PhD Defense

William de Vazelhes¹, PhD Candidate
Dr. Bin Gu¹, Supervisor
Dr. Xiaotong Yuan², External Examiner
Dr. Chih-Jen Lin¹, Internal Member
Dr. Karthik Nandakumar¹, Internal Member
Dr. Zhiqiang Xu¹, Internal Member

April 17, 2024



Mohamed bin Zayed University of Artificial Intelligence, ² Nanjing University

- 1 Iterative Hard Thresholding
 - Introduction
 - Convergence Rate
- 2 Zeroth-Order Hard-Thresholding
 - Introduction
 - Convergence Rate
- 3 Additional Constraints
 - Introduction
 - Convergence Rate
- 4 A Dual Perspective on IHT
 - Interlude
 - IRKSN
 - Conditions for recovery
- 5 QA

Iterative Hard Thresholding

Iterative Hard Thresholding

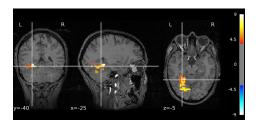
Introduction

Sparse Optimization:

$$\min_{\boldsymbol{x} \in \mathbb{R}^d: \|\boldsymbol{x}\|_0 \le k} f(\boldsymbol{x})$$

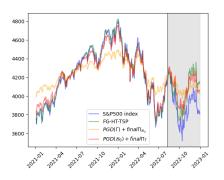
└─Introduction

Application: fMRI



- x: map of functional region of the brain (d = number of voxels)
- $f(\mathbf{x}) := \|\mathbf{y} \mathbf{A}\mathbf{x}\|^2$ with $y_i \in \{-1, 1\}$ standing for $\{'face', 'house'\}$ and $\mathbf{A}_{i, \cdot}$ being the recorded activation map at time i.

Application: Index Tracking



- **x**: amount invested in each of d stocks
- $f(\mathbf{x}) := \|\mathbf{y} \mathbf{A}\mathbf{x}\|^2$ with \mathbf{y}_i : S&P returns for day i, $\mathbf{A}_{i,j}$: return of stock j on day i

Application: Sparse Adversarial Attacks



Perturbation x



'bird'



'dog'

- x: perturbation of an image z
- $f(x) = \max\{F_y(\text{clip}(z+x)) \max_{j \neq y} F_j(\text{clip}(z+x)), 0\}$ with y: true class of the image, F_j : prediction score for class j

└─ Introduction

The Iterative Hard Thresholding (IHT) algorithm

```
Algorithm 1: Iterative Hard-Thresholding (IHT)

Initialization: \mathbf{x}_0

for t = 0, ..., T do

\mathbf{x}_{t+1} := \mathcal{H}_k(\mathbf{x}_t - \eta \nabla f(\mathbf{x}_t))

end

Output: \hat{\mathbf{x}_T} := \text{e.g. } \mathbf{x}_T \text{ or arg min}_{\mathbf{x} \in \{\mathbf{x}_t\}_{t=1}^T} f(\mathbf{x}_t)
```

$$\mathcal{H}_k(\mathbf{x}) := \min_{\mathbf{y} \in \mathcal{B}_0(k)} \|\mathbf{y} - \mathbf{x}\|_2$$

 $\mathcal{B}_0(k) := \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_0 \le k\}$

Goal: Convergence Rate

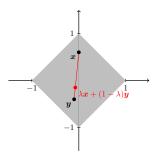
Goal: Prove Convergence Rate Why?

- To make sure it does not diverge.
- To have an estimate of how feasible it is for a large scale task.
- To set the hyperparameters of the algorithm properly (e.g. η).

Warm Up: Convex Case

$$\min_{\mathbf{x}\in\mathcal{C}}f(\mathbf{x})$$

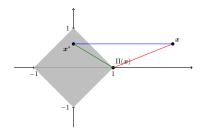
with C convex : $\forall (x,y) \in (C)^2$: $\lambda x + (1-\lambda)y \in C$.



Projection onto ${\cal C}$

3 Point Lemma:

$$\|\mathbf{x} - \mathbf{x}^*\|^2 \ge \|\Pi_{\mathcal{C}}(\mathbf{x}) - \mathbf{x}\|^2 + \|\Pi_{\mathcal{C}}(\mathbf{x}) - \mathbf{x}^*\|^2.$$

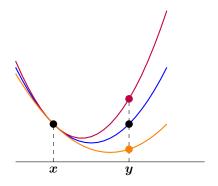


Proj. onto the ℓ_1 unit ball.

Strong Convexity and Smoothness

Assumptions: strong convexity and smoothness. $\forall (x, y) \in \mathcal{C}^2$:

$$f(x) + \langle \nabla f(x), y - x \rangle + \frac{\nu}{2} ||x - y||^2 \le f(y) \le f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} ||x - y||^2$$



Proof of Convergence (Convex Case)

Take
$$\eta := \frac{1}{L}$$
.

$$f(\mathbf{x}_{t}) \leq f(\mathbf{x}_{t-1}) + \langle \nabla f(\mathbf{x}_{t-1}), \mathbf{x}_{t} - \mathbf{x}_{t-1} \rangle + \frac{L}{2} \| \mathbf{x}_{t} - \mathbf{x}_{t-1} \|^{2}$$

$$= f(\mathbf{x}_{t-1}) + \frac{L}{2} \| \mathbf{x}_{t} - \mathbf{x}_{t-1} + \frac{1}{L} \nabla f(\mathbf{x}_{t-1}) \|^{2} - \frac{1}{2L} \| \nabla f(\mathbf{x}_{t-1}) \|^{2}$$

$$\leq f(\mathbf{x}_{t-1}) + \frac{L}{2} \| \mathbf{x}^{*} - \mathbf{x}_{t-1} + \frac{1}{L} \nabla f(\mathbf{x}_{t-1}) \|^{2} - \frac{L}{2} \| \mathbf{x}_{t} - \mathbf{x}^{*} \|^{2} - \frac{1}{2L} \| \nabla f(\mathbf{x}_{t-1}) \|^{2}$$

$$= f(\mathbf{x}_{t-1}) + \langle \nabla f(\mathbf{x}_{t-1}), \mathbf{x}^{*} - \mathbf{x}_{t-1} \rangle + \frac{L}{2} \| \mathbf{x}_{t-1} - \mathbf{x}^{*} \|^{2} - \frac{L}{2} \| \mathbf{x}_{t} - \mathbf{x}^{*} \|^{2}$$

$$\leq f(\mathbf{x}^{*}) + \frac{L - \nu}{2} \| \mathbf{x}_{t-1} - \mathbf{x}^{*} \|^{2} - \frac{L}{2} \| \mathbf{x}_{t} - \mathbf{x}^{*} \|^{2}$$

Proof of Convergence (Convex Case)

$$\left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}}\right)^{T-t} [f(\mathbf{x}_t) - f(\mathbf{x}^*)] \le \left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}}\right)^{T-t} \frac{L-\nu}{2} \|\mathbf{x}_{t-1} - \mathbf{x}^*\|^2 - \left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}}\right)^{T-t} \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2$$

$$\left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}}\right)^{T-2} \left[f(x_2) - f(x^*)\right] \le \left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}}\right)^{T-2} \frac{L-\nu}{2} \|x_1 - x^*\|^2 - \left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}}\right)^{T-2} \frac{L}{2} \|x_2 - x^*\|^2$$

$$\left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}}\right)^{T-1} \left[f(\mathbf{x}_1) - f(\mathbf{x}^*)\right] \le \left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}}\right)^{T-1} \frac{L-\nu}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}}\right)^{T-1} \frac{L}{2} \|\mathbf{x}_1 - \mathbf{x}^*\|^2$$

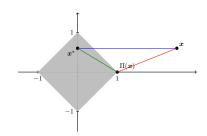
$$\sum_{t=1}^{T} \left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}} \right)^{T-t} \left[f(\mathbf{x}_t) - f(\mathbf{x}^*) \right] \leq \left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}} \right)^{T-1} \frac{L-\nu}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}} \right)^{T-t} \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2$$

$$f(\mathbf{x}_{\hat{T}}) - f(\mathbf{x}^*) \le C\omega^T$$

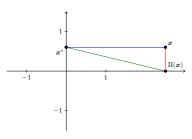
Non-Convex case: C is the ℓ_0 pseudo-ball

$$\|\mathbf{x} - \mathbf{x}^*\|^2 \ge \|\mathcal{H}_k(\mathbf{x}) - \mathbf{x}\|^2 + \left(1 - \sqrt{\frac{k^*}{k}}\right) \|\mathcal{H}_k(\mathbf{x}) - \mathbf{x}^*\|^2.$$

$$\mathbf{x}^* \in \mathcal{B}_0(k^*), \quad k^* \le k$$



Proj. onto the ℓ_1 unit ball.



Proj. onto the ℓ_0 unit pseudo-ball.

Non-convex case: Assumptions

Assumptions: restricted strong convexity and restricted smoothness. $\forall (x, y) \in \mathbb{R}^d$ s.t. $||x - y||_0 \le s$ (s := 3k).

$$f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\nu_s}{2} \|\mathbf{x} - \mathbf{y}\|^2 \le f(\mathbf{y}) \le f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L_s}{2} \|\mathbf{x} - \mathbf{y}\|^2$$

Proof of Convergence (IHT)

We take
$$\eta=rac{1}{L_s}$$
, and $k\geq 4\kappa_s^2k^*$, with $\kappa_s:=rac{L_s}{
u_s}\implies \sqrt{eta}\leq rac{
u_s}{2L_s}.$

$$f(\mathbf{x}_{t}) \leq f(\mathbf{x}_{t-1}) + \langle \nabla f(\mathbf{x}_{t-1}), \mathbf{x}_{t} - \mathbf{x}_{t-1} \rangle + \frac{L_{s}}{2} \| \mathbf{x}_{t} - \mathbf{x}_{t-1} \|^{2}$$

$$= f(\mathbf{x}_{t-1}) + \frac{L_{s}}{2} \| \mathbf{x}_{t} - \mathbf{x}_{t-1} + \frac{1}{L_{s}} \nabla f(\mathbf{x}_{t-1}) \|^{2} - \frac{1}{2L_{s}} \| \nabla f(\mathbf{x}_{t-1}) \|^{2}$$

$$\leq f(\mathbf{x}_{t-1}) + \frac{L_{s}}{2} \| \mathbf{x}^{*} - \mathbf{x}_{t-1} + \frac{1}{L_{s}} \nabla f(\mathbf{x}_{t-1}) \|^{2} - \frac{L_{s}}{2} (1 - \sqrt{\beta}) \| \mathbf{x}_{t} - \mathbf{x}^{*} \|^{2} - \frac{1}{2L_{s}} \| \nabla f(\mathbf{x}_{t-1}) \|^{2}$$

$$= f(\mathbf{x}_{t-1}) + \langle \nabla f(\mathbf{x}_{t-1}), \mathbf{x}^{*} - \mathbf{x}_{t-1} \rangle + \frac{L_{s}}{2} \| \mathbf{x}_{t-1} - \mathbf{x}^{*} \|^{2} - \frac{L_{s}}{2} (1 - \sqrt{\beta}) \| \mathbf{x}_{t} - \mathbf{x}^{*} \|^{2}$$

$$\leq f(\mathbf{x}^{*}) + \frac{L_{s} - \nu_{s}}{2} \| \mathbf{x}_{t-1} - \mathbf{x}^{*} \|^{2} - \frac{2L_{s} - \nu_{s}}{4} \| \mathbf{x}_{t} - \mathbf{x}^{*} \|^{2}$$

$$\leq f(\mathbf{x}^{*}) + \frac{L_{s} - \nu_{s}}{2} \| \mathbf{x}_{t-1} - \mathbf{x}^{*} \|^{2} - \frac{2L_{s} - \nu_{s}}{4} \| \mathbf{x}_{t} - \mathbf{x}^{*} \|^{2}$$

Literative Hard Thresholding

Convergence Rate

Proof of Convergence (IHT)

$$\left(\frac{\frac{L_{s}-\nu_{s}}{2}}{\frac{2L_{s}-\nu_{s}}{4}}\right)^{T-t} [f(\mathbf{x}_{t})-f(\mathbf{x}^{*})] \leq \left(\frac{\frac{L_{s}-\nu_{s}}{2}}{\frac{2L_{s}-\nu_{s}}{4}}\right)^{T-t} \frac{L_{s}-\nu_{s}}{2} \|\mathbf{x}_{t-1}-\mathbf{x}^{*}\|^{2} - \left(\frac{\frac{L_{s}-\nu_{s}}{2}}{\frac{2L_{s}-\nu_{s}}{4}}\right)^{T-t} \frac{2L_{s}-\nu_{s}}{4} \|\mathbf{x}_{t}-\mathbf{x}^{*}\|^{2}$$

٠.

$$\left(\frac{\frac{L_{s}-\nu_{s}}{2}}{\frac{2L_{s}-\nu_{s}}{4}}\right)^{T-2} [f(x_{2}) - f(x^{*})] \leq \left(\frac{\frac{L_{s}-\nu_{s}}{2}}{\frac{2L_{s}-\nu_{s}}{4}}\right)^{T-2} \frac{L_{s}-\nu_{s}}{2} \|x_{1}-x^{*}\|^{2} - \left(\frac{\frac{L_{s}-\nu_{s}}{2}}{\frac{2L_{s}-\nu_{s}}{4}}\right)^{T-2} \frac{2L_{s}-\nu_{s}}{4} \|x_{2}-x^{*}\|^{2}$$

$$\left(\frac{\frac{L_{s}-\nu_{s}}{2}}{\frac{2L_{s}-\nu_{s}}{4}}\right)^{T-1}\left[f(x_{1})-f(x^{*})\right] \leq \left(\frac{\frac{L_{s}-\nu_{s}}{2}}{\frac{2L_{s}-\nu_{s}}{4}}\right)^{T-1}\frac{L_{s}-\nu_{s}}{2}\|x_{0}-x^{*}\|^{2}-\left(\frac{\frac{L_{s}-\nu_{s}}{2}}{\frac{2L_{s}-\nu_{s}}{4}}\right)^{T-1}\frac{2L_{s}-\nu_{s}}{4}\|x_{1}-x^{*}\|^{2}$$

$$\sum_{t=1}^{T} \left(\frac{\frac{L_s - \nu_s}{2}}{\frac{2L_s - \nu_s}{4}} \right)^{T-t} \left[f(x_t) - f(x^*) \right] \leq \left(\frac{\frac{L_s - \nu_s}{2}}{\frac{2L_s - \nu_s}{4}} \right)^{T-1} \frac{L_s - \nu_s}{2} \left\| x_0 - x^* \right\|^2 - \left(\frac{\frac{L_s - \nu_s}{2}}{\frac{2L_s - \nu_s}{4}} \right)^{T-t} \frac{2L_s - \nu_s}{4} \left\| x_t - x^* \right\|^2$$

 $f(\mathbf{x}_{\hat{T}}) - f(\mathbf{x}^*) \leq C\omega^T$

Zeroth-Order Hard-Thresholding

Zeroth-Order Hard-Thresholding

Zeroth-Order Hard-Thresholding (ZOHT)

Algorithm 2: Hard-Thresholding

Initialization:
$$x_0$$

for $t = 0, ..., T$ do
$$x_{t+1} := \mathcal{H}_k(x_t - \eta \nabla f(x_t))$$

end

Output: $\hat{x_T} := \text{e.g. } x_T \text{ or } \arg\min_{x \in \{x_i\}_{t=1}^T} f(x_t)$

What if we don't know $\nabla f(\mathbf{x}_t)$? e.g. for privacy or computational reasons.

Approximating $\nabla f(x)$: two points approximation [1] [2]:

One random direction u:

$$oldsymbol{g}_t = drac{f(oldsymbol{x}_t + \mu oldsymbol{u}) - f(oldsymbol{x}_t)}{\mu} oldsymbol{u} \quad ext{with} \quad oldsymbol{u} \sim ext{Uni}(\mathbb{S}_d)$$

q random directions $\{u_i\}_{i=1}^q$:

$$m{g}_t = rac{d}{q} \sum_{i=1}^q rac{f(m{x}_t + \mu m{u}_i) - f(m{x}_t)}{\mu} m{u}_i \; \; ext{with} \; \; \{m{u}_i\}_{i=1}^q \stackrel{ ext{i.i.d.}}{\sim} \mathsf{Uni}(\mathbb{S}_d)$$

Curse of dimensionality: An impossibility result [5]

Under standard assumptions (strongly cvx, smooth, noisy obs.):

"\forall algorithm, $\exists f_{adv} \ s.t.$ we need more than $O(d/\varepsilon^2)$ queries to achieve $\mathbb{E}[f_{adv}(\hat{\mathbf{x}}_T) - f_{adv}(\mathbf{x}_*)] \leq \varepsilon$ "

Solutions in litterature: more assumptions on f:

- f(x) = g(Ax) with rank $(A) \ll d$ [3]
- sparse/compressible gradients [4]
- What happens in our non-convex case ?

Key Insight: Error of g_t on a Support F

 $F := \operatorname{supp}(\mathbf{x}_t) \cup \operatorname{supp}(\mathbf{x}_{t-1}) \cup \operatorname{supp}(\mathbf{x}^*) \implies |F| = O(k).$

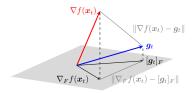
Bias:

$$\|[\mathbb{E}\boldsymbol{g}_t]_F - [\nabla f(\boldsymbol{x}_t)]_F\|^2 \le L^2 \epsilon_\mu \mu^2$$

Variance:

$$\mathbb{E}\|[\mathbf{g}_t]_F - \mathbb{E}[\mathbf{g}_t]_F\|^2 \leq \frac{\varepsilon_F}{q} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{\varepsilon_{abs}}{q} \mu^2, \text{ with } \varepsilon_F = O(k)$$

⇒ Dimension Independent! (Note: we assume full smoothness here for simplicity)





ZOHT: Convergence Analysis

Proof is similar as before, except that we:

- "extract" out the error terms
- keep the constants free at the beginning, and later choose them to make things work

$$\begin{split} f(x_t) &\leq f(x_{t-1}) + \frac{1}{2\eta} \|x^* - x_{t-1}\|^2 - \langle \nabla f(x_{t-1}), x_{t-1} - x^* \rangle + \langle [\nabla f(x_{t-1}) - \mathbf{g}_{t-1}]_F, x_{t-1} - x^* \rangle \\ &- \frac{1}{2\eta} (1 - \sqrt{\beta}) \|x_t - x^*\|^2 + \left[\frac{L - \frac{1}{\eta} + C}{2} \right] \|x_t - x_{t-1}\|^2 + \frac{1}{2C} \|[\nabla f(x_{t-1}) - \mathbf{g}_{t-1}]_F\|^2 \end{split}$$

ZOHT: Convergence Analysis

Choose $\eta:=\frac{1}{L+C}=\frac{1}{\alpha L}$, $k\geq 16\alpha^2\kappa_s^2k^*$ $q_t:=\left\lceil\frac{\tau}{\omega^t}\right\rceil$ with $\omega:=1-\frac{1}{8\alpha\kappa_s}$ and $\tau:=16\kappa_s\frac{\varepsilon_F}{(\alpha-1)}$. Use algebraic manipulations, RSC, expression of bias and variance, and smoothness again:

$$\mathbb{E}f(\mathbf{x}_{t}) - f(\mathbf{x}^{*}) \leq \frac{1}{2\eta} \left[\left(1 - \frac{1}{\alpha' \kappa_{s}} \right) \mathbb{E} \|\mathbf{x}^{*} - \mathbf{x}_{t-1}\|^{2} - (1 - \sqrt{\beta}) \mathbb{E} \|\mathbf{x}_{t} - \mathbf{x}^{*}\|^{2} + 2\eta \left(\frac{G}{2} C_{3} + \frac{1}{C} \left(2C_{1} \|\nabla f(\mathbf{x}^{*})\|^{2} + C_{2}\mu^{2} + C_{3} \right) \right) \right]$$

ZOHT: Convergence Analysis

$$\mathbb{E}f(\hat{\mathbf{x}}_{T}) - f(\mathbf{x}^{*}) \leq F\omega^{T} + H\mu^{2}$$

Query Complexity =
$$\mathcal{O}\left(\frac{\varepsilon_F \kappa_s^3 L}{\varepsilon}\right) = \mathcal{O}\left(\frac{k \kappa_s^3 L}{\varepsilon}\right)$$

Dimension Independent!

L-Additional Constraints

IHT with Additional Constraints

IHT + Additional Constraints

We now consider the following problem:

$$\min_{\boldsymbol{x} \in \mathbb{R}^d: \|\boldsymbol{x}\|_0 \le k, \ \boldsymbol{x} \in \Gamma} f(\boldsymbol{x})$$

Application: e.g. Index Tracking with sector constraints.

 $\Gamma = \{ \boldsymbol{x} \in \mathbb{R}^d : \forall i \in [c], \|\boldsymbol{x}_{G_i}\|_1 \leq D \}$, where \boldsymbol{x}_{G_i} is the restriction of \boldsymbol{x} to group G_i (i.e. for $j \in [d]$, $\boldsymbol{x}_{G_ij} = \boldsymbol{x}_j$ if $j \in G_i$ and 0 otherwise).

IHT + Additional Constraints

Assumption (k-support-preserving set)

 Γ is convex and for any $\mathbf{x} \in \mathbb{R}^d$ s.t. $\|\mathbf{x}\|_0 \le k$: $\sup_{\mathbf{x} \in \mathbb{R}^d} \sup_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{x}\|_0 \le k$:

Algorithm 3: IHT with Two-Step Proj. (TSP)

Initialization: x_0 for t = 0, ..., T do $v_t := \mathcal{H}_k(x_t - \eta \nabla f(x_t))$ $x_{t+1} := \Pi_{\Gamma}(v_t)$

end

Output: $\hat{x_T} := \text{e.g. } x_T \text{ or arg min}_{x \in \{x_i\}_{t=1}^T} f(x_t)$

Support Preserving Set and TSP

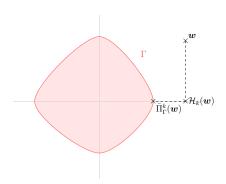


Figure: Support-preserving set and two-step projection (d = 2, k = 1).

$$\bar{\Pi}_{\Gamma}^{k}(\boldsymbol{w}) := \Pi_{\Gamma}(\mathcal{H}_{k}(\boldsymbol{w}))$$

3 Point Lemma with Extra Constraint

New Three (Four) - Point Lemma:

$$\|\bar{\Pi}_{\Gamma}^{k}(\mathbf{x}) - \mathbf{x}\|^{2} \leq \|\mathbf{x} - \mathbf{x}^{*}\|^{2} - \|\bar{\Pi}_{\Gamma}^{k}(\mathbf{x}) - \mathbf{x}^{*}\|^{2} + \sqrt{\beta}\|\mathcal{H}_{k}(\mathbf{x}) - \mathbf{x}^{*}\|^{2}$$

Proof of Convergence

With $\rho \in (0, \frac{1}{2}]$ and $k \ge \frac{4(1-\rho)^2 L_s^2}{\rho^2 \nu_s^2} k^*$:

$$(f(\mathbf{x}_{t}) \leq f(\mathbf{x}^{*}) + \frac{L_{s} - \nu_{s}}{2} \|\mathbf{x}_{t-1} - \mathbf{x}^{*}\|^{2} - \frac{L_{s}}{2} \|\mathbf{x}_{t} - \mathbf{x}^{*}\|^{2} + \frac{L_{s}}{2} \sqrt{\beta} \|\mathbf{v}_{t} - \mathbf{x}^{*}\|^{2}) \times (1 - \rho)$$

$$(f(\mathbf{v}_{t}) \leq f(\mathbf{x}^{*}) + \frac{L_{s} - \nu_{s}}{2} \|\mathbf{x}_{t-1} - \mathbf{x}^{*}\|^{2} - \frac{2L_{s} - \nu_{s}}{4} \|\mathbf{v}_{t} - \mathbf{x}^{*}\|^{2}) \times \rho$$

$$\begin{split} (1 - \rho)f(\mathbf{x}_{t}) + \rho f(\mathbf{v}_{t}) &\leq f(\mathbf{x}^{*}) + \frac{L_{s} - \nu_{s}}{2} \|\mathbf{x}_{t-1} - \mathbf{x}^{*}\|^{2} - \frac{(1 - \rho)L_{s}}{2} \|\mathbf{x}_{t} - \mathbf{x}^{*}\|^{2} - \frac{\rho(L_{s} - \nu_{s})}{2} \|\mathbf{v}_{t} - \mathbf{x}^{*}\|^{2} \\ &= f(\mathbf{x}^{*}) + \frac{L_{s} - \nu_{s}}{2} \|\mathbf{x}_{t-1} - \mathbf{x}^{*}\|^{2} - \frac{L_{s} - \rho\nu_{s}}{2} \|\mathbf{x}_{t} - \mathbf{x}^{*}\|^{2} \\ &\boxed{\min_{t \in [T]} f\left(\mathbf{x}_{t}\right) \leq \left(1 + 2\rho\right) f\left(\mathbf{x}^{*}\right) + \varepsilon} \end{split}$$

$$\text{with} \quad T \geq \left\lceil \frac{L_s}{\nu_s} \log \left(\frac{(L_s - \nu_s) \|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2\varepsilon(1-\rho)} \right) \right\rceil + 1 = \mathcal{O}(\kappa_s \log(\frac{1}{\varepsilon}))$$

Proof of Convergence

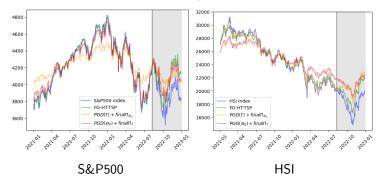
Further, if \mathbf{x}^* is a global minimizer of f over $\mathcal{B}_0(k)$, with $\rho = 0.5$ in the expressions of k and T above:

$$\min_{t\in[T]}f(\mathbf{x}_t)\leq f(\mathbf{x}^*)+\varepsilon.$$

Application: Index Tracking

$$\min_{\boldsymbol{x} \in \mathcal{B}_0(k) \cap \Gamma} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{y}\|^2$$

 $\Gamma = \{ \boldsymbol{x} \in \mathbb{R}^d : \forall i \in [c], \|\boldsymbol{x}_{G_i}\|_1 \leq D \}, \text{ where } \boldsymbol{x}_{G_i} \text{ is the restriction of } \boldsymbol{x} \text{ to group } G_i \text{ (i.e. for } j \in [d], \ \boldsymbol{x}_{G_{ij}} = \boldsymbol{x}_j \text{ if } j \in G_i \text{ and 0 otherwise)}.$



Dual Perspective on IHT

A dual perspective on IHT: Iterative Regularization with *k*-Support Norm (IRKSN)

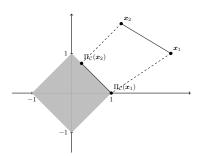
Dual Perspective on IHT

Variant of Projected Gradient Descent:

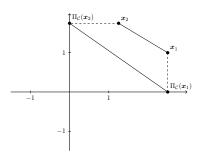
Dual Averaging (DA)[6]/(Lazy) Mirror Descent (MD)[7]/Lazy OCO[8]/Bregman Iterations [9]:

$$egin{aligned} oldsymbol{y}_{t+1} &= oldsymbol{y}_t - \eta_t
abla f(oldsymbol{x}_t) \ oldsymbol{x}_{t+1} &= oldsymbol{\mathcal{H}}_{oldsymbol{k}}(oldsymbol{y}_{t+1}) \end{aligned}$$

Dual Perspective on IHT



Projection onto the ℓ_1 unit ball.

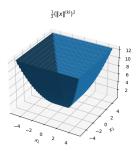


Projection onto the ℓ_0 unit pseudo-ball.

Figure: For projection onto the ℓ_1 ball, we have $\|\Pi_{\mathcal{C}}(\mathbf{x}_1) - \Pi_{\mathcal{C}}(\mathbf{x}_2)\| \leq \|\mathbf{x}_1 - \mathbf{x}_2\|$ (contractivity), but this is not true if \mathcal{C} is the ℓ_0 pseudo-ball.

Projection and Mirror Map

Contractivity of Π = Smoothness of some function ϕ $\mathcal{H}_k(\cdot) = \partial \phi(\cdot)$ with $\phi(\cdot) = \frac{1}{2}(\|\cdot\|^{(k)})^2$ (top-k norm): but ϕ not smooth.



But we can take the δ -Moreau smoothing:

$$\phi_{\delta}(\cdot) = (\underbrace{-\frac{1}{2}}_{k\text{-support norm (KSN)}})^2 + \frac{1}{2}(\|\cdot\|_2^2) -)^*$$

Note on the k-support norm (KSN)

■ KSN ball is tightest convex relaxation of ℓ_0 and ℓ_2 ball:

$$\{ \mathbf{x} : \|\mathbf{x}\|_k^{sp} \le D \} = \operatorname{conv}(\{ \mathbf{x} : \|\mathbf{x}\|_0 \le k \} \cap \{ \mathbf{x} : \|\mathbf{x}\|_2 \le D \})$$

■ The proximal operator for the squared KSN is known [10].



Figure: k-support norm ball (source: [11])

Dual Perspective on IHT

Algorithm becomes:

$$egin{aligned} oldsymbol{y}_{t+1} &= oldsymbol{y}_t - \eta_t
abla f(oldsymbol{x}_t) \ oldsymbol{x}_{t+1} &= \operatorname{prox}_{rac{1}{2\delta}(\|\cdot\|_k^{sp})^2} \left(rac{oldsymbol{y}_{t+1}}{\delta}
ight) \end{aligned}$$

Some properties:

- MD/DA Converges to x^* (not sparse in general)
- For overparam. linear models: implicit bias towards min KSN² $(+\delta \ell_2^2)$ solution
- BUT: may still not be k-sparse in general

IRKSN

We consider the **sparse recovery** problem:

$$egin{aligned} oldsymbol{y}^\delta &= oldsymbol{X} oldsymbol{w}^* + oldsymbol{\epsilon} \ \|oldsymbol{\epsilon}\| \leq \delta \end{aligned}$$

Solved by ADGD [12], solving, with early stopping:

$$\min_{\boldsymbol{w}} f(\boldsymbol{w}) \text{ s.t. } \boldsymbol{X} \boldsymbol{w} = \boldsymbol{y}^{\delta}$$

with
$$f(w) = F(w) + \frac{\alpha}{2} ||w||_2^2$$
 with $F(w) = \frac{1-\alpha}{2} (||w||_k^{sp})^2$

IRKSN

Algorithm 4: IRKSN

$$\begin{split} & \textbf{Initialization:} \ \hat{\pmb{v}}_0 = \hat{\pmb{z}}_{-1} = \hat{\pmb{z}}_0 \in \mathbb{R}^d, \gamma = \alpha \|\pmb{X}\|^{-2}, \pmb{x}_0 = 1 \\ & \textbf{for} \ t = 0, ..., T \ \textbf{do} \\ & \begin{vmatrix} \hat{\pmb{w}}_t \leftarrow \operatorname{prox}_{\alpha^{-1}F} \left(-\alpha^{-1}\pmb{X}^T\hat{\pmb{z}}_t \right) \\ \hat{\pmb{r}}_t \leftarrow \operatorname{prox}_{\alpha^{-1}F} \left(-\alpha^{-1}\pmb{X}^T\hat{\pmb{v}}_t \right) \\ \hat{\pmb{z}}_t \leftarrow \hat{\pmb{v}}_t + \gamma \left(\pmb{X}\hat{\pmb{r}}_t - \pmb{y}^\delta \right) \\ & \theta_{t+1} \leftarrow \left(1 + \sqrt{1 + 4\theta_t^2} \right)/2 \\ & \hat{\pmb{v}}_{t+1} = \hat{\pmb{z}}_t + \frac{\theta_{t-1}}{\theta_{t+1}} \left(\hat{\pmb{z}}_t - \hat{\pmb{z}}_{t-1} \right) \\ & \textbf{end} \end{aligned}$$

Notations

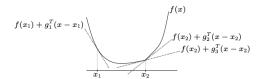
- For $S \subseteq [d]$, $\bar{S} := [d] \setminus S$
- *M*[†]: Moore-Penrose pseudo-inverse [13]
- M_S column-restriction of M to support $S \subseteq [d]$, i.e. the $n \times |S|$ matrix composed of the |S| columns of M of indices in S
- supp(w): support of w (coordinates of the non-zero components of w)
- $\mathbf{w}_S \in \mathbb{R}^k$ restriction of \mathbf{w}_S to a support S of size k, i.e. the sub-vector of size k formed by extracting only the components w_i with $i \in S$
- \blacksquare sgn(\boldsymbol{w}) vector of signs of \boldsymbol{w}

Conditions for Recovery

Метнор	Condition on X
IHT [14]	RESTRICTED ISOMETRY PROPERTY (RIP)
Lasso [15]	$\max_{\ell \in \bar{S}} \langle \pmb{X}_S^{\dagger} \pmb{x}_{\ell}, \operatorname{sgn}(\pmb{w}_S^*) \rangle < 1 \& \pmb{X}_S \text{ INJECTIVE}$
ElasticNet [16]	-
KSN PEN. [11]	-
OMP [17]	RIP
SRDI [18]	$\left\{egin{array}{l} \exists \gamma \in (0,1]: \; extbf{\textit{X}}_{S}^{ op} extbf{\textit{X}}_{S} \geq n \gamma I_{d,d} \ \exists \eta \in (0,1): \; \ extbf{\textit{X}}_{S}^{ op} extbf{\textit{X}}_{S}^{\dagger}\ _{\infty} \leq 1-\eta \end{array} ight.$
IROSR [19]	RIP
IRCR [20]	$\max_{\ell \in ar{\mathcal{S}}} \langle oldsymbol{\mathcal{X}}_{\mathcal{S}}^{\dagger} oldsymbol{x}_{\ell}, \operatorname{sgn}(oldsymbol{w}_{\mathcal{S}}^{*}) angle < 1 \ \& \ oldsymbol{\mathcal{X}}_{\mathcal{S}} \ ext{INJECTIVE}$
IRKSN (ours)	$\max_{\ell \in \bar{S}} \langle \boldsymbol{X}_{S}^{\dagger} \boldsymbol{x}_{\ell}, \boldsymbol{w}_{S}^{*} \rangle < \min_{j \in S} \langle \boldsymbol{X}_{S}^{\dagger} \boldsymbol{x}_{j}, \boldsymbol{w}_{S}^{*} \rangle $

Conditions for recovery

Finding Sufficient Conditions: Proof Technique



Subdifferential of the (half-squared) top-k norm:

$$\partial \left[\frac{1}{2} \left(\|\cdot\|^{(k)}\right)^{2}\right] = \mathsf{conv}(\mathcal{H}_{k}(\cdot))$$

Example with k = 1:

$$\partial \left[\frac{1}{2} \left(\| [-1.2, 1] \|^{(k)} \right)^2 \right] = \{ [-1.2, 0] \}$$

$$\partial \left[\frac{1}{2} \left(\|[-1.2, 1.2]\|^{(k)}\right)^2\right] = \operatorname{conv}(\{[-1.2, 0], [0, 1.2]\}) = \{[-1.2\lambda, 1.2(1-\lambda)], \lambda \in [0, 1]\}$$

Sufficient conditions for recovery: comparison with ℓ_1 norm

Assumption (Conditions for recovery with ℓ_1 norm-based algorithms)

Let \mathbf{w}^* be supported on a support $S \subset [d]$. \mathbf{w}^* is such that:

- $\mathbf{Z} X_S$ is injective
- $oxed{3} \; \mathsf{max}_{\ell \in ar{S}} \, |\langle oldsymbol{X}_S^\dagger oldsymbol{x}_\ell, \mathsf{sgn}(oldsymbol{w}_S^*)
 angle| < 1$

Assumption (Conditions for recovery with IRKSN)

- \mathbf{w}^* k-sparse, supp(\mathbf{w}^*) = $S \subset [d]$, $\mathbf{X}\mathbf{w}^* = \mathbf{y}$
- $\mathbf{w}_{S}^{*} = \operatorname{arg\,min}_{\mathbf{z} \in \mathbb{R}^{k}: \mathbf{X}_{S}\mathbf{z} = \mathbf{v}} \|\mathbf{z}\|_{2}$
- Does not need X_S to be injective !

Conditions for recovery, case where X_S is injective

If X_S is injective and $Xw^* = y$, the conditions become:

- lacksquare (A) (ℓ_1 -norm based): $\max_{\ell \in \bar{S}} |\langle \pmb{X}_S^\dagger \pmb{x}_\ell, \operatorname{sgn}(\pmb{w}_S^*)
 angle| < 1$
- $\blacksquare \text{ (B) (IRKSN): } \max\nolimits_{\ell \in \bar{\mathcal{S}}} |\langle \textbf{\textit{X}}_{\mathcal{S}}^{\dagger} \textbf{\textit{x}}_{\ell}, \frac{\textbf{\textit{w}}_{\mathcal{S}}^{*}}{\min_{j \in \mathcal{S}} |\textbf{\textit{w}}_{\mathcal{S}}^{*}|} \rangle| < 1$

It is possible to find examples of design matrix \boldsymbol{X} and vector \boldsymbol{w}^* which verify (B) but not (A): IRKSN is ensured to recover \boldsymbol{w}^* there, contrary to ℓ_1 norm-based algorithms.

Conditions for recovery

Experiments: Synthetic design matrix X

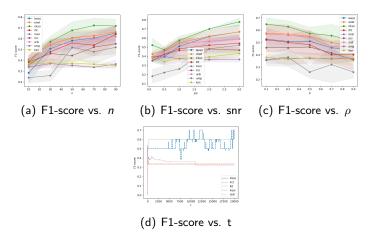


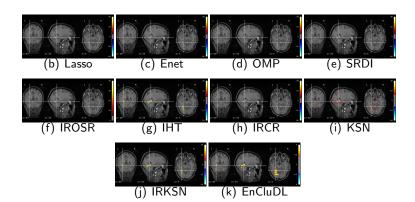
Figure: F1-score of support recovery for a correlated design matrix [20] ρ : correlation, snr: signal/noise ratio, n: num. samples.

Conditions for recovery

Experiments: fMRI decoding

	Lasso	ElasticNet	OMP	IHT	KSN	IRKSN	IRCR	IROSR	SRDI
face'/'house'	.425	.349	.938	.2441	.247	.2440	.341	.381	.314
'house'/'shoe'	.528	.500	.938	.2968	.299	.2965	.407	.502	.357

Model estimation $\| \mathbf{w} - \mathbf{w}^* \|$ (\mathbf{w}^* : obtained by EnCluDL).



QA

QΑ

References I

- [1] S. Liu, P.-Y. Chen, B. Kailkhura, G. Zhang, A. O. Hero III, and P. K. Varshney, "A primer on zeroth-order optimization in signal processing and machine learning: Principals, recent advances, and applications," *IEEE Signal Processing Magazine*, vol. 37, no. 5, pp. 43–54, 2020.
- [2] Y. Nesterov and V. Spokoiny, "Random gradient-free minimization of convex functions," Foundations of Computational Mathematics, vol. 17, no. 2, pp. 527–566, 2017.

References II

- [3] D. Golovin, J. Karro, G. Kochanski, C. Lee, X. Song, and Q. Zhang, "Gradientless descent: High-dimensional zeroth-order optimization," in *International Conference on Learning Representations*, 2019.
- [4] H. Cai, D. McKenzie, W. Yin, and Z. Zhang, "Zeroth-order regularized optimization (zoro): Approximately sparse gradients and adaptive sampling," *SIAM Journal on Optimization*, vol. 32, no. 2, pp. 687–714, 2022.
- [5] K. G. Jamieson, R. Nowak, and B. Recht, "Query complexity of derivative-free optimization," in Advances in Neural Information Processing Systems, vol. 25, 2012.

References III

- [6] Y. Nesterov, "Primal-dual subgradient methods for convex problems," *Mathematical programming*, vol. 120, no. 1, pp. 221–259, 2009.
- [7] S. Bubeck *et al.*, "Convex optimization: Algorithms and complexity," *Foundations and Trends®* in Machine Learning, vol. 8, no. 3-4, pp. 231–357, 2015.
- [8] M. Zinkevich, "Online convex programming and generalized infinitesimal gradient ascent," in *Proceedings of the 20th* international conference on machine learning (icml-03), 2003, pp. 928–936.

References IV

- [9] L. Bungert, T. Roith, D. Tenbrinck, and M. Burger, "A bregman learning framework for sparse neural networks," *Journal of Machine Learning Research*, vol. 23, no. 192, pp. 1–43, 2022.
- [10] A. McDonald, M. Pontil, and D. Stamos, "Fitting spectral decay with the k-support norm," in *Artificial Intelligence and Statistics*, PMLR, 2016, pp. 1061–1069.
- [11] A. Argyriou, R. Foygel, and N. Srebro, "Sparse prediction with the *k*-support norm," *Advances in Neural Information Processing Systems*, vol. 25, 2012.

References V

- [12] S. Matet, L. Rosasco, S. Villa, and B. L. Vu, "Don't relax: Early stopping for convex regularization," arXiv preprint arXiv:1707.05422, 2017.
- [13] G. H. Golub and C. F. Van Loan, *Matrix computations*. JHU press, 2013.
- [14] T. Blumensath and M. E. Davies, "Iterative hard thresholding for compressed sensing," Applied and computational harmonic analysis, vol. 27, no. 3, pp. 265–274, 2009.
- [15] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society: Series B* (Methodological), vol. 58, no. 1, pp. 267–288, 1996.

References VI

- [16] H. Zou and T. Hastie, "Regularization and variable selection via the elastic net," *Journal of the royal statistical society:* series B (statistical methodology), vol. 67, no. 2, pp. 301–320, 2005.
- [17] J. A. Tropp and A. C. Gilbert, "Signal recovery from random measurements via orthogonal matching pursuit," *IEEE Transactions on information theory*, vol. 53, no. 12, pp. 4655–4666, 2007.
- [18] S. Osher, F. Ruan, J. Xiong, Y. Yao, and W. Yin, "Sparse recovery via differential inclusions," *Applied and Computational Harmonic Analysis*, vol. 41, no. 2, pp. 436–469, 2016.

References VII

- [19] T. Vaskevicius, V. Kanade, and P. Rebeschini, "Implicit regularization for optimal sparse recovery," Advances in Neural Information Processing Systems, vol. 32, 2019.
- [20] C. Molinari, M. Massias, L. Rosasco, and S. Villa, "Iterative regularization for convex regularizers," in *International* conference on artificial intelligence and statistics, PMLR, 2021, pp. 1684–1692.

Some images were taken from the MTH702 course at MBZUAI.