

Department of Industrial and Systems Engineering  
University of Southern California

# Comparing Solution Paths of Sparse Quadratic Minimization with a Stieltjes Matrix

May 7, 2024

**Ziyu He**, Shaoning Han, Andrés Gómez, Ying Cui, Jong-Shi Pang

# Hyperparameter selection for sparse estimation

## Basic Sparse Estimation:

Constraints  $\ell \in \mathbb{R}_-^n, u \in \mathbb{R}^+$

$$\underset{\ell \leq x \leq u}{\text{minimize}} \quad \underbrace{\frac{1}{2} x^\top Q x + q^\top x}_{\text{Stieltjes matrix}} \quad + \quad \gamma \underbrace{\Phi(x)}_{\text{Sparsity inducing regularizer}} \quad (1)$$

$Q \in \mathbb{R}^{n \times n}$ : Stieltjes matrix

Applications e.g., Markov random field

**Hyperparameter  $\gamma \geq 0$**

Sparsity inducing  
regularizer

# Hyperparameter selection for sparse estimation

## Basic Sparse Estimation:

Constraints  $l \in \mathbb{R}_-^n, u \in \mathbb{R}^+$

$$\begin{array}{ccc} \downarrow & & \downarrow \\ \text{minimize}_{l \leq x \leq u} & \underbrace{\frac{1}{2} x^\top Q x + q^\top x}_{Q \in \mathbb{R}^{n \times n}: \text{Stieltjes matrix}} & + \gamma \underbrace{\Phi(x)}_{\text{Sparsity inducing regularizer}} \end{array} \quad (1)$$

Applications e.g., Markov random field

**Hyperparameter  $\gamma \geq 0$**

## Hyperparameter Selection:

- Select the best  $\gamma$  by some criteria for (1)'s solution, e.g., test error.

# Hyperparameter selection for sparse estimation

## Basic Sparse Estimation:

Constraints  $\ell \in \mathbb{R}_-^n, u \in \mathbb{R}^+$

**Hyperparameter  $\gamma \geq 0$**

$$\begin{array}{ccc} \downarrow & & \downarrow \\ \text{minimize}_{\ell \leq x \leq u} & \underbrace{\frac{1}{2} x^\top Q x + q^\top x}_{Q \in \mathbb{R}^{n \times n}: \text{Stieltjes matrix}} & + \gamma \underbrace{\Phi(x)}_{\text{Sparsity inducing regularizer}} \end{array} \quad (1)$$

**Applications e.g., Markov random field**

## Hyperparameter Selection:

- Select the best  $\gamma$  by some criteria for (1)'s solution, e.g., test error.

## Parametric Programming:

- A whole path of (1)'s solution as a function of  $\gamma$ .

# Choice of $\Phi$ : a dilemma

## The ideal choice

- Weighted  $\ell_0$  :  $\sum_{i=1}^n p_i |x_i|_0$
- Mixed integer hence can be computationally prohibitive.

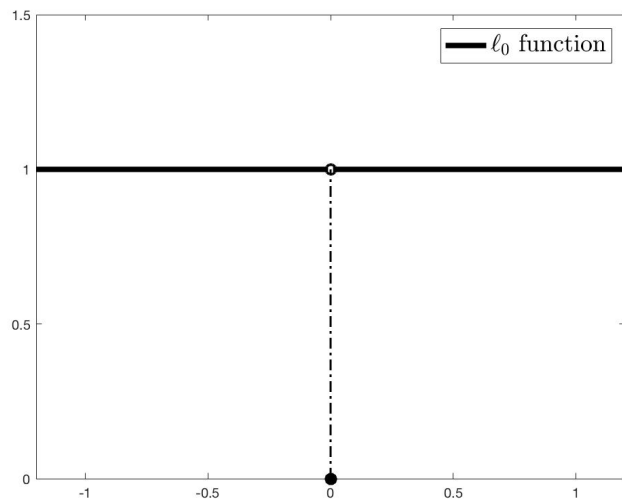
## Convex relaxation

- Weighted  $\ell_1$  :  $\sum_{i=1}^n p_i |x_i|$
- Easier to compute but can give us undesirable results.

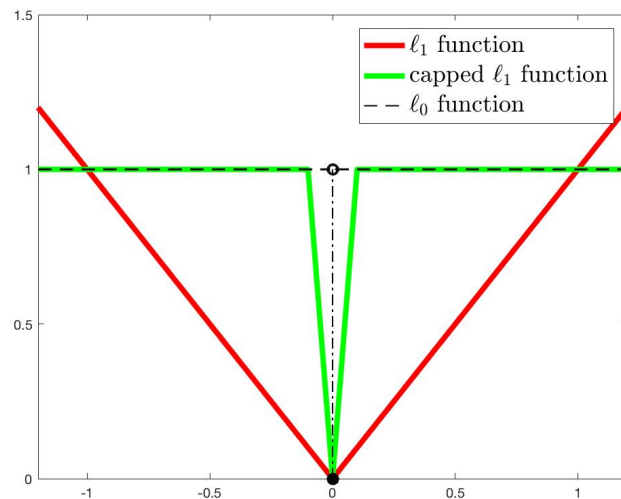
## Nonconvex surrogate

- Weighted capped  $\ell_1$ :  $\sum_{i=1}^n p_i \min(\frac{|x_i|}{\delta}, 1), \delta > 0$

# $\ell_0$ , $\ell_1$ and capped $\ell_1$



(a)  $\ell_0$  function



(b)  $\ell_0$ ,  $\ell_1$  and Capped  $\ell_1$

Figure 1: Capped  $\ell_1$  is a better approximation to  $\ell_0$  than  $\ell_1$

# Choice of $\Phi$ : a dilemma

## The ideal choice

- Weighted  $\ell_0$  :  $\sum_{i=1}^n p_i |x_i|_0$
- Mixed integer hence can be computationally prohibitive.

## Convex relaxation

- Weighted  $\ell_1$  :  $\sum_{i=1}^n p_i |x_i|$
- Easier to compute but can us give undesirable results.

## Nonconvex surrogate

- Weighted capped  $\ell_1$ :  $\sum_{i=1}^n p_i \min(\frac{|x_i|}{\delta}, 1), \delta > 0$
- A better approximation to  $\ell_0$  but nonconvex and nonsmooth.
- Analytical properties? How to compute?
- D(irectional)-stationary [ $\Leftrightarrow$  strongly local] solution path.

# The overall goals

- Studying and comparing paths from  $\ell_0, \ell_1$  and capped  $\ell_1$ .
- Emphasizing on the d-stationary (**d-stat.**) path of capped  $\ell_1$ 
  - Analytical (**e.g., number of pieces**)
  - Computational (**the first rigorous computational study for nonconvex paths**).
- Highlighting the benefit of nonconvex approaches in balancing
  - Computational effort
  - Statistical and optimization performances



# Previous studies

## $\ell_1$ -path ( $Q$ is PD)

- Continuous piecewise affine [Efron et al., 2004].
- In worst case exponentially many pieces [Mairal and Yu, 2012].
- **Cases guaranteed to have polynomially many pieces?**

# Previous studies

## $\ell_1$ -path ( $Q$ is PD)

- Continuous piecewise affine [Efron et al., 2004].
- In worst case exponentially many pieces [Mairal and Yu, 2012].
- **Cases guaranteed to have polynomial pieces?**

## $\ell_0$ -path (equal weights $p_i \equiv 1$ )

- Discontinuous piecewise affine ( $n + 1$  pieces) [Soussen et al., 2015].
- **Unequal weights?**

# Previous studies

## $\ell_1$ -path ( $Q$ is PD)

- Continuous piecewise affine [Efron et al., 2004].
- In worst case exponentially many pieces [Mairal and Yu, 2012].
- **Cases guaranteed to have polynomial pieces?**

## $\ell_0$ -path (equal weights $p_i \equiv 1$ )

- Discontinuous piecewise affine ( $n + 1$  pieces) [Soussen et al., 2015].
- **Unequal weights?**

## Other nonconvex surrogates (e.g., $\ell_p$ , MCP)

- Parametric nonlinear systems, e.g.,  $\gamma$  in quadratic terms.
- Piecewise smooth paths, cannot be exactly traced in finite time.
- Either heuristic [Yukawa and Amari, 2015].
- Or fails to approximate the exact  $\ell_0$ -path [Zhang, 2010].
- **Capped  $\ell_1$  doesn't have these issues, but how to compute?**

# Our contributions: analytical

Deriving special classes:

- Guaranteed to have polynomial many pieces.
- Worst case exponential complexity.

# Our contributions: analytical

Deriving special classes:

- Guaranteed to have polynomial many pieces.
- Worst case exponential complexity.

Regularizer	Pieces	Optimality	Class
$\ell_1$	$2n + 1$	global	$\ell = 0$ , Stieltjes $Q$
$\ell_0$	exponential	global	Non-Stieltjes $Q$
(unequal weights)	$n + 1$	global	$\ell = 0$ , Stieltjes $Q$
Capped $\ell_1$	$2n^2 + 3n + 1$	global	$\ell = 0$ , Stieltjes $Q$
	$n + 1$	d-stat.	$\ell = 0$ , Stieltjes $Q$

Table 1: Summary of some analytical results

# Our contributions: analytical

Deriving special classes:

- Guaranteed to have polynomial many pieces.
- Worst case exponential complexity.

Regularizer	Pieces	Optimality	Class
$\ell_1$	$2n + 1$	global	$\ell = 0$ , Stieltjes $Q$
$\ell_0$	exponential	global	Non-Stieltjes $Q$
(unequal weights)	$n + 1$	global	$\ell = 0$ , Stieltjes $Q$
Capped $\ell_1$	$2n^2 + 3n + 1$	global	$\ell = 0$ , Stieltjes $Q$
	$n + 1$	d-stat.	$\ell = 0$ , Stieltjes $Q$

Table 1: Summary of some analytical results

- **They are all piecewise affine.**
- **Stieltjes structure is the key for polynomial complexity.**

# Our contributions: computational

- A rigorous method to compute d-stat. paths for capped  $\ell_1$ .
  - **Can be discontinuous!**

# Our contributions: computational

- A rigorous method to compute d-stat. paths for capped  $\ell_1$ .
  - **Can be discontinuous!**
- Efficient algorithm (GHP) to restore discontinuity.
  - **Complexity is strongly polynomial (Stieltjes  $Q$ ).**



# Our contributions: computational

- A rigorous method to compute d-stat. paths for capped  $\ell_1$ .
  - **Can be discontinuous!**
- Efficient algorithm (GHP) to restore discontinuity.
  - **Complexity is strongly polynomial (Stieltjes  $Q$ ).**
- Benefits of capped  $\ell_1$  d-stat. path supported by numerical results:
  - **Way faster than computing  $\ell_0$ .**
  - **Superior optimization and statistical performance than  $\ell_1$ .**

## Pivoting method: high level ideas

$$\underset{\ell \leq x \leq u}{\text{minimize}} \quad \frac{1}{2} x^\top Q x + q^\top x + \gamma \sum_{i=1}^n p_i \min \left\{ \frac{|x_i|}{\delta}, 1 \right\} \quad (2)$$

- D-stat. path of (2) is piecewise affine in  $\gamma$ .

# Pivoting method: high level ideas

$$\underset{\ell \leq x \leq u}{\text{minimize}} \quad \frac{1}{2} x^\top Q x + q^\top x + \gamma \sum_{i=1}^n p_i \min \left\{ \frac{|x_i|}{\delta}, 1 \right\} \quad (2)$$

- D-stat. path of (2) is piecewise affine in  $\gamma$ .
- To trace the path is to compute these pieces one by one.
  - **In the direction of  $\gamma \downarrow 0$**

# Pivoting method: high level ideas

$$\underset{\ell \leq x \leq u}{\text{minimize}} \quad \frac{1}{2} x^\top Q x + q^\top x + \gamma \sum_{i=1}^n p_i \min \left\{ \frac{|x_i|}{\delta}, 1 \right\} \quad (2)$$

- D-stat. path of (2) is piecewise affine in  $\gamma$ .
- To trace the path is to compute these pieces one by one.
- Each piece is associated with a “basis”.
  - **Tuple of index sets to restrict the values of solution for (2).**

# Pivoting method: high level ideas

$$\underset{\ell \leq x \leq u}{\text{minimize}} \quad \frac{1}{2} x^\top Q x + q^\top x + \gamma \sum_{i=1}^n p_i \min \left\{ \frac{|x_i|}{\delta}, 1 \right\} \quad (2)$$

- D-stat. path of (2) is piecewise affine in  $\gamma$ .
- To trace the path is to compute these pieces one by one.
- Each piece is associated with a “basis”.
- Fixed basis  $\Rightarrow$  reduced problem  $\Rightarrow$  reduced solution.
  - **Not necessarily d-stat. for (2).**

# Pivoting method: high level ideas

$$\underset{\ell \leq x \leq u}{\text{minimize}} \quad \frac{1}{2} x^\top Q x + q^\top x + \gamma \sum_{i=1}^n p_i \min \left\{ \frac{|x_i|}{\delta}, 1 \right\} \quad (2)$$

- D-stat. path of (2) is piecewise affine in  $\gamma$ .
- To trace the path is to compute these pieces one by one.
- Each piece is associated with a “basis”.
- Fixed basis  $\Rightarrow$  reduced problem  $\Rightarrow$  reduced solution.
- Conditions for the reduced solution to be d-stat. of (2).
  - **Linear inequalities in  $\gamma$ .**
  - **Ratio test: the smallest  $\gamma$  for the current basis to be d-stat.**

## Pivoting method: high level ideas

$$\underset{\ell \leq x \leq u}{\text{minimize}} \quad \frac{1}{2} x^\top Q x + q^\top x + \gamma \sum_{i=1}^n p_i \min \left\{ \frac{|x_i|}{\delta}, 1 \right\} \quad (2)$$

### Pivoting method (informal)

- Compute the pieces/bases from right to the left.
- At the current piece/basis, do ratio test to get  $\gamma^*$ .
- When we go beyond  $\gamma^*$ , change the basis accordingly.

# Pivoting method: high level ideas

$$\underset{\ell \leq x \leq u}{\text{minimize}} \quad \frac{1}{2} x^\top Q x + q^\top x + \gamma \sum_{i=1}^n p_i \min \left\{ \frac{|x_i|}{\delta}, 1 \right\} \quad (2)$$

## Pivoting method (informal)

- Compute the pieces/bases from right to the left.
- At the current piece/basis, do ratio test to get  $\gamma^*$ .
- When we go beyond  $\gamma^*$ , change the basis accordingly.
- **Discontinuous at  $\gamma^*$  if taking  $\pm\delta$  (not allowed for (2)'s d-stat.).**



# Pivoting method: high level ideas

$$\underset{\ell \leq x \leq u}{\text{minimize}} \quad \frac{1}{2} x^\top Q x + q^\top x + \gamma \sum_{i=1}^n p_i \min \left\{ \frac{|x_i|}{\delta}, 1 \right\} \quad (2)$$

## Pivoting method (informal)

- Compute the pieces/bases from right to the left.
- At the current piece/basis, do ratio test to get  $\gamma^*$ .
- When we go beyond  $\gamma^*$ , change the basis accordingly.
- Discontinuous at  $\gamma^*$  if taking  $\pm\delta$  (not allowed for (2)'s d-stat.).
- We need an algorithm to restore a d-stat. solution at  $\gamma^*$  that:
  - **Doesn't need to compute from scratch.**
  - **Leverages the Stieltjes  $Q$  for faster computation.**

# The GHP method (informal)

$$\underset{\ell \leq x \leq u}{\text{minimize}} \quad \frac{1}{2} x^\top Q x + q^\top x + \gamma \sum_{i=1}^n p_i \min \left\{ \frac{|x_i|}{\delta}, 1 \right\} \quad (2)$$

## High level ideas

- Basis enumeration method, just for an alternative basis.
- Fixed basis  $\Rightarrow$  reduced problem  $\Rightarrow$  reduced solution.

# The GHP method (informal)

$$\underset{\ell \leq x \leq u}{\text{minimize}} \quad \frac{1}{2} x^\top Q x + q^\top x + \gamma \sum_{i=1}^n p_i \min \left\{ \frac{|x_i|}{\delta}, 1 \right\} \quad (2)$$

## High level ideas

- Basis enumeration method, just for an alternative basis.
- Fixed basis  $\Rightarrow$  reduced problem  $\Rightarrow$  reduced solution.
- Conditions for the reduced solution to be d-stat. for (2).
  - **Certain index sets should be empty.**
  - **If not, form a new basis by moving them accordingly.**
  - **Proceed to the next steps until such index sets are empty.**

# The GHP method

## Back to the pivoting method

- Suppose we are at a discontinuous  $\gamma^*$  where we need to restore.
- Initialize the GHP method with the current “non-d-stationary” solution, we can prove the following theorem.

# The GHP method

## Back to the pivoting method

- Suppose we are at a discontinuous  $\gamma^*$  which we need to restore.
- Initialize the GHP method with the current “non-d-stationary” solution, we can prove the following theorem.

## Main theorem of GHP (informal)

GHP with this specialized initialization will terminate in  $3n$  steps with a d-stationary solution of (2) at  $\gamma^*$ .

# The GHP method

## Back to the pivoting method

- Suppose we are at a discontinuous  $\gamma^*$  which we need to restore.
- Initialize the GHP method with the current “non-d-stationary” solution, we can prove the following theorem.

## Main theorem of GHP (informal)

GHP with this specialized initialization will terminate in  $3n$  steps with a d-stationary solution of (2) at  $\gamma^*$ .

## Remark

- [At most  $3n$  GHP subproblems] + [each  $\mathcal{O}(n^3)$  by  $Q$  Stieltjes]  $\implies$  [total complexity  $\mathcal{O}(n^4)$ ].
- It is inductively proved by leveraging a key property of Stieltjes matrices named the “least element property”. [Pang, 1979].

# Numerical experiments: GMRF

## Gaussian Markov Random Field

The maximum a posteriori estimation of Gaussian Markov random field (GMRF) naturally gives rise to the Stieltjes structure.

Given graph  $(V, E)$ :

$$\underset{x}{\text{minimize}} \quad \sum_{i \in V} \frac{1}{\sigma_i^2} (y_i - x_i)^2 + \sum_{(i,j) \in E} \frac{1}{d_{ij}} (x_i - x_j)^2$$

# Numerical experiments: GMRF

## Summary: capped $\ell_1$ vs. $\ell_1$ vs. $\ell_0$

- **Settings:**

- $n \in \{100, 10000\}$ , noise level  $\in (0, 1]$
- $\ell_0$  (only for  $n = 100$ ),  $\ell_1$ , capped  $\ell_1$
- $\delta \in \{10, 1, 10^{-1}, 10^{-4}\}$  for capped  $\ell_1$  (**controls its approx. to  $\ell_0$** )



# Numerical experiments: GMRF

## Summary: capped $\ell_1$ vs. $\ell_1$ vs. $\ell_0$

### • Settings:

- $n \in \{100, 10000\}$ , noise level  $\in (0, 1]$
- $\ell_0$  (only for  $n = 100$ ),  $\ell_1$ , capped  $\ell_1$
- $\delta \in \{10, 1, 10^{-1}, 10^{-4}\}$  for capped  $\ell_1$

### • Computation time:

- Capped  $\ell_1$  d-stat. path can be **60 - 3,000 times faster** than  $\ell_0$ .
- When  $\delta$  is large, capped  $\ell_1$  is basically the same as  $\ell_1$ .
- When  $\delta$  is small, capped  $\ell_1$  needs more effort  
**(can be 10 times slower, no free lunch).**

## Theorem (no free lunch)

When  $\delta < |\bar{x}_i^0|, \forall i \in [n]$  where  $\bar{x}^0$  is the unique solution at  $\gamma = 0$ , then the unique continuous d-stat. path is  $\bar{x}(\gamma) = \bar{x}^0, \forall \gamma \geq 0$ .

# Numerical experiments: GMRF

## Summary: capped $\ell_1$ vs. $\ell_1$ vs. $\ell_0$

- **Settings:**

- $n \in \{100, 10000\}$ , noise level  $\in (0, 1]$
- $\ell_0$  (only for  $n = 100$ ),  $\ell_1$ , capped  $\ell_1$
- $\delta \in \{10, 1, 10^{-1}, 10^{-4}\}$  for capped  $\ell_1$

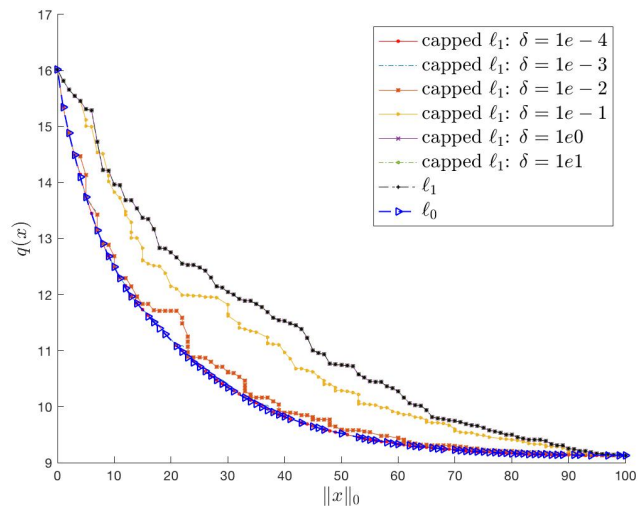
- **Computation time:**

- Capped  $\ell_1$  d-stat. path can be 60 - 3,000 times faster than  $\ell_0$ .
- When  $\delta$  is large, capped  $\ell_1$  is basically the same as  $\ell_1$ .
- When  $\delta$  is small, capped  $\ell_1$  needs more effort (can be 10 times slower, no free lunch).

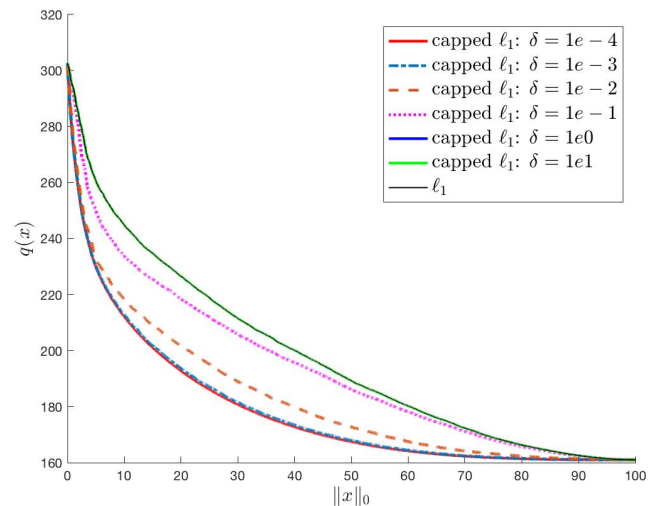
- **Loss value:**

- Capped  $\ell_1$  achieves better loss value ( $\frac{1}{2}x^\top Qx + q^\top x$ ) than  $\ell_1$  when the solution sparsity is the same!

# GMRF: loss vs. sparsity



(a)  $n = 100$



(b)  $n = 10000$

- When  $\delta$  is small, capped  $\ell_1$  behaves like  $\ell_0$  (blue curve in the bottom).
- When  $\delta$  is large, capped  $\ell_1$  behaves like  $\ell_1$  (black curve on the top).
- Capped  $\ell_1$  trade-off (acceptable) computation time to gain superior optimization and statistical (to be shown) performance.

# Hyperparameter selection

- We care about the following quantities ( $X$  is some ground truth)

**Signal recovery:**  $\sum_{i=1}^p \sum_{j=1}^p (x_{ij}^*(\gamma) - X_{ij})^2$

**Support recovery:**  $\sum_{i=1}^p \sum_{j=1}^p \left| |x_{ij}^*(\gamma)|_0 - |X_{ij}|_0 \right|$

# Hyperparameter selection

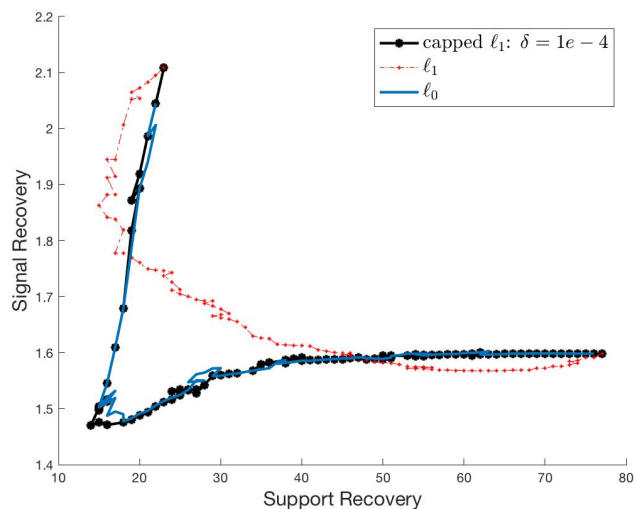
- We care about the following quantities ( $X$  is some ground truth)

$$\text{Signal recovery: } \sum_{i=1}^p \sum_{j=1}^p (x_{ij}^*(\gamma) - X_{ij})^2$$

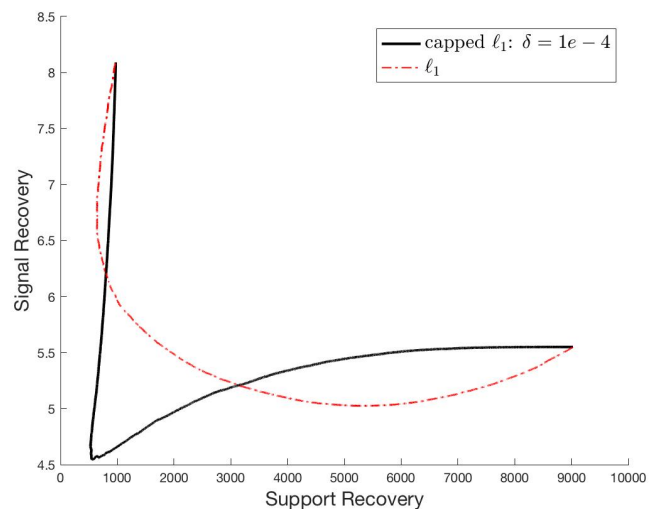
$$\text{Support recovery: } \sum_{i=1}^p \sum_{j=1}^p \left| |x_{ij}^*(\gamma)|_0 - |X_{ij}|_0 \right|$$

- Capped  $\ell_1$  d-stat. path is superior in hyperparameter selection
  - The best capped  $\ell_1$  solution dominates  $\ell_1$ 's in both quantities.
  - Capped  $\ell_1$  achieves the minimum of both quantities at the same  $\gamma$ .
  - $\ell_1$  cannot achieve this: we always have to sacrifice one of them.

# Signal vs. support recovery



(a)  $n = 100$

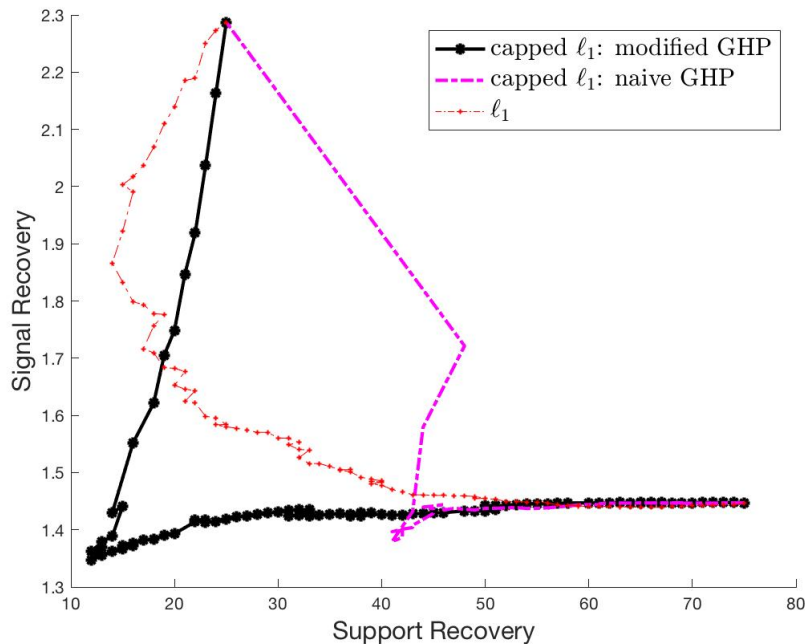


(b)  $n = 10000$

- For  $\ell_0$ , capped  $\ell_1$ , good signal and support recovery are highly correlated.
- But for  $\ell_1$ , we always have to sacrifice one quantity.

# GHP restoration

Our specialized initialization for GHP restoration leverages the most recent basis, which turns out to be the key for all the nice properties.



**Hyperparameter selection from the naïve initialization is worse than  $\ell_1$ .**

# Numerical experiments: summary

- **Computation time:**

- Capped  $\ell_1$  path is much more scalable than  $\ell_0$  path.
- Capped  $\ell_1$  path (small  $\delta$ ) needs more (acceptable) effort than  $\ell_1$ .



# Numerical experiments: summary

- **Computation time:**

- Capped  $\ell_1$  path is much more scalable than  $\ell_0$  path.
- Capped  $\ell_1$  path (small  $\delta$ ) needs more (acceptable) effort than  $\ell_1$ .

- **Optimization performance:**

- Capped  $\ell_1$  path (small  $\delta$ ) is a better approximation to the  $\ell_0$  path.
- For the same sparsity, capped  $\ell_1$  achieves better loss than  $\ell_1$ .

# Numerical experiments: summary

- **Computation time:**

- Capped  $\ell_1$  path is much more scalable than  $\ell_0$  path.
- Capped  $\ell_1$  path (small  $\delta$ ) needs more (acceptable) effort than  $\ell_1$ .

- **Optimization performance:**

- Capped  $\ell_1$  path (small  $\delta$ ) is a better approximation to the  $\ell_0$  path.
- For the same sparsity, capped  $\ell_1$  achieves better loss than  $\ell_1$ .

- **Hyperparameter selection:**

- Capped  $\ell_1$ : minimal signal and support recovery at the same  $\gamma$ .
- $\ell_1$  path does not have such  $\gamma$ .

# Numerical experiments: summary

- **Computation time:**

- Capped  $\ell_1$  path is much more scalable than  $\ell_0$  path.
- Capped  $\ell_1$  path (small  $\delta$ ) needs more (acceptable) effort than  $\ell_1$ .

- **Optimization performance:**








- Capped  $\ell_1$  path (small  $\delta$ ) is a better approximation to the  $\ell_0$  path.
- For the same sparsity, capped  $\ell_1$  always achieves better loss than  $\ell_1$ .

- **Hyperparameter selection:**

- Capped  $\ell_1$ : minimal signal and support recovery at the same  $\gamma$ .
- $\ell_1$  path does not have such  $\gamma$ .

- The **GHP** restoration is critical for all the nice practical properties.

# Reference

-  Efron, B., Hastie, T., Johnstone, I., & Tibshirani, R. (2004). Least angle regression. *The Annals of statistics*, 32(2), 407-499.
-  Mairal, J., & Yu, B. (2012). Complexity analysis of the lasso regularization path. *arXiv preprint arXiv:1205.0079*.
-  Soussen, C., Idier, J., Duan, J., & Brie, D. (2015). Homotopy Based Algorithms for  $\ell_0$ -Regularized Least-Squares. *IEEE Transactions on Signal Processing*, 63(13), 3301-3316.
-  Yukawa, M., & Amari, S. I. (2015).  $\ell_p$ -Regularized Least Squares ( $0 < p < 1$ ) and Critical Path. *IEEE Transactions on Information Theory*, 62(1), 488-502.
-  Zhang, C. H. (2010). Nearly Unbiased Variable Selection Under Minimax Concave Penalty. *The Annals of statistics*, 38(2), 894-942.
-  Pang, J. S. (1979). On a class of least-element complementarity problems. *Mathematical Programming*, 16(1), 111-126.
-  Pang, J. S., & Chandrasekaran, R. (1985). Linear complementarity problems solvable by a polynomially bounded pivoting algorithm. In *Mathematical Programming Essays in Honor of George B. Dantzig Part II* (pp. 13-27). Springer, Berlin, Heidelberg.