# The Trimmed Lasso:

## Sparse recovery guarantees and practical optimization by the Generalized Soft-Min Penalty

Boaz Nadler

Weizmann Institute of Science

Joint work Tal Amir and Ronen Basri

Statistics in the Big Data Era

June 2022

# Peter's Non-Sparse Influence on My Work

- Some theory for Fisher's LDA... when there are many more variables than observations, 2004'

- Some theory for Fisher's LDA... when there are many more variables than observations, 2004'

$\rightarrow$ The prediction error in PLS and CLS, 05'

# Peter's Non-Sparse Influence on My Work

- Some theory for Fisher's LDA... when there are many more variables than observations, 2004'

$\rightarrow$ The prediction error in PLS and CLS, 05'

- Covariance Regularization by Thresholding, 08'

$\rightarrow$ Minimax bounds on sparse PCA,

- Some theory for Fisher's LDA... when there are many more variables than observations, 2004'

$\rightarrow$ The prediction error in PLS and CLS, 05'

- Covariance Regularization by Thresholding, 08'

$\rightarrow$ Minimax bounds on sparse PCA,

- Simultaneous analysis of Lasso and Dantzig, 09'

- Some theory for Fisher's LDA... when there are many more variables than observations, 2004'

$\rightarrow$ The prediction error in PLS and CLS, 05'

- Covariance Regularization by Thresholding, 08'

$\rightarrow$ Minimax bounds on sparse PCA,

- Simultaneous analysis of Lasso and Dantzig, 09'

Today's talk: Sparse Linear Regression

# Sparse Approximation / Best subset selection

**Problem setup:**

Observe

(i) $n \times d$ matrix $A$

(ii) response vector $\boldsymbol{y} \in \mathbb{R}^n$

# Sparse Approximation / Best subset selection

**Problem setup:**

Observe
  (i) $n \times d$ matrix $A$
 (ii) response vector $\boldsymbol{y} \in \mathbb{R}^n$

Given sparsity parameter $k$

# Sparse Approximation / Best subset selection

**Problem setup:**

Observe

(i) $n \times d$ matrix $A$

(ii) response vector $\boldsymbol{y} \in \mathbb{R}^n$

Given sparsity parameter $k$

solve

$$\min_{\boldsymbol{x}} \|A\boldsymbol{x} - \boldsymbol{y}\|_2 \quad \text{subject to} \quad \|\boldsymbol{x}\|_0 \leq k \qquad \text{(P0)}$$

# Sparse Approximation

$$\min_{\boldsymbol{x}} \|A\boldsymbol{x} - \boldsymbol{y}\|_2 \quad \text{subject to} \quad \|\boldsymbol{x}\|_0 \leq k \qquad \text{(P0)}$$

# Sparse Approximation

$$\min_{\boldsymbol{x}} \|A\boldsymbol{x} - \boldsymbol{y}\|_2 \quad \text{subject to} \quad \|\boldsymbol{x}\|_0 \leq k \qquad \text{(P0)}$$

**Signal/Image processing**:

$\boldsymbol{y} = (y_1, \ldots, y_n)$ are $n$ samples of unknown function

$A =$ dictionary, whose columns are basic signals / atoms

$$\min_{\boldsymbol{x}} \|A\boldsymbol{x} - \boldsymbol{y}\|_2 \quad \text{subject to} \quad \|\boldsymbol{x}\|_0 \leq k \qquad \text{(P0)}$$

**Signal/Image processing**:

$\boldsymbol{y} = (y_1, \ldots, y_n)$ are $n$ samples of unknown function

$A =$ dictionary, whose columns are basic signals / atoms

Seek best representation of $\boldsymbol{y}$ by at most $k$ *dictionary atoms*.

$$\min_{\boldsymbol{x}} \|A\boldsymbol{x} - \boldsymbol{y}\|_2 \quad \text{subject to } \|\boldsymbol{x}\|_0 \leq k \qquad \text{(P0)}$$

**Signal/Image processing**:

$\boldsymbol{y} = (y_1, \ldots, y_n)$ are $n$ samples of unknown function

$A =$ dictionary, whose columns are basic signals / atoms

Seek best representation of $\boldsymbol{y}$ by at most $k$ *dictionary atoms*.

**Compressed sensing**:

Wish to recover unknown signal $\boldsymbol{x} \in \mathbb{R}^d$, from $n$ noisy observations

$$y_i = \mathbf{w}_i^\top \boldsymbol{x} + \sigma \xi_i$$

Assume that $\boldsymbol{x}$ is (approximately) $k$-sparse

# Sparse Approximation

**Statistics:** sparse linear regression

given $n$ observations $(X_i, y_i)$, assumed of the form

$$y = X^\top \beta + \varepsilon$$

$y$ is a response variable that we wish to predict from an explanatory vector $X \in \mathbb{R}^d$

# Sparse Approximation

**Statistics:** sparse linear regression

given $n$ observations $(X_i, y_i)$, assumed of the form

$$y = X^\top \beta + \varepsilon$$

$y$ is a response variable that we wish to predict from an explanatory vector $X \in \mathbb{R}^d$

...using at most $k$ explanatory variables.

Often $k$ is unknown and needs to be estimated

Often $k$ is unknown and needs to be estimated

A common approach: Solve (P0) for several values of $k$ and apply:

- Cross validation
- Model selection criterion

# Sparsity parameter $k$

Often $k$ is unknown and needs to be estimated

A common approach: Solve (P0) for several values of $k$ and apply:

- Cross validation
- Model selection criterion

In rest of talk: Assume $k$ is given

Often $k$ is unknown and needs to be estimated

A common approach: Solve (P0) for several values of $k$ and apply:

- Cross validation
- Model selection criterion

In rest of talk: Assume $k$ is given

Focus on solving (P0) for a *given* value of $k$

The key challenge in solving (P0) is *support detection*, finding the optimal $k$ columns of $A$ to include in the solution

# Support Detection

The key challenge in solving (P0) is *support detection*, finding the optimal $k$ columns of $A$ to include in the solution

Once support has been found, problem reduces to solving least squares on these $k$ columns.

# Support Detection

The key challenge in solving (P0) is *support detection*, finding the optimal $k$ columns of $A$ to include in the solution

Once support has been found, problem reduces to solving least squares on these $k$ columns.

[Natarajan 95', Davis et al 97']

Unfortunately, this problem is NP-hard...

# Support Detection

The key challenge in solving (P0) is *support detection*, finding the optimal $k$ columns of $A$ to include in the solution

Once support has been found, problem reduces to solving least squares on these $k$ columns.

[Natarajan 95', Davis et al 97']

Unfortunately, this problem is NP-hard...

Yet, extensive *prior work*, on algorithms, theory, lower bounds, etc.

# Support Detection

The key challenge in solving (P0) is *support detection*, finding the optimal $k$ columns of $A$ to include in the solution

Once support has been found, problem reduces to solving least squares on these $k$ columns.

[Natarajan 95', Davis et al 97']

Unfortunately, this problem is NP-hard...

Yet, extensive *prior work*, on algorithms, theory, lower bounds, etc.

Over a hundred methods to approximately solve (P0)

# Support Detection

The key challenge in solving (P0) is *support detection*, finding the optimal $k$ columns of $A$ to include in the solution

Once support has been found, problem reduces to solving least squares on these $k$ columns.

[Natarajan 95', Davis et al 97']

Unfortunately, this problem is NP-hard...

Yet, extensive *prior work*, on algorithms, theory, lower bounds, etc.

Over a hundred methods to approximately solve (P0)
lots of theoretical results, recovery guarantees, etc.

(Almost) all prior work on (P0) in 3 slides...

# Previous Work

**Greedy methods:**

- ○ Matching Pursuit algorithms
    - Orthogonal Matching Pursuit (OMP), CoSaMP [Needell, Tropp, ACHA 2009] and more
- ○ Iterative Hard Thresholding [Blumensath, Davies, ACHA 2009]
- ○ Iterative Support Detection (ISD) [Wang, Yin, Im. Sc. 2010]
- ○ Forward stepwise linear regression (1960's), etc.

# Previous Work

**Greedy methods:**

- Matching Pursuit algorithms
    - Orthogonal Matching Pursuit (OMP), CoSaMP [Needell, Tropp, ACHA 2009] and more
- Iterative Hard Thresholding [Blumensath, Davies, ACHA 2009]
- Iterative Support Detection (ISD) [Wang, Yin, Im. Sc. 2010]
- Forward stepwise linear regression (1960's), etc.

Advantages: Easy to program, run very fast.

# Previous Work

**Greedy methods:**

- ○ Matching Pursuit algorithms
  - Orthogonal Matching Pursuit (OMP), CoSaMP [Needell, Tropp, ACHA 2009] and more
- ○ Iterative Hard Thresholding [Blumensath, Davies, ACHA 2009]
- ○ Iterative Support Detection (ISD) [Wang, Yin, Im. Sc. 2010]
- ○ Forward stepwise linear regression (1960's), etc.

Advantages: Easy to program, run very fast.

Limitation: May yield suboptimal solutions.

# Penalty Methods

Replace constraint $\|x\|_0 \le k$ by a penalty $\rho(x)$:

$$\min_{x} \ \tfrac{1}{2}\|Ax - y\|^2 + \lambda\rho(x).$$

# Penalty Methods

Replace constraint $\|x\|_0 \leq k$ by a penalty $\rho(x)$:

$$\min_{x} \ \tfrac{1}{2}\|Ax - y\|^2 + \lambda\rho(x).$$

- To obtain a $k$-sparse solution, $\lambda$ needs to be tuned.

## Penalty Methods

Replace constraint $\|x\|_0 \leq k$ by a penalty $\rho(x)$:

$$\min_{x} \tfrac{1}{2}\|Ax - y\|^2 + \lambda\rho(x).$$

- To obtain a $k$-sparse solution, $\lambda$ needs to be tuned.

The most popular penalty is the convex *lasso*: $\rho(x) = \|x\|_1$

# Penalty Methods

Replace constraint $\|x\|_0 \leq k$ by a penalty $\rho(x)$:

$$\min_x \; \tfrac{1}{2}\|Ax - y\|^2 + \lambda\rho(x).$$

- To obtain a $k$-sparse solution, $\lambda$ needs to be tuned.

The most popular penalty is the convex *lasso*: $\rho(x) = \|x\|_1$

**Lasso:**

- Recovery guarantees under various conditions (Incoherence, RIP, Restricted Eigenvalue, ...)
- Fast optimization schemes developed

# Penalty Methods

Replace constraint $\|\boldsymbol{x}\|_0 \leq k$ by a penalty $\rho(\boldsymbol{x})$:

$$\min_{\boldsymbol{x}} \ \tfrac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|^2 + \lambda\rho(\boldsymbol{x}).$$

- To obtain a $k$-sparse solution, $\lambda$ needs to be tuned.

The most popular penalty is the convex *lasso*: $\rho(\boldsymbol{x}) = \|\boldsymbol{x}\|_1$

**Lasso:**

- Recovery guarantees under various conditions (Incoherence, RIP, Restricted Eigenvalue, ...)
- Fast optimization schemes developed
- May yield suboptimal solutions

# Exact / Approximate Mixed Integer Programming

- During optimization, calculate lower bound for objective
- If current objective equals lower bound, terminate with a global optimality certificate.

[Bertsimas, King, Mazumder, AoS '16]

# Exact / Approximate Mixed Integer Programming

- During optimization, calculate lower bound for objective
- If current objective equals lower bound, terminate with a global optimality certificate.

  [Bertsimas, King, Mazumder, AoS '16]

- MIP solves (P0) *globally*
- Applicable with $d = O(100)$, much faster than exhaustive search

# Exact / Approximate Mixed Integer Programming

- During optimization, calculate lower bound for objective
- If current objective equals lower bound, terminate with a global optimality certificate.

  [Bertsimas, King, Mazumder, AoS '16]

- MIP solves (P0) *globally*
- Applicable with $d = O(100)$, much faster than exhaustive search

Limitation: May be very slow

# Exact / Approximate Mixed Integer Programming

○ During optimization, calculate lower bound for objective

○ If current objective equals lower bound, terminate with a global optimality certificate.

[Bertsimas, King, Mazumder, AoS '16]

○ MIP solves (P0) *globally*

○ Applicable with $d = O(100)$, much faster than exhaustive search

Limitation: May be very slow

○ On $30 \times 180$ matrix $A$ and $k = 15$, may take several days

# Exact / Approximate Mixed Integer Programming

- During optimization, calculate lower bound for objective
- If current objective equals lower bound, terminate with a global optimality certificate.

  [Bertsimas, King, Mazumder, AoS '16]

- MIP solves (P0) *globally*
- Applicable with $d = O(100)$, much faster than exhaustive search

Limitation: May be very slow

- On $30 \times 180$ matrix $A$ and $k = 15$, may take several days

  [Bertsimas, Van Parys, AoS '20]

**Cutting plane method**

  globally solve $d = 15000$, $n = 200$, $k = 10$ in minutes

# Approximate MIP

[Hazimeh & Mazumder, Oper. Res. '20]

Greedy coordinate descent + local combinatorial search

# Approximate MIP

[Hazimeh & Mazumder, Oper. Res. '20]

Greedy coordinate descent + local combinatorial search

- − No optimality certificate
- − Extremely fast, can handle $d = 10^6$ in less than a minute
- − state of the art performance

In addition to algorithm development, substantial body of literature on conditions for perfect recovery (noiseless setting), accurate and stable recovery in presence of noise.

# Theoretical Guarantees

In addition to algorithm development, substantial body of literature on conditions for perfect recovery (noiseless setting), accurate and stable recovery in presence of noise.

**Key notions:** Coherence of dictionary, restricted isometry property, etc.

# Theoretical Guarantees

In addition to algorithm development, substantial body of
literature on conditions for perfect recovery (noiseless setting),
accurate and stable recovery in presence of noise.

**Key notions:** Coherence of dictionary, restricted isometry
property, etc.

Under some conditions, current methods are *optimal*

In addition to algorithm development, substantial body of literature on conditions for perfect recovery (noiseless setting), accurate and stable recovery in presence of noise.

**Key notions:** Coherence of dictionary, restricted isometry property, etc.

Under some conditions, current methods are *optimal*

Has the problem not been solved yet?

In addition to algorithm development, substantial body of literature on conditions for perfect recovery (noiseless setting), accurate and stable recovery in presence of noise.

**Key notions:** Coherence of dictionary, restricted isometry property, etc.

Under some conditions, current methods are *optimal*

Has the problem not been solved yet?

**No !**

In addition to algorithm development, substantial body of literature on conditions for perfect recovery (noiseless setting), accurate and stable recovery in presence of noise.

**Key notions:** Coherence of dictionary, restricted isometry property, etc.

Under some conditions, current methods are *optimal*

Has the problem not been solved yet?

**No !**

Key limitation of above methods:
with few observations $n \ll d$,
higher values of $k$ (not so sparse vectors)
nearly all prior methods either compute far from optimal solutions
or run essentially forever...

## Example

Matrix $A$ of size $100 \times 800$, random i.i.d. $\mathcal{N}(0,1)$ entries followed by column normalization.

## Example

Matrix $A$ of size $100 \times 800$, random i.i.d. $\mathcal{N}(0,1)$ entries followed by column normalization.

For various sparsity values $k$, generate random $k$-sparse vector $\boldsymbol{x}_0$. Its non-zero entries are i.i.d. $\mathcal{N}(0,1)$.

## Example

Matrix $A$ of size $100 \times 800$, random i.i.d. $\mathcal{N}(0, 1)$ entries followed by column normalization.

For various sparsity values $k$, generate random $k$-sparse vector $\mathbf{x}_0$. Its non-zero entries are i.i.d. $\mathcal{N}(0, 1)$.
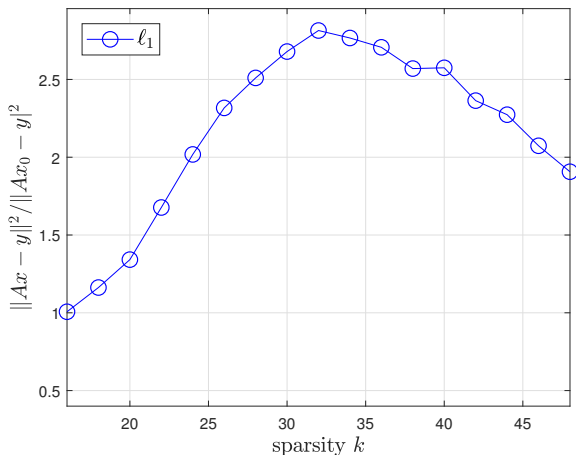
Generate

$$\mathbf{y} = A\mathbf{x}_0 + \mathbf{e}$$

where vector $\mathbf{e} \sim \sigma \mathcal{N}(\mathbf{0}, I_n)$, with $\mathbb{E}\|\mathbf{e}\|^2 = (0.05)^2 \cdot \mathbb{E}\|A\mathbf{x}_0\|^2$.

## Example

Matrix $A$ of size $100 \times 800$, random i.i.d. $\mathcal{N}(0,1)$ entries followed by column normalization.

For various sparsity values $k$, generate random $k$-sparse vector $\boldsymbol{x}_0$. Its non-zero entries are i.i.d. $\mathcal{N}(0,1)$.

Generate

$$\boldsymbol{y} = A\boldsymbol{x}_0 + \mathbf{e}$$

where vector $\mathbf{e} \sim \sigma\mathcal{N}(\mathbf{0}, I_n)$, with $\mathbb{E}\|\mathbf{e}\|^2 = (0.05)^2 \cdot \mathbb{E}\|A\boldsymbol{x}_0\|^2$.

**Measure of optimization success:**

$$\frac{\|A\hat{\boldsymbol{x}} - \boldsymbol{y}\|}{\|A\boldsymbol{x}_0 - \boldsymbol{y}\|}.$$

## Example

Matrix $A$ of size $100 \times 800$, random i.i.d. $\mathcal{N}(0,1)$ entries followed by column normalization.

For various sparsity values $k$, generate random $k$-sparse vector $\mathbf{x}_0$. Its non-zero entries are i.i.d. $\mathcal{N}(0,1)$.

Generate

$$\mathbf{y} = A\mathbf{x}_0 + \mathbf{e}$$

where vector $\mathbf{e} \sim \sigma\mathcal{N}(\mathbf{0}, I_n)$, with $\mathbb{E}\|\mathbf{e}\|^2 = (0.05)^2 \cdot \mathbb{E}\|A\mathbf{x}_0\|^2$.

**Measure of optimization success:**

$$\frac{\|A\hat{\mathbf{x}} - \mathbf{y}\|}{\|A\mathbf{x}_0 - \mathbf{y}\|}.$$

If ratio $\leq 1$ then $\hat{\mathbf{x}}$ is *potentially* accurate estimate of $\mathbf{x}_0$

# An Example



In our setting, $\ell_1$ penalty (Lasso / Basis Pursuit) essentially works only up to sparsity levels $k \leq 16$.
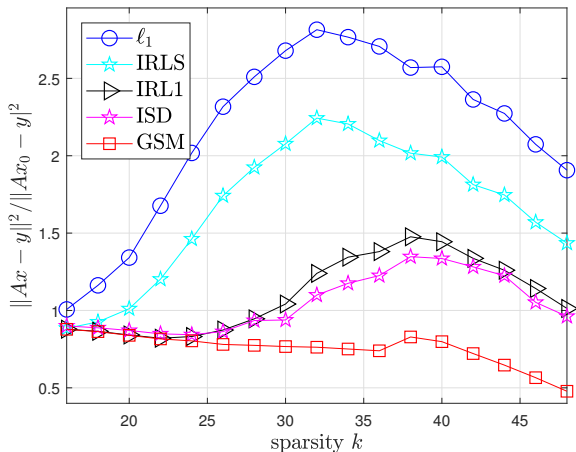
# An Example



IRLS and IRL-1 solve $\ell_q$ penalized objectives with $q < 1$. Solved with 10 values of $q < 1$ and took solution with minimal $\|A\boldsymbol{x} - \boldsymbol{y}\|$.

# An Example



ISD=Iterative Support Detection [Wang & Yin 2010'].
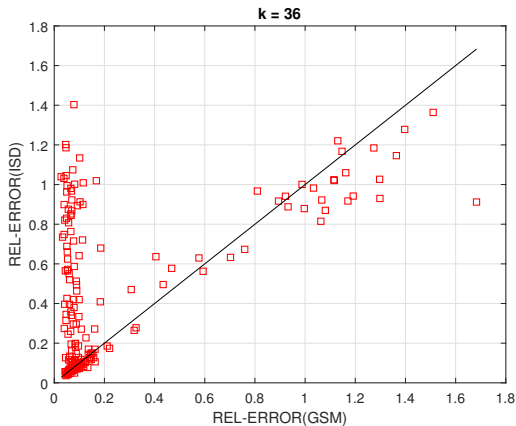Sophisticated greedy support-detection strategy.

GSM= our proposed method. Superior at the more challenging settings with larger values of $k$ and/or correlated dictionaries

# An Example

Successful optimization often (but not always) translates into better recovery



Showing $\|\hat{x} - x_0\|_1 / \|x_0\|_1$

**Desired properties for a penalty function:**

**Desired properties for a penalty function:**

(i) A penalty $\rho(\boldsymbol{x}) = \rho_k(\boldsymbol{x})$ that *explicitly* takes into account the sparsity level $k$

**Desired properties for a penalty function:**

(i) A penalty $\rho(x) = \rho_k(x)$ that *explicitly* takes into account the sparsity level $k$

(ii) For large $\lambda$, solutions of

$$\min \|Ax - y\|_2^2 + \lambda \rho_k(x)$$

are close to those of (P0).

○ Better yet - they *coincide*

# Solving (P0) by a Penalized Objective

**Desired properties for a penalty function:**

(i) A penalty $\rho(\boldsymbol{x}) = \rho_k(\boldsymbol{x})$ that *explicitly* takes into account the sparsity level $k$

(ii) For large $\lambda$, solutions of

$$\min \|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\rho_k(\boldsymbol{x})$$

are close to those of (P0).

○ Better yet - they *coincide*

(iii) Objective would be easy to optimize

## The Trimmed Lasso

A penalty that satisfies (i) and (ii) above: (Not our contribution)

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

where $|x|_{(1)} \geq |x|_{(2)} \geq \ldots \geq |x|_{(d)}$ are the entries of $\boldsymbol{x}$ in absolute value, sorted in decreasing order

# The Trimmed Lasso

A penalty that satisfies (i) and (ii) above: (Not our contribution)

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

where $|x|_{(1)} \geq |x|_{(2)} \geq \ldots \geq |x|_{(d)}$ are the entries of $\boldsymbol{x}$ in absolute value, sorted in decreasing order

Penalize "tail" of $\boldsymbol{x}$: the $\ell_1$ distance to the nearest $k$-sparse vector

# The Trimmed Lasso

A penalty that satisfies (i) and (ii) above: (Not our contribution)

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

where $|x|_{(1)} \geq |x|_{(2)} \geq \ldots \geq |x|_{(d)}$ are the entries of $\boldsymbol{x}$ in absolute value, sorted in decreasing order

Penalize "tail" of $\boldsymbol{x}$: the $\ell_1$ distance to the nearest $k$-sparse vector

Early related works:
- [Cohen, Dahmen, DeVore, *JAMS* '08]
- [Huang, Liu, Shi, Van Huffel, Suykens, *Sig. Proc.* '15]

# The Trimmed Lasso

A penalty that satisfies (i) and (ii) above: (Not our contribution)

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

where $|x|_{(1)} \geq |x|_{(2)} \geq \ldots \geq |x|_{(d)}$ are the entries of $\boldsymbol{x}$ in absolute value, sorted in decreasing order

Penalize "tail" of $\boldsymbol{x}$: the $\ell_1$ distance to the nearest $k$-sparse vector

Early related works:

- [Cohen, Dahmen, DeVore, *JAMS* '08]
- [Huang, Liu, Shi, Van Huffel, Suykens, *Sig. Proc.* '15]

Penalty studied by:

- [Gotoh, Takeda, Tono, *Math. Prog.* '18]
- [Bertsimas, Copenhaver, Mazumder, '17], who coined the term *trimmed Lasso*

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

**Theoretical questions:**

# The Trimmed Lasso

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

**Theoretical questions:**

1. Relation to original problem (P0)?

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

**Theoretical questions:**

1. Relation to original problem (P0)?
2. What value to use for $\lambda$?

# The Trimmed Lasso

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

**Theoretical questions:**

1. Relation to original problem (P0)?
2. What value to use for $\lambda$?
3. Can we recover $\boldsymbol{x}$ using $\tau_k(\boldsymbol{x})$?

# The Trimmed Lasso

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

**Theoretical questions:**

1. Relation to original problem (P0)?
2. What value to use for $\lambda$?
3. Can we recover $\boldsymbol{x}$ using $\tau_k(\boldsymbol{x})$?

**Practical question:** How to optimize an objective with $\tau_k(\boldsymbol{x})$?

# The Trimmed Lasso

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

**Theoretical questions:**

1. Relation to original problem (P0)?
2. What value to use for $\lambda$?
3. Can we recover $\boldsymbol{x}$ using $\tau_k(\boldsymbol{x})$?

**Practical question:** How to optimize an objective with $\tau_k(\boldsymbol{x})$?

**Our contribution:**

# The Trimmed Lasso

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

**Theoretical questions:**

1. Relation to original problem (P0)?
2. What value to use for $\lambda$?
3. Can we recover $\boldsymbol{x}$ using $\tau_k(\boldsymbol{x})$?

**Practical question:** How to optimize an objective with $\tau_k(\boldsymbol{x})$?

**Our contribution:**

1. Theoretical study of $\tau_k(\boldsymbol{x})$, addressing questions 1-3

# The Trimmed Lasso

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

**Theoretical questions:**

1. Relation to original problem (P0)?
2. What value to use for $\lambda$?
3. Can we recover $\boldsymbol{x}$ using $\tau_k(\boldsymbol{x})$?

**Practical question:** How to optimize an objective with $\tau_k(\boldsymbol{x})$?

**Our contribution:**

1. Theoretical study of $\tau_k(\boldsymbol{x})$, addressing questions 1-3
   $\rightarrow$ $\tau_k(\boldsymbol{x})$ is a good candidate for solving (P0)

# The Trimmed Lasso

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

**Theoretical questions:**

1. Relation to original problem (P0)?
2. What value to use for $\lambda$?
3. Can we recover $\boldsymbol{x}$ using $\tau_k(\boldsymbol{x})$?

**Practical question:** How to optimize an objective with $\tau_k(\boldsymbol{x})$?

**Our contribution:**

1. Theoretical study of $\tau_k(\boldsymbol{x})$, addressing questions 1-3
   $\rightarrow \tau_k(\boldsymbol{x})$ is a good candidate for solving (P0)
2. Novel surrogate penalty that satisfies (i)-(iii)

## The Trimmed Lasso

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

**Theoretical questions:**

1. Relation to original problem (P0)?
2. What value to use for $\lambda$?
3. Can we recover $\boldsymbol{x}$ using $\tau_k(\boldsymbol{x})$?

**Practical question:** How to optimize an objective with $\tau_k(\boldsymbol{x})$?

**Our contribution:**

1. Theoretical study of $\tau_k(\boldsymbol{x})$, addressing questions 1-3
   $\rightarrow \tau_k(\boldsymbol{x})$ is a good candidate for solving (P0)
2. Novel surrogate penalty that satisfies (i)-(iii)
3. Practical optimization method, state-of-the-art results

$$\min_{\boldsymbol{x}} \ \mathsf{F}_\lambda(\boldsymbol{x}) := \tfrac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_k(\boldsymbol{x}) \qquad\qquad (\mathsf{P}_\lambda)$$

$$\min_{\boldsymbol{x}} \; F_\lambda(\boldsymbol{x}) := \tfrac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_k(\boldsymbol{x}) \qquad (P_\lambda)$$

**How to choose $\lambda$?**

# The Trimmed Lasso: Choosing $\lambda$

Define $\beta = \max_{i=1,\dots,d} \|\boldsymbol{a}_i\|_2$, where $\boldsymbol{a}_i$ are the columns of $A$.

### Lemma

*If $\lambda > \bar{\lambda} = \beta \|\boldsymbol{y}\|_2$, then any local minimum of $(P_\lambda)$ is k-sparse.*

Define $\beta = \max_{i=1,\ldots,d} \|\boldsymbol{a}_i\|_2$, where $\boldsymbol{a}_i$ are the columns of $A$.

### Lemma

*If $\lambda > \bar{\lambda} = \beta \|\boldsymbol{y}\|_2$, then any local minimum of $(P_\lambda)$ is k-sparse.*

Define $\beta = \max_{i=1,\ldots,d} \|\boldsymbol{a}_i\|_2$, where $\boldsymbol{a}_i$ are the columns of $A$.

### Lemma

*If $\lambda > \bar{\lambda} = \beta\|\boldsymbol{y}\|_2$, then any local minimum of $(P_\lambda)$ is k-sparse.*

○ For large enough $\lambda$, optimal solutions of $(P_\lambda)$ coincide with those of (P0).

Define $\beta = \max_{i=1,\ldots,d} \|\boldsymbol{a}_i\|_2$, where $\boldsymbol{a}_i$ are the columns of $A$.

### Lemma

If $\lambda > \bar{\lambda} = \beta \|\boldsymbol{y}\|_2$, then any local minimum of $(P_\lambda)$ is k-sparse.

- For large enough $\lambda$, optimal solutions of $(P_\lambda)$ coincide with those of (P0).
- Strategy: Solve with increasing values of $\lambda$, until a *k*-sparse solution is obtained.
  - $\rightarrow$ Guaranteed to happen when $\lambda$ surpasses the threshold.

# Sparse Signal Recovery Guarantees

Suppose that

$$y = Ax_0 + e \in \mathbb{R}^n$$

$x_0 \in \mathbb{R}^d$ = unknown vector to be recovered
$e$ = measurement error

# Sparse Signal Recovery Guarantees

Suppose that

$$y = Ax_0 + e \in \mathbb{R}^n$$

$x_0 \in \mathbb{R}^d$ = unknown vector to be recovered
$e$ = measurement error

## Assumptions:

$x_0$ is approximately $k$-sparse $(\tau_k(x_0) \ll \|x_0\|_1)$

$\|e\|_2$ is small

# Sparse Signal Recovery Guarantees

Suppose that

$$y = Ax_0 + e \in \mathbb{R}^n$$

$x_0 \in \mathbb{R}^d =$ unknown vector to be recovered
$e =$ measurement error

**Assumptions:**

$x_0$ is approximately $k$-sparse $(\tau_k(x_0) \ll \|x_0\|_1)$

$\|e\|_2$ is small

**Goal:** Recover $x_0$ given $A, y$ and $k$.

Suppose that

$$y = Ax_0 + e \in \mathbb{R}^n$$

$x_0 \in \mathbb{R}^d =$ unknown vector to be recovered
$e =$ measurement error

**Assumptions:**

$x_0$ is approximately $k$-sparse $(\tau_k(x_0) \ll \|x_0\|_1)$

$\|e\|_2$ is small

**Goal:** Recover $x_0$ given $A, y$ and $k$.

**Question:**

Can one accurately recover $x_0$ by solving problem $(P_\lambda)$ ?

# Sparse Signal Recovery

Without additional assumptions on $A$, this problem is ill posed

# Sparse Signal Recovery

Without additional assumptions on $A$, this problem is ill posed

- Even in the absence of noise, to be able to recover $x_0$, any $2k$ columns of $A$ must be linearly independent

# Sparse Signal Recovery

Without additional assumptions on $A$, this problem is ill posed

- Even in the absence of noise, to be able to recover $\boldsymbol{x}_0$, any $2k$ columns of $A$ must be linearly independent

### Assumption

*There exists a constant $\alpha_{2k} > 0$ such that for all $\boldsymbol{x} \in \mathbb{R}^d$ with $\|\boldsymbol{x}\|_0 \leq 2k$,*

$$\|A\boldsymbol{x}\|_2 \geq \alpha_{2k}\|\boldsymbol{x}\|_1$$

# Sparse Signal Recovery

Without additional assumptions on $A$, this problem is ill posed

- Even in the absence of noise, to be able to recover $x_0$, any $2k$ columns of $A$ must be linearly independent

## Assumption

*There exists a constant $\alpha_{2k} > 0$ such that for all $x \in \mathbb{R}^d$ with $\|x\|_0 \leq 2k$,*

$$\|Ax\|_2 \geq \alpha_{2k}\|x\|_1$$

Variant of the *Restricted Isometry Property*: One-sided, with mixed norms

# Sparse Signal Recovery

Without additional assumptions on $A$, this problem is ill posed

- Even in the absence of noise, to be able to recover $\boldsymbol{x}_0$, any $2k$ columns of $A$ must be linearly independent

### Assumption

There exists a constant $\alpha_{2k} > 0$ such that for all $\boldsymbol{x} \in \mathbb{R}^d$ with $\|\boldsymbol{x}\|_0 \leq 2k$,
$$\|A\boldsymbol{x}\|_2 \geq \alpha_{2k}\|\boldsymbol{x}\|_1$$

Variant of the *Restricted Isometry Property*: One-sided, with mixed norms

**Notation**:
For a vector $\boldsymbol{x} \in \mathbb{R}^d$, denote by $\Pi_k(\boldsymbol{x})$ the *k*-sparse *projection* of $\boldsymbol{x}$, namely the nearest *k*-sparse vector to $\boldsymbol{x}$

# The Trimmed Lasso: Sparse Recovery Guarantees

## Theorem

*Suppose that for some $\lambda > 0$, an optimization algorithm outputs a solution $\hat{x}$ such that*

$$F_\lambda(\hat{x}) \leq F_\lambda(\Pi_k(x_0)).$$

# The Trimmed Lasso: Sparse Recovery Guarantees

## Theorem

*Suppose that for some $\lambda > 0$, an optimization algorithm outputs a solution $\hat{x}$ such that*

$$F_\lambda(\hat{x}) \leq F_\lambda(\Pi_k(\boldsymbol{x}_0)).$$

*Let $\xi = \|\mathbf{e}\|_2 + \beta\tau_k(\boldsymbol{x}_0)$. Then,*

# The Trimmed Lasso: Sparse Recovery Guarantees

## Theorem

*Suppose that for some $\lambda > 0$, an optimization algorithm outputs a solution $\hat{x}$ such that*

$$F_\lambda(\hat{x}) \leq F_\lambda(\Pi_k(x_0)).$$

*Let $\xi = \|e\|_2 + \beta\tau_k(x_0)$. Then,*

*1. The projected solution $\Pi_k(\hat{x})$ is close to $x_0$,*

$$\|\Pi_k(\hat{x}) - x_0\|_1 \leq \tau_k(x_0) + \tfrac{2}{\alpha_{2k}}\xi + \tfrac{1}{2\lambda\alpha_{2k}}\xi^2$$

# The Trimmed Lasso: Sparse Recovery Guarantees

## Theorem

*Suppose that for some $\lambda > 0$, an optimization algorithm outputs a solution $\hat{x}$ such that*

$$F_\lambda(\hat{x}) \leq F_\lambda(\Pi_k(x_0)).$$

*Let $\xi = \|e\|_2 + \beta\tau_k(x_0)$. Then,*

*1. The projected solution $\Pi_k(\hat{x})$ is close to $x_0$,*

$$\|\Pi_k(\hat{x}) - x_0\|_1 \leq \tau_k(x_0) + \frac{2}{\alpha_{2k}}\xi + \frac{1}{2\lambda\alpha_{2k}}\xi^2$$

*2. If $\hat{x}$ itself is k-sparse, then the following tighter bound holds,*

$$\|\hat{x} - x_0\|_1 \leq \tau_k(x_0) + \frac{2}{\alpha_{2k}}\xi$$

**Implication:** We can well-approximate $x_0$ by solving $(P_\lambda)$ with $\lambda$ *smaller* than $\bar{\lambda}$

**Implication:** We can well-approximate $x_0$ by solving ($P_\lambda$) with $\lambda$ *smaller* than $\bar{\lambda}$

- ○ We don't need the optimal solutions of ($P_\lambda$) to coincide with those of (P0)

**Implication:** We can well-approximate $\boldsymbol{x}_0$ by solving ($P_\lambda$) with $\lambda$ *smaller* than $\bar{\lambda}$

- We don't need the optimal solutions of ($P_\lambda$) to coincide with those of (P0)
- Potentially, solving ($P_\lambda$) with smaller $\lambda$ is easier

**Implication:** We can well-approximate $x_0$ by solving $(P_\lambda)$ with $\lambda$ *smaller* than $\bar{\lambda}$

- We don't need the optimal solutions of $(P_\lambda)$ to coincide with those of $(P0)$
- Potentially, solving $(P_\lambda)$ with smaller $\lambda$ is easier
- Recovery is stable w.r.t. measurement error $\|e\|_2$ and inexactness of sparsity $\tau_k(x_0)$

**Note:** Theoretical guarantee for Lasso has better dependence on $\tau_k(\boldsymbol{x}_0)$, by a factor of $\mathcal{O}\left(\sqrt{k}\right)$.

**Note:** Theoretical guarantee for Lasso has better dependence on $\tau_k(\boldsymbol{x}_0)$, by a factor of $\mathcal{O}\left(\sqrt{k}\right)$.

- However, it requires the RIP constant to be bounded away from zero.

**Note:** Theoretical guarantee for Lasso has better dependence on $\tau_k(\boldsymbol{x}_0)$, by a factor of $\mathcal{O}\left(\sqrt{k}\right)$.

- However, it requires the RIP constant to be bounded away from zero.
  Even w/out noise, Lasso/BP requires $\alpha_{2k}$ to be bounded away from zero for recovery guarantees.

# The Trimmed Lasso: Sparse Recovery Guarantees

**Note:** Theoretical guarantee for Lasso has better dependence on $\tau_k(\boldsymbol{x}_0)$, by a factor of $\mathcal{O}\left(\sqrt{k}\right)$.

- However, it requires the RIP constant to be bounded away from zero.
  Even w/out noise, Lasso/BP requires $\alpha_{2k}$ to be bounded away from zero for recovery guarantees.
- Our guarantee only requires $\alpha_{2k} > 0$.

**Note:** Theoretical guarantee for Lasso has better dependence on $\tau_k(\boldsymbol{x}_0)$, by a factor of $\mathcal{O}\left(\sqrt{k}\right)$.

- However, it requires the RIP constant to be bounded away from zero.
  Even w/out noise, Lasso/BP requires $\alpha_{2k}$ to be bounded away from zero for recovery guarantees.

- Our guarantee only requires $\alpha_{2k} > 0$.
  - $\rightarrow$ a necessary condition for successful recovery by *any* algorithm

# The Trimmed Lasso: Sparse Recovery Guarantees

**Note:** Theoretical guarantee for Lasso has better dependence on $\tau_k(\boldsymbol{x}_0)$, by a factor of $\mathcal{O}\left(\sqrt{k}\right)$.

- However, it requires the RIP constant to be bounded away from zero.
  Even w/out noise, Lasso/BP requires $\alpha_{2k}$ to be bounded away from zero for recovery guarantees.
- Our guarantee only requires $\alpha_{2k} > 0$.
  - $\rightarrow$ a necessary condition for successful recovery by *any* algorithm

**In conclusion:**

Optimizing trimmed-lasso penalized objectives is a promising approach to (P0).

**Reminder:**

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

**Reminder:**

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

**Goal:**

$$\min_{\boldsymbol{x}} \frac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_k(\boldsymbol{x})$$

# The Trimmed Lasso: Practical Optimization

**Reminder:**

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

**Goal:**

$$\boxed{\min_{\boldsymbol{x}} \frac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_k(\boldsymbol{x})}$$

**Previous Optimization Methods:**

○ Difference of Convex Programming (DCP)

[Gotoh, Takeda, Tono, *Math. Prog.* '18]

○ Alternating Direction Method of Multipliers (ADMM)

[Bertsimas, Copenhaver, Mazumder, '17]

**Comparison to DC-programming**

# The Trimmed Lasso: Practical Optimization



Comparison to ADMM

# The Trimmed Lasso: Practical Optimization

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

Alternative formula:

$$\tau_k(\boldsymbol{x}) = \min_{|\Lambda|=d-k} \sum_{i\in\Lambda} |x_i|$$

# The Trimmed Lasso: Practical Optimization

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

Alternative formula:

$$\tau_k(\boldsymbol{x}) = \min_{|\Lambda|=d-k} \sum_{i\in\Lambda} |x_i|$$

Trimmed Lasso as a *hard* minimum:
Out of all $\binom{d}{k}$ subsets of $\{1, \ldots, d\}$, choose one with minimal $\ell_1$-norm.

$$\tau_k(\boldsymbol{x}) = \sum_{j=k+1}^{d} |x|_{(j)}$$

Alternative formula:

$$\tau_k(\boldsymbol{x}) = \min_{|\Lambda|=d-k} \sum_{i\in\Lambda} |x_i|$$

Trimmed Lasso as a *hard* minimum:

Out of all $\binom{d}{k}$ subsets of $\{1, \ldots, d\}$, choose one with minimal $\ell_1$-norm.

**Our Key Idea**: Replace the hard minimum by a *soft* minimum.

Let $z \in \mathbb{R}^m$ with $m = \binom{d}{k}$, whose entries consist of the $\ell_1$-norms of all subvectors of $x$ of size $d - k$. Formally:

# Surrogate for Trimmed Lasso

Let $z \in \mathbb{R}^m$ with $m = \binom{d}{k}$, whose entries consist of the $\ell_1$-norms of all subvectors of $x$ of size $d - k$. Formally:

$z$ is indexed by subsets $\Lambda \subset \{1, \ldots, d\}$ of size $d - k$:

$$z = (z_\Lambda), \quad |\Lambda| = d - k$$

# Surrogate for Trimmed Lasso

Let $z \in \mathbb{R}^m$ with $m = \binom{d}{k}$, whose entries consist of the $\ell_1$-norms of all subvectors of $x$ of size $d - k$. Formally:

$z$ is indexed by subsets $\Lambda \subset \{1, \ldots, d\}$ of size $d - k$:

$$z = (z_\Lambda), \quad |\Lambda| = d - k$$

Each entry of $z$ is given by

$$z_\Lambda = \sum_{i \in \Lambda} |x_i|$$

# Surrogate for Trimmed Lasso

Let $z \in \mathbb{R}^m$ with $m = \binom{d}{k}$, whose entries consist of the $\ell_1$-norms of all subvectors of $x$ of size $d - k$. Formally:

$z$ is indexed by subsets $\Lambda \subset \{1, \ldots, d\}$ of size $d - k$:

$$z = (z_\Lambda), \quad |\Lambda| = d - k$$

Each entry of $z$ is given by

$$z_\Lambda = \sum_{i \in \Lambda} |x_i|$$

Note that

$$\tau_k(x) = \min_{|\Lambda| = d - k} z_\Lambda$$

# Surrogate for Trimmed Lasso

Let $z \in \mathbb{R}^m$ with $m = \binom{d}{k}$, whose entries consist of the $\ell_1$-norms of all subvectors of $x$ of size $d - k$. Formally:

$z$ is indexed by subsets $\Lambda \subset \{1, \ldots, d\}$ of size $d - k$:

$$z = (z_\Lambda), \quad |\Lambda| = d - k$$

Each entry of $z$ is given by

$$z_\Lambda = \sum_{i \in \Lambda} |x_i|$$

We wish:

$$\rho(x) = \underset{|\Lambda| = d - k}{\text{soft min}} \; z_\Lambda$$

# Surrogate for Trimmed Lasso

Let $z \in \mathbb{R}^m$ with $m = \binom{d}{k}$, whose entries consist of the $\ell_1$-norms of all subvectors of $x$ of size $d - k$. Formally:

$z$ is indexed by subsets $\Lambda \subset \{1, \ldots, d\}$ of size $d - k$:

$$z = (z_\Lambda), \quad |\Lambda| = d - k$$

Each entry of $z$ is given by

$$z_\Lambda = \sum_{i \in \Lambda} |x_i|$$

We wish:

$$\rho(x) = \underset{|\Lambda|=d-k}{\text{soft min}}\ z_\Lambda$$

- As in the *softmax* function in multi-class classification.

# Surrogate for Trimmed Lasso

Soft maximum of $\mathbf{z} = (z_1, \ldots, z_m)$:

$$\log \left( \sum_{j=1}^m \exp(z_j) \right)$$

# Surrogate for Trimmed Lasso

Soft minimum of $z$:

$$-\log\left(\sum_{j=1}^{m}\exp\left(-z_j\right)\right)$$

# Surrogate for Trimmed Lasso

Add a smoothness parameter $\gamma$:

$$-\frac{1}{\gamma} \log \left( \sum_{j=1}^{m} \exp\left(-\gamma z_j\right) \right)$$

# Surrogate for Trimmed Lasso

Add averaging:

$$-\frac{1}{\gamma} \log \left( \frac{1}{m} \sum_{j=1}^{m} \exp\left(-\gamma z_j\right) \right)$$

# Surrogate for Trimmed Lasso

Plug in the original definition of $z$:

$$-\frac{1}{\gamma} \log \left( \frac{1}{\binom{d}{k}} \sum_{|\Lambda|=d-k} \exp\left( -\gamma \sum_{i \in \Lambda} |x_i| \right) \right)$$

# Surrogate for Trimmed Lasso

$$-\frac{1}{\gamma} \log \left( \frac{1}{\binom{d}{k}} \sum_{|\Lambda| = d-k} \exp \left( -\gamma \sum_{i \in \Lambda} |x_i| \right) \right)$$

$$\tau_{k,\gamma}(\boldsymbol{x}) = -\frac{1}{\gamma} \log \left( \frac{1}{\binom{d}{k}} \sum_{|\Lambda|=d-k} \exp\left( -\gamma \sum_{i \in \Lambda} |x_i| \right) \right)$$

*Generalized Soft-Min Penalty*

# Surrogate for Trimmed Lasso

$$\tau_{k,\gamma}(\boldsymbol{x}) = -\frac{1}{\gamma} \log \left( \frac{1}{\binom{d}{k}} \sum_{|\Lambda|=d-k} \exp \left( -\gamma \sum_{i \in \Lambda} |x_i| \right) \right)$$

*Generalized Soft-Min Penalty*

○ Infinitely differentiable as a function of $|\boldsymbol{x}|$

# Surrogate for Trimmed Lasso

$$\tau_{k,\gamma}(\boldsymbol{x}) = -\frac{1}{\gamma} \log \left( \frac{1}{\binom{d}{k}} \sum_{|\Lambda|=d-k} \exp\left( -\gamma \sum_{i \in \Lambda} |x_i| \right) \right)$$

*Generalized Soft-Min Penalty*

○ Infinitely differentiable as a function of $|\boldsymbol{x}|$
- Parameter $\gamma$ controls level of smoothness

# Surrogate for Trimmed Lasso

$$\tau_{k,\gamma}(\boldsymbol{x}) = -\frac{1}{\gamma} \log \left( \frac{1}{\binom{d}{k}} \sum_{|\Lambda|=d-k} \exp\left( -\gamma \sum_{i\in\Lambda} |x_i| \right) \right)$$

*Generalized Soft-Min Penalty*

- Infinitely differentiable as a function of $|\boldsymbol{x}|$
    - Parameter $\gamma$ controls level of smoothness
- Takes into account all possible $\binom{d}{k}$ sparsity patterns of $\boldsymbol{x}$

# Surrogate for Trimmed Lasso

$$\tau_{k,\gamma}(\boldsymbol{x}) = -\frac{1}{\gamma} \log \left( \frac{1}{\binom{d}{k}} \sum_{|\Lambda|=d-k} \exp\left( -\gamma \sum_{i \in \Lambda} |x_i| \right) \right)$$

*Generalized Soft-Min Penalty*

○ Infinitely differentiable as a function of $|\boldsymbol{x}|$
- Parameter $\gamma$ controls level of smoothness
○ Takes into account all possible $\binom{d}{k}$ sparsity patterns of $\boldsymbol{x}$
○ Significantly easier to optimize

## Lemma

*For any $\boldsymbol{x} \in \mathbb{R}^d$, the function $\tau_{k,\gamma}(\boldsymbol{x})$ is monotone-decreasing with respect to $\gamma$. Moreover,*

# Generalized Soft-Min Properties

## Lemma

*For any $\boldsymbol{x} \in \mathbb{R}^d$, the function $\tau_{k,\gamma}(\boldsymbol{x})$ is monotone-decreasing with respect to $\gamma$. Moreover,*

$$\lim_{\gamma \to 0} \tau_{k,\gamma}(\boldsymbol{x}) = \frac{d-k}{d} \|\boldsymbol{x}\|_1$$

## Lemma

For any $\boldsymbol{x} \in \mathbb{R}^d$, the function $\tau_{k,\gamma}(\boldsymbol{x})$ is monotone-decreasing with respect to $\gamma$. Moreover,

$$\lim_{\gamma \to 0} \tau_{k,\gamma}(\boldsymbol{x}) = \frac{d-k}{d} \|\boldsymbol{x}\|_1$$

$$\lim_{\gamma \to \infty} \tau_{k,\gamma}(\boldsymbol{x}) = \tau_k(\boldsymbol{x})$$

## Lemma

For any $\boldsymbol{x} \in \mathbb{R}^d$, the function $\tau_{k,\gamma}(\boldsymbol{x})$ is monotone-decreasing with respect to $\gamma$. Moreover,

$$\lim_{\gamma \to 0} \tau_{k,\gamma}(\boldsymbol{x}) = \frac{d-k}{d} \|\boldsymbol{x}\|_1$$

$$\lim_{\gamma \to \infty} \tau_{k,\gamma}(\boldsymbol{x}) = \tau_k(\boldsymbol{x})$$

# A Homotopy Scheme

Instead of directly minimizing

$$\tfrac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_k(\boldsymbol{x})$$

# A Homotopy Scheme

Instead of directly minimizing

$$\tfrac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_k(\boldsymbol{x})$$

Solve a sequence of problems

$$\min_{\boldsymbol{x}} \mathsf{F}_{\lambda,\gamma}(\boldsymbol{x}) = \tfrac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_{k,\gamma}(\boldsymbol{x})$$

with an increasing sequence $\gamma_0 < \gamma_1 < \ldots$, while tracing path of solutions.

# A Homotopy Scheme

Instead of directly minimizing

$$\tfrac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_k(\boldsymbol{x})$$

Solve a sequence of problems

$$\min_{\boldsymbol{x}} \mathsf{F}_{\lambda,\gamma}(\boldsymbol{x}) = \tfrac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_{k,\gamma}(\boldsymbol{x})$$

with an increasing sequence $\gamma_0 < \gamma_1 < \ldots$, while tracing path of solutions.

- Start at $\gamma = 0$: $\tau_{k,0}(\boldsymbol{x})$ is the convex $\ell_1$ norm (Lasso problem).

# A Homotopy Scheme

Instead of directly minimizing

$$\tfrac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_k(\boldsymbol{x})$$

Solve a sequence of problems

$$\min_{\boldsymbol{x}} F_{\lambda,\gamma}(\boldsymbol{x}) = \tfrac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_{k,\gamma}(\boldsymbol{x})$$

with an increasing sequence $\gamma_0 < \gamma_1 < \ldots$, while tracing path of solutions.

○ Start at $\gamma = 0$: $\tau_{k,0}(\boldsymbol{x})$ is the convex $\ell_1$ norm (Lasso problem).

○ Slowly increase $\gamma$. At iteration $t$ with $\gamma = \gamma_t$, initialize optimization method with previous solution $\hat{\boldsymbol{x}}_{t-1}$.

**Problem:** How to minimize each nonconvex objective

$$F_{\lambda,\gamma}(\boldsymbol{x}) = \tfrac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_{k,\gamma}(\boldsymbol{x})?$$

**Problem:** How to minimize each nonconvex objective

$$F_{\lambda,\gamma}(\boldsymbol{x}) = \tfrac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_{k,\gamma}(\boldsymbol{x})?$$

**Approach:** Majorization-Minimization

## Majorization Minimization Scheme

**Problem:** How to minimize each nonconvex objective

$$F_{\lambda,\gamma}(\boldsymbol{x}) = \tfrac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_{k,\gamma}(\boldsymbol{x})?$$

**Approach:** Majorization-Minimization

Construct a function $G_{\lambda,\gamma}(\boldsymbol{x}, \tilde{\boldsymbol{x}})$ such that

$$G_{\lambda,\gamma}(\boldsymbol{x}, \tilde{\boldsymbol{x}}) \geq F_{\lambda,\gamma}(\boldsymbol{x}), \quad G_{\lambda,\gamma}(\boldsymbol{x}, \boldsymbol{x}) = F_{\lambda,\gamma}(\boldsymbol{x}).$$

Iterate:

$$\boldsymbol{x}^t = \arg\min_{x} G_{\lambda,\gamma}\big(\boldsymbol{x}, \boldsymbol{x}^{t-1}\big).$$

## Majorization Minimization Scheme

**Problem:** How to minimize each nonconvex objective

$$F_{\lambda,\gamma}(\boldsymbol{x}) = \tfrac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_{k,\gamma}(\boldsymbol{x})?$$

**Approach:** Majorization-Minimization

Construct a function $G_{\lambda,\gamma}(\boldsymbol{x}, \tilde{\boldsymbol{x}})$ such that

$$G_{\lambda,\gamma}(\boldsymbol{x}, \tilde{\boldsymbol{x}}) \geq F_{\lambda,\gamma}(\boldsymbol{x}), \quad G_{\lambda,\gamma}(\boldsymbol{x}, \boldsymbol{x}) = F_{\lambda,\gamma}(\boldsymbol{x}).$$

Iterate:

$$\boldsymbol{x}^t = \arg\min_x G_{\lambda,\gamma}\big(\boldsymbol{x}, \boldsymbol{x}^{t-1}\big).$$

- Objective is guaranteed to decrease monotonically.

## Majorization Minimization Scheme

**Problem:** How to minimize each nonconvex objective

$$F_{\lambda,\gamma}(\boldsymbol{x}) = \tfrac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda \tau_{k,\gamma}(\boldsymbol{x})?$$

**Approach:** Majorization-Minimization

Construct a function $G_{\lambda,\gamma}(\boldsymbol{x}, \tilde{\boldsymbol{x}})$ such that

$$G_{\lambda,\gamma}(\boldsymbol{x}, \tilde{\boldsymbol{x}}) \geq F_{\lambda,\gamma}(\boldsymbol{x}), \quad G_{\lambda,\gamma}(\boldsymbol{x}, \boldsymbol{x}) = F_{\lambda,\gamma}(\boldsymbol{x}).$$

Iterate:

$$\boldsymbol{x}^t = \arg\min_x G_{\lambda,\gamma}\big(\boldsymbol{x}, \boldsymbol{x}^{t-1}\big).$$

- Objective is guaranteed to decrease monotonically.
- Under some assumptions, guaranteed to converge to a stationary point.

**Constructing a majorizer for $F_{\lambda,\gamma}(x)$:**

## Majorization Minimization Scheme

**Constructing a majorizer for $F_{\lambda,\gamma}(x)$:**

Define $w_{k,\gamma} : \mathbb{R}^d \to \mathbb{R}^d$ for $0 \le \gamma < \infty$ by

$$w_{k,\gamma}^i(x) = \frac{\sum_{|\Lambda|=d-k,\, i \in \Lambda} \exp\left(-\gamma \sum_{j \in \Lambda} |x_j|\right)}{\sum_{|\Lambda|=d-k} \exp\left(-\gamma \sum_{j \in \Lambda} |x_j|\right)}$$

**Constructing a majorizer for $F_{\lambda,\gamma}(x)$:**

Define $w_{k,\gamma} : \mathbb{R}^d \to \mathbb{R}^d$ for $0 \leq \gamma < \infty$ by

$$w_{k,\gamma}^i(x) = \frac{\sum_{|\Lambda|=d-k, i\in\Lambda} \exp\left(-\gamma \sum_{j\in\Lambda} |x_j|\right)}{\sum_{|\Lambda|=d-k} \exp\left(-\gamma \sum_{j\in\Lambda} |x_j|\right)}$$

**Lemma:** The following function is a majorizer of $F_{\lambda,\gamma}(x)$:

$$G_{\lambda,\gamma}(x, \tilde{x}) = \frac{1}{2}\|Ax - y\|^2 + \lambda\tau_{k,\gamma}(\tilde{x}) + \lambda\langle w_{k,\gamma}(\tilde{x}), |x| - |\tilde{x}|\rangle$$

# Majorization Minimization Scheme

**Constructing a majorizer for $\mathbf{F}_{\lambda,\gamma}(\boldsymbol{x})$:**

Define $w_{k,\gamma} : \mathbb{R}^d \to \mathbb{R}^d$ for $0 \le \gamma < \infty$ by

$$w_{k,\gamma}^i(\boldsymbol{x}) = \frac{\sum_{|\Lambda|=d-k, i \in \Lambda} \exp\left(-\gamma \sum_{j \in \Lambda} |x_j|\right)}{\sum_{|\Lambda|=d-k} \exp\left(-\gamma \sum_{j \in \Lambda} |x_j|\right)}$$

**Lemma:** The following function is a majorizer of $\mathsf{F}_{\lambda,\gamma}(\boldsymbol{x})$:

$$\mathsf{G}_{\lambda,\gamma}(\boldsymbol{x}, \tilde{\boldsymbol{x}}) = \frac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|^2 + \lambda\tau_{k,\gamma}(\tilde{\boldsymbol{x}}) + \lambda\langle w_{k,\gamma}(\tilde{\boldsymbol{x}}), |\boldsymbol{x}| - |\tilde{\boldsymbol{x}}|\rangle$$

constant w.r.t. $\boldsymbol{x}$

# Majorization Minimization Scheme

**MM scheme to minimize $F_{\lambda,\gamma}(x)$:**

$$\mathbf{w}^t = \mathbf{w}_{k,\gamma}(\mathbf{x}^{t-1})$$
$$\mathbf{x}^t = \arg\min_x \frac{1}{2}\|A\mathbf{x} - \mathbf{y}\|^2 + \lambda\langle\mathbf{w}^t, |\mathbf{x}|\rangle$$

# Majorization Minimization Scheme

**MM scheme to minimize $F_{\lambda,\gamma}(x)$:**

$$\mathbf{w}^t = \mathbf{w}_{k,\gamma}(\mathbf{x}^{t-1})$$
$$\mathbf{x}^t = \arg\min_x \frac{1}{2}\|A\mathbf{x} - \mathbf{y}\|^2 + \lambda\langle\mathbf{w}^t, |\mathbf{x}|\rangle$$

Each subproblem is a *convex* weighted $\ell_1$ problem.

## Majorization Minimization Scheme

**MM scheme to minimize $F_{\lambda,\gamma}(x)$:**

$$\mathbf{w}^t = \mathbf{w}_{k,\gamma}(\mathbf{x}^{t-1})$$
$$\mathbf{x}^t = \arg\min_x \frac{1}{2}\|A\mathbf{x} - \mathbf{y}\|^2 + \lambda\langle\mathbf{w}^t, |\mathbf{x}|\rangle$$

Each subproblem is a *convex* weighted $\ell_1$ problem.
Similar to IRL1...

# Majorization Minimization Scheme

**MM scheme to minimize $F_{\lambda,\gamma}(x)$:**

$$w^t = w_{k,\gamma}(x^{t-1})$$
$$x^t = \arg\min_x \frac{1}{2}\|Ax - y\|^2 + \lambda\langle w^t, |x|\rangle$$

Each subproblem is a *convex* weighted $\ell_1$ problem.
Similar to IRL1... with a key difference:

### Lemma

*For any $x \in \mathbb{R}^d$, $k$, $\gamma$,*

1. *All weights $w_{k,\gamma}^i(x) \in [0,1]$*
2. *$\sum_{i=1}^d w_{k,\gamma}^i(x) = d - k$*

# Majorization Minimization Scheme

**MM scheme to minimize $F_{\lambda,\gamma}(x)$:**

$$w^t = w_{k,\gamma}(x^{t-1})$$
$$x^t = \arg\min_x \frac{1}{2}\|Ax - y\|^2 + \lambda\langle w^t, |x|\rangle$$

Each subproblem is a *convex* weighted $\ell_1$ problem.
Similar to IRL1... with a key difference:

### Lemma

For any $x \in \mathbb{R}^d$, $k$, $\gamma$,

1. All weights $w_{k,\gamma}^i(x) \in [0,1]$

2. $\sum_{i=1}^d w_{k,\gamma}^i(x) = d - k$

Since all weights are in [0,1], and their sum is constant, they do
not require regularization.

**Problem:**   How to compute $\tau_{k,\gamma}(\mathbf{x})$ and $\mathbf{w}_{k,\gamma}(\mathbf{x})$?

Their formulas involve sums of $\binom{d}{k}$ terms.

**Problem:** How to compute $\tau_{k,\gamma}(\mathbf{x})$ and $\mathbf{w}_{k,\gamma}(\mathbf{x})$?

Their formulas involve sums of $\binom{d}{k}$ terms.

Naïve calculation would be...

# Computing $\tau_{k,\gamma}$ and $\mathbf{w}_{k,\gamma}$

**Problem:** How to compute $\tau_{k,\gamma}(\boldsymbol{x})$ and $\mathbf{w}_{k,\gamma}(\boldsymbol{x})$?

Their formulas involve sums of $\binom{d}{k}$ terms.

Naïve calculation would be...

- *prohibitively* slow.

# Computing $\tau_{k,\gamma}$ and $\mathbf{w}_{k,\gamma}$

**Problem:** How to compute $\tau_{k,\gamma}(\boldsymbol{x})$ and $\mathbf{w}_{k,\gamma}(\boldsymbol{x})$?

Their formulas involve sums of $\binom{d}{k}$ terms.

Naïve calculation would be...

- *prohibitively* slow.
- highly prone to numerical corruption by arithmetic overflow and underflow, due to the log and exp operations.

# Computing $\tau_{k,\gamma}$ and $\mathbf{w}_{k,\gamma}$

**Problem:** How to compute $\tau_{k,\gamma}(\boldsymbol{x})$ and $\mathbf{w}_{k,\gamma}(\boldsymbol{x})$?

Their formulas involve sums of $\binom{d}{k}$ terms.

Naïve calculation would be...

- *prohibitively* slow.
- highly prone to numerical corruption by arithmetic overflow and underflow, due to the log and exp operations.

Developed numerical scheme to accurately compute $\tau_{k,\gamma}(\boldsymbol{x})$ and $\mathbf{w}_{k,\gamma}(\boldsymbol{x})$

# Computing $\tau_{k,\gamma}$ and $\mathbf{w}_{k,\gamma}$

**Problem:** How to compute $\tau_{k,\gamma}(\boldsymbol{x})$ and $\mathbf{w}_{k,\gamma}(\boldsymbol{x})$?

Their formulas involve sums of $\binom{d}{k}$ terms.

Naïve calculation would be...

- *prohibitively* slow.
- highly prone to numerical corruption by arithmetic overflow and underflow, due to the log and exp operations.

Developed numerical scheme to accurately compute $\tau_{k,\gamma}(\boldsymbol{x})$ and $\mathbf{w}_{k,\gamma}(\boldsymbol{x})$

- Recursive, takes $\mathcal{O}(kd)$ operations

# Computing $\tau_{k,\gamma}$ and $\mathbf{w}_{k,\gamma}$

**Problem:** How to compute $\tau_{k,\gamma}(\boldsymbol{x})$ and $\mathbf{w}_{k,\gamma}(\boldsymbol{x})$?

Their formulas involve sums of $\binom{d}{k}$ terms.

Naïve calculation would be...

- *prohibitively* slow.
- highly prone to numerical corruption by arithmetic overflow and underflow, due to the log and exp operations.

Developed numerical scheme to accurately compute $\tau_{k,\gamma}(\boldsymbol{x})$ and $\mathbf{w}_{k,\gamma}(\boldsymbol{x})$

- Recursive, takes $\mathcal{O}(kd)$ operations

Approach also relevant for top-$k$ classification. Method to compute similar functions for small $k$ was proposed by [Berrada, Zisserman, Kumar, *ICLR* '18].

## Outline of our method

(a) We seek a solution of (P0) by solving

$$\frac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda \tau_k(\boldsymbol{x})$$

for increasing values of $\lambda < \bar{\lambda}$, till a $k$-sparse solution found.

## Outline of our method

(a) We seek a solution of (P0) by solving

$$\frac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_k(\boldsymbol{x})$$

for increasing values of $\lambda < \bar{\lambda}$, till a $k$-sparse solution found.

(b) Each such problem solved by homotopy: Minimize

$$\frac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda\tau_{k,\gamma}(\boldsymbol{x})$$

for increasing sequence of values of $\gamma$.

## Outline of our method

(a) We seek a solution of (P0) by solving

$$\frac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda \tau_k(\boldsymbol{x})$$

for increasing values of $\lambda < \bar{\lambda}$, till a $k$-sparse solution found.

(b) Each such problem solved by homotopy: Minimize

$$\frac{1}{2}\|A\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \lambda \tau_{k,\gamma}(\boldsymbol{x})$$

for increasing sequence of values of $\gamma$.

(c) Each such problem solved by MM, requiring solution of several weighted $\ell_1$ problems.

(a) We seek a solution of (P0) by solving

$$\frac{1}{2}\|A\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \tau_k(\mathbf{x})$$

for increasing values of $\lambda < \bar{\lambda}$, till a $k$-sparse solution found.

(b) Each such problem solved by homotopy: Minimize

$$\frac{1}{2}\|A\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \tau_{k,\gamma}(\mathbf{x})$$

for increasing sequence of values of $\gamma$.

(c) Each such problem solved by MM, requiring solution of several weighted $\ell_1$ problems.

Running time for one $\lambda$: $\approx 500\times$ slower than single $\ell_1$ problem.

# Comparison to current state of the art

(As in [Bertsimas and Van Parys, 2020])

- $x_0 \in \mathbb{R}^d$ is $k$-sparse, $d = 15000$, $k = 10$, with entries $\pm 1$
- $A \in \mathbb{R}^{n \times d}$ with uncorrelated $\mathcal{N}(0, 1)$ entries

# Comparison to current state of the art

(As in [Bertsimas and Van Parys, 2020])

- $x_0 \in \mathbb{R}^d$ is $k$-sparse, $d = 15000$, $k = 10$, with entries $\pm 1$
- $A \in \mathbb{R}^{n \times d}$ with uncorrelated $\mathcal{N}(0,1)$ entries
- Observation: $y = Ax_0 + e$, with 5% noise (SNR=400)
- True $k$ is known to all methods

# Comparison to current state of the art

(As in [Bertsimas and Van Parys, 2020])

- $x_0 \in \mathbb{R}^d$ is $k$-sparse, $d = 15000$, $k = 10$, with entries $\pm 1$
- $A \in \mathbb{R}^{n \times d}$ with uncorrelated $\mathcal{N}(0,1)$ entries
- Observation: $\mathbf{y} = A x_0 + \mathbf{e}$, with 5% noise (SNR=400)
- True $k$ is known to all methods
- Coordinate descent returns multiple solutions

  Chose the one whose support is closest to the true support

# Comparison to current state of the art

(As in [Bertsimas and Van Parys, 2020])

- $x_0 \in \mathbb{R}^d$ is $k$-sparse, $d = 15000$, $k = 10$, with entries $\pm 1$
- $A \in \mathbb{R}^{n \times d}$ with uncorrelated $\mathcal{N}(0, 1)$ entries
- Observation: $\mathbf{y} = A x_0 + \mathbf{e}$, with 5% noise (SNR=400)
- True $k$ is known to all methods
- Coordinate descent returns multiple solutions

    Chose the one whose support is closest to the true support

**Measure of success:**

- Support accuracy: $\frac{|\hat{S} \cap S_0|}{k}$

# Comparison to current state of the art



$k = 10$, $d = 15000$, 5% noise (SNR=400)

# Comparison to current state of the art



$k = 10$, $d = 15000$, 5% noise (SNR=400)

# Comparison to current state of the art



$k = 10$, $d = 15000$, $5\%$ noise (SNR=400)

$k = 10$, $d = 15000$, 5% noise (SNR=400)

(As in [Hazimeh, Mazumder 2020])

- $k$-sparse signal $\boldsymbol{x}_0 \in \mathbb{R}^d$, $k = 50$, $d = 20000$
- Entries $\pm 1$

# Comparison to current state of the art

(As in [Hazimeh, Mazumder 2020])

- $k$-sparse signal $\boldsymbol{x}_0 \in \mathbb{R}^d$, $k = 50$, $d = 20000$
- Entries $\pm 1$
- $A \in \mathbb{R}^{n \times d}$, $\mathcal{N}(0, \Sigma)$, $\Sigma_{i,j} = 0.5^{|i-j|}$
- Observation: $\mathbf{y} = A\boldsymbol{x}_0 + \mathbf{e}$, varying noise levels

# Comparison to current state of the art

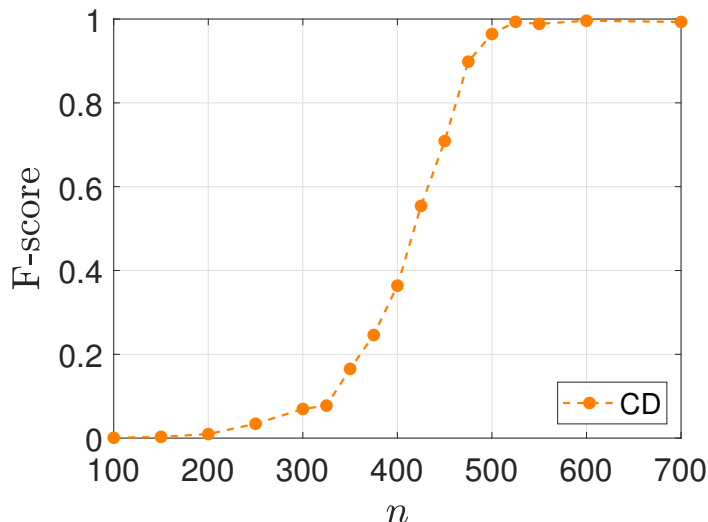(As in [Hazimeh, Mazumder 2020])

- $k$-sparse signal $\mathbf{x}_0 \in \mathbb{R}^d$, $k = 50$, $d = 20000$
- Entries $\pm 1$
- $A \in \mathbb{R}^{n \times d}$, $\mathcal{N}(0, \Sigma)$, $\Sigma_{i,j} = 0.5^{|i-j|}$
- Observation: $\mathbf{y} = A\mathbf{x}_0 + \mathbf{e}$, varying noise levels
- Each method chooses $k$ using a separate validation set: $\tilde{\mathbf{y}} = \tilde{A}\mathbf{x}_0 + \tilde{\mathbf{e}}$

# Comparison to current state of the art

(As in [Hazimeh, Mazumder 2020])

- $k$-sparse signal $\boldsymbol{x}_0 \in \mathbb{R}^d$, $k = 50$, $d = 20000$
- Entries $\pm 1$
- $A \in \mathbb{R}^{n \times d}$, $\mathcal{N}(0, \Sigma)$, $\Sigma_{i,j} = 0.5^{|i-j|}$
- Observation: $\boldsymbol{y} = A\boldsymbol{x}_0 + \mathbf{e}$, varying noise levels
- Each method chooses $k$ using a separate validation set:
  $\tilde{\mathbf{y}} = \tilde{A}\boldsymbol{x}_0 + \tilde{\mathbf{e}}$

Measures of success:

# Comparison to current state of the art

(As in [Hazimeh, Mazumder 2020])

- $k$-sparse signal $\boldsymbol{x}_0 \in \mathbb{R}^d$, $k = 50$, $d = 20000$
- Entries $\pm 1$
- $A \in \mathbb{R}^{n \times d}$, $\mathcal{N}(0, \Sigma)$, $\Sigma_{i,j} = 0.5^{|i-j|}$
- Observation: $\boldsymbol{y} = A\boldsymbol{x}_0 + \boldsymbol{e}$, varying noise levels
- Each method chooses $k$ using a separate validation set: $\tilde{\boldsymbol{y}} = \tilde{A}\boldsymbol{x}_0 + \tilde{\boldsymbol{e}}$

Measures of success:

- F-score: $2\dfrac{|\hat{S} \cap S_0|}{|\hat{S}| + |S_0|}$

# Comparison to current state of the art

(As in [Hazimeh, Mazumder 2020])

- $k$-sparse signal $\boldsymbol{x}_0 \in \mathbb{R}^d$, $k = 50$, $d = 20000$
- Entries $\pm 1$
- $A \in \mathbb{R}^{n \times d}$, $\mathcal{N}(0, \Sigma)$, $\Sigma_{i,j} = 0.5^{|i-j|}$
- Observation: $\boldsymbol{y} = A\boldsymbol{x}_0 + \mathbf{e}$, varying noise levels
- Each method chooses $k$ using a separate validation set:
  $\tilde{\boldsymbol{y}} = \tilde{A}\boldsymbol{x}_0 + \tilde{\mathbf{e}}$

Measures of success:

- F-score: $2\dfrac{|\hat{S} \cap S_0|}{|\hat{S}| + |S_0|}$

- Expected prediction error: $\sqrt{\dfrac{\mathbb{E}_{\mathbf{A},\mathbf{y}}\left[\left\|A\hat{\boldsymbol{x}} - y\right\|^2\right]}{\mathbb{E}_{\mathbf{y}}\left[\|y\|^2\right]}}$

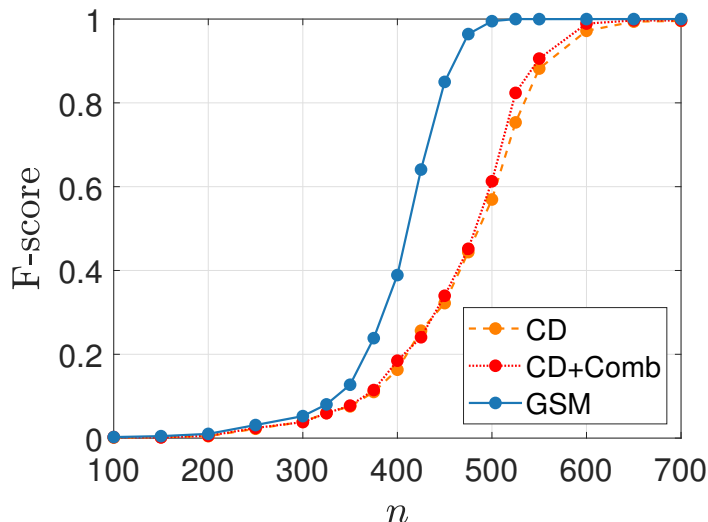# Comparison to current state of the art



$d = 20000$   $k = 50$   Entries: $\pm 1$   5% noise (SNR=400)

# Comparison to current state of the art



$d = 20000$   $k = 50$   Entries: $\pm 1$   5% noise (SNR=400)
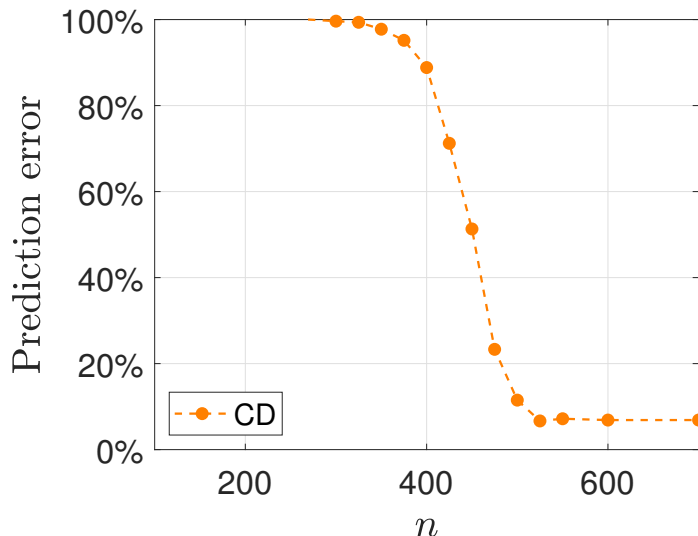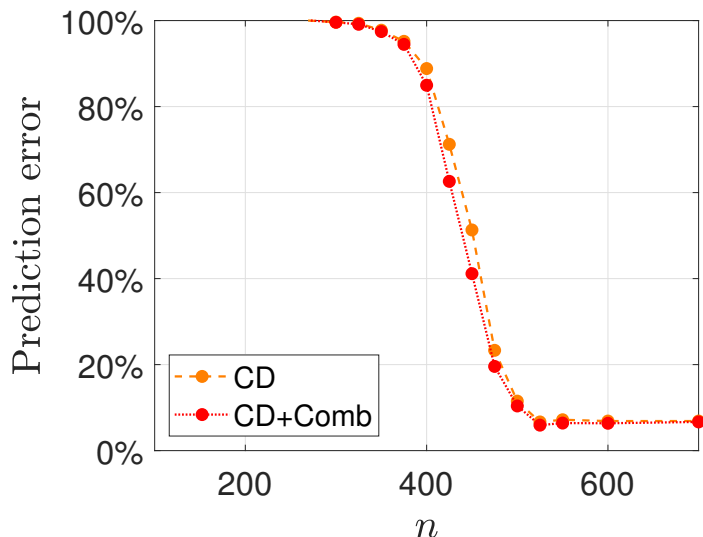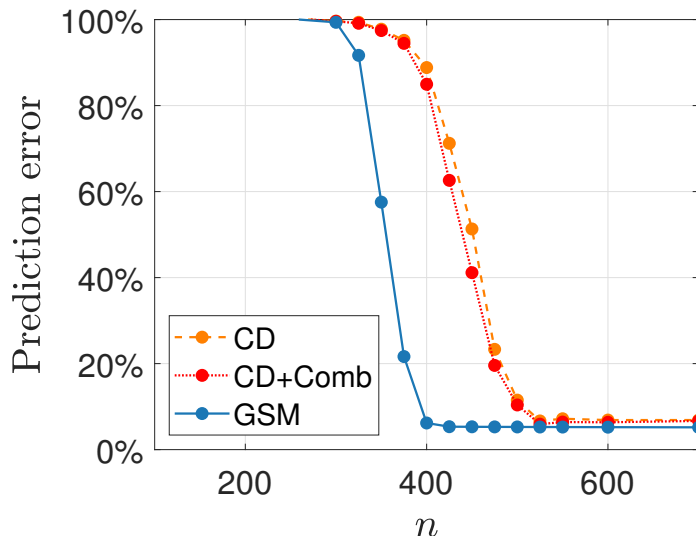
# Comparison to current state of the art



$d = 20000$ $k = 50$ Entries: $\pm 1$ 5% noise (SNR=400)

$d = 20000 \quad k = 50 \quad$ Entries: $\pm 1 \quad 33.3\%$ noise (SNR=9)

# Comparison to current state of the art



$d = 20000$  $k = 50$  Entries: $\pm 1$  5% noise (SNR=400)

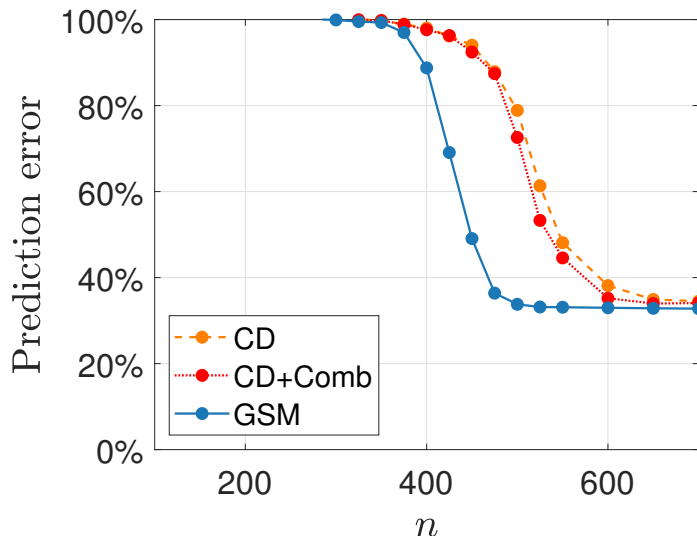# Comparison to current state of the art



$d = 20000$   $k = 50$   Entries: $\pm 1$   5% noise (SNR=400)

$d = 20000$  $k = 50$  Entries: $\pm 1$  5% noise (SNR=400)

$d = 20000$   $k = 50$   Entries: $\pm 1$   33.3% noise (SNR=9)

# Conclusion

- Problem (P0) plays a key role in multiple applications.
- Still room for improvements for challenging instances of (P0)

# Conclusion

- Problem (P0) plays a key role in multiple applications.
- Still room for improvements for challenging instances of (P0)
- Trimmed Lasso - desirable theoretical properties to solve (P0)

# Conclusion

- Problem (P0) plays a key role in multiple applications.
- Still room for improvements for challenging instances of (P0)
- Trimmed Lasso - desirable theoretical properties to solve (P0)
- Practical optimization method for Trimmed-Lasso penalty

# Conclusion

- Problem (P0) plays a key role in multiple applications.
- Still room for improvements for challenging instances of (P0)
- Trimmed Lasso - desirable theoretical properties to solve (P0)
- Practical optimization method for Trimmed-Lasso penalty
    - Novel surrogate penalty (GSM)
    - Accurate numerical scheme
    - Accompanying optimization algorithm
- Approach potentially applicable to other sparse combinatorial search problems

# Conclusion

- Problem (P0) plays a key role in multiple applications.
- Still room for improvements for challenging instances of (P0)
- Trimmed Lasso - desirable theoretical properties to solve (P0)
- Practical optimization method for Trimmed-Lasso penalty
    - Novel surrogate penalty (GSM)
    - Accurate numerical scheme
    - Accompanying optimization algorithm
- Approach potentially applicable to other sparse combinatorial search problems

code on GitHub.

Amir, T., Basri, R. and Nadler, B., The Trimmed Lasso: Sparse Recovery
Guarantees and Practical Optimization by the Generalized Soft-Min Penalty.
*SIAM J. Math. Data Science, 2021*

**Thank You**

**The End**