

Trimming the ℓ_1 Regularizer:

Statistical Analysis, Optimization, and Applications to Deep Learning

Jihun Yun¹, Peng Zheng², Eunho Yang^{1,3}, Aurélie C. Lozano⁴, Aleksandr Aravkin²

¹KAIST ²University of Washington ³AITRICS ⁴IBM T.J. Watson Research Center

arcprime@kaist.ac.kr

International Conference on Machine Learning
June 12, 2019

Table of Contents

- 1 Introduction and Setup
- 2 Statistical Analysis
- 3 Optimization
- 4 Experiments & Applications to Deep Learning

Table of Contents

1 Introduction and Setup

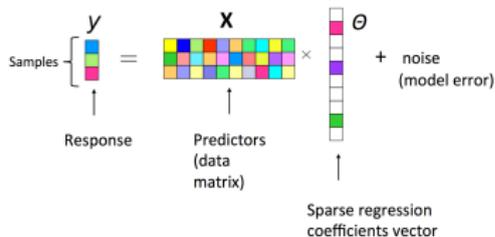
2 Statistical Analysis

3 Optimization

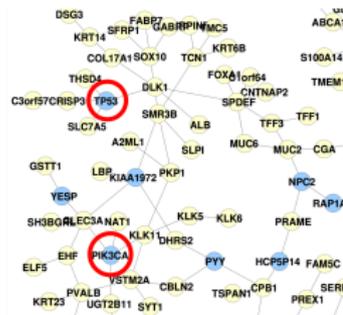
4 Experiments & Applications to Deep Learning

ℓ_1 Regularization is Popular

- High-dimensional data with ℓ_1 regularization ($n \ll p$)
 - Genomic Data, Matrix Completion, Deep Learning, etc.



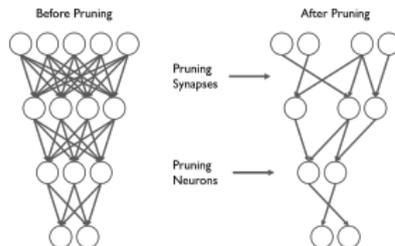
(a) Sparse linear models



(b) Sparse graphical models

	Movies									
	1	2	3	4	5	6	7	8	9	10
1		4		2	4					
2	3		3	1			3		3	
3		3	2	3						4
4			2			4		1	2	5
5	3			3				1		
6			2						3	2

(c) Matrix Completion



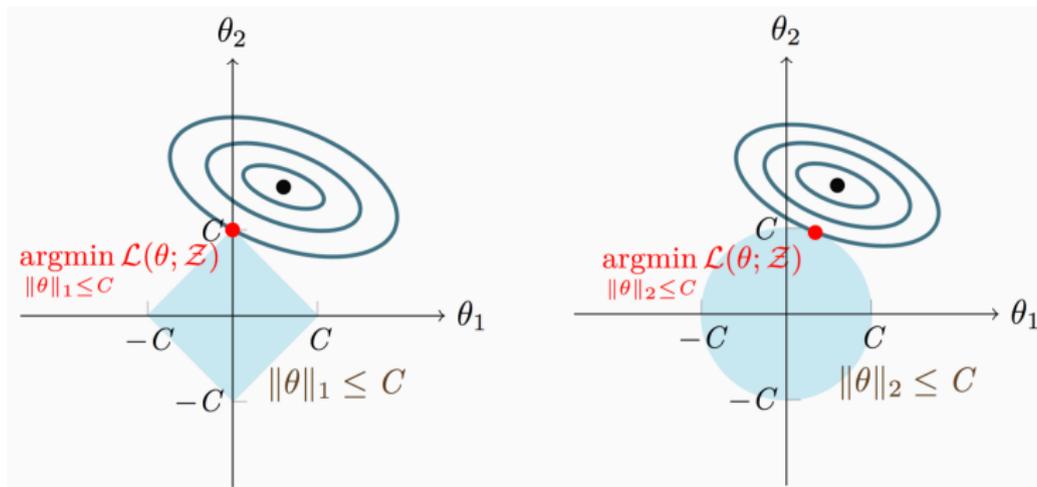
(d) Sparse neural networks

Concrete Example 1

Lasso

Example 1: Lasso* (Sparse Linear Regression)

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \Omega} \frac{1}{2n} \|y - X\theta\|_2^2 + \lambda_n \|\theta\|_1$$



*R. Tibshirani. Regression shrinkage and selection via the lasso. JRSS, Series B, 1996.

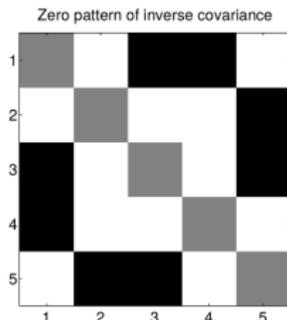
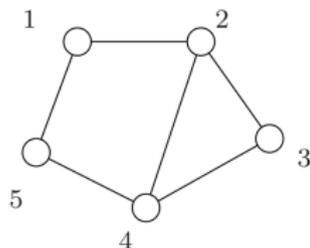
Concrete Example 2

Graphical Lasso

Example 2: Graphical Lasso* (Sparse Concentration Matrix)

$$\hat{\Theta} \in \operatorname{argmin}_{\Theta \in \mathcal{S}_{++}^p} \operatorname{trace}(\hat{\Sigma}\Theta) - \log \det(\Theta) + \lambda_n \|\Theta\|_{1,\text{off}}$$

where $\hat{\Sigma}$ is a sample covariance matrix, \mathcal{S}_{++}^p the symmetric and strictly positive definite matrices, and $\|\Theta\|_{1,\text{off}}$ the ℓ_1 -norm on the off-diagonal elements of Θ .



*P. Ravikumar, M. J. Wainwright, G. Raskutti, and B. Yu. High-dimensional covariance estimation by minimizing l1-penalized log-determinant divergence. EJS, 2011

Concrete Example 3

Group ℓ_1 on Network Pruning Task

Example 3: Group ℓ_1^* (Structured Sparsity of Weight Parameters)

$$\hat{\theta} \in \underset{\theta \in \Omega}{\operatorname{argmin}} \mathcal{L}(\theta; \mathcal{D}) + \lambda_n \|\theta\|_g$$

where $\hat{\theta}$ is a collection of weight parameters of neural networks, \mathcal{L} the neural network loss (ex. softmax), and $\|\theta\|_g$ the group sparsity regularizer.

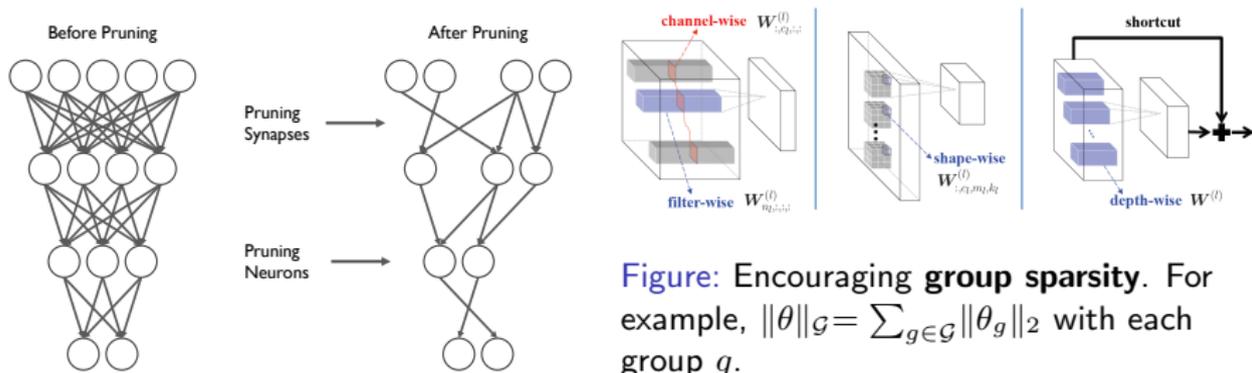


Figure: Encouraging group sparsity. For example, $\|\theta\|_g = \sum_{g \in \mathcal{G}} \|\theta_g\|_2$ with each group g .

*W. Wen, C. Wu, Y. Wang, Y. Chen, and H. Li. Learning Structured Sparsity in Deep Neural Networks. NIPS, 2016

Shrinkage Bias of Standard ℓ_1 Penalty

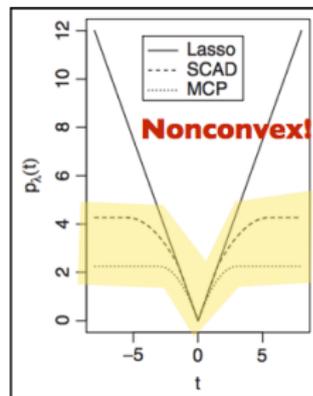
- As parameter size gets larger, the shrinkage bias effect also tends to be larger.
 - The ℓ_1 penalty is proportional to the size of parameters.

Despite the popularity of ℓ_1 penalty
(and also strong statistical guarantees),
Is it **really good enough?**

Non-convex Regularizers

Previous Work

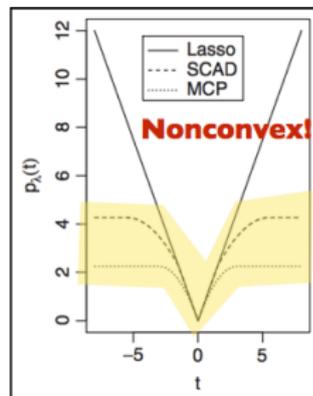
- For *amenable* non-convex regularizers (such as SCAD* and MCP**),
 - ▷ *Amenable regularizer*: Resembles ℓ_1 at the origin and has vanishing derivatives at the tail.
→ *coordinate-wise decomposable*.
 - ▷ (Loh & Wainwright)*** provide the statistical analysis on amenable regularizers.



Non-convex Regularizers

Previous Work

- For *amenable* non-convex regularizers (such as SCAD* and MCP**),
 - ▷ *Amenable regularizer*: Resembles ℓ_1 at the origin and has vanishing derivatives at the tail. \rightarrow *coordinate-wise decomposable*.
 - ▷ (Loh & Wainwright)*** provide the statistical analysis on amenable regularizers.



What about **more structurally complex** regularizer?

* J. Fan and R. Li. Variable selection via non-concave penalized likelihood and its oracle properties. *Jour. Amer. Stat. Ass.*, 96(456):1348-1360, December 2001.

** Cun-Hui Zhang et al. Nearly unbiased variable selection under minimax concave penalty. *The Annals of statistics*, 38(2):894-942, 2010.

*** P. Loh and M. J. Wainwright. Regularized M -estimators with non-convexity: statistical and algorithmic theory for local optima and algorithmic. *JMLR*, 2015.

*** P. Loh and M. J. Wainwright. Support recovery without incoherence: A case for nonconvex regularization. *The Annals of Statistics*, 2017.

Trimmed ℓ_1 Penalty

Definition

- In this paper, we study the **Trimmed ℓ_1 penalty**.
 - New class of regularizers.

Trimmed ℓ_1 Penalty

Definition

- In this paper, we study the **Trimmed ℓ_1 penalty**.
 - New class of regularizers.
- **Definition:**
For a parameter vector $\theta \in \mathbb{R}^p$, we only ℓ_1 -**penalize** each entry **except largest h entries** (We call h the trimming parameter).

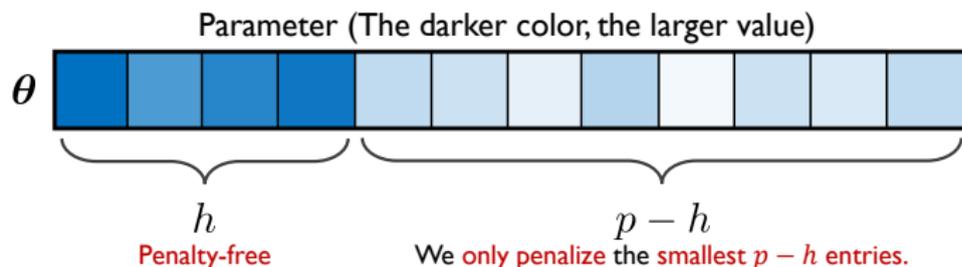
Trimmed ℓ_1 Penalty

Definition

- In this paper, we study the **Trimmed ℓ_1 penalty**.
 - New class of regularizers.

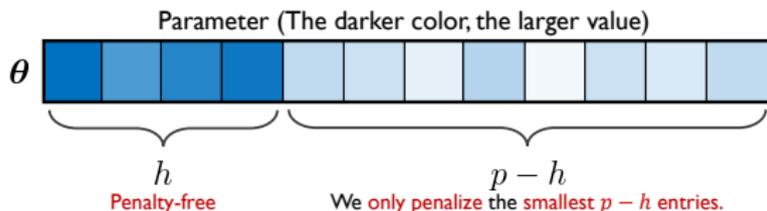
- **Definition:**

For a parameter vector $\theta \in \mathbb{R}^p$, we only ℓ_1 -**penalize** each entry **except largest h entries** (We call h the trimming parameter).



Trimmed ℓ_1 Penalty

First Formulation



- We can formalize by defining the order statistics of the parameter vector $|\theta_{(1)}| > |\theta_{(2)}| > \dots > |\theta_{(p)}|$, the M -estimation with the Trimmed ℓ_1 penalty is

$$\underset{\theta \in \Omega}{\text{minimize}} \mathcal{L}(\theta; \mathcal{D}) + \lambda_n \mathcal{R}(\theta; h)$$

where the regularizer $\mathcal{R}(\theta; h) = \sum_{j=h+1}^p |\theta_{(j)}|$ (sum of smallest $p - h$ entries in absolute values).

- Importantly, the Trimmed ℓ_1 is **not amenable nor coordinate-wise separable**.

M -estimation with the Trimmed ℓ_1 penalty

Second Formulation

- We can rewrite the M -estimation with the Trimmed ℓ_1 penalty by introducing additional variable \mathbf{w} :

$$\underset{\boldsymbol{\theta} \in \Omega, \mathbf{w} \in [0,1]^p}{\text{minimize}} \quad \mathcal{F}(\boldsymbol{\theta}, \mathbf{w}) := \mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) + \lambda_n \sum_{j=1}^p w_j |\theta_j|$$

$$\text{such that } \mathbf{1}^T \mathbf{w} \geq p - h$$

M -estimation with the Trimmed ℓ_1 penalty

Second Formulation

- We can rewrite the M -estimation with the Trimmed ℓ_1 penalty by introducing additional variable \mathbf{w} :

$$\underset{\boldsymbol{\theta} \in \Omega, \mathbf{w} \in [0,1]^p}{\text{minimize}} \quad \mathcal{F}(\boldsymbol{\theta}, \mathbf{w}) := \mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) + \lambda_n \sum_{j=1}^p w_j |\theta_j|$$

such that $\mathbf{1}^T \mathbf{w} \geq p - h$

- The variable \mathbf{w} encodes the **sparsity pattern** and **order information** of $\boldsymbol{\theta}$.
As an ideal case,
 - $w_j = 0$ for largest h entries
 - $w_j = 1$ for smallest $p - h$ entries

M -estimation with the Trimmed ℓ_1 penalty

Second Formulation

- We can rewrite the M -estimation with the Trimmed ℓ_1 penalty by introducing additional variable \mathbf{w} :

$$\underset{\boldsymbol{\theta} \in \Omega, \mathbf{w} \in [0,1]^p}{\text{minimize}} \quad \mathcal{F}(\boldsymbol{\theta}, \mathbf{w}) := \mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) + \lambda_n \sum_{j=1}^p w_j |\theta_j|$$

such that $\mathbf{1}^T \mathbf{w} \geq p - h$

- The variable \mathbf{w} encodes the **sparsity pattern** and **order information** of $\boldsymbol{\theta}$.
As an ideal case,
 - $w_j = 0$ for largest h entries
 - $w_j = 1$ for smallest $p - h$ entries
- If we set the trimming parameter $h = 0$, it is just a standard ℓ_1 .

M-estimation with the Trimmed ℓ_1 penalty

Second Formulation: Important Properties

$$\begin{aligned} & \underset{\boldsymbol{\theta} \in \Omega, \mathbf{w} \in [0,1]^p}{\text{minimize}} \quad \mathcal{F}(\boldsymbol{\theta}, \mathbf{w}) := \mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) + \lambda_n \sum_{j=1}^p w_j |\theta_j| \\ & \text{such that } \mathbf{1}^T \mathbf{w} \geq p - h \end{aligned}$$

- The objective function \mathcal{F} is
 - Weighted ℓ_1 -regularized if we fix \mathbf{w} .
 - Linear in \mathbf{w} with fixing $\boldsymbol{\theta}$.
 - However, \mathcal{F} is **non-convex** in jointly $(\boldsymbol{\theta}, \mathbf{w})$ because of **coupling of $\boldsymbol{\theta}$ and \mathbf{w}** .
- We use this second formulation for an **optimization**.
 - Since we don't need to sort the parameter.

Trimmed ℓ_1 Penalty

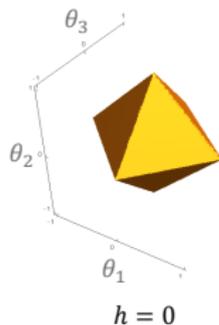
Unit Balls Visualization

- **Trimmed ℓ_1 Unit balls** of $\theta = (\theta_1, \theta_2, \theta_3)$ in the 3-dimensional space.

Trimmed ℓ_1 Penalty

Unit Balls Visualization

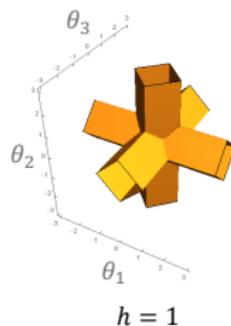
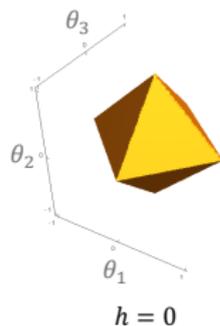
- **Trimmed ℓ_1 Unit balls** of $\theta = (\theta_1, \theta_2, \theta_3)$ in the 3-dimensional space.



Trimmed ℓ_1 Penalty

Unit Balls Visualization

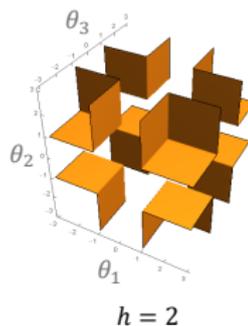
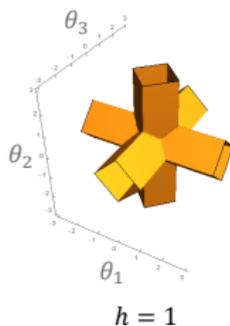
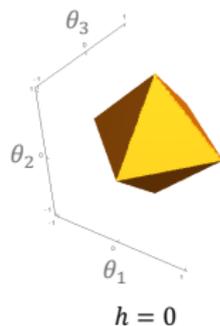
- **Trimmed ℓ_1 Unit balls** of $\theta = (\theta_1, \theta_2, \theta_3)$ in the 3-dimensional space.



Trimmed ℓ_1 Penalty

Unit Balls Visualization

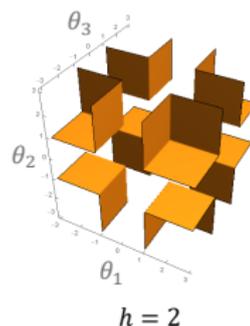
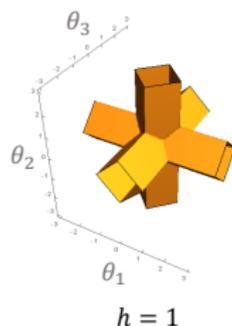
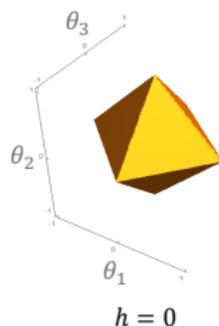
- **Trimmed ℓ_1 Unit balls** of $\theta = (\theta_1, \theta_2, \theta_3)$ in the 3-dimensional space.



Trimmed ℓ_1 Penalty

Unit Balls Visualization

- **Trimmed ℓ_1 Unit balls** of $\theta = (\theta_1, \theta_2, \theta_3)$ in the 3-dimensional space.



- For $h = 0$, the shape is the same as standard ℓ_1 unit ball.
- For $h > 0$, the penalty could be unbounded.
 - Since the largest h entries are not penalized, the unit ball could extend to infinity in these directions.
 - As h increases, the penalty would be more complicated.

Table of Contents

1 Introduction and Setup

2 Statistical Analysis

3 Optimization

4 Experiments & Applications to Deep Learning

Statistical Analysis: Key Assumptions and Quantity

Assumptions:

(C1) The loss \mathcal{L} is differentiable and convex.

(C2) **Restricted Strong Convexity on θ** : Let \mathbb{D} be the set of all possible error vectors for θ . Then, for all $\theta - \theta^* \in \mathbb{D}$,

$$\langle \nabla \mathcal{L}(\theta^*, \Delta) - \nabla \mathcal{L}(\theta^*), \Delta \rangle \geq \kappa_l \|\Delta\|_2^2 - \tau_1 \frac{\log p}{n} \|\Delta\|_1^2,$$

where κ_l is a “curvature” parameter, and τ_1 a “tolerance”.

- Allowing a small loss difference to be translated to a small error $\theta - \theta^*$.
- **RSC condition** is a **standard one** in this line of work.

Quantity:

- Let $\hat{Q} = \int_0^1 \nabla^2 \mathcal{L}(\theta^* + t(\hat{\theta} - \theta^*)) dt$.

Statistical Analysis

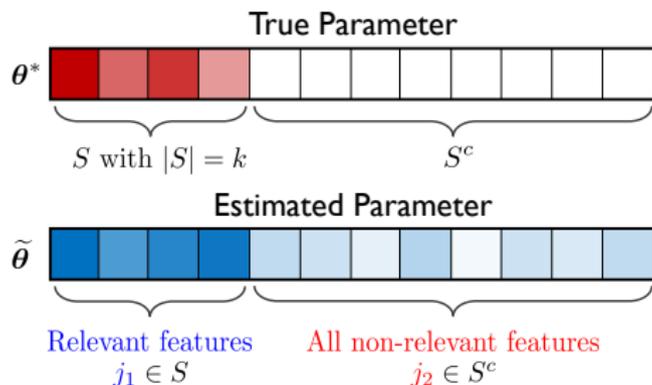
Theorem 1: General ℓ_∞ -error Bound and Variable Selection

- Consider an M -estimation problem with the Trimmed ℓ_1 penalty.
- Under (C1)&(C2) and standard conditions, for *any* local minimum $\tilde{\theta}$, we have
 - ① For every pair $j_1 \in S$, $j_2 \in S^c$, we have $|\tilde{\theta}_{j_1}| > |\tilde{\theta}_{j_2}|$

Statistical Analysis

Theorem 1: General ℓ_∞ -error Bound and Variable Selection

- Consider an M -estimation problem with the Trimmed ℓ_1 penalty.
- Under (C1)&(C2) and standard conditions, for *any* local minimum $\tilde{\theta}$, we have
 - 1 For every pair $j_1 \in S$, $j_2 \in S^c$, we have $|\tilde{\theta}_{j_1}| > |\tilde{\theta}_{j_2}|$



Statistical Analysis

Theorem 1: General ℓ_∞ -error Bound and Variable Selection

- 1 For every pair $j_1 \in S$, $j_2 \in S^c$, we have $|\tilde{\theta}_{j_1}| > |\tilde{\theta}_{j_2}|$
- 2 If $h < k$, all $j \in S^c$ are successfully estimated as zero and

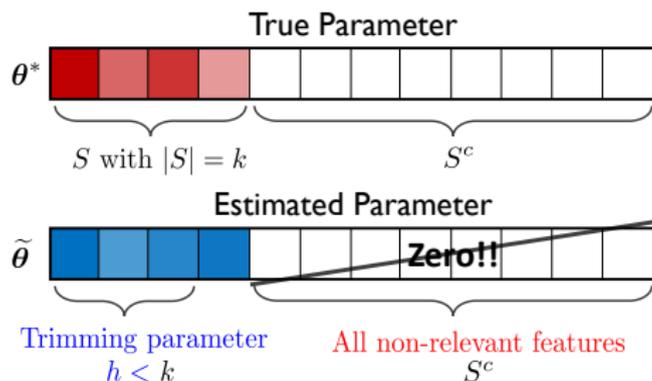
$$\|\tilde{\theta} - \theta^*\|_\infty \leq \left\| (\hat{Q}_{SS})^{-1} \nabla \mathcal{L}(\theta^*)_S \right\|_\infty + \lambda_n \left\| (\hat{Q}_{SS})^{-1} \right\|_\infty$$

Statistical Analysis

Theorem 1: General ℓ_∞ -error Bound and Variable Selection

- 1 For every pair $j_1 \in S$, $j_2 \in S^c$, we have $|\tilde{\theta}_{j_1}| > |\tilde{\theta}_{j_2}|$
- 2 If $h < k$, all $j \in S^c$ are successfully estimated as zero and

$$\|\tilde{\theta} - \theta^*\|_\infty \leq \left\| (\hat{Q}_{SS})^{-1} \nabla \mathcal{L}(\theta^*)_S \right\|_\infty + \lambda_n \left\| (\hat{Q}_{SS})^{-1} \right\|_\infty$$



Statistical Analysis

Theorem 1: General ℓ_∞ -error Bound and Variable Selection

- 1 For every pair $j_1 \in S$, $j_2 \in S^c$, we have $|\tilde{\theta}_{j_1}| > |\tilde{\theta}_{j_2}|$
- 2 If $h < k$, all $j \in S^c$ are successfully estimated as zero and

$$\|\tilde{\theta} - \theta^*\|_\infty \leq \left\| (\hat{Q}_{SS})^{-1} \nabla \mathcal{L}(\theta^*)_S \right\|_\infty + \lambda_n \left\| (\hat{Q}_{SS})^{-1} \right\|_\infty$$

- 3 If $h \geq k$, at least the smallest (in absolute) $p - h$ entries in S^c are exactly zero and $\|\tilde{\theta} - \theta\|_\infty \leq \left\| (\hat{Q}_{\hat{U}\hat{U}})^{-1} \nabla \mathcal{L}(\theta^*)_{\hat{U}} \right\|_\infty$ where \hat{U} is defined as the h largest absolute entries of $\tilde{\theta}$ including S .

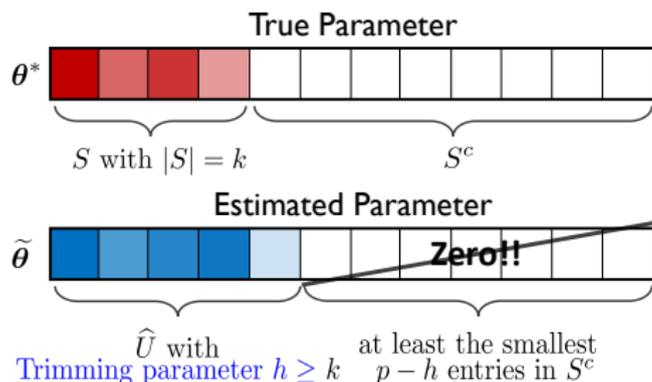
Statistical Analysis

Theorem 1: General ℓ_∞ -error Bound and Variable Selection

- 1 For every pair $j_1 \in S$, $j_2 \in S^c$, we have $|\tilde{\theta}_{j_1}| > |\tilde{\theta}_{j_2}|$
- 2 If $h < k$, all $j \in S^c$ are successfully estimated as zero and

$$\|\tilde{\theta} - \theta^*\|_\infty \leq \left\| (\hat{Q}_{SS})^{-1} \nabla \mathcal{L}(\theta^*)_S \right\|_\infty + \lambda_n \left\| (\hat{Q}_{SS})^{-1} \right\|_\infty$$

- 3 If $h \geq k$, at least the smallest (in absolute) $p - h$ entries in S^c are exactly zero and $\|\tilde{\theta} - \theta^*\|_\infty \leq \left\| (\hat{Q}_{\hat{U}\hat{U}})^{-1} \nabla \mathcal{L}(\theta^*)_{\hat{U}} \right\|_\infty$ where \hat{U} is defined as the h largest absolute entries of $\tilde{\theta}$ including S .



Statistical Analysis

Theorem 2: General ℓ_2 -error Bound

Theorem 2

- Consider an M -estimation problem with Trimmed ℓ_1 regularization where all conditions in Theorem 1 hold.
- For any local minimum $\tilde{\boldsymbol{\theta}}$, the parameter estimation error in terms of ℓ_2 -norm is upper bounded as:

$$\|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_2 = \begin{cases} C\lambda_n \left(\sqrt{k}/2 + \sqrt{k-h} \right) & \text{if } h < k \\ C\lambda_n \sqrt{h}/2 & \text{otherwise} \end{cases}$$

Statistical Analysis

Theorem 2: General ℓ_2 -error Bound

Theorem 2

- Consider an M -estimation problem with Trimmed ℓ_1 regularization where all conditions in Theorem 1 hold.
- For any local minimum $\tilde{\theta}$, the parameter estimation error in terms of ℓ_2 -norm is upper bounded as:

$$\|\tilde{\theta} - \theta^*\|_2 = \begin{cases} C\lambda_n \left(\sqrt{k}/2 + \sqrt{k-h} \right) & \text{if } h < k \\ C\lambda_n \sqrt{h}/2 & \text{otherwise} \end{cases}$$

- From our bound, $h = k$ is the **best case!**
 - We can choose $h \asymp k$ via **cross-validation**.

Table: ℓ_2 -error bound for different h values.

	$h < k$	$h = k$	$h > k$
$\ \tilde{\theta} - \theta^*\ _2$	$C\lambda_n \left(\frac{\sqrt{k}}{2} + \sqrt{k-h} \right)$	$C\lambda_n \frac{\sqrt{k}}{2}$	$C\lambda_n \frac{\sqrt{h}}{2}$

Statistical Analysis

Remarks: Other alternative penalties vs. Trimmed ℓ_1

$$\|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_2 = \begin{cases} C\lambda_n \left(\sqrt{k}/2 + \sqrt{k-h} \right) & \text{if } h < k \\ C\lambda_n \sqrt{h}/2 & \text{otherwise} \end{cases}$$

- $\rho_\lambda(t)$: $(\boldsymbol{\mu}, \gamma)$ -amenable
 - $\rho_\lambda(t) + \frac{1}{2}\boldsymbol{\mu}t^2$ is convex.
 - $\rho'_\lambda(t) = 0$ for $|t| > \gamma$.

Statistical Analysis

Remarks: Other alternative penalties vs. Trimmed ℓ_1

$$\|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_2 = \begin{cases} C\lambda_n \left(\sqrt{k}/2 + \sqrt{k-h} \right) & \text{if } h < k \\ C\lambda_n \sqrt{h}/2 & \text{otherwise} \end{cases}$$

- $\rho_\lambda(t)$: (μ, γ) -amenable
 - $\rho_\lambda(t) + \frac{1}{2}\mu t^2$ is convex.
 - $\rho'_\lambda(t) = 0$ for $|t| > \gamma$.

Table: ℓ_2 -error bound comparison with universal constant c_0 in sub-Gaussian tail bounds.

	Standard ℓ_1 ($h = 0$)	(μ, γ) -amenable	Trimmed ℓ_1 ($h = k$)
$\ \tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\ _2$	$\frac{3c_0}{\kappa_l} \frac{\lambda_n \sqrt{k}}{2}$	$\frac{c_0}{\kappa_l - \frac{3}{2}\mu} \frac{\lambda_n \sqrt{k}}{2}$	$\frac{c_0}{\kappa_l} \frac{\lambda_n \sqrt{k}}{2}$

Statistical Analysis

Remarks: Other alternative penalties vs. Trimmed ℓ_1

$$\|\tilde{\theta} - \theta^*\|_2 = \begin{cases} C\lambda_n \left(\sqrt{k}/2 + \sqrt{k-h} \right) & \text{if } h < k \\ C\lambda_n \sqrt{h}/2 & \text{otherwise} \end{cases}$$

- $\rho_\lambda(t)$: (μ, γ) -amenable
 - $\rho_\lambda(t) + \frac{1}{2}\mu t^2$ is convex.
 - $\rho'_\lambda(t) = 0$ for $|t| > \gamma$.

Table: ℓ_2 -error bound comparison with universal constant c_0 in sub-Gaussian tail bounds.

	Standard ℓ_1 ($h = 0$)	(μ, γ) -amenable	Trimmed ℓ_1 ($h = k$)
$\ \tilde{\theta} - \theta^*\ _2$	$\frac{3c_0}{\kappa_l} \frac{\lambda_n \sqrt{k}}{2}$	$\frac{c_0}{\kappa_l - \frac{3}{2}\mu} \frac{\lambda_n \sqrt{k}}{2}$	$\frac{c_0}{\kappa_l} \frac{\lambda_n \sqrt{k}}{2}$

- Trimmed ℓ_1 can achieve **three times smaller bound** than standard one.

Statistical Analysis

Remarks: Other alternative penalties vs. Trimmed ℓ_1

$$\|\tilde{\theta} - \theta^*\|_2 = \begin{cases} C\lambda_n \left(\sqrt{k}/2 + \sqrt{k-h} \right) & \text{if } h < k \\ C\lambda_n \sqrt{h}/2 & \text{otherwise} \end{cases}$$

- $\rho_\lambda(t)$: (μ, γ) -amenable
 - $\rho_\lambda(t) + \frac{1}{2}\mu t^2$ is convex.
 - $\rho'_\lambda(t) = 0$ for $|t| > \gamma$.

Table: ℓ_2 -error bound comparison with universal constant c_0 in sub-Gaussian tail bounds.

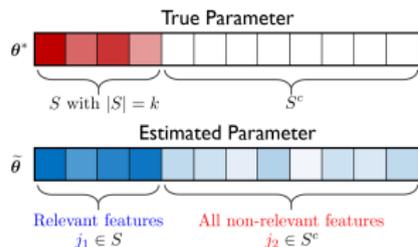
	Standard ℓ_1 ($h = 0$)	(μ, γ) -amenable	Trimmed ℓ_1 ($h = k$)
$\ \tilde{\theta} - \theta^*\ _2$	$\frac{3c_0}{\kappa_l} \frac{\lambda_n \sqrt{k}}{2}$	$\frac{c_0}{\kappa_l - \frac{3}{2}\mu} \frac{\lambda_n \sqrt{k}}{2}$	$\frac{c_0}{\kappa_l} \frac{\lambda_n \sqrt{k}}{2}$

- Also, we have a smaller bound than non-convex regularizers since (μ, γ) -amenable regularizers have (possibly large) μ in the denominator.

Statistical Analysis

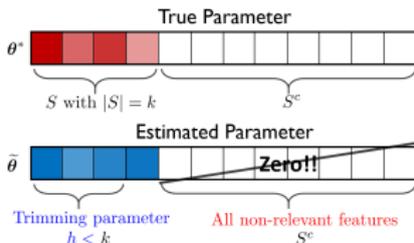
Corollary 1: General ℓ_∞ -error Bound for Linear Regression

- Consider a linear regression problem with sub-Gaussian error ϵ .
- Under standard conditions as in **Theorem 1** and **incoherence condition** on sample covariance, with high probability, *any* local minimum $\tilde{\theta}$ satisfies



① The absolute value of relevant features is always larger than non-relevant features.

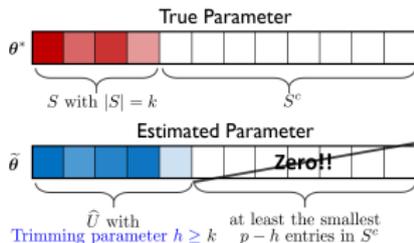
θ^* : True parameter
 $\tilde{\theta}$: Estimated parameter



② If we set the trimming parameter h smaller than true sparsity level k , all non-relevant parameters are estimated as zero.

$$\|\tilde{\theta} - \theta^*\|_\infty \leq c_1 \sqrt{\frac{\log p}{n}} + \lambda_n c_{\infty}$$

$$\|\tilde{\theta} - \theta^*\|_2 \leq c_4 \sqrt{\frac{\log p}{n}} (\sqrt{k/2} + \sqrt{k-h})$$



③ If we set h larger than true sparsity level k , at least the smallest $p-h$ entries are estimated as zero.

$$\|\tilde{\theta} - \theta^*\|_\infty \leq c_1 \sqrt{\frac{\log p}{n}}$$

$$\|\tilde{\theta} - \theta^*\|_2 \leq \frac{c_4}{2} \sqrt{\frac{h \log p}{n}}$$

Table of Contents

1 Introduction and Setup

2 Statistical Analysis

3 Optimization

4 Experiments & Applications to Deep Learning

Optimization for Trimmed ℓ_1 Regularized Program

- For an optimization, we use our **second formulation** of trimmed regularization problem

$$\underset{\boldsymbol{\theta} \in \Omega, \mathbf{w} \in [0,1]^p}{\text{minimize}} \mathcal{F}(\boldsymbol{\theta}, \mathbf{w}) := \mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) + \lambda \sum_{j=1}^p w_j |\theta_j| \quad \text{s.t.} \quad \mathbf{1}^T \mathbf{w} \geq p - h$$

Optimization for Trimmed ℓ_1 Regularized Program

- For an optimization, we use our **second formulation** of trimmed regularization problem

$$\underset{\boldsymbol{\theta} \in \Omega, \mathbf{w} \in [0,1]^p}{\text{minimize}} \mathcal{F}(\boldsymbol{\theta}, \mathbf{w}) := \mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) + \lambda \sum_{j=1}^p w_j |\theta_j| \quad \text{s.t.} \quad \mathbf{1}^T \mathbf{w} \geq p - h$$

- We update $(\boldsymbol{\theta}, \mathbf{w})$ in an alternating manner.

$$\begin{aligned} \boldsymbol{\theta}^{k+1} &\leftarrow \text{prox}_{\eta\lambda\mathcal{R}(\cdot, \mathbf{w}^k)}[\boldsymbol{\theta}^k - \eta\nabla_{\boldsymbol{\theta}}\mathcal{L}(\boldsymbol{\theta}^k)] \\ \mathbf{w}^{k+1} &\leftarrow \text{proj}_{\mathcal{S}}[\mathbf{w}^k - \tau\mathbf{r}(\boldsymbol{\theta}^{k+1})] \end{aligned}$$

- Fixing \mathbf{w} , prox operator is weighted ℓ_1 norm.
- By fixing $\boldsymbol{\theta}$, the objective function \mathcal{F} is linear in \mathbf{w} .
- $\text{proj}_{\mathcal{S}}$ is a projection onto the constraint set $\mathcal{S} = \{\mathbf{w} \in [0,1]^p \mid \mathbf{1}^T \mathbf{w} = p - h\}$.

Optimization: Comparison with DC-based Approach

- Convergence history **our algorithm** vs. **Algorithm 2 of (Khamaru & Wainwright, 2018)***.
 - **Algorithm 2 of (Khamaru & Wainwright, 2018)** is an optimization method for (non-convex and non-smooth) objective functions of the form **difference of convex functions** ($f := g + \phi - h$).
 - Trimmed regularized problem can be formulated as a DC.

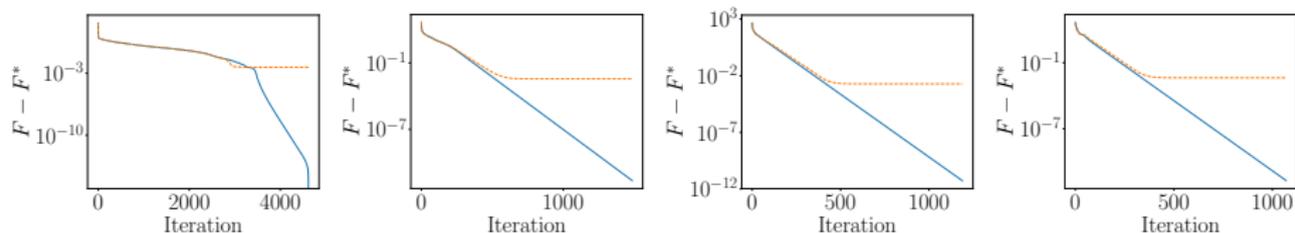


Figure: Algorithm comparison with $\lambda \in \{0.5, 5, 10, 20\}$.

*K. Khamaru and M. J. Wainwright. Convergence guarantees for a class of non-convex and non-smooth optimization problems. ICML, 2018

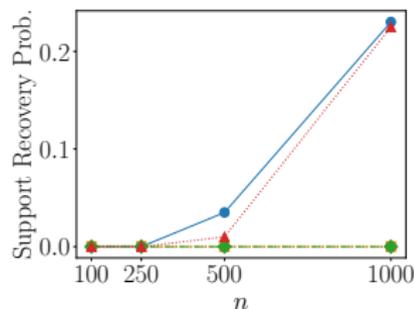
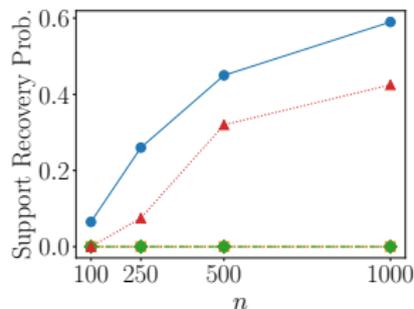
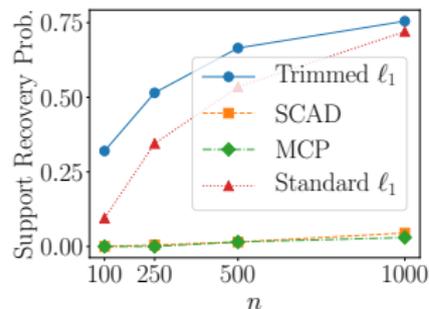
Table of Contents

- 1 Introduction and Setup
- 2 Statistical Analysis
- 3 Optimization
- 4 Experiments & Applications to Deep Learning

Simulation Experiments

Incoherent Case: Support Recovery

- Scenario 1: Incoherence condition is satisfied

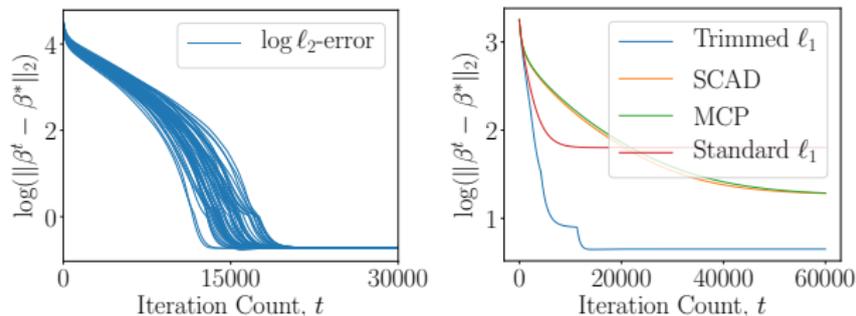


- Probability of successful support recovery for Trimmed Lasso, SCAD, MCP, and standard Lasso with $(p, k) = (128, 8), (256, 16), (512, 32)$.

Simulation Experiments

Incoherent Case: Stationary & $\log \ell_2$ -error Comparison

- Scenario 1: Incoherence condition is satisfied

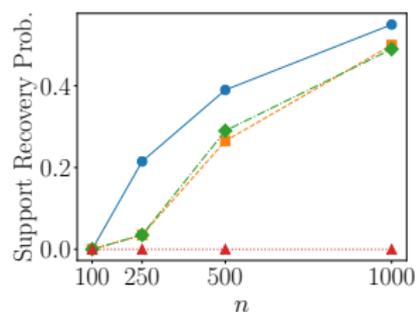
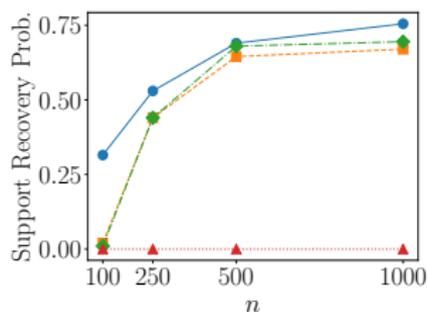
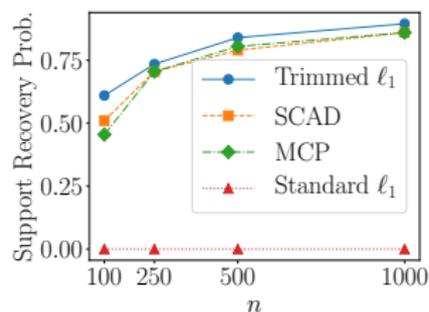


- (Left)** 50 random initializations for a setting with $(n, p, k) = (160, 256, 16)$.
- (Right)** $\log \ell_2$ -error comparison.

Simulation Experiments

Nonincoherent Case: Support Recovery

- Scenario 2: Incoherence condition violated
 - Note that we need an **incoherence condition** in our Corollary 1.
 - Interestingly, the Trimmed Lasso outperforms all the other comparison regularizers even in this regime.

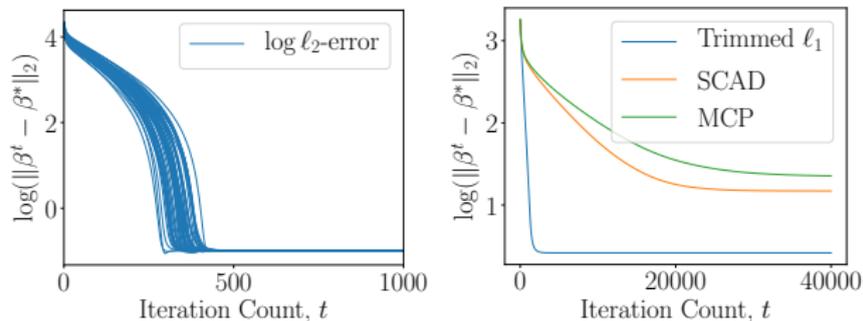


- Probability of successful support recovery for Trimmed Lasso, SCAD, MCP, and standard Lasso with $(p, k) = (128, 8), (256, 16), (512, 32)$.

Simulation Experiments

Nonincoherent Case: Stationary & $\log \ell_2$ -error Comparison

- Scenario 2: Incoherence condition violated



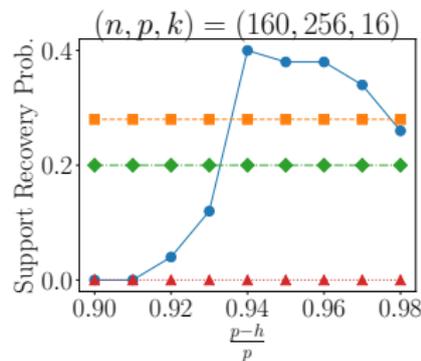
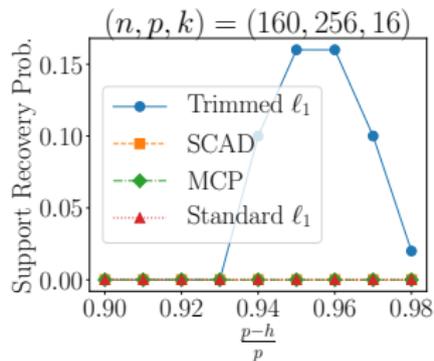
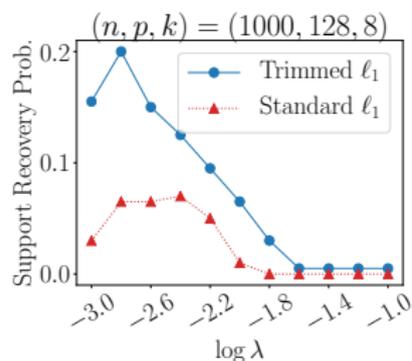
- (Left)** 50 random initializations for a setting with $(n, p, k) = (160, 256, 16)$.
- (Right)** $\log \ell_2$ -error comparison.

Simulation Experiments

Nonincoherent Case: Stationary

• Scenario 3

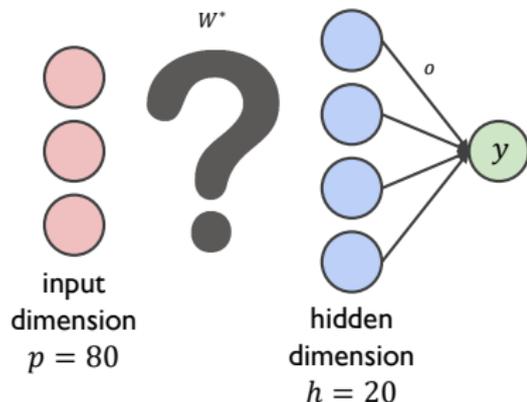
- **(Left)** True signals and regularization parameter λ are both small (Small regime)
- Investigating the choice of the trimming parameter h (**Middle:** Incoherent case, **Right:** Non-incoherent case).



Applications to Deep Learning 1

Input Structure Recovery of Compact Neural Networks

- We apply trimmed regularization to recover the weight structure of neural networks as parameter support recovery.
- Motivated by the recent work of [Oymak \(2018\)*](#), we consider



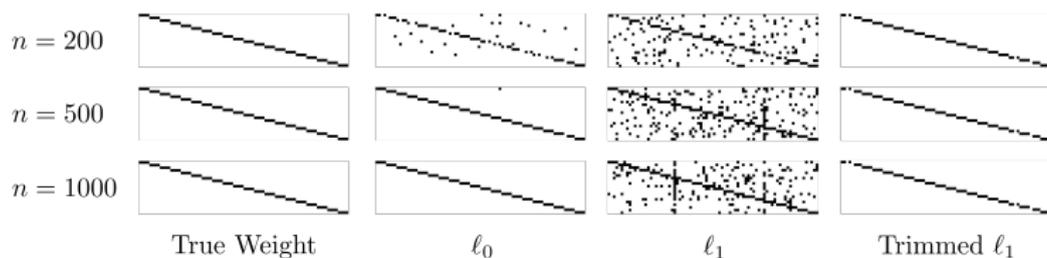
- The regression model, $y_i = \boldsymbol{o}^T \text{ReLU}(\boldsymbol{W}^* \boldsymbol{x}_i)$ with $\boldsymbol{o} = \mathbf{1}$.
- Each hidden node is connected to only 4 input features.

* Samet Oymak. Learning Compact Neural Networks with Regularization. ICML, 2018.

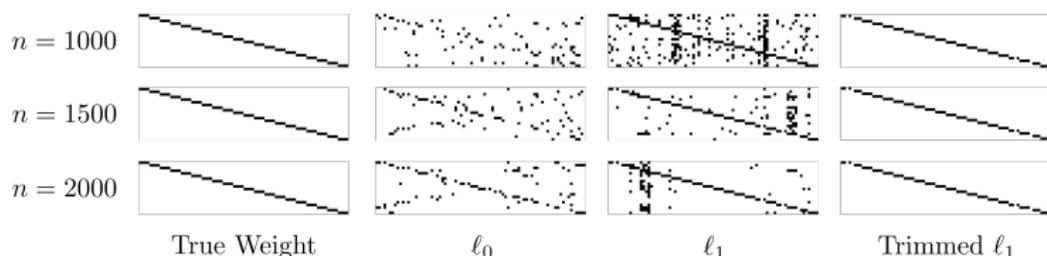
Applications to Deep Learning 1

Input Structure Recovery of Compact Neural Networks: Results

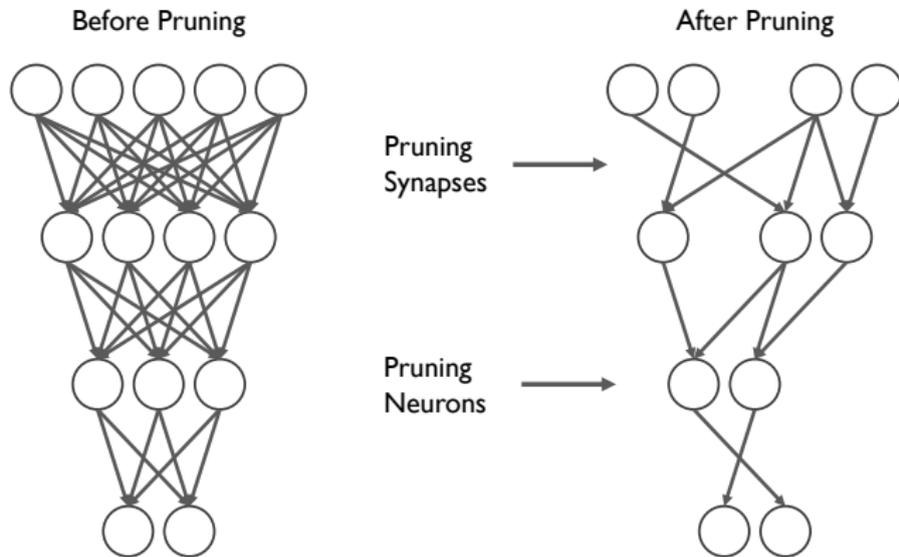
- With good initialization (**small perturbation** from true weight)



- With random initialization



Applications to Deep Learning 2: Pruning Deep Networks



- **Pruning neurons** is more computationally efficient than edge-wise pruning.

Applications to Deep Learning: Pruning Deep Networks

Trimmed Group ℓ_1 Regularization on Deep Networks

To encourage group sparsity on neural networks, we consider two cases:

- **Neuron sparsity (for fully-connected layers)**

- Let $\theta_l \in \mathbb{R}^{n_{\text{in}} \times n_{\text{out}}}$ be a weight parameter, then we can enforce group-wise sparsity via **Trimmed group ℓ_1 penalty** as

$$\mathcal{R}_l(\theta_l, \mathbf{w}) = \lambda_l \sum_{j=1}^{n_{\text{in}}} w_j \sqrt{\theta_{j,1}^2 + \theta_{j,2}^2 + \cdots + \theta_{j,n_{\text{out}}}^2}$$

- **Activation map sparsity (for convolutional layers)**

- Similarly, let $\theta_l \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}} \times H \times W}$ be a weight parameter, then

$$\mathcal{R}_l(\theta, \mathbf{w}) = \lambda_l \sum_{j=1}^{C_{\text{out}}} w_j \sqrt{\sum_{m,n,k} \theta_{j,m,n,k}^2}$$

for all possible indices (m, n, k) .

with the constraint $\mathbf{1}^T \mathbf{w} = n_{\text{in}} - h_l$ or $C_{\text{out}} - h_l$ respectively.

Applications to Deep Learning: Pruning Deep Networks

Results on MNIST dataset

- Comparison with vanilla group ℓ_1 penalty vs. Trimmed group ℓ_1 penalty on LeNet-300-100 structure

Method	Pruned Model	Error (%)
No Regularization	784-300-100	1.6
grp ℓ_1	784-241-67	1.7
grp $\ell_{1\text{trim}}$, $h = \text{half of original}$	392-150-50	1.6

Applications to Deep Learning: Pruning Deep Networks

Bayesian Neural Networks with Trimmed ℓ_1 Regularization

- Most modern algorithms for network pruning are based on **Bayesian variational framework**. We propose a Bayesian neural network with Trimmed ℓ_1 regularization regarding only θ as Bayesian.

Applications to Deep Learning: Pruning Deep Networks

Bayesian Neural Networks with Trimmed ℓ_1 Regularization

- Most modern algorithms for network pruning are based on **Bayesian variational framework**. We propose a Bayesian neural network with Trimmed ℓ_1 regularization regarding only θ as Bayesian.
- By relationship between **Bayesian neural networks** and **variational dropout**, we choose $q_{\theta, \alpha}(\theta_{i,j}) = \mathcal{N}(\phi_{i,j}, \alpha_{i,j} \phi_{i,j}^2)$ as a variational distribution.

Applications to Deep Learning: Pruning Deep Networks

Bayesian Neural Networks with Trimmed ℓ_1 Regularization

- Most modern algorithms for network pruning are based on **Bayesian variational framework**. We propose a Bayesian neural network with Trimmed ℓ_1 regularization **regarding only θ as Bayesian**.
- By relationship between **Bayesian neural networks** and **variational dropout**, we choose $q_{\theta, \alpha}(\theta_{i,j}) = \mathcal{N}(\phi_{i,j}, \alpha_{i,j} \phi_{i,j}^2)$ as a variational distribution.
- Combined with Trimmed ℓ_1 regularization, the objective is

$$\underbrace{\mathbb{E}_{q_{\phi, \alpha}(\theta)} \left[-\mathcal{L}(\mathcal{W}; \mathcal{D}) \right] + \mathbb{KL}(q_{\phi, \alpha}(\mathcal{W}) \| p(\mathcal{W}))}_{\text{ELBO}} + \underbrace{\mathbb{E}_{q_{\phi, \alpha}(\theta)} \left[\sum_{l=1}^{L+1} \lambda_l \mathcal{R}_l(\theta_l, \mathbf{w}_l) \right]}_{\text{Expected Trimmed group } \ell_1 \text{ penalty}}$$

Applications to Deep Learning: Pruning Deep Networks

Results on MNIST dataset (Cont' d)

- With Bayesian extensions on LeNet-300-100
 - We compare with a smoothed ℓ_0 -norm under Bayesian variational framework proposed by Louizos et al. (2018)*

Method	Pruned Model	Error (%)
ℓ_0 (Louizos et al., 2018)	219-214-100	1.4
ℓ_0, λ sep. (Louizos et al., 2018)	266-88-33	1.8
Bayes grp $\ell_{1\text{trim}}, h = \ell_0$	219-214-100	1.4
Bayes grp $\ell_{1\text{trim}}, h = \ell_0, \lambda$ sep.	266-88-33	1.6
Bayes grp $\ell_{1\text{trim}}, h < \ell_0, \lambda$ sep.	245-75-25	1.7

- With Bayesian extensions on LeNet-5-Caffe

Method	Pruned Model	Error (%)
ℓ_0 (Louizos et al., 2018)	20-25-45-462	0.9
ℓ_0, λ sep. (Louizos et al., 2018)	9-18-65-25	1.0
Bayes grp $\ell_{1\text{trim}}, h < \ell_0$	20-25-45-150	0.9
Bayes grp $\ell_{1\text{trim}}, h = \ell_0, \lambda$ sep.	9-18-65-25	1.0
Bayes grp $\ell_{1\text{trim}}, h < \ell_0, \lambda$ sep.	8-17-53-19	1.0

* Louizos et al. Learning Sparse Neural Networks through ℓ_0 Regularization. ICLR, 2018

Concluding Remarks

- **High-dimensional M -estimators with Trimmed ℓ_1 penalty:** Alleviate the bias incurred by the vanilla ℓ_1 penalty by leaving the h largest parameter entries penalty-free.
- **Theoretical Results** on support recovery and ℓ_2 -error hold for any local optima and are competitive with other non-convex regularizers.
- **Simulation experiments** demonstrated the value of approach compared to Lasso and non-convex penalties.
- Future work:
 - Trimming for other standard regularizers beyond sparsity
 - Bypassing incoherence condition in corollaries
 - More experiments and theories when RSC does not hold
 - Investigating the use of trimmed regularization in deep models.

THANK YOU!
Any Questions?

Poster Session at Pacific Ballroom #186
6:30pm – 9:00pm