta^{\text{relax}}$  satisfying

$$\min_{\xi \in \partial \|\beta^{\text{relax}}\|_1} \left\| \nabla \mathcal{L}(\beta^{\text{relax}}) + \lambda_0 \xi \right\|_{\infty} \leq \frac{\lambda_0}{8}.$$

- **Step 2 – PICASSO:** Solve the optimization problem with PICASSO, and use  $\beta^{\text{relax}}$  as initialization for  $\lambda_0$ .

**Remark:** The low precision makes Step 1 very efficient.

**Remark:** The restricted strong convexity holds for  $\|\beta - \beta^*\|_2 \leq R$  (e.g. logistic loss, huber loss), where  $R$  is a constant and does not scale with  $(n, d, s^*)$ . All previous theoretical results hold.

# Nearly Unbiased Estimation

Suppose that Assumptions A and B. We have

$$\left\| \hat{\beta}^{\{N\}} - \beta^* \right\|_2 = \mathcal{O} \left( \underbrace{\frac{\|\nabla_{\mathcal{S}_1} \mathcal{L}(\beta^*)\|_2}{\tilde{\rho}_-(s^* + 2\tilde{s})}}_{\text{Strong Signals}} + \underbrace{\frac{\lambda_N \sqrt{|\mathcal{S}_2|}}{\tilde{\rho}_-(s^* + \tilde{s})}}_{\text{Weak Signals}} \right),$$

where  $\mathcal{S}_1 = \{j \mid |\beta_j^*| \geq \gamma \lambda_N\}$  and  $\mathcal{S}_2 = \{j \mid 0 < |\beta_j^*| < \gamma \lambda_N\}$ .

**Clarification:** To establish the theoretical analysis for each individual problem, we need to assume that the model is **CORRECTLY** specified. This is a very common assumption in high dimensional statistical theories.



# Model Specification

**Sparse Linear Regression:** We consider a linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\epsilon},$$

where  $\boldsymbol{\epsilon} \sim N(0, \sigma^2 \mathbf{I}_n)$  is the observational noise vector.

**Sparse Logistic Regression:** We consider a logistic model

$$y_i \sim \text{Bernoulli} \left( \frac{\exp(\mathbf{X}_{i*}^T \boldsymbol{\beta}^*)}{1 + \exp(\mathbf{X}_{i*}^T \boldsymbol{\beta}^*)} \right)$$

for  $i = 1, \dots, n$ .

# Application to Sparse Linear Regression

**Verify Assumption A:** Given  $\lambda_N = 8\sigma \sqrt{\log d/n}$ , with **HIGH PROBABILITY**, we have

$$\lambda_N \geq 4 \|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_\infty.$$

**Verify Assumption B:** Suppose that each column of  $\mathbf{X}$  is independently sampled from a sub-Gaussian distribution with mean  $\mathbf{0}$  and covariance  $\boldsymbol{\Sigma}$ , where  $\Lambda_{\min}(\boldsymbol{\Sigma}) \geq \psi_{\min}$  and  $\Lambda_{\max}(\boldsymbol{\Sigma}) \leq \psi_{\max}$ . Given  $\alpha = \psi_{\min}/4$ , there exists an  $\tilde{s}$  such that for large enough  $n$ , **HIGH PROBABILITY**, we have

- (1)  $\tilde{s} \geq [484\kappa^2 + 100\kappa] \cdot s^*$ ,
- (2)  $\tilde{\rho}_-(s^* + 2\tilde{s} + 2) \geq \psi_{\min}/4$ ,
- (3)  $\rho_+(s^* + 2\tilde{s} + 2) \leq 3\psi_{\max}/2$ .

# Application to Sparse Linear Regression

## Parameter Estimation:

Given  $\alpha = \psi_{\min}/4$  and  $\lambda_N = 8\sigma \sqrt{\log d/n}$ , we have

$$\left\| \widehat{\boldsymbol{\beta}}^{\{N\}} - \boldsymbol{\beta}^* \right\|_2 = \mathcal{O}_P \left( \underbrace{\sigma \sqrt{\frac{s_1^*}{n}}}_{\text{Strong Signals}} + \underbrace{\sigma \sqrt{\frac{s_2^* \log d}{n}}}_{\text{Weak Signals}} \right),$$

where  $s_1^* = \{j \mid |\beta_j^*| \geq \gamma \lambda_N\}$  and  $s_2^* = \{j \mid 0 < |\beta_j^*| < \gamma \lambda_N\}$ .

## MCP v.s. $\ell_1$ :

$$\left\| \widehat{\boldsymbol{\beta}}^{\ell_1} - \boldsymbol{\beta}^* \right\|_2 = \mathcal{O}_P \left( \sigma \sqrt{\frac{s^* \log d}{n}} \right).$$

# Application to Sparse Linear Regression

**Minimum Signal Strength:**  $\min_{j \in \mathcal{S}} |\beta_j^*| \geq \frac{C' \sigma}{\psi_{\min}} \sqrt{\frac{\log d}{n}}$ .

**Support Recovery:**

Given  $\alpha = \psi_{\min}/4$  and  $\lambda_N = 8\sigma \sqrt{\log d/n}$ , we have

$$\bar{\beta}^{\lambda_N} = \arg \min_{\beta} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \quad \text{subject to } \beta_{\bar{\mathcal{S}}} = \mathbf{0}$$

with high probability.

**MCP v.s.  $\ell_1$ :** Restricted Strong Convexity v.s. Irrepresentability.

# Application to Sparse Logistic Regression

**Verify Assumption A:** Given  $\lambda_N = 8 \sqrt{\log d/n}$ , with high probability we have

$$\lambda_N \geq 4 \|\nabla \mathcal{L}(\beta^*)\|_\infty.$$

**Verify Assumption B:** Suppose that each column of  $\mathbf{X}$  is independently sampled from a sub-Gaussian distribution with mean  $\mathbf{0}$  and covariance  $\Sigma$ , where  $\Lambda_{\min}(\Sigma) \geq \psi_{\min}$  and  $\Lambda_{\max}(\Sigma) \leq \psi_{\max}$ . Given  $\alpha = \psi_{\min}/4$ , there exists an  $\tilde{s}$  such that for large enough  $n$  and any  $\|\beta - \beta^*\| \leq R$ , with high probability, we have

- (1)  $\tilde{s} \geq [484\kappa^2 + 100\kappa] \cdot s^*$ ,
- (2)  $\tilde{\rho}_-(s^* + 2\tilde{s} + 2) \geq \psi_{\min}/4$ ,
- (3)  $\rho_+(s^* + 2\tilde{s} + 2) \leq 3\psi_{\max}/2$ .

# Application to Sparse Logistic Regression

## Parameter Estimation:

Given  $\alpha = \psi_{\min}/4$  and  $\lambda_N = 8\sqrt{\log d/n}$ , we have

$$\left\| \widehat{\boldsymbol{\beta}}^{\{N\}} - \boldsymbol{\beta}^* \right\|_2 = \mathcal{O}_P \left( \underbrace{\sqrt{\frac{s_1^*}{n}}}_{\text{Strong Signals}} + \underbrace{\sqrt{\frac{s_2^* \log d}{n}}}_{\text{Weak Signals}} \right),$$

where  $s_1^* = \{j \mid |\beta_j^*| \geq \gamma \lambda_N\}$  and  $s_2^* = \{j \mid 0 < |\beta_j^*| < \gamma \lambda_N\}$ .

## MCP v.s. $\ell_1$ :

$$\left\| \widehat{\boldsymbol{\beta}}^{\ell_1} - \boldsymbol{\beta}^* \right\|_2 = \mathcal{O}_P \left( \sqrt{\frac{s^* \log d}{n}} \right).$$

# Numerical Simulations

# Numerical Simulations

- PICASSO with Greedy selection, denoted by “G-PICASSO”.
- PICASSO with Randomized selection, denoted by “R-PICASSO”.
- PICASSO with Truncated Cyclic selection, denoted by “TC-PICASSO”.
- SPARSENET proposed in Mazumder et al. 2011.
- PISTA proposed in Wang et al. 2014.



# Numerical Simulations

**Table:** 1. Quantitative comparison on sparse linear regression  
 ( $N = 100, n = 60, d = 1000, \sigma = 1, \lambda_N = 0.25 \sqrt{\log d/n}, \gamma = 1.05$ ).

Method	$\ \hat{\beta} - \beta^*\ _2$	$\ \hat{\beta}_S\ _0$	$\ \hat{\beta}_{S^c}\ _0$	Correct Selection	Timing
G-PICASSO	<b>0.8003</b> (0.8908)	<b>2.812</b> (0.4997)	<b>0.844</b> (2.066)	<b>666/1000</b>	0.0169(0.0027)
R-PICASSO	0.8102(0.9663)	2.791(0.5355)	0.902(2.353)	653/1000	0.0186(0.0034)
TC-PICASSO	0.8057(0.8374)	2.800(0.4839)	0.888(2.038)	645/1000	0.0167(0.0024)
SPARSENET	1.1260(1.2708)	2.669(0.6942)	1.678(3.191)	514/1000	0.0171(0.0025)
PISTA	0.8135(0.8998)	2.797(0.5115)	0.881(2.112)	664/1000	2.1771(0.3805)

## Conclusions

# Conclusions

- Multistage convex relaxation (Zhang, 2010; Zhang 2012):  
No theoretical guarantee on the iteration complexity;  
Needs to be combined with an efficient solver for each subproblem.
- One-step convex relaxation method (Zou and Li, 2008; Wang and Li, 2013; Fan et al. 2014): Attains suboptimal statistical rates of convergence; Requires a stronger minimum signal strength assumption; Needs to be combined with an efficient solver for each subproblem.
- Path-following proximal gradient algorithm (Wang et al. 2014): Worse empirical computational performance; Requires  $\|\beta^*\|_2 \leq R/2$  for sparse generalized linear model estimation.

# Conclusions

- Proximal gradient algorithm (Loh and Wainwright, 2013):  
Solves

$$\min_{\beta \in \mathbb{R}^d} \mathcal{L}(\beta) + \mathcal{R}_\lambda(\beta) \quad \text{subject to } \|\beta\|_1 \leq R/2. \quad (1)$$

Sophisticated parameter tuning; Inexact convergence;

Slower parameter estimation rates of convergence;

Requires  $\|\beta^*\|_1 \leq R/2$  for all nonconvex sparse learning problems.

- Pathwise Calibrated Sparse Shooting Algorithm: Concrete theoretical guarantees; Empirically very efficient; Weaker Assumptions.

**Thank You! Questions?**