

# Two-timescale Extragradient for Finding Local Minimax Points

Jiseok Chae<sup>1</sup>   Kyuwon Kim<sup>1</sup>   Donghwan Kim<sup>1</sup>

<sup>1</sup>Department of Mathematical Sciences, KAIST

2023 SAARC CoLLabor Workshop

# Contents

## 1 Introduction

## 2 The Two-timescale Extragradient Method

# Problem Setting

We consider minimax problems:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \max_{\mathbf{y} \in \mathbb{R}^m} f(\mathbf{x}, \mathbf{y}).$$

Many ML applications of minimax problems:

- Generative Adversarial Networks
- Adversarial Learning
- Fair Classification
- Sharpness-aware Minimization

# Optimality

Minimax problems have a hierarchical structure

$$\min_{\mathbf{x} \in \mathbb{R}^n} \left( \Phi(\mathbf{x}) := \max_{\mathbf{y} \in \mathbb{R}^m} f(\mathbf{x}, \mathbf{y}) \right).$$

The definition of (local) optimal point  $(\mathbf{x}^*, \mathbf{y}^*)$  reflects this hierarchy.

## Definition 1 (Informal; Jin et al. (2020))

There exists  $h : \mathbb{R} \rightarrow \mathbb{R}$  such that for any small  $\delta$ ,

$$f(\mathbf{x}^*, \mathbf{y}) \leq f(\mathbf{x}^*, \mathbf{y}^*) \leq \max_{\|\tilde{\mathbf{y}} - \mathbf{y}^*\| \leq h(\delta)} f(\mathbf{x}, \tilde{\mathbf{y}})$$

holds for any  $\|\mathbf{x} - \mathbf{x}^*\| \leq \delta$  and  $\|\mathbf{y} - \mathbf{y}^*\| \leq \delta$ .

# Contents

1 Introduction

**2 The Two-timescale Extragradient Method**

## Two-timescale Methods

Gradient descent ascent (GDA) is a widely used simple modification of gradient descent

$$\begin{aligned}\mathbf{x}_{k+1} &= \mathbf{x}_k - \eta \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k) \\ \mathbf{y}_{k+1} &= \mathbf{y}_k + \eta \nabla_{\mathbf{y}} f(\mathbf{x}_k, \mathbf{y}_k)\end{aligned}$$

but the hierarchy of minimax problems is not well incorporated.

To put more emphasis on the maximization, one can introduce a timescale parameter  $\tau \geq 1$  and consider the *two-timescale* GDA

$$\begin{aligned}\mathbf{x}_{k+1} &= \mathbf{x}_k - \eta/\tau \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k) \\ \mathbf{y}_{k+1} &= \mathbf{y}_k + \eta \nabla_{\mathbf{y}} f(\mathbf{x}_k, \mathbf{y}_k)\end{aligned}$$

## Two-timescale GDA is Not Sufficient

Does the two-timescale GDA work well?

Yes and no...

### Theorem 2 (Fiez and Ratliff (2021) and Jin et al. (2020))

If  $\nabla_{yy}^2 f$  is negative definite, then for a sufficiently large  $\tau$ , the limit points of the two-timescale GDA are the local optimum.

Local optimal points with  $\nabla_{yy}^2 f \prec \mathbf{0}$  are called *strict* local minimax points. For *non-strict* local optimal points, no convergence guarantees were known.

# The Extragradient Method

Maybe using GDA is too naïve?

GDA already lacks convergence guarantees even for convex-concave problems. (Mescheder et al., 2018)

The extragradient(EG) method

$$\mathbf{x}_{k+1/2} = \mathbf{x}_k - \eta \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k)$$

$$\mathbf{y}_{k+1/2} = \mathbf{y}_k + \eta \nabla_{\mathbf{y}} f(\mathbf{x}_k, \mathbf{y}_k)$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \nabla_{\mathbf{x}} f(\mathbf{x}_{k+1/2}, \mathbf{y}_{k+1/2})$$

$$\mathbf{y}_{k+1} = \mathbf{y}_k + \eta \nabla_{\mathbf{y}} f(\mathbf{x}_{k+1/2}, \mathbf{y}_{k+1/2})$$

is well known for its better convergence.



# The Extragradient Method

Maybe using GDA is too naïve?

GDA already lacks convergence guarantees even for convex-concave problems. (Mescheder et al., 2018)

The *two-timescale* extragradient(EG) method

$$\mathbf{x}_{k+1/2} = \mathbf{x}_k - \eta/\tau \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k)$$

$$\mathbf{y}_{k+1/2} = \mathbf{y}_k + \eta \nabla_{\mathbf{y}} f(\mathbf{x}_k, \mathbf{y}_k)$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta/\tau \nabla_{\mathbf{x}} f(\mathbf{x}_{k+1/2}, \mathbf{y}_{k+1/2})$$

$$\mathbf{y}_{k+1} = \mathbf{y}_k + \eta \nabla_{\mathbf{y}} f(\mathbf{x}_{k+1/2}, \mathbf{y}_{k+1/2})$$

should do a better job in finding local optimal points!

# Dynamical Systems Approach

Using some notations...

$$\Lambda_\tau := \begin{bmatrix} 1/\tau \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}, \quad \mathbf{F}(\cdot) := \begin{bmatrix} \nabla_{\mathbf{x}} f(\cdot) \\ -\nabla_{\mathbf{y}} f(\cdot) \end{bmatrix}, \quad \mathbf{z} := \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}$$

... two-timescale EG can be viewed as a dynamical system :

$$\mathbf{z}_{k+1} = \mathbf{w}_\tau(\mathbf{z}_k) := \mathbf{z}_k - \eta \Lambda_\tau \mathbf{F}(\mathbf{z}_k - \eta \Lambda_\tau \mathbf{F}(\mathbf{z}_k)).$$

## Fact 3 (Galor (2007, Theorem 4.8))

For an equilibrium point  $\mathbf{z}^*$  of a dynamical system  $\mathbf{z}_{k+1} = \mathbf{w}(\mathbf{z}_k)$ , the iterates locally converge to  $\mathbf{z}^*$  if the Jacobian matrix  $D\mathbf{w}(\mathbf{z}^*)$  has spectral radius smaller than 1, *i.e.*,  $\rho(D\mathbf{w}(\mathbf{z}^*)) < 1$ .

# Asymptotic Behavior of Eigenvalues

## Fact 4 (Chae, Kim, and Kim (2024))

The condition for  $\rho(D\mathbf{w}_\tau(\mathbf{z}^*)) < 1$  can be characterized by the spectrum of  $\mathbf{\Lambda}_\tau D\mathbf{F}(\mathbf{z}^*)$ .

## Theorem 5 (Chae, Kim, and Kim (2024))

For  $\epsilon := 1/\tau$ , the complex eigenvalues  $\lambda_j(\epsilon)$  of  $\mathbf{\Lambda}_\tau D\mathbf{F}(\mathbf{x}^*, \mathbf{y}^*)$  have one of the following three asymptotics as  $\tau \rightarrow +\infty$ :

- (i)  $|\lambda_j(\epsilon) \pm i\sigma_j\sqrt{\epsilon}| = o(\sqrt{\epsilon})$ ,
- (ii)  $|\lambda_j(\epsilon) - \epsilon\mu_j| = o(\epsilon)$ ,
- (iii)  $|\lambda_j(\epsilon) - \nu_j| = o(1)$ .

Here,  $\sigma_j$ ,  $\nu_j$ , and  $\mu_j$  are some values determined by  $D\mathbf{F}(\mathbf{x}^*, \mathbf{y}^*)$ .

# Limit Points of the Two-timescale EG

## Theorem 6

For an equilibrium point  $z^*$  of the two-timescale EG

$$z_{k+1} = z_k - \eta \Lambda_\tau \mathbf{F}(z_k - \eta \Lambda_\tau \mathbf{F}(z_k)),$$

TFAE:

- $\mathbf{S}_{\text{res}}(D\mathbf{F}(z^*)) \succeq \mathbf{0}$ ,  $\nabla_{\mathbf{y}\mathbf{y}}^2 f \preceq \mathbf{0}$ , and  $s_0 < \frac{1}{2L}$ .
- There exists a number  $0 < \eta^* < \frac{1}{L}$  such that the two-timescale EG locally converges to  $z^*$  for all sufficiently large  $\tau$  and  $\eta^* < \eta < \frac{1}{L}$ .

Here,  $s_0$  is a scalar depending on  $D\mathbf{F}(z^*)$ , and  $\mathbf{S}_{\text{res}}$  is a matrix-valued function. (Details omitted.)

## Relating to Local Optimality

What points satisfy  $\mathbf{S}_{\text{res}}(DF(\mathbf{z}^*)) \succeq \mathbf{0}$  and  $\nabla_{\mathbf{y}\mathbf{y}}^2 f \preceq \mathbf{0}$ ?

### Proposition 7 (Chae, Kim, and Kim (2024))

If  $f \in \mathcal{C}^2$ , then any local optimal point satisfies  $\nabla_{\mathbf{y}\mathbf{y}}^2 f \preceq \mathbf{0}$ .

If we further assume that  $\limsup_{\delta \rightarrow 0} h(\delta)/\delta < \infty$ , then any local optimal point satisfies  $\mathbf{S}_{\text{res}}(DF(\mathbf{z}^*)) \succeq \mathbf{0}$ .

Therefore, under those mild conditions...

*Two-timescale EG can find non-strict local optimal points.*

This is the first ever known convergence result of a first-order method to non-strict local optimal points.

Thank you for your attention.

# References I

- Chae, Jiseok, Kyuwon Kim, and Donghwan Kim (2024). “Two-timescale Extragradient for Finding Local Minimax Points”. In: *The Twelfth International Conference on Learning Representations*.
- Fiez, Tanner and Lillian J. Ratliff (2021). “Local convergence analysis of gradient descent ascent with finite timescale separation”. In: *International Conference on Learning Representations*.
- Galor, Oded (2007). *Discrete dynamical systems*. Springer-Verlag.
- Jin, Chi, Praneeth Netrapalli, and Michael Jordan (2020). “What is local optimality in nonconvex-nonconcave minimax optimization?” In: *International Conference on Machine Learning*. PMLR, pp. 4880–4889.
- Mescheder, Lars, Andreas Geiger, and Sebastian Nowozin (2018). “Which training methods for GANs do actually converge?” In: *International Conference on Machine Learning*. PMLR, pp. 3481–3490.