

Candidacy Exam

William de Vazelhes ¹ ,	PhD Candidate
Dr. Bin Gu ¹ ,	Supervisor
Dr. Xiaotong Yuan ² ,	External Examiner
Dr. Chih-Jen Lin ¹ ,	Internal Member
Dr. Karthik Nandakumar ¹ ,	Internal Member
Dr. Zhiqiang Xu ¹ ,	Internal Member

¹ Mohamed bin Zayed University of Artificial Intelligence,

² Nanjing University of Information Science & Technology

May 11, 2023



1 Introduction

2 Research Progress

- ZOHT
- IRKSN

3 Future research

- Variance Reduction
- Additional Constraints
- Structural Sparsity
- Reinforcement learning
- Others

Introduction

Sparse Optimization:

$$\min_{\|\mathbf{x}\|_0 \leq k} f(\mathbf{x})$$

Applications:

- Sparse regression/classification (e.g. gene array data)
- Sparse recovery

TL;DR

Main contributions:

- ZOHT: condition on number of random directions q
- IRKSN: new conditions for linear sparse recovery

Zeroth-Order

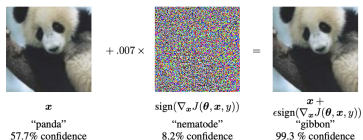
$$\min_{\mathbf{x} \in \mathbb{R}} f(\mathbf{x})$$

Gradient descent:

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

What if we don't know $\nabla f(\cdot)$, but only $f(\cdot)$?

- Black-Box Adversarial attacks [1]



- Reinforcement learning [2]



Approximate $\nabla f(\mathbf{x})$: two points approximation [3] [4]:

- One random direction \mathbf{u} :

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

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

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

Curse of dimensionality: An impossibility result [7]

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 $\text{rank}(\mathbf{A}) \ll d$ [5]
- sparse/compressible gradients [6]

Vanilla ZO

Most ZO algorithms [4], [8], [9]:

Algorithm 1: Vanilla ZO

Initialization: $\eta, T, \mathbf{x}^{(0)}$

Output: \mathbf{x}_T .

for $t = 1, \dots, T$ **do**

 Sample $\mathbf{u} \sim \text{Uni}(\mathbb{S}_d)$

$\hat{\nabla} f(\mathbf{x}_{t-1}) \leftarrow \frac{d}{\mu} (f(\mathbf{x} + \mu \mathbf{u}) - f(\mathbf{x})) \mathbf{u};$

$\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} - \eta \hat{\nabla} f(\mathbf{x}_{t-1});$

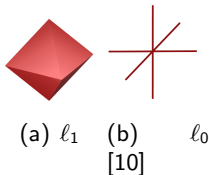
end

Note: just **1 random direction** \mathbf{u} is sufficient (with proper η)

Zeroth-Order Hard-Thresholding: Our approach

Consider the **non-convex** ℓ_0 “ball”

$$\min_{\mathbf{x}} \text{ s.t. } \|\mathbf{x}\|_0 \leq k \quad f(\mathbf{x})$$



Why not ℓ_1 ? ℓ_1 is convex (impossibility result)

ZOHT: Zeroth-Order Hard-Thresholding

Algorithm 2: SZOHT (simplified)

Initialization: $\eta, T, q, k = \mathcal{O}(\kappa^4 k^*), \mathbf{x}^{(0)}$ **Output:** \mathbf{x}_T .**for** $t = 1, \dots, T$ **do** **for** $i = 1, \dots, q$ **do** Sample $\mathbf{u}_i \sim \text{Uni}(\mathbb{S}_d)$ $\hat{\nabla} f(\mathbf{x}_{t-1}; \mathbf{u}_i) \leftarrow \frac{d}{\mu} (f(\mathbf{x} + \mu \mathbf{u}_i) - f(\mathbf{x})) \mathbf{u}_i;$ **end** $\hat{\nabla} f(\mathbf{x}_{t-1}) \leftarrow \frac{1}{q} \sum_{i=1}^q \hat{\nabla} f(\mathbf{x}_{t-1}; \mathbf{u}_i)$ $\mathbf{x}_t \leftarrow \Phi_k(\mathbf{x}_{t-1} - \eta \hat{\nabla} f(\mathbf{x}_{t-1}));$ hard-thresholding**end**

Dimension independence

- **vanilla ZO:** $O(d)$ QC comes from $\|\hat{\nabla}f(\mathbf{x}) - \nabla f(\mathbf{x})\|^2 \leq O(d)$
- **ZOHT:** $\|\hat{\nabla}_F f(\mathbf{x}) - \nabla_F f(\mathbf{x})\|^2 \leq O(k)$ (props. of projections [11])

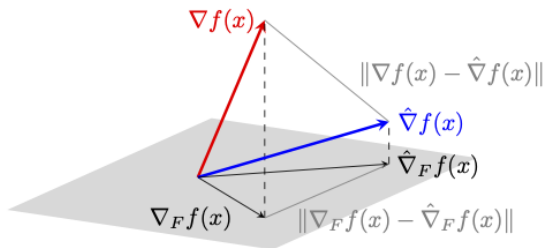


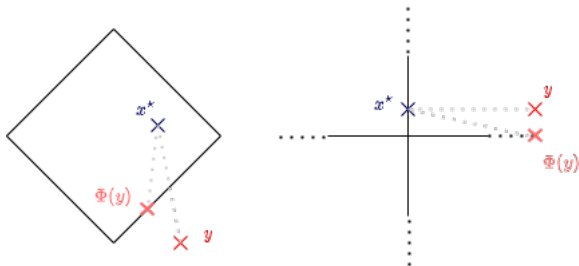
Figure: Gradient estimate and its projections

Tuning q : gradient error vs. expansivity

Main difference with vanilla ZO: projection on the ℓ_0 ball ($\mathcal{B}_{\ell_0, k}$) is not non-expansive:

$$\forall \mathbf{y} \in \mathbb{R}^d, \mathbf{x}^* \in \mathcal{B}_{\ell_0, k} : \|\Phi_k(\mathbf{y}) - \mathbf{x}^*\| \leq \gamma \|\mathbf{y} - \mathbf{x}^*\|$$

with $\gamma > 1$



Tuning q : gradient error vs. expansivity

Convergence rate:

$$\mathbb{E} \|\mathbf{x}_t - \mathbf{x}^*\| \leq (\rho\gamma)^t \|\mathbf{x}_0 - \mathbf{x}^*\| + (\cdot)\sigma + (\cdot)\mu$$

With $\eta = \frac{\nu}{(4\epsilon_{err}+1)L^2}$, $\rho = 1 - \frac{\nu^2}{(4\epsilon_{err}+1)L^2}$, $\epsilon_{err} \leq O(\frac{k}{q})$, $k \geq \frac{\rho^2 k^*}{(1-\rho^2)^2}$

Sufficient: (by inspection) valid q s.t. $\rho\gamma < 1$: $q \geq 2(3k + 2)$

Necessary: minimum q so that \exists valid k s.t. $\rho\gamma < 1$:

$$q \geq 4\kappa^2 \sqrt{\frac{k^* d}{s_2}} > 1$$

Tuning q : gradient error vs. expansivity

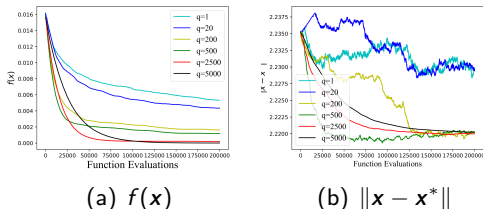


Figure: Sensitivity analysis

Toy XP: $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{a} \odot (\mathbf{x} - \mathbf{b})\|^2$, (\mathbf{a} and \mathbf{b} chosen to have $\|\nabla f(\mathbf{x}^*)\|$ small enough). $\eta = \frac{1}{(4\varepsilon_F + 1)}$. For $q = 1$ and 20, $\|\mathbf{x}^{(t)} - \mathbf{x}^*\|$ does not converge.

Improvements: sampling along a random support

$\mathbf{u}_i \sim \text{Uni}(\mathbb{S}_S)$ with $S \sim \text{Uni}(\binom{[d]}{s_2})$

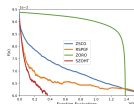
- memory efficiency (if distributed learning)
- allows to work with “restricted smoothness” only
- can improve the condition number ν/L

Experiments

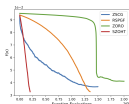
- Asset management [12], (a), (b), (c) :

$$\min_{\mathbf{x} \in \mathbb{R}^d} \frac{\mathbf{x}^\top \mathbf{C} \mathbf{x}}{2 \left(\sum_{i=1}^d x_i \right)^2} + \lambda \left(\min \left\{ \frac{\sum_{i=1}^d m_i x_i}{\sum_{i=1}^d x_i} - r, 0 \right\} \right)^2 \quad \text{s.t.} \quad \|\mathbf{x}\|_0 \leq k$$

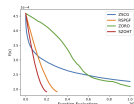
- Few pixels adv. attacks [13], (d), (e), (f) : $\min_{\delta} f(\mathbf{x} + \delta)$ such that $\|\delta\|_0 \leq k$
- Comparison with ZSCG [14], RSPGF [15], ZORO [6] ($f(\mathbf{x})$ vs QC)



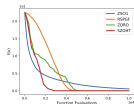
(a) port3



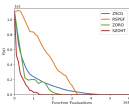
(b) port4



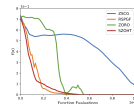
(c) port5



(d) MNIST



(e) CIFAR



(f) Imagenet

Iterative Regularization with k -support norm: A Dual Perspective on Hard-Thresholding

Original Goal: Online ℓ_0 optimization.

Attempt 1: Modify Online Convex Optimization ([16])

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

Problem: Could not get sublinear regret

Dual Perspective on IHT (contd.)

Attempt 2: Dual Averaging[17]/(Lazy) Mirror Descent[18]/Lazy OCO[16]/Bregman Iterations [19]:

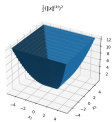
$$\mathbf{y}_{t+1} = \mathbf{y}_t - \eta_t \nabla f(\mathbf{x}_t)$$

$$\mathbf{x}_{t+1} = \Phi_k(\mathbf{y}_{t+1})$$

Problem: $\Phi_k(\cdot) = \partial\phi(\cdot)$ with $\phi(\cdot) = \frac{1}{2}(\|\cdot\|^{(k)})^2$ (top- k norm). ϕ not smooth...(proof cannot work)

But we can take the δ -Moreau smoothing, to get:

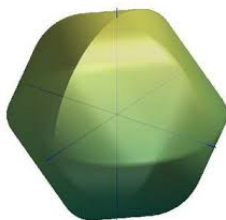
$$\phi_\delta(\cdot) = \left(\frac{1}{2} \left(\underbrace{\|\cdot\|_k^{sp}}_{k\text{-support norm (KSN)}} \right)^2 + \frac{1}{2} (\|\cdot\|_2^2) \right)^*$$



Note on the k -support norm

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

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



Dual Perspective on IHT (contd.)

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: (not just online)

- Not really new now (MD/DA/BI, with just new use of KSN)
- \mathbf{x}_t empirically "almost" sparse (prelim. XPs)
- BUT: MD, so convergence to \mathbf{x}^* (maybe not sparse)
- **For overparam. linear models: implicit bias** towards min $\text{KSN}^2 (+\delta\ell_2^2)$ solution
- BUT: may still not be k -sparse

IRKSN

We simplify the problem, to make proofs easier → **sparse recovery**:

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

,

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

Solved by a tweaked version of ADGD [20], solving, **with early stopping**

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

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

Algorithm

Algorithm 3: IRKSN

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

Output: $\hat{\mathbf{w}}_T$

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

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

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

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

$$\mathbf{x}_{t+1} \leftarrow \left(1 + \sqrt{1 + 4\mathbf{x}_t^2} \right) / 2$$

$$\hat{\mathbf{v}}_{t+1} = \hat{\mathbf{z}}_t + \frac{\mathbf{x}_t - 1}{\mathbf{x}_{t+1}} \left(\hat{\mathbf{z}}_t - \hat{\mathbf{z}}_{t-1} \right)$$

end

Recovery

Assumption

\mathbf{w}^* k -sparse, $\text{supp}(\mathbf{w}^*) = S \subset [d]$, $\mathbf{X}\mathbf{w}^* = \mathbf{y}$,

$\mathbf{w}_S^* = \arg \min_{\mathbf{z} \in \mathbb{R}^k: \mathbf{X}_S \mathbf{z} = \mathbf{y}} \|\mathbf{z}\|_2$

$$\max_{\ell \in \bar{S}} |\langle \mathbf{X}_S^\dagger \mathbf{x}_\ell, \mathbf{w}_S^* \rangle| < \min_{j \in S} |\langle \mathbf{X}_S^\dagger \mathbf{x}_j, \mathbf{w}_S^* \rangle|$$

Theorem (Early Stopping Bound)

$$\|\hat{\mathbf{w}}_t - \mathbf{w}^*\|_2 \leq at\delta + bt^{-1}$$

$$\text{with } a = 4\|\mathbf{X}\|^{-1} \quad \text{and} \quad b = \frac{2\|\mathbf{X}\| \|(\mathbf{X}_S^\top)^\dagger \mathbf{w}_S^*\|}{\alpha}$$

Comparison with ℓ_1 -based recovery

Assumption (Recovery with ℓ_1 norm.)

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

If \mathbf{X}_S is injective

IRKSN

$$\max_{\ell \in \bar{S}} \left| \langle \mathbf{X}_S^\dagger \mathbf{x}_\ell, \frac{\mathbf{w}_S^*}{\min_{j \in S} |\mathbf{w}_S^*|} \rangle \right| < 1$$

l1 iter. reg.

$$\max_{\ell \in \bar{S}} \left| \langle \mathbf{X}_S^\dagger \mathbf{x}_\ell, \text{sgn}(\mathbf{w}_S^*) \rangle \right| < 1$$

Variance Reduction

VR-SZHT introduced by Xinzhe Yuan:

Algorithm 4: Stochastic variance reduced zeroth-order Hard-Thresholding (VR-SZHT)

Initialization: η, T, \mathbf{x}^0 , SVRG update frequency m, q, k .

Output: \mathbf{x}^T .

for $r = 1, \dots, T$ **do**

$\mathbf{x}^{(0)} = \mathbf{x}^{r-1}$; $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \hat{\nabla} f_i(\mathbf{x}^{(0)})$; **for** $t = 0, 1, \dots, m - 1$

do

Randomly sample $i_t \in \{1, 2, \dots, n\}$; Compute ZO

estimate $\hat{\nabla} f_{i_t}(\mathbf{x}^{(t)})$, $\hat{\nabla} f_{i_t}(\mathbf{x}^{(0)})$;

$\bar{\mathbf{x}}^{(t+1)} = \mathbf{x}^{(t)} - \eta(\hat{\nabla} f_{i_t}(\mathbf{x}^{(t)}) - \hat{\nabla} f_{i_t}(\mathbf{x}^{(0)}) + \hat{\mu})$;

$\mathbf{x}^{(t+1)} = \phi_k(\bar{\mathbf{x}}^{(t+1)})$;

end

$\mathbf{x}^r = \mathbf{x}^{(m)}$;

end

VR-SZHT

Removes need for minimum q : VR can compensate the variance of gradient estimate

Additional Constraints

$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{F}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}), \quad s.t. \|\mathbf{x}\|_0 \leq k \text{ and } \mathbf{x} \in \mathcal{S}. \quad (1)$$

Useful e.g. in **adversarial attacks**.

Structural sparsity

We may want to enforce constraints of the form:

$$\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}_{\mathcal{G}_1}\|_0 < D_1 \wedge \|\mathbf{x}_{\mathcal{G}_2}\|_0 < D_2\}$$

$\mathbf{x}_{\mathcal{G}_1}$ and $\mathbf{x}_{\mathcal{G}_2}$: a partition of \mathbf{x} into coordinates from a group \mathcal{G}_1 and a group \mathcal{G}_2 .

Reinforcement learning

Salimans et al. [21]: Evolution Strategies (\approx ZO) very efficient for RL (esp. distributed)

BUT: dependence on d ($d \gg 1$ for DNNs).

\implies Could using ZOHT reduce the dependence in d ?

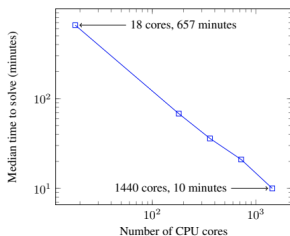


Figure 1: Time to reach a score of 6000 on 3D Humanoid with different number of CPU cores. Experiments are repeated 7 times and median time is reported.

Others

Other potential ideas:

- **Low-rank Matrices**
- **Sparse graphs**
- **Acceleration of ZOHT [22]**
- **Relaxed Assumptions (non-RSC)**
- **Lower bound for ZOO with sparse optima**
- **Non-convex regularization: $h(\mathbf{x}) = \lambda(\frac{1}{2}(\|\mathbf{x}\|_k^{sp})^2 - \frac{1}{2}\|\mathbf{x}\|_2^2)$**

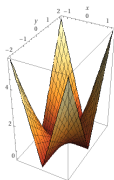


Figure: Nonconvex penalty based on the k -support norm.

QA

References I

- [1] I. J. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” *arXiv preprint arXiv:1412.6572*, 2014.
- [2] V. Mnih, K. Kavukcuoglu, D. Silver, *et al.*, “Playing atari with deep reinforcement learning,” *arXiv preprint arXiv:1312.5602*, 2013.
- [3] 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.

References II

- [4] Y. Nesterov and V. Spokoiny, “Random gradient-free minimization of convex functions,” *Foundations of Computational Mathematics*, vol. 17, no. 2, pp. 527–566, 2017.
- [5] 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.
- [6] 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.

References III

- [7] K. G. Jamieson, R. Nowak, and B. Recht, “Query complexity of derivative-free optimization,” in *Advances in Neural Information Processing Systems*, vol. 25, 2012.
- [8] J. C. Duchi, M. I. Jordan, M. J. Wainwright, and A. Wibisono, “Optimal rates for zero-order convex optimization: The power of two function evaluations,” *IEEE Transactions on Information Theory*, vol. 61, no. 5, pp. 2788–2806, 2015.
- [9] K. Balasubramanian and S. Ghadimi, “Zeroth-order nonconvex stochastic optimization: Handling constraints, high dimensionality, and saddle points,” *Foundations of Computational Mathematics*, pp. 1–42, 2021.

References IV

- [10] J. Francis, B. Madathil, S. N. George, and S. George, “A robust tensor-based submodule clustering for imaging data using $l_1 l_2$ regularization and simultaneous noise recovery via sparse and low rank decomposition approach,” *Journal of Imaging*, vol. 7, no. 12, p. 279, 2021.
- [11] S. Sykora, “Surface integrals over n-dimensional spheres,” *Stan’s Library*, no. Volume I, May 2005. DOI: [10.3247/sl1math05.002](https://doi.org/10.3247/sl1math05.002).
- [12] T.-J. Chang, N. Meade, J. E. Beasley, and Y. M. Sharaiha, “Heuristics for cardinality constrained portfolio optimisation,” *Computers & Operations Research*, vol. 27, no. 13, pp. 1271–1302, 2000.

References V

- [13] P.-Y. Chen, H. Zhang, Y. Sharma, J. Yi, and C.-J. Hsieh, “Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models,” in *Proceedings of the 10th ACM workshop on artificial intelligence and security*, 2017, pp. 15–26.
- [14] K. Balasubramanian and S. Ghadimi, “Zeroth-order (non)-convex stochastic optimization via conditional gradient and gradient updates,” in *Advances in Neural Information Processing Systems*, vol. 31, 2018.

References VI

- [15] S. Ghadimi, G. Lan, and H. Zhang, “Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization,” *Mathematical Programming*, vol. 155, no. 1, pp. 267–305, 2016.
- [16] 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.
- [17] Y. Nesterov, “Primal-dual subgradient methods for convex problems,” *Mathematical programming*, vol. 120, no. 1, pp. 221–259, 2009.

References VII

- [18] S. Bubeck *et al.*, “Convex optimization: Algorithms and complexity,” *Foundations and Trends® in Machine Learning*, vol. 8, no. 3-4, pp. 231–357, 2015.
- [19] 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.
- [20] S. Matet, L. Rosasco, S. Villa, and B. L. Vu, “Don’t relax: Early stopping for convex regularization,” *arXiv preprint arXiv:1707.05422*, 2017.

References VIII

- [21] T. Salimans, J. Ho, X. Chen, S. Sidor, and I. Sutskever, “Evolution strategies as a scalable alternative to reinforcement learning,” *arXiv preprint arXiv:1703.03864*, 2017.
- [22] K. Axiotis and M. Sviridenko, “Iterative hard thresholding with adaptive regularization: Sparser solutions without sacrificing runtime,” in *International Conference on Machine Learning*, PMLR, 2022, pp. 1175–1197.