

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

¹ Mohamed bin Zayed University of Artificial Intelligence,

² Nanjing University

April 17, 2024



- 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

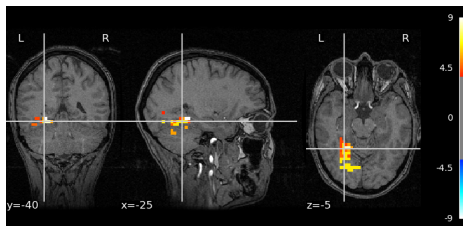
Iterative Hard Thresholding

Introduction

Sparse Optimization:

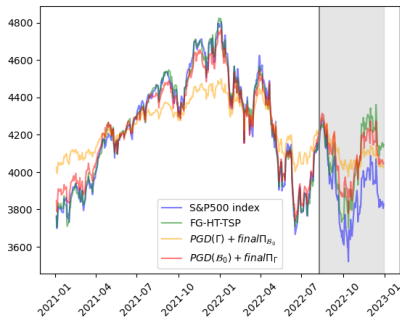
$$\min_{\mathbf{x} \in \mathbb{R}^d: \|\mathbf{x}\|_0 \leq k} f(\mathbf{x})$$

Application: fMRI



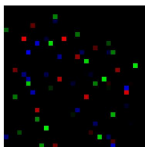
- \mathbf{x} : map of functional region of the brain (d = number of voxels)
- $f(\mathbf{x}) := \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2$ with $y_i \in \{-1, 1\}$ standing for $\{'face', 'house'\}$ and $\mathbf{A}_{i,\cdot}$ being the recorded activation map at time i .

Application: Index Tracking

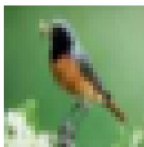


- \mathbf{x} : amount invested in each of d stocks
- $f(\mathbf{x}) := \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2$ with \mathbf{y}_i : S&P returns for day i , $\mathbf{A}_{i,j}$: return of stock j on day i

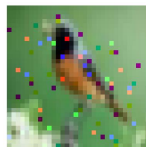
Application: Sparse Adversarial Attacks



Perturbation \mathbf{x}



'bird'



'dog'

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

The Iterative Hard Thresholding (IHT) algorithm

Algorithm 1: Iterative Hard-Thresholding (IHT)

Initialization: \mathbf{x}_0 **for** $t = 0, \dots, T$ **do**| $\mathbf{x}_{t+1} := \mathcal{H}_k(\mathbf{x}_t - \eta \nabla f(\mathbf{x}_t))$ **end****Output:** $\hat{\mathbf{x}}_T := \text{e.g. } \mathbf{x}_T \text{ or } \arg \min_{\mathbf{x} \in \{\mathbf{x}_i\}_{i=1}^T} f(\mathbf{x}_t)$

It is a **Projected Gradient Descent** algorithm:

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

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

Goal: Convergence Rate

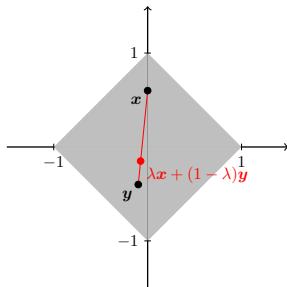
Goal: Prove Convergence Rate Why ?

- 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

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

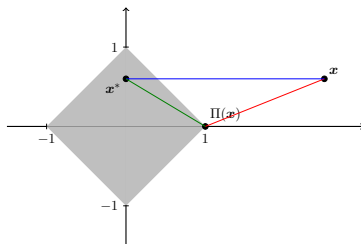
with \mathcal{C} **convex** : $\forall (\mathbf{x}, \mathbf{y}) \in (\mathcal{C})^2 : \lambda \mathbf{x} + (1 - \lambda) \mathbf{y} \in \mathcal{C}$.



Projection onto \mathcal{C}

3 Point Lemma:

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

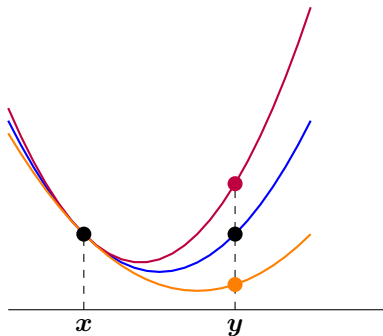


Proj. onto the ℓ_1 unit ball.

Strong Convexity and Smoothness

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

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



Proof of Convergence (Convex Case)

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

$$\begin{aligned} f(\mathbf{x}_t) &\leq f(\mathbf{x}_{t-1}) + \langle \nabla f(\mathbf{x}_{t-1}), \mathbf{x}_t - \mathbf{x}_{t-1} \rangle + \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}_{t-1}\|^2 \\ &= f(\mathbf{x}_{t-1}) + \frac{L}{2} \left\| \mathbf{x}_t - \mathbf{x}_{t-1} + \frac{1}{L} \nabla f(\mathbf{x}_{t-1}) \right\|^2 - \frac{1}{2L} \|\nabla f(\mathbf{x}_{t-1})\|^2 \\ &\leq f(\mathbf{x}_{t-1}) + \frac{L}{2} \left\| \mathbf{x}^* - \mathbf{x}_{t-1} + \frac{1}{L} \nabla f(\mathbf{x}_{t-1}) \right\|^2 - \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \frac{1}{2L} \|\nabla f(\mathbf{x}_{t-1})\|^2 \\ &= f(\mathbf{x}_{t-1}) + \langle \nabla f(\mathbf{x}_{t-1}), \mathbf{x}^* - \mathbf{x}_{t-1} \rangle + \frac{L}{2} \|\mathbf{x}_{t-1} - \mathbf{x}^*\|^2 - \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2 \\ &\leq f(\mathbf{x}^*) + \frac{L - \nu}{2} \|\mathbf{x}_{t-1} - \mathbf{x}^*\|^2 - \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2 \end{aligned}$$

Proof of Convergence (Convex Case)

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

...

$$\left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}}\right)^{T-2} [f(\mathbf{x}_2) - f(\mathbf{x}^*)] \leq \left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}}\right)^{T-2} \frac{L-\nu}{2} \|\mathbf{x}_1 - \mathbf{x}^*\|^2 - \left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}}\right)^{T-2} \frac{L}{2} \|\mathbf{x}_2 - \mathbf{x}^*\|^2$$

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

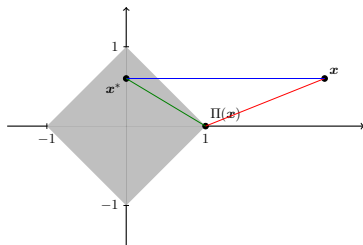
$$\sum_{t=1}^T \left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}}\right)^{T-t} [f(\mathbf{x}_t) - f(\mathbf{x}^*)] \leq \left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}}\right)^{T-1} \frac{L-\nu}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \left(\frac{\frac{L-\nu}{2}}{\frac{L}{2}}\right)^{T-t} \frac{L}{2} \|\mathbf{x}_t -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

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

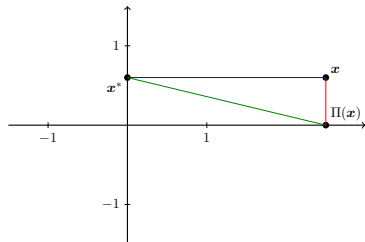
Non-Convex case: \mathcal{C} is the ℓ_0 pseudo-ball

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

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



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**. $\forall (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^d$ s.t. $\|\mathbf{x} - \mathbf{y}\|_0 \leq s$ ($s := 3k$).

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

Proof of Convergence (IHT)

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

$$\begin{aligned}
 f(\mathbf{x}_t) &\leq f(\mathbf{x}_{t-1}) + \langle \nabla f(\mathbf{x}_{t-1}), \mathbf{x}_t - \mathbf{x}_{t-1} \rangle + \frac{L_s}{2} \|\mathbf{x}_t - \mathbf{x}_{t-1}\|^2 \\
 &= f(\mathbf{x}_{t-1}) + \frac{L_s}{2} \left\| \mathbf{x}_t - \mathbf{x}_{t-1} + \frac{1}{L_s} \nabla f(\mathbf{x}_{t-1}) \right\|^2 - \frac{1}{2L_s} \|\nabla f(\mathbf{x}_{t-1})\|^2 \\
 &\leq f(\mathbf{x}_{t-1}) + \frac{L_s}{2} \left\| \mathbf{x}^* - \mathbf{x}_{t-1} + \frac{1}{L_s} \nabla f(\mathbf{x}_{t-1}) \right\|^2 - \frac{L_s}{2} (1 - \sqrt{\beta}) \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \frac{1}{2L_s} \|\nabla f(\mathbf{x}_{t-1})\|^2 \\
 &= f(\mathbf{x}_{t-1}) + \langle \nabla f(\mathbf{x}_{t-1}), \mathbf{x}^* - \mathbf{x}_{t-1} \rangle + \frac{L_s}{2} \|\mathbf{x}_{t-1} - \mathbf{x}^*\|^2 - \frac{L_s}{2} (1 - \sqrt{\beta}) \|\mathbf{x}_t - \mathbf{x}^*\|^2 \\
 &\leq f(\mathbf{x}^*) + \frac{L_s - \nu_s}{2} \|\mathbf{x}_{t-1} - \mathbf{x}^*\|^2 - \frac{L_s}{2} (1 - \sqrt{\beta}) \|\mathbf{x}_t - \mathbf{x}^*\|^2 \\
 &\leq f(\mathbf{x}^*) + \frac{L_s - \nu_s}{2} \|\mathbf{x}_{t-1} - \mathbf{x}^*\|^2 - \frac{2L_s - \nu_s}{4} \|\mathbf{x}_t - \mathbf{x}^*\|^2
 \end{aligned}$$

Proof of Convergence (IHT)

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

...

$$\left(\frac{\frac{L_S - \nu_S}{2}}{\frac{2L_S - \nu_S}{4}} \right)^{T-2} [f(\mathbf{x}_2) - f(\mathbf{x}^*)] \leq \left(\frac{\frac{L_S - \nu_S}{2}}{\frac{2L_S - \nu_S}{4}} \right)^{T-2} \frac{L_S - \nu_S}{2} \|\mathbf{x}_1 - \mathbf{x}^*\|^2 - \left(\frac{\frac{L_S - \nu_S}{2}}{\frac{2L_S - \nu_S}{4}} \right)^{T-2} \frac{2L_S - \nu_S}{4} \|\mathbf{x}_2 - \mathbf{x}^*\|^2$$

$$\left(\frac{\frac{L_S - \nu_S}{2}}{\frac{2L_S - \nu_S}{4}} \right)^{T-1} [f(\mathbf{x}_1) - f(\mathbf{x}^*)] \leq \left(\frac{\frac{L_S - \nu_S}{2}}{\frac{2L_S - \nu_S}{4}} \right)^{T-1} \frac{L_S - \nu_S}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \left(\frac{\frac{L_S - \nu_S}{2}}{\frac{2L_S - \nu_S}{4}} \right)^{T-1} \frac{2L_S - \nu_S}{4} \|\mathbf{x}_1 - \mathbf{x}^*\|^2$$

$$\sum_{t=1}^T \left(\frac{\frac{L_S - \nu_S}{2}}{\frac{2L_S - \nu_S}{4}} \right)^{T-t} [f(\mathbf{x}_t) - f(\mathbf{x}^*)] \leq \left(\frac{\frac{L_S - \nu_S}{2}}{\frac{2L_S - \nu_S}{4}} \right)^{T-1} \frac{L_S - \nu_S}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \left(\frac{\frac{L_S - \nu_S}{2}}{\frac{2L_S - \nu_S}{4}} \right)^{T-t} \frac{2L_S - \nu_S}{4} \|\mathbf{x}_t - \mathbf{x}^*\|^2$$

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

Zeroth-Order Hard-Thresholding

Zeroth-Order Hard-Thresholding

Zeroth-Order Hard-Thresholding (ZOHT)

Algorithm 2: Hard-Thresholding

Initialization: \mathbf{x}_0 **for** $t = 0, \dots, T$ **do** |
 $\mathbf{x}_{t+1} := \mathcal{H}_k(\mathbf{x}_t - \eta \nabla f(\mathbf{x}_t))$ **end****Output:** $\hat{\mathbf{x}}_T := \text{e.g. } \mathbf{x}_T \text{ or } \arg \min_{\mathbf{x} \in \{\mathbf{x}_i\}_{i=1}^T} f(\mathbf{x}_t)$

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

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

- One random direction \mathbf{u} :

$$\mathbf{g}_t = d \frac{f(\mathbf{x}_t + \mu \mathbf{u}) - f(\mathbf{x}_t)}{\mu} \mathbf{u} \quad \text{with} \quad \mathbf{u} \sim \text{Uni}(\mathbb{S}_d)$$

- q random directions $\{\mathbf{u}_i\}_{i=1}^q$:

$$\mathbf{g}_t = \frac{d}{q} \sum_{i=1}^q \frac{f(\mathbf{x}_t + \mu \mathbf{u}_i) - f(\mathbf{x}_t)}{\mu} \mathbf{u}_i \quad \text{with} \quad \{\mathbf{u}_i\}_{i=1}^q \stackrel{\text{i.i.d.}}{\sim} \text{Uni}(\mathbb{S}_d)$$

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 literature: more assumptions on f :

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

Key Insight: Error of \mathbf{g}_t on a Support F

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

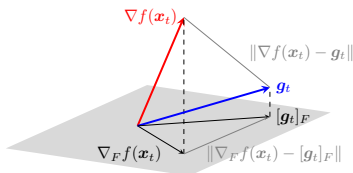
■ Bias:

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

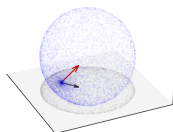
■ Variance:

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

\implies **Dimension Independent** ! (Note: we assume full smoothness here for simplicity)



$q = 1$



$q = 10^6$

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

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

ZOHT: Convergence Analysis

Choose $\eta := \frac{1}{L+C} = \frac{1}{\alpha L}$, $k \geq 16\alpha^2\kappa_s^2 k^*$ $q_t := \lceil \frac{\tau}{\omega^t} \rceil$ 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:

$$\begin{aligned} \mathbb{E}f(\mathbf{x}_t) - f(\mathbf{x}^*) &\leq \frac{1}{2\eta} \left[\left(1 - \frac{1}{\alpha'\kappa_s}\right) \mathbb{E}\|\mathbf{x}^* - \mathbf{x}_{t-1}\|^2 - (1 - \sqrt{\beta})\mathbb{E}\|\mathbf{x}_t - \mathbf{x}^*\|^2 \right. \\ &\quad \left. + 2\eta \left(\frac{G}{2} C_3 + \frac{1}{C} (2C_1\|\nabla f(\mathbf{x}^*)\|^2 + C_2\mu^2 + C_3) \right) \right] \end{aligned}$$

ZOHT: Convergence Analysis

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

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

Dimension Independent !

IHT with Additional Constraints

IHT + Additional Constraints

We now consider the following problem:

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

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

$\Gamma = \{\mathbf{x} \in \mathbb{R}^d : \forall i \in [c], \|\mathbf{x}_{G_i}\|_1 \leq D\}$, where \mathbf{x}_{G_i} is the restriction of \mathbf{x} to group G_i (i.e. for $j \in [d]$, $\mathbf{x}_{G_i j} = \mathbf{x}_j$ 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. $\|\mathbf{x}\|_0 \leq k$:
 $\text{supp}(\Pi_\Gamma(\mathbf{x})) \subseteq \text{supp}(\mathbf{x})$.

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

Initialization: \mathbf{x}_0

for $t = 0, \dots, T$ **do**

$\mathbf{v}_t := \mathcal{H}_k(\mathbf{x}_t - \eta \nabla f(\mathbf{x}_t))$

$\mathbf{x}_{t+1} := \Pi_\Gamma(\mathbf{v}_t)$

end

Output: $\hat{\mathbf{x}}_T := \text{e.g. } \mathbf{x}_T \text{ or } \arg \min_{\mathbf{x} \in \{\mathbf{x}_i\}_{i=1}^T} f(\mathbf{x}_t)$

Support Preserving Set and TSP

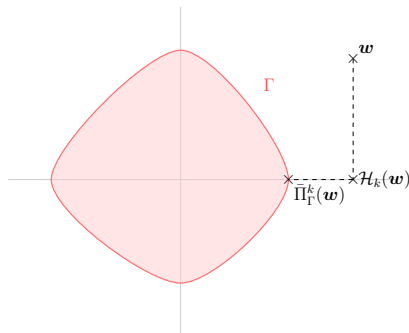


Figure: Support-preserving set and two-step projection ($d = 2, k = 1$).

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

3 Point Lemma with Extra Constraint

New Three (Four) - Point Lemma:

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

Proof of Convergence

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

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

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

$$\begin{aligned} (1 - \rho)f(\mathbf{x}_t) + \rho f(\mathbf{v}_t) &\leq f(\mathbf{x}^*) + \frac{L_s - \nu_s}{2} \|\mathbf{x}_{t-1} - \mathbf{x}^*\|^2 - \frac{(1 - \rho)L_s}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \frac{\rho(L_s - \nu_s)}{2} \|\mathbf{v}_t - \mathbf{x}^*\|^2 \\ &= f(\mathbf{x}^*) + \frac{L_s - \nu_s}{2} \|\mathbf{x}_{t-1} - \mathbf{x}^*\|^2 - \frac{L_s - \rho\nu_s}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2 \end{aligned}$$

$$\boxed{\min_{t \in [T]} f(\mathbf{x}_t) \leq (1 + 2\rho)f(\mathbf{x}^*) + \varepsilon}$$

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

Proof of Convergence

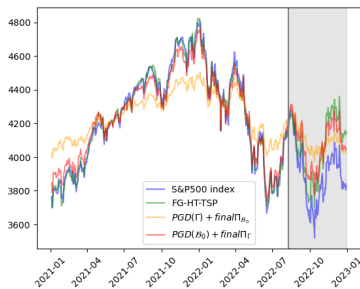
Further, if \mathbf{x}^* is a global minimizer of f over $\mathcal{B}_0(k)$, with $\rho = 0.5$ in the expressions of k and T above:

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

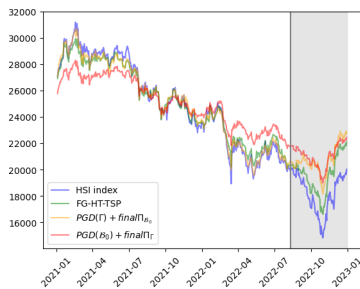
Application: Index Tracking

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

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



S&P500



HSI

Dual Perspective on IHT

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

Dual Perspective on IHT

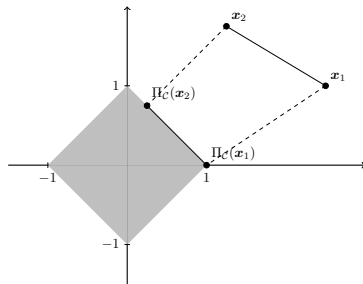
Variant of Projected Gradient Descent:

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

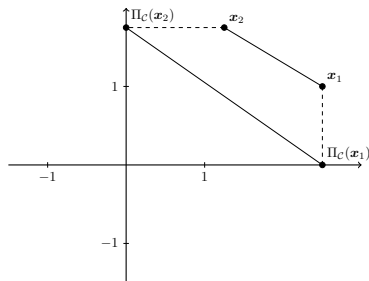
$$\mathbf{y}_{t+1} = \mathbf{y}_t - \eta_t \nabla f(\mathbf{x}_t)$$

$$\mathbf{x}_{t+1} = \mathcal{H}_k(\mathbf{y}_{t+1})$$

Dual Perspective on IHT



Projection onto the ℓ_1 unit ball.



Projection onto the ℓ_0 unit pseudo-ball.

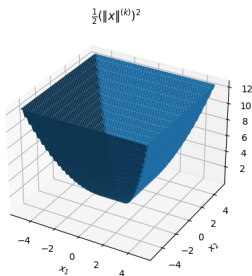
Figure: For projection onto the ℓ_1 ball, we have

$\|\Pi_C(x_1) - \Pi_C(x_2)\| \leq \|x_1 - x_2\|$ (contractivity), but this is not true if C is the ℓ_0 pseudo-ball.

Projection and Mirror Map

Contractivity of Π = Smoothness 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:

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

Note on the k -support norm (KSN)

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

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

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

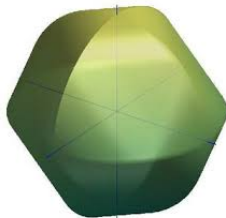


Figure: k -support norm ball (source: [11])

Dual Perspective on IHT

Algorithm becomes:

$$\mathbf{y}_{t+1} = \mathbf{y}_t - \eta_t \nabla f(\mathbf{x}_t)$$

$$\mathbf{x}_{t+1} = \text{prox}_{\frac{1}{2\delta}(\|\cdot\|_k^{sp})^2} \left(\frac{\mathbf{y}_{t+1}}{\delta} \right)$$

Some properties:

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

IRKSN

We consider the **sparse recovery** problem:

$$\mathbf{y}^\delta = \mathbf{X}\mathbf{w}^* + \boldsymbol{\epsilon}$$

$$\|\boldsymbol{\epsilon}\| \leq \delta$$

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

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

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

IRKSN

Algorithm 4: IRKSN

Initialization: $\hat{\mathbf{v}}_0 = \hat{\mathbf{z}}_{-1} = \hat{\mathbf{z}}_0 \in \mathbb{R}^d, \gamma = \alpha \|\mathbf{X}\|^{-2}, \mathbf{x}_0 = 1$

for $t = 0, \dots, T$ **do**

$$\hat{\mathbf{w}}_t \leftarrow \text{prox}_{\alpha^{-1}F}(-\alpha^{-1}\mathbf{X}^T \hat{\mathbf{z}}_t)$$

$$\hat{\mathbf{r}}_t \leftarrow \text{prox}_{\alpha^{-1}F}(-\alpha^{-1}\mathbf{X}^T \hat{\mathbf{w}}_t)$$

$$\hat{\mathbf{z}}_t \leftarrow \hat{\mathbf{v}}_t + \gamma (\mathbf{X} \hat{\mathbf{r}}_t - \mathbf{y}^\delta)$$

$$\theta_{t+1} \leftarrow (1 + \sqrt{1 + 4\theta_t^2}) / 2$$

$$\hat{\mathbf{v}}_{t+1} = \hat{\mathbf{z}}_t + \frac{\theta_t - 1}{\theta_{t+1}} (\hat{\mathbf{z}}_t - \hat{\mathbf{z}}_{t-1})$$

end

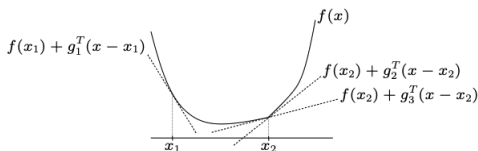
Notations

- For $S \subseteq [d]$, $\bar{S} := [d] \setminus S$
- \mathbf{M}^\dagger : Moore-Penrose pseudo-inverse [13]
- \mathbf{M}_S column-restriction of \mathbf{M} to support $S \subseteq [d]$, i.e. the $n \times |S|$ matrix composed of the $|S|$ columns of \mathbf{M} of indices in S
- $\text{supp}(\mathbf{w})$: support of \mathbf{w} (coordinates of the non-zero components of \mathbf{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 $i \in S$
- $\text{sgn}(\mathbf{w})$ vector of signs of \mathbf{w}

Conditions for Recovery

METHOD	CONDITION ON \mathbf{X}
IHT [14]	RESTRICTED ISOMETRY PROPERTY (RIP)
LASSO [15]	$\max_{\ell \in \bar{S}} \langle \mathbf{X}_S^\dagger \mathbf{x}_\ell, \text{sgn}(\mathbf{w}_S^*) \rangle < 1$ & \mathbf{X}_S INJECTIVE
ELASTICNET [16]	-
KSN PEN. [11]	-
OMP [17]	RIP
SRDI [18]	$\begin{cases} \exists \gamma \in (0, 1] : \mathbf{X}_S^\top \mathbf{X}_S \geq n\gamma I_{d,d} \\ \exists \eta \in (0, 1) : \ \mathbf{X}_{\bar{S}} \mathbf{X}_S^\dagger\ _\infty \leq 1 - \eta \end{cases}$
IROSR [19]	RIP
IRCR [20]	$\max_{\ell \in \bar{S}} \langle \mathbf{X}_S^\dagger \mathbf{x}_\ell, \text{sgn}(\mathbf{w}_S^*) \rangle < 1$ & \mathbf{X}_S INJECTIVE
IRKSN (ours)	$\max_{\ell \in \bar{S}} \langle \mathbf{X}_S^\dagger \mathbf{x}_\ell, \mathbf{w}_S^* \rangle < \min_{j \in S} \langle \mathbf{X}_S^\dagger \mathbf{x}_j, \mathbf{w}_S^* \rangle $

Finding Sufficient Conditions: Proof Technique



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

$$\partial \left[\frac{1}{2} (\|\cdot\|_k^{sp})^2 \right] = \text{conv}(\mathcal{H}_k(\cdot))$$

Example with $k = 1$:

$$\partial \left[\frac{1}{2} (\|[-1.2, 1]\|_1^{sp})^2 \right] = \{[-1.2, 0]\}$$

$$\partial \left[\frac{1}{2} (\|[-1.2, 1.2]\|_1^{sp})^2 \right] = \text{conv}(\{[-1.2, 0], [0, 1.2]\}) = \{[-1.2\lambda, 1.2(1-\lambda)], \lambda \in [0, 1]\}$$

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:

- 1 $\mathbf{X}\mathbf{w}^* = \mathbf{y}$
- 2 \mathbf{X}_S is injective
- 3 $\max_{\ell \in \bar{S}} |\langle \mathbf{X}_S^\dagger \mathbf{x}_\ell, \text{sgn}(\mathbf{w}_S^*) \rangle| < 1$

Assumption (Conditions for recovery with IRKSN)

- \mathbf{w}^* k -sparse, $\text{supp}(\mathbf{w}^*) = S \subset [d]$, $\mathbf{X}\mathbf{w}^* = \mathbf{y}$
- $\mathbf{w}_S^* = \arg \min_{\mathbf{z} \in \mathbb{R}^k: \mathbf{X}_S \mathbf{z} = \mathbf{y}} \|\mathbf{z}\|_2$
- $\max_{\ell \in \bar{S}} |\langle \mathbf{X}_S^\dagger \mathbf{x}_\ell, \mathbf{w}_S^* \rangle| < \min_{j \in S} |\langle \mathbf{X}_S^\dagger \mathbf{x}_j, \mathbf{w}_S^* \rangle|$
- Does not need \mathbf{X}_S to be injective !

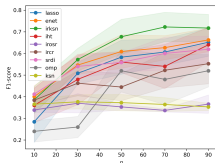
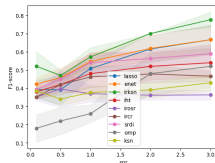
Conditions for recovery, case where \mathbf{X}_S is injective

If \mathbf{X}_S is injective and $\mathbf{X}\mathbf{w}^* = \mathbf{y}$, the conditions become:

- (A) (ℓ_1 -norm based): $\max_{\ell \in \bar{S}} |\langle \mathbf{X}_S^\dagger \mathbf{x}_\ell, \text{sgn}(\mathbf{w}_S^*) \rangle| < 1$
- (B) (IRKSN): $\max_{\ell \in \bar{S}} |\langle \mathbf{X}_S^\dagger \mathbf{x}_\ell, \frac{\mathbf{w}_S^*}{\min_{j \in S} |\mathbf{w}_S^*|} \rangle| < 1$

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

Experiments: Synthetic design matrix X

(a) F1-score vs. n 

(b) F1-score vs. snr

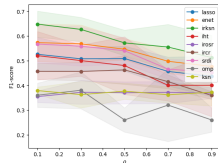
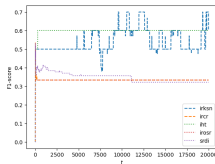
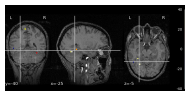
(c) F1-score vs. ρ (d) F1-score vs. t

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

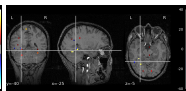
Experiments: fMRI decoding

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

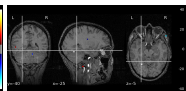
Model estimation $\|\mathbf{w} - \mathbf{w}^*\|$ (\mathbf{w}^* : obtained by EnCluDL).



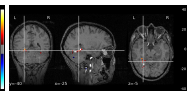
(b) Lasso



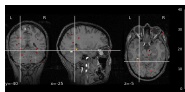
(c) Enet



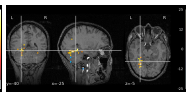
(d) OMP



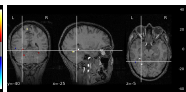
(e) SRDI



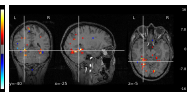
(f) IROS



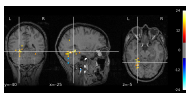
(g) IHT



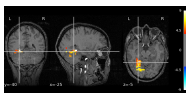
(h) IRCR



(i) KSN



(j) IRKSN



(k) EnCluDL

QA

QA

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.