

# Recipes for Hard Thresholding Methods

A. Kyrillidis<sup>1</sup> and V. Cevher<sup>1,2</sup>



<sup>1</sup>Laboratory for Information and Inference Systems (EPFL) <sup>2</sup>Idiap Research Institute  
{anastasios.kyrillidis,volkan.cevher}@epfl.ch

<http://lions.epfl.ch/>



Grant Acknowledgements:

KEEP project through DARPA KeCoM, Marie Curie reintegration.  
VC also acknowledges his Faculty Fellowship position at Rice University.



# Restricted Isometry Property

- Sparsity: not enough by itself.
- Conditions on  $\mathcal{A}$ : stable embedding, null space property, spark, unique representation property, exact recovery condition, etc.
- In this work  $\rightarrow$  stable embedding.
- Let  $\Sigma_K^N$ : union-of-subspaces with at most  $K$ -nonzero entries in  $N$ -dimensions.

## Restricted Isometry Property (RIP)

Let  $\mathcal{A} \in \mathbb{R}^{M \times N}$  and  $K < N$  be an integer number. Then,  $\mathcal{A}$  satisfies the restricted isometry property with constant  $\delta_K$  iff

$$(1 - \delta_K) \leq \frac{\|\mathcal{A}\mathbf{x}\|_2^2}{\|\mathbf{x}\|_2^2} \leq (1 + \delta_K),$$

is satisfied for any  $\mathbf{x} \in \Sigma_K^N$ .

- Real data example -  $\Phi_{300 \times 1000}$  i.i.d.  $\sim \mathcal{N}(0, 1/M)$ :  $\frac{\|\mathcal{A}\mathbf{x}\|_2^2}{\|\mathbf{x}\|_2^2} \in [\sim 0.5, \sim 2.5] \Rightarrow$

Non-Symmetric

# Restricted Isometry Property

- Sparsity: not enough by itself.
- Conditions on  $\mathcal{A}$ : stable embedding, null space property, spark, unique representation property, exact recovery condition, etc.
- In this work  $\rightarrow$  stable embedding.
- Let  $\Sigma_K^N$ : union-of-subspaces with at most  $K$ -nonzero entries in  $N$ -dimensions.

## Non-Symmetric RIP

Let  $\mathcal{A} \in \mathbb{R}^{M \times N}$  and  $\alpha_K, \beta_K$  be two positive numbers. Then,  $\mathcal{A}$  satisfies the non-symmetric restricted isometry property with constants  $\alpha_K, \beta_K$  iff

$$\alpha_K \leq \frac{\|\mathcal{A}\mathbf{x}\|_2^2}{\|\mathbf{x}\|_2^2} \leq \beta_K,$$

is satisfied for any  $\mathbf{x} \in \Sigma_K^N$ .

# Hard Thresholding Methods

$l_0$ -“norm” minimization	$l_1$ -norm minimization
$\hat{x} = \arg \min_{x: x \in \Sigma_K^N} f(x)$	$\hat{x} = \arg \min_{x: \ x\ _1 \leq \lambda} f(x)$

$$\leftarrow f(x) \triangleq \|y - \mathcal{A}x\|_2^2$$

- Both approaches are computationally attractive  $\rightarrow l_0$ -“norm” minimization allows the use of model-CS framework (e.g., incorporate structured sparsity models) [BaraniukCevherDuarteHedge10].
- We focus on: *Iterative Hard Thresholding (IHT)* algorithm [NowakFigueiredo98, KingsburyReeves03, DaubechiesDefriseDeMol04; BlumensathDavies08, ...]
- Let  $H_K(y) = \arg \min_{x: x \in \Sigma_K^N} \|x - y\|_2^2$ ,  $\mathcal{X}_{(\cdot)} = \text{supp}(x_{(\cdot)})$  and  $\nabla f(x) = -2\mathcal{A}^*(y - \mathcal{A}x)$ .

---

## Algorithm 1: IHT Algorithm

---

**Input:**  $y, \mathcal{A}, K$ , Tolerance, MaxIterations

**Initialize:**  $\mathbf{X}_0 \leftarrow \mathbf{X}_{\text{init}}, \mathcal{X}_0 \leftarrow \mathcal{X}_{\text{init}}, i \leftarrow 0$

**repeat**

$$b \leftarrow x_i - \frac{\mu_i}{2} \nabla f(x_i)$$

$$x_{i+1} \leftarrow H_K(b)$$

$$i \leftarrow i + 1$$

(Update current estimate)

(Best  $K$ -term approximation)

**until**  $\|\mathbf{X}_i - \mathbf{X}_{i-1}\|_2 \leq \text{Tolerance}\|\mathbf{X}_i\|_2$  or MaxIterations.

# Convergence Guarantees - prior work

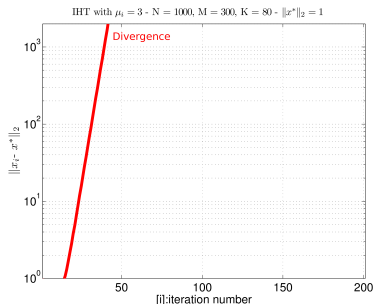
- General convergence formula:

$$\|x_{i+1} - x^*\|_2 \leq \rho \|x_i - x^*\|_2 + \gamma \|\varepsilon\|_2,$$

where

- 1  $\rho \rightarrow$  convergence rate,
- 2  $\gamma \rightarrow$  approximation guarantee.

Reference	Convergence Guarantee ( $\rho < 1$ )	Assumptions
[BlumensathDavies09]	$\delta_{3K} < 1/\sqrt{8}$	$\ \mathcal{A}\ _{2 \rightarrow 2}^2 < 1, \mu_i = 1, \forall i$
[Foucart10]	$\delta_{3K} < 1/2$ or $\delta_{2K} < 1/4$	$\ \mathcal{A}\ _{2 \rightarrow 2}^2 < 1, \mu_i = 1, \forall i$



- Trade-off: the faster you converge, the worse approximation guarantee you get.
- Majority of prior work  $\rightarrow$  optimizing convergence speed, not approximation guarantee.

# Convergence Guarantees - prior work

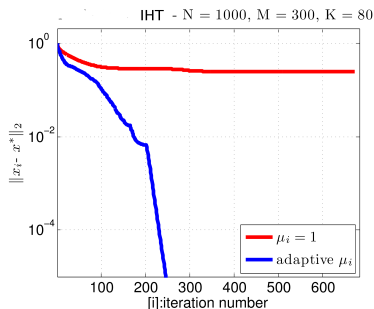
- General convergence formula:

$$\|x_{i+1} - x^*\|_2 \leq \rho \|x_i - x^*\|_2 + \gamma \|\varepsilon\|_2,$$

where

- 1  $\rho \rightarrow$  convergence rate,
- 2  $\gamma \rightarrow$  approximation guarantee.

Reference	Convergence Guarantee ( $\rho < 1$ )	Assumptions
[BlumensathDavies09]	$\delta_{3K} < 1/\sqrt{8}$	$\ \mathcal{A}\ _{2 \rightarrow 2}^2 < 1, \mu_i = 1, \forall i$
[Foucart10]	$\delta_{3K} < 1/2$ or $\delta_{2K} < 1/4$	$\ \mathcal{A}\ _{2 \rightarrow 2}^2 < 1, \mu_i = 1, \forall i$



- Trade-off: the faster you converge, the worse approximation guarantee you get.
- Majority of prior work  $\rightarrow$  optimizing convergence speed, not approximation guarantee.

We present:

- Basic “ingredients” of Hard Thresholding methods.
- Optimal/efficient strategies under various problem assumptions.
- General IHT template (ALgebraic PursuitS, dubbed as ALPS) that “mixes” the ingredients into one framework.



## Step Size Selection



$$x_{i+1} = H_K \left( \underbrace{x_i - \frac{\mu_i}{2} \nabla f(x_i)}_{=b} \right)$$

- Key observation:

“ $x_{i+1}$  is the *best*  $K$ -sparse approximation to  $b$ ”

- Convergence proof of IHT:

$$\|x_{i+1} - b\|_2^2 \leq \|x^* - b\|_2^2 \Rightarrow$$

$$\|x_{i+1} - x^*\|_2 \leq 2\|\mathbb{I} - \mu_i \mathcal{A}_T^* \mathcal{A}_T\|_{2 \rightarrow 2} \|x_i - x^*\|_2 + 2\mu_i \sqrt{1 + \delta_{2K}} \|\varepsilon\|_2 \quad (*)$$

where  $T = \text{supp}(x^*) \cup \text{supp}(x_{i+1}) \cup \text{supp}(x_i)$ .

- Fundamental operator in IHT convergence rate proofs:

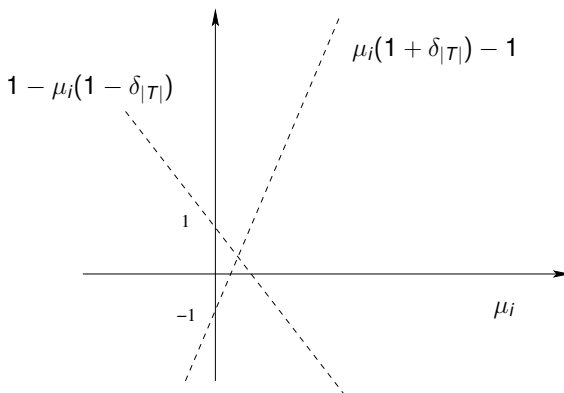
$$\|\mathbb{I} - \mu_i \mathcal{A}_T^* \mathcal{A}_T\|_{2 \rightarrow 2} \leq \max \left\{ \mu_i \lambda_{\max}(\mathcal{A}_T^* \mathcal{A}_T) - 1, \quad 1 - \mu_i \lambda_{\min}(\mathcal{A}_T^* \mathcal{A}_T) \right\}$$

# Constant Step Size Selection

- Symmetric RIP:

$$\lambda(\mathcal{A}_T^* \mathcal{A}_T) \in [1 - \delta_{|T|}, 1 + \delta_{|T|}]$$

- $\min_{\mu_i} \|\mathbb{I} - \mu_i \mathcal{A}_T^* \mathcal{A}_T\|_{2 \rightarrow 2} \leq \min_{\mu_i} \max \left\{ \mu_i(1 + \delta_{|T|}) - 1, 1 - \mu_i(1 - \delta_{|T|}) \right\}$



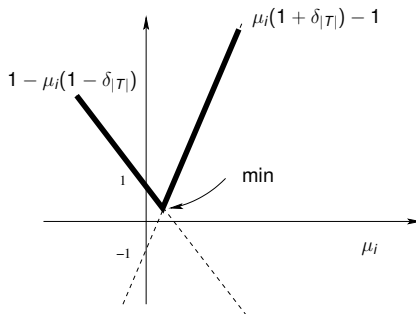
# Constant Step Size Selection

## Lemma 1

Given the symmetric RIP assumption:  $(1 - \delta_K) \leq \frac{\|\mathbf{Ax}\|_2^2}{\|\mathbf{x}\|_2^2} \leq (1 + \delta_K), \forall \mathbf{x} \in \Sigma_K^N$ , the step size  $\mu_i$  that implies the fastest convergence rate in (\*) amounts to

$$\mu_i = 1, \forall i = \{1, 2, \dots, \},$$

where  $\rho = 2\delta_{3K} < 1 \Rightarrow \delta_{3K} < 1/2$  and  $\gamma = 2\sqrt{1 + \delta_{2K}}$ .



## Constant Step Size Selection

- For non-symmetric RIP with known upper/lower bounds, (\*) becomes:

$$\|x_{i+1} - x^*\|_2 \leq 2\|\mathbb{I} - \mu_i \mathbf{A}_T^* \mathbf{A}_T\|_{2 \rightarrow 2} \|x_i - x^*\|_2 + 2\mu_i \sqrt{\beta_{2K}} \|\epsilon\|_2. \quad (**)$$

### Corollary 1

Given non-symmetric RIP with known upper/lower bounds:

$$\alpha_K \leq \frac{\|\mathbf{A}\mathbf{X}\|_2^2}{\|\mathbf{X}\|_2^2} \leq \beta_K, \quad \forall \mathbf{X} \in \Sigma_K^N,$$

the step size  $\mu_i$  that implies the fastest convergence rate in (\*\*) amounts to

$$\mu_i = \frac{2}{\alpha_{3K} + \beta_{3K}}, \quad \forall i = \{1, 2, \dots\},$$

where  $\rho = \frac{2(\beta_{3K} - \alpha_{3K})}{\alpha_{3K} + \beta_{3K}} < 1 \Rightarrow \beta_{3K} < 3\alpha_{3K}$  and  $\gamma = \frac{2\sqrt{\beta_{2K}}}{\alpha_{3K} + \beta_{3K}}$ .

# Potential pitfalls in step size selection

- No knowledge about RIP condition.
- Can we leverage the ideas from convex optimization?
  - ① At each iteration, pick a conservative (small) value for  $\mu_j \rightarrow$  premature termination of the algorithm.
- Perform binary search over step size  $\mu_j$ :
  - ① No knowledge about RIP bounds  $\rightarrow$  we may miss the “sweet”  $\mu_j$  range that leads to convergence near the true vector.

## Adaptive Step Size Selection

- Let  $\bar{\mathcal{X}}_i = \text{supp}(H_K(\nabla_{\mathcal{X}_i^c} f(x_i)))$ .
- Key observation:  $x_{i+1}$  contains non-zero elements at positions from the set  $\mathcal{S}_i = \mathcal{X}_i \cup \bar{\mathcal{X}}_i$ ,  $|\mathcal{S}_i| \leq 2K$ .

$$x_{i+1} = \mathcal{H}_K \left( x_i - \frac{\mu_i}{2} \nabla f(x_i) \right)$$

## Adaptive Step Size Selection

- Let  $\bar{\mathcal{X}}_i = \text{supp}(H_K(\nabla_{\mathcal{X}_i^c} f(x_i)))$ .
- Key observation:  $x_{i+1}$  contains non-zero elements at positions from the set  $\mathcal{S}_i = \mathcal{X}_i \cup \bar{\mathcal{X}}_i$ ,  $|\mathcal{S}_i| \leq 2K$ .

$$x_{i+1} = H_K \left( x_i - \frac{\mu_i}{2} \nabla f(x_i) \right)$$

- $H_K(\cdot) \implies \mathcal{O}(K \log K)$  vs.  $\mathcal{O}(N \log N)$  complexity.
- More sophisticated method: median method with  $\mathcal{O}(K)$  amortized complexity.



# Adaptive Step Size Selection

- Let  $\bar{\mathcal{X}}_i = \text{supp}(H_K(\nabla_{\mathcal{X}_i^c} f(x_i)))$ .
- Key observation:  $x_{i+1}$  contains non-zero elements at positions from the set  $\mathcal{S}_i = \mathcal{X}_i \cup \bar{\mathcal{X}}_i$ ,  $|\mathcal{S}_i| \leq 2K$ .

- Thus:

$$x_{i+1} = H_K \left( \underbrace{x_i - \frac{\mu_i}{2} \nabla_{\mathcal{S}_i} f(x_i)}_{=b} \right).$$

- Observe  $b \in \Sigma_{2K}^N$  since  $|\mathcal{S}_i| \leq 2K$ .
- Calculate step size that minimizes the objective value, i.e.

$$\mu_i = \arg \min_{\mu} \|y - \mathcal{A}b\|_2^2 = \frac{\|\nabla_{\mathcal{S}_i} f(x_i)\|_2^2}{\|\underbrace{\mathcal{A} \nabla_{\mathcal{S}_i} f(x_i)}_{2K\text{-sparse}}\|_2^2}.$$

- RIP:

$$\frac{1}{1 + \delta_{2K}} \leq \mu_i \leq \frac{1}{1 - \delta_{2K}}$$

- Non-symmetric RIP:

$$\frac{1}{\beta_{2K}} \leq \mu_i \leq \frac{1}{\alpha_{2K}}$$

# Convergence guarantees

## Theorem 1 [Iteration Invariant]

Assume  $\mathcal{A} \in \mathbb{R}^{M \times N}$  satisfies:

$$\alpha_K \leq \frac{\|\mathcal{A}\mathbf{X}\|_2^2}{\|\mathbf{X}\|_2^2} \leq \beta_K, \quad \forall \mathbf{X} \in \Sigma_K^N,$$

where  $\alpha_K, \beta_K$  are unknown. Then, in the worst case scenario, IHT with adaptive step size selection approximates the true  $K$ -sparse signal  $\mathbf{X}^*$  with convergence rate  $\rho$  according to:

$$\|\mathbf{X}_{i+1} - \mathbf{X}^*\|_2 \leq \rho \|\mathbf{X}_i - \mathbf{X}^*\|_2 + \gamma \|\varepsilon\|_2,$$

where  $\rho = 2 \max\left\{\frac{\beta_{3K}}{\alpha_{2K}} - 1, 1 - \frac{\alpha_{3K}}{\beta_{2K}}\right\} < 1$  and  $\gamma = \frac{2\sqrt{\beta_{2K}}}{\alpha_{2K}}$ .

- If symmetric RIP holds:  $(1 - \delta_K) \leq \frac{\|\mathcal{A}\mathbf{X}\|_2^2}{\|\mathbf{X}\|_2^2} \leq (1 + \delta_K)$ ,  $\forall \mathbf{X} \in \Sigma_K^N$ , then

$$\|\mathbf{X}_{i+1} - \mathbf{X}^*\|_2 \leq \rho \|\mathbf{X}_i - \mathbf{X}^*\|_2 + \frac{2\sqrt{1 + \delta_{2K}}}{1 - \delta_{2K}} \|\varepsilon\|_2,$$

where  $\rho \triangleq \frac{\delta_{3K} + \delta_{2K}}{1 - \delta_{2K}} < 1 \Rightarrow \delta_{3K} < 1/5$  (worst-case scenario -  $\mu_i = \frac{1}{1 - \delta_{2K}}$ ).

Memory



$$x_{i+1} = H_K(x_i - \frac{\mu_i}{2} \nabla_{S_i} f(x_i))$$

- Idea: why not use information from previous estimates ( $x_{i-1}$ ,  $x_{i-2}$ , etc.)?
- Nesterov's one-memory utilization scheme over convex sets:

$$\tau_i = \frac{\alpha_i(1 - \alpha_i)}{\alpha_i^2 + \alpha_{i+1}}, \quad (1)$$

where  $\alpha_0 \in (0, 1)$  and  $\alpha_{i+1} \in (0, 1)$  is computed as the root of

$$\alpha_{i+1}^2 = (1 - \alpha_{i+1})\alpha_i^2 + q\alpha_{i+1}, \text{ for } q \triangleq \frac{\lambda_{\min}(\mathcal{A}^* \mathcal{A})}{\lambda_{\max}(\mathcal{A}^* \mathcal{A})}. \quad (2)$$

- In our case:

$$x_i = H_K(y_i - \frac{\mu_i}{2} \nabla_{S_i} f(y_i)), \quad y_{i+1} = x_i + \tau_i(x_i - x_{i-1})$$

where  $\mathcal{Y}_i = \text{supp}(y_i)$  and  $S_i = \mathcal{Y}_i \cup \text{supp}(H_K(\nabla_{\mathcal{Y}_i^c} f(y_i)))$  with  $|S_i| \leq 3K$ .

# Nesterov's scheme

- Nesterov's scheme is not optimal when the algorithm is linear convergent.

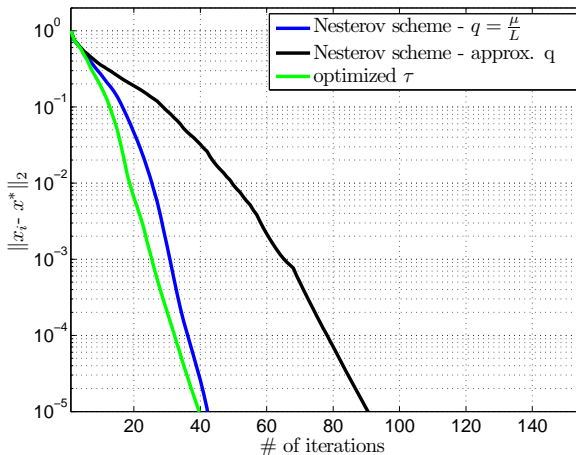


Figure: ALPS convergence rate example using memory.  $N = 2000$ ,  $M = 600$ ,  $K = 120$ . Blue and black lines represent Nesterov's  $\tau_i$  selection scheme with  $q = \frac{\lambda_{\min}(\mathcal{A}^* \mathcal{A})}{\lambda_{\max}(\mathcal{A}^* \mathcal{A})}$  and  $q \sim \frac{\mu_{i \min}}{\mu_{i \max}}$ , respectively; green line represents the proposed momentum step size selection.

# Momentum step size selection

- Momentum step size selection:

- (i) Constant  $\tau_i$ , e.g.  $\tau_i = 1/2, \forall i \rightarrow$  No additional computational cost,
- (ii) Nesterov's scheme  $\tau_i = (a_i - 1)/(a_{i+1})$  where  $a_{i+1} = \frac{1 + \sqrt{4a_i^2 + 1}}{2} \rightarrow$  No additional computational cost,
- (iii) Objective minimizer  $\tau_i$ :  $\tau_i = \arg \min_{\tau} \|\mathbf{y} - \mathcal{A}\mathbf{y}_{i+1}\|_2^2 \rightarrow$  **No additional computational cost (!!!):**

$$\begin{aligned}\tau_i &= \arg \min_{\tau} \|\mathbf{y} - \mathcal{A}\mathbf{y}_{i+1}\|_2^2 \\ &= \arg \min_{\tau} \|(\mathbf{y} - \mathcal{A}\mathbf{x}_i) - \tau \mathcal{A}(\mathbf{x}_i - \mathbf{x}_{i-1})\|_2^2 \\ &= \frac{\langle \mathbf{y} - \mathcal{A}\mathbf{x}_i, \mathcal{A}\mathbf{x}_i - \mathcal{A}\mathbf{x}_{i-1} \rangle}{\|\mathcal{A}\mathbf{x}_i - \mathcal{A}\mathbf{x}_{i-1}\|_2^2}\end{aligned}$$

where  $\mathcal{A}\mathbf{x}_i, \mathcal{A}\mathbf{x}_{i-1}$  are previously computed - similar memory-based “tricks” in [Blumensath10].

Gradient updates over restricted support sets



# Gradient updates over restricted support sets

- Proxy vector  $b = x_i - \frac{\mu_i}{2} \nabla_{\mathcal{S}_i} f(x_i)$ : gradient descent over support set  $\mathcal{S}_i$ .
- Alternatively,  $b$  can be computed as the minimizer:

$$b = \arg \min_{x: \text{supp}(x) \subseteq \mathcal{S}_i} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2.$$

- Based on Hard Thresholding Pursuit [Foucart10] and Subspace Pursuit [WeiMilenkovic09], we can further refine  $x_{i+1} = H_K(b)$  by

$$x_{i+1} = x_{i+1} - \frac{\bar{\mu}_i}{2} \nabla_{\mathcal{X}_{i+1}} f(x_{i+1}), \quad \text{where } \bar{\mu}_i = \frac{\|\nabla_{\mathcal{X}_{i+1}} f(x_{i+1})\|_2^2}{\|\mathbf{A} \nabla_{\mathcal{X}_{i+1}} f(x_{i+1})\|_2^2},$$

or

$$x_{i+1} = \arg \min_{x: \text{supp}(x) \subseteq \mathcal{X}_{i+1}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2.$$



## Experiments

(Cookbook: ALgebraic PursuitS (ALPS))  
Please check <http://lions.epfl.ch/ALPS/>



# Execution time

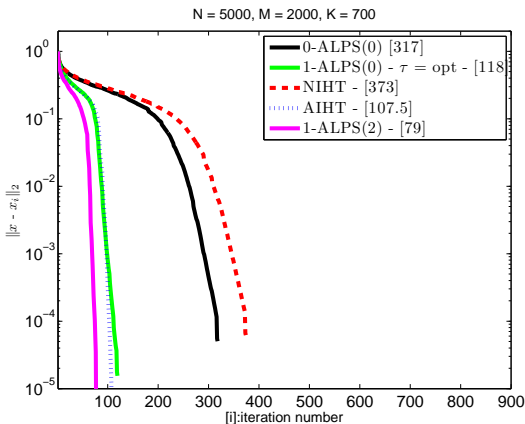


Figure: Median error per iter. - [median # of iter.]  
 #-ALPS(0): adaptive  $\mu_i$  with # memory,  
 NIHT: Normalized IHT, AIHT: NIHT with  
 Double Relaxation, 1-ALPS(2): adaptive  $\mu_i$   
 and additional gradient update.

- Remark:** model-based projection  $\mathcal{M}_K(\cdot) \Rightarrow$  # of thresholding operations **does matter!**

## Complexity per iter.

0-ALPS(0)	$\mathcal{O}(MN) + 3\mathcal{O}(MK)$
*NIHT	$\mathcal{O}(MN) + 2\mathcal{O}(MK)$
*AIHT	$\mathcal{O}(MN) + 3\mathcal{O}(MK)$
1-ALPS(0)	$\mathcal{O}(MN) + 3\mathcal{O}(MK)$
1-ALPS(2)	$2\mathcal{O}(MN) + 5\mathcal{O}(MK)$

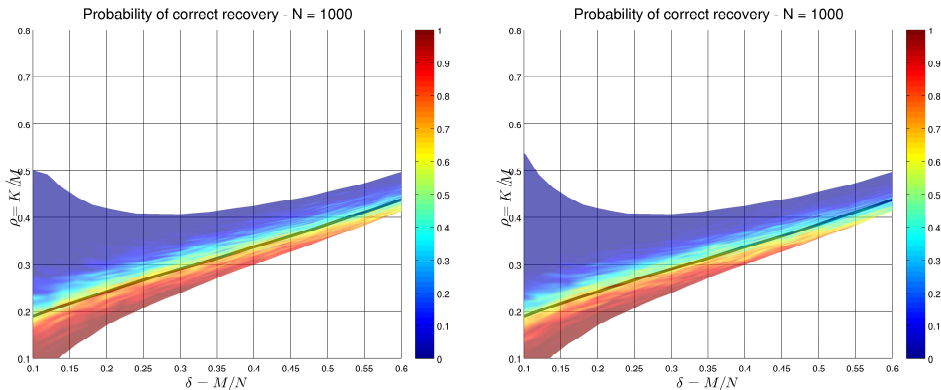
## # of $H_K(\cdot)$ per iter.

0-ALPS(0)	2
*NIHT	2
*AIHT	3
1-ALPS(0)	2
1-ALPS(2)	2

Table: (\*) Best case scenario

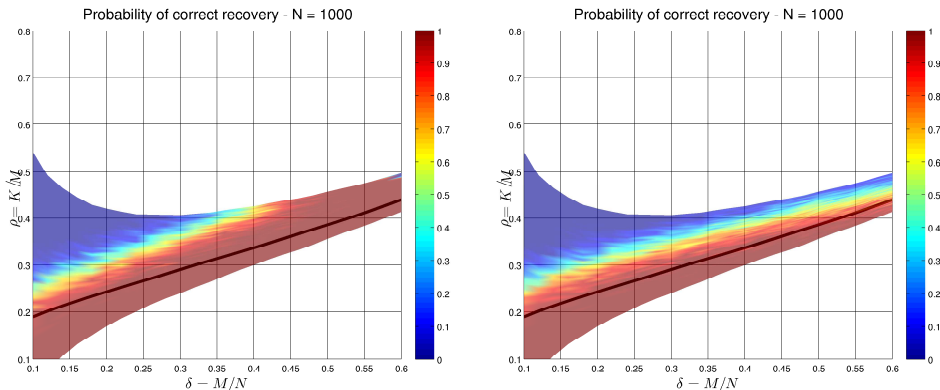
- Less # of iterations  $\Rightarrow$  faster algorithm.

# Phase transitions



**Figure:** Empirical phase transition performance of 1-ALPS(0) (left column) and AIHT (right column) algorithms.  $\rho = K/M$  and  $\delta = M/N$  where  $0 \leq \rho, \delta \leq 1$ . A signal recovery with solution  $\hat{\mathbf{X}}$  is considered successful provided that  $\|\hat{\mathbf{X}} - \mathbf{X}^*\|_2 < 10^{-6}$ . Solid black line denotes the theoretical  $l_1$  minimization phase transition curve.

# Phase transitions



**Figure:** Empirical phase transition performance of HTP with the proposed step size selection (left column) and HTP with NIHT step size selection (right column).  $\rho = K/M$  and  $\delta = M/N$  where  $0 \leq \rho, \delta \leq 1$ . A signal recovery with solution  $\hat{\mathbf{X}}$  is considered successful provided that  $\|\hat{\mathbf{X}} - \mathbf{X}^*\|_2 < 10^{-6}$ . Solid black line denotes the theoretical  $l_1$  minimization phase transition curve.

- Basic “ingredients” of IHT method:
  - ① Step size  $\mu_i$ ,
  - ② Memory,
  - ③ Additional gradient updates on restricted support sets.
- Step size selection  $\mu_i$ : different strategies for different problem assumptions.
- Memory: usage of memory leads to faster convergence rate with (almost) no additional computational cost.
- Additional gradient updates over restricted support sets  $\rightarrow$  better phase transition performance.

Thank you

Bon Appetit!



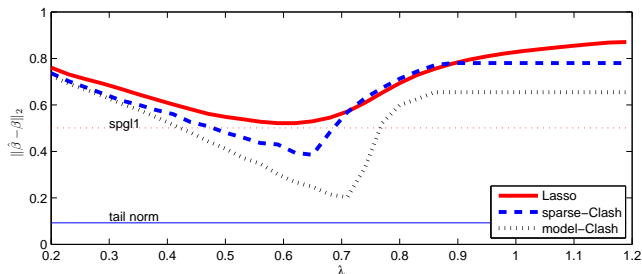
# References

- Please check out the ALPS homepage for further information and codes:

<http://lions.epfl.ch/ALPS>

- Also, check out the CLASH (Combinatorial selection and Least Absolute SHrinkage) algorithm:

<http://lions.epfl.ch/CLASH>



# Adaptive Step Size Selection - prior work

- Normalized Iterative Hard Thresholding [BlumensathDavies10].

---

## Algorithm 2: NIHT Algorithm

---

Input:  $u, \Phi, K$ , Stopping Criteria

Initialize:  $\mathbf{x}_0 \leftarrow \mathbf{x}_{\text{init}}, \mathcal{X}_0 \leftarrow \mathcal{X}_{\text{init}}, i \leftarrow 0$

repeat

$$\mu_i = \frac{\|\nabla_{\mathcal{X}_i} f(\mathbf{x}_i)\|_2^2}{\|\Phi \nabla_{\mathcal{X}_i} f(\mathbf{x}_i)\|_2^2} \quad (\text{Step size selection})$$

$$\mathbf{b} \leftarrow \mathbf{x}_i - \frac{\mu_i}{2} \nabla f(\mathbf{x}_i) \quad (\text{Update current estimate})$$

$$\hat{\mathbf{x}}_{i+1} \leftarrow H_K(\mathbf{b}) \quad (\text{Best } K\text{-term approximation})$$

if  $\text{supp}(\hat{\mathbf{x}}_{i+1}) = \mathcal{X}_i$  then

$$\mathcal{X}_{i+1} = \mathcal{X}_i$$

else

Iterate by decreasing  $\mu_i$  until specific step size criteria are met.

end if

$$i \leftarrow i + 1$$

until Stopping criteria are met.

---

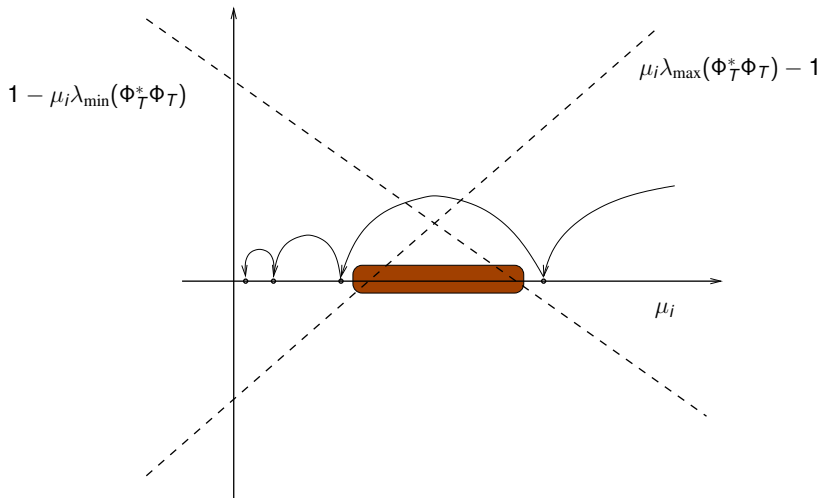


## Adaptive Step Size Selection - prior work

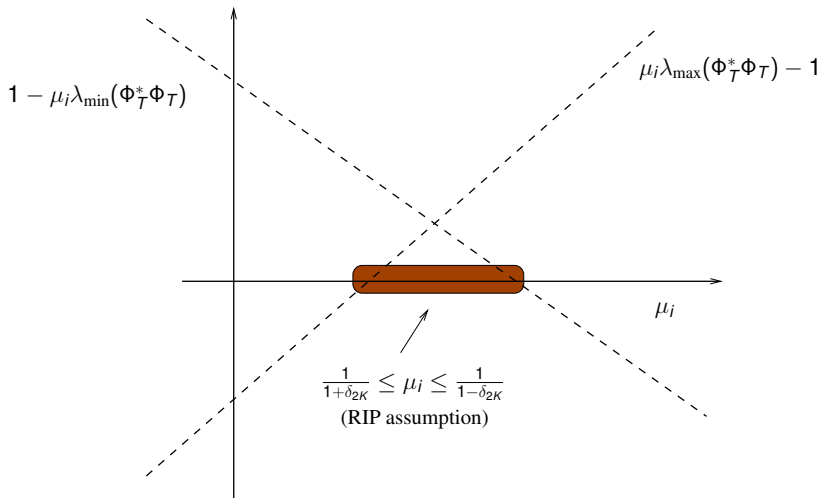
- Normalized Iterative Hard Thresholding [BlumensathDavies10].

## Adaptive Step Size Selection - prior work

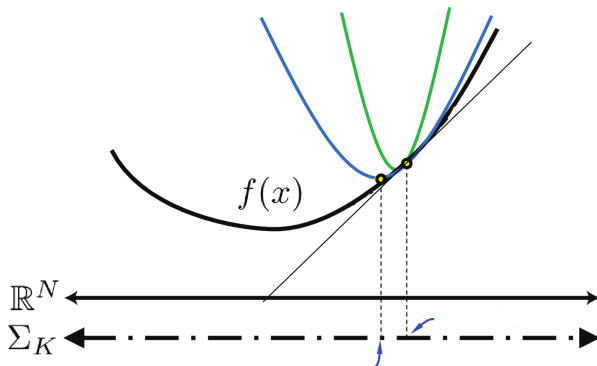
- Why NIHT with binary search over  $\mu_i$  may not work?



# Adaptive step size selection



# Adaptive step size selection - convex optimization approach



	<b>Convex-based</b>
Solution via	Convex relaxation $\ \cdot\ _* + \ \cdot\ _1, \dots$
Criteria example	$\min_{\ \mathbf{y} - \mathcal{A}(\mathbf{L} + \mathbf{M})\ _2 \leq \sigma} \ \mathbf{L}\ _* + \lambda \ \mathbf{M}\ _1$
Algorithms	(S)PCP <sup>1,2,3</sup> , CPCP <sup>1,2,3,4</sup> , SVT <sup>1,3</sup> , ...
	<b>Greedy-based</b>
Solution via	Non-convex projections, ...
Criteria example	$\min_{\text{rank}(\mathbf{L}) \leq k, \ \mathbf{M}\ _0 \leq s} \ \mathbf{y} - \mathcal{A}(\mathbf{L} + \mathbf{M})\ _2^2$
Algorithms	SpaRCS <sup>1,2,3,4</sup> , GoDec <sup>1,2</sup> , SVP <sup>1,3</sup> , ...
	<b>Manifold-based</b>
Solution via	Manifold Trust regions, subspace identification, ...
Criteria example	$\min_{\text{rank}(\mathbf{US}) \leq k, \ \mathbf{M}\ _0 \leq s} \ \mathbf{y} - \mathcal{A}(\mathbf{US} + \mathbf{M})\ _2^2$
Algorithms	RTRMC <sup>1</sup> , GROUSE <sup>1</sup> , GRASTA <sup>1</sup> , ...

<sup>1</sup>MC, <sup>2</sup>RPCA, <sup>3</sup>ARM, <sup>4</sup>handles CS data