# Convex-Concave Minmax Optimization

## Applications and Methods

School of Data Science

August 5, 2022

Yilin Gu

# Outline

Motivation

Background
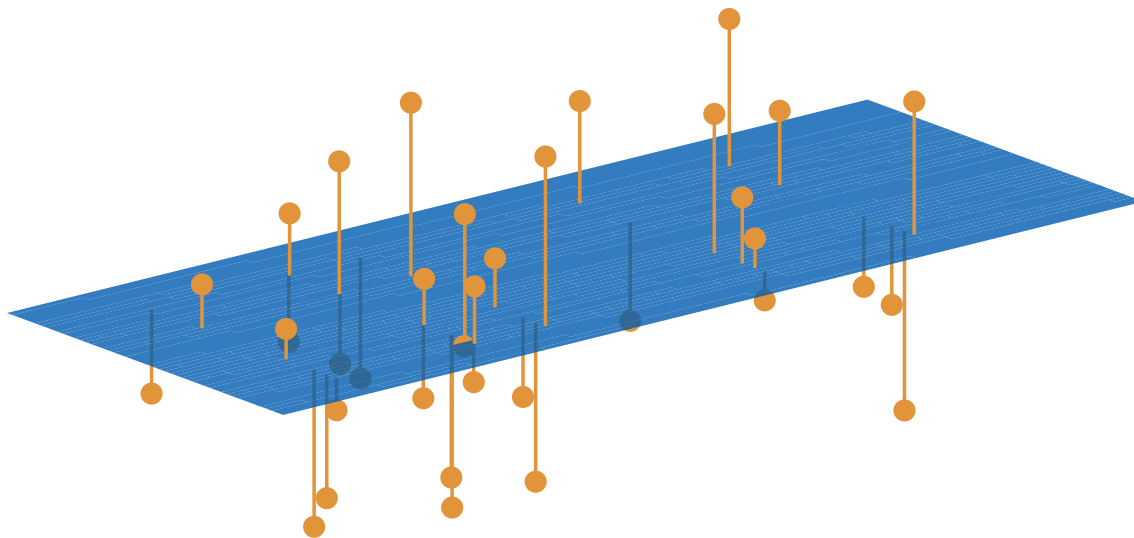
Algorithms

Last Iterate

# Robustness in Learning (I)

**Standard training:** Minimize empirical loss by selecting parameters $\boldsymbol{x}$

$$L(\boldsymbol{x}) := \frac{1}{N} \sum_{i=1}^{N} \ell(a_i, b_i | \boldsymbol{x})$$

$(a_i, b_i)$ is a training sample, $a_i$ is the input and $b_i$ is the expected output



**Linear regression:** Consider $\ell(a_i, b_i | \boldsymbol{x}) = \|a_i^\mathsf{T} \boldsymbol{x} - b_i\|^2$

$$L(\boldsymbol{x}) := \frac{1}{N} \sum_{i=1}^{N} \|a_i^\mathsf{T} \boldsymbol{x} - b_i\|^2 = \frac{1}{N} \|A\boldsymbol{x} - b\|^2$$
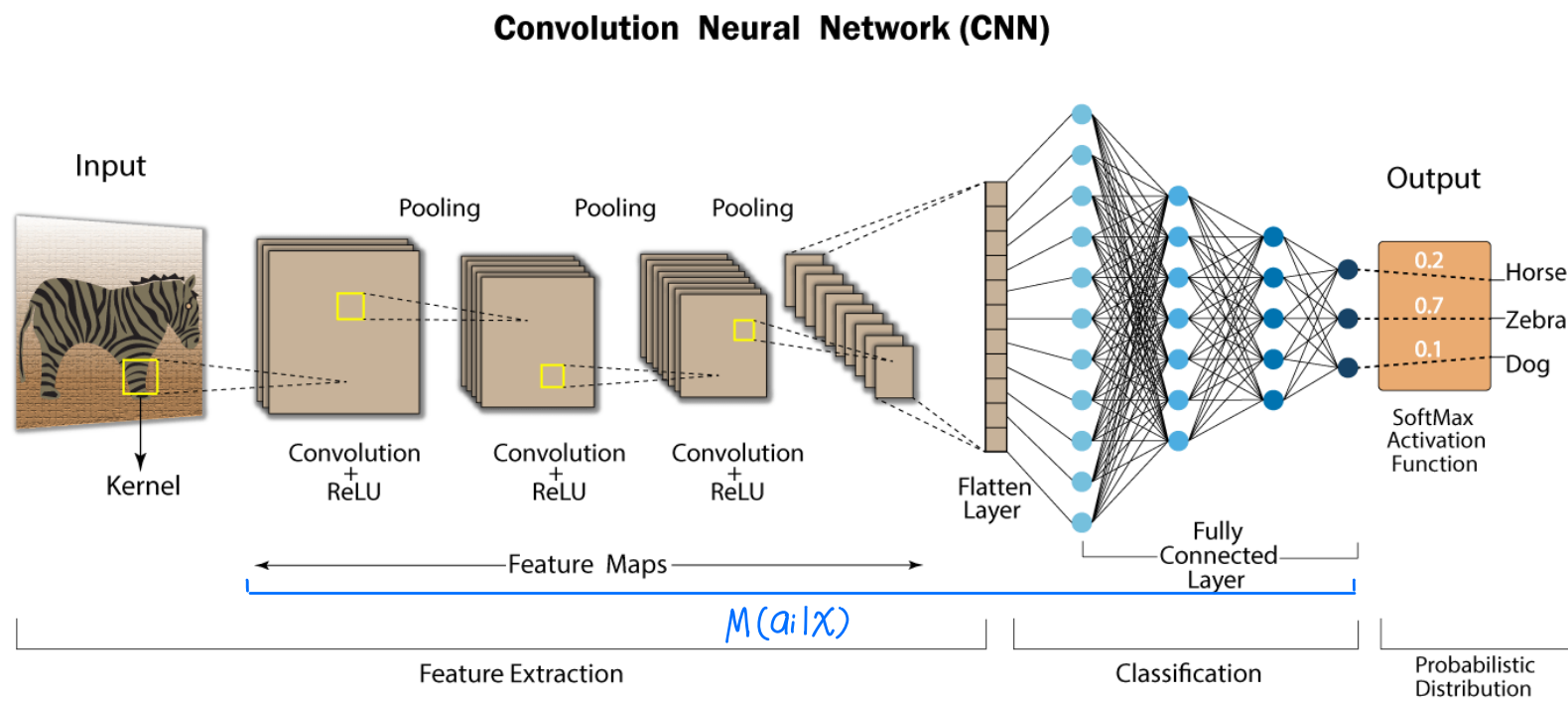
# Robustness in Learning (II)

**Neural network:** Consider $\ell(a_i, b_i | \boldsymbol{x}) = \| \mathcal{M}(a_i | \boldsymbol{x}) - \underline{b_i} \|^2$

$$L(\boldsymbol{x}) := \frac{1}{N} \sum_{i=1}^{N} \| \mathcal{M}(a_i | \boldsymbol{x}) - b_i \|^2$$

where $\mathcal{M}(\cdot | \boldsymbol{x})$ denotes the model with parameters $\boldsymbol{x}$

**Convolution  Neural  Network (CNN)**

# Robustness in Learning (III)



+0.01 X

"Zebra" 89%          =          "Horse" 95%

▶ **Robust training:** Consider inputs with modifications represented as perturbations $\boldsymbol{y}$ of data.

▶ It amounts to choosing $\boldsymbol{x}$ to solve the **minmax problem:**

$$\min_{\boldsymbol{x}\in\mathbb{R}^m} \frac{1}{N}\sum_{i=1}^{N} \underbrace{\max_{\boldsymbol{y}\in\mathcal{S}} \overbrace{\ell(a_i+\boldsymbol{y}, b_i|\boldsymbol{x})}^{\text{worst case of loss}}}_{\text{noise}\uparrow} \xrightarrow{\text{minimize}} \text{robustness}\uparrow$$

where $\mathcal{S}$ denotes allowable perturbations

# Outline

# Minmax Problems

Consider the following minmax problem:

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

**Applications:**

▶ Worst-case design (robust optimization): Minimize over $\boldsymbol{x}$ the loss function with the worst possible value of $\boldsymbol{y}$

# Minmax Problems

Consider the following minmax problem:

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

**Applications:**

▶ Worst-case design (robust optimization): Minimize over $\boldsymbol{x}$ the loss function with the worst possible value of $\boldsymbol{y}$

▶ Duality theory for constrained optimization:

   ▶ Primal problem

   $$\min_{\boldsymbol{x} \in \mathbb{R}^m} f(\boldsymbol{x}), \quad \text{s.t. } g(\boldsymbol{x}) \leq 0$$

   $$\max_{y} f(x) + y^\top g(x) \rightarrow \text{close to } \min_{x} f(x)$$

   ▶ Lagarangian function

   $$\mathcal{L}(\boldsymbol{x}, y) = f(\boldsymbol{x}) + yg(\boldsymbol{x}), \quad y \geq 0$$

   ▶ Dual problem is a minmax problem

   $$\max_{y \geq 0} \min_{\boldsymbol{x} \in \mathbb{R}^m} \mathcal{L}(\boldsymbol{x}, y) \quad \Longleftrightarrow \quad -\min_{y \geq 0} \max_{\boldsymbol{x} \in \mathbb{R}^m} -\mathcal{L}(\boldsymbol{x}, y)$$

# Convex-concave Functions

(1) $\max\limits_{y \geq 0} \min\limits_{x} L(x,y)$

$= \max\limits_{y} \min\limits_{x} f(x) + y^T g(x) \cdots$ ①

suppose $x^*$ minimizes $L(x,y)$. then.

① $= \max\limits_{y} f(x^*) + y^T g(x^*) \cdots$ ②

$\because g(x^*) \leq 0$. $y \geq 0$

To maximize ②. $y = 0$.

$\therefore \max\limits_{y} \min\limits_{x} L(x,y) = f(x^*)$.

(2) $-\min\limits_{y \geq 0} \max\limits_{x} -L(x,y)$

$= -\min\limits_{y} \max\limits_{x} -f(x) - y^T g(x) \cdots$ ③

suppose $x^{**}$ maximize $-L(x,y)$. then.

③ $= -\min\limits_{y} -f(x^{**}) - y^T g(x^{**}) \cdots$ ④

$\because g(x^{**}) \leq 0$. $y \geq 0$

To minimize ④. $y = 0$.

$\therefore -\min\limits_{y} \max\limits_{x} -L(x,y) = f(x^{**})$.

$f(x,y) = x^2 - y^2$

High

Low

function w.r.t. $x$

function w.r.t. $y$

$x^* = x^{**}$.

min $L(x^*,y)$   $L(x,y)$

$x^* = x^{**}$

$-L(x,y)$

max $-L(x^*,y)$

# Convex-concave Functions



$$f(x, y) = x^2 - y^2$$

function w.r.t. $x$

function w.r.t. $y$

## Definition: Convex-concave Function

The function $f(\boldsymbol{x}, \boldsymbol{y})$ is convex-concave if

- for any $\boldsymbol{y} \in \mathbb{R}^n$, the function $f(\boldsymbol{x}, \boldsymbol{y})$ is a convex function of $\boldsymbol{x}$; and
- for any $\boldsymbol{x} \in \mathbb{R}^m$, the function $f(\boldsymbol{x}, \boldsymbol{y})$ is a concave function of $\boldsymbol{y}$

# Saddle Points



$f(x, y) = x^2 - y^2$ with saddle point $(0, 0)$

$$f(0, y) \leqslant f(0, 0) \leqslant f(x, 0)$$

# Saddle Points



$f(x, y) = x^2 - y^2$ with saddle point $(0, 0)$

## Definition: Saddle Points

A saddle point of the minmax problem is a pair $(\boldsymbol{x}^*, \boldsymbol{y}^*) \in \mathbb{R}^m \times \mathbb{R}^n$ that

$$f(\boldsymbol{x}^*, \boldsymbol{y}) \leq f(\boldsymbol{x}^*, \boldsymbol{y}^*) \leq f(\boldsymbol{x}, \boldsymbol{y}^*)$$

for all $\boldsymbol{x} \in \mathbb{R}^m$ and $\boldsymbol{y} \in \mathbb{R}^n$

# Outline

Motivation

Background
- ▷ Minmax Problems
- ▷ Convergence Measure

Algorithms

Last Iterate

# Primal-dual Gap

Define the constant $D$ and the **neighborhood** $\mathcal{S}$ of saddle point $(\boldsymbol{x}^*, \boldsymbol{y}^*)$

$$D := \|\boldsymbol{x}_0 - \boldsymbol{x}^*\|^2 + \|\boldsymbol{y}_0 - \boldsymbol{y}^*\|^2 \longrightarrow \|\text{ initial } - \text{ saddle}\|^2$$

$$\mathcal{S} := \left\{ (\boldsymbol{x}, \boldsymbol{y}) : \|\boldsymbol{x} - \boldsymbol{x}^*\|^2 + \|\boldsymbol{y} - \boldsymbol{y}^*\|^2 \leq 2D \right\}$$

all the iterations will in $\mathcal{S}$.



$D$

$D$

$(x_0, y_0)$

$(x^*, y^*)$

# Primal-dual Gap

Define the constant $D$ and the neighborhood $\mathcal{S}$ of saddle point $(\boldsymbol{x}^*, \boldsymbol{y}^*)$

$$D := \|\boldsymbol{x}_0 - \boldsymbol{x}^*\|^2 + \|\boldsymbol{y}_0 - \boldsymbol{y}^*\|^2$$

$$\mathcal{S} := \left\{ (\boldsymbol{x}, \boldsymbol{y}) : \|\boldsymbol{x} - \boldsymbol{x}^*\|^2 + \|\boldsymbol{y} - \boldsymbol{y}^*\|^2 \leq 2D \right\}$$

## Definition: Primal-dual Gap

For fixed $\bar{\boldsymbol{x}}$ and $\bar{\boldsymbol{y}}$, the primal-dual gap is

$$\left| f(\bar{x}, \bar{y}) - f(x^*, y^*) \right| < \varepsilon \quad \Leftarrow \quad \max_{\boldsymbol{y} : (\bar{\boldsymbol{x}}, \boldsymbol{y}) \in \mathcal{S}} f(\bar{\boldsymbol{x}}, \boldsymbol{y}) - \min_{\boldsymbol{x} : (\boldsymbol{x}, \bar{\boldsymbol{y}}) \in \mathcal{S}} f(\boldsymbol{x}, \bar{\boldsymbol{y}}) \; < \; \varepsilon$$

$$\max_y f(\bar{x}, y) \geqslant f(\bar{x}, y^*) \geqslant \underbrace{f(x^*, y^*)}_{\text{Saddle point}} \geqslant \min_x f(x, \bar{y})$$

**Remark:**

▶ The primal-dual gap is zero iff $(\bar{\boldsymbol{x}}, \bar{\boldsymbol{y}})$ is a saddle point

▶ We also write the primal dual gap as

$$\underbrace{\left[ \max_{\boldsymbol{y} : (\bar{\boldsymbol{x}}, \boldsymbol{y}) \in \mathcal{S}} f(\bar{\boldsymbol{x}}, \boldsymbol{y}) - f(\boldsymbol{x}^*, \boldsymbol{y}^*) \right]}_{>0} + \underbrace{\left[ f(\boldsymbol{x}^*, \boldsymbol{y}^*) - \min_{\boldsymbol{x} : (\boldsymbol{x}, \bar{\boldsymbol{y}}) \in \mathcal{S}} f(\boldsymbol{x}, \bar{\boldsymbol{y}}) \right]}_{>0}$$

# Monotone Operator

Consider the minmax problem with convex-concave objective function

- Saddle point satisfies the first-order optimality condition

$$\nabla_x f(x^*, y^*) = 0 \quad \text{and} \quad \nabla_y f(x^*, y^*) = 0$$

- Define $z := [x; y] \in \mathbb{R}^{m+n}$ and the monotone operator

$$F(z) := [\nabla_x f(x, y); \ominus \nabla_y f(x, y)] \quad \Longrightarrow \quad F(z^*) = 0$$

*y concave.*

## Definition: Monotone Operator

$F$ is a monotone operator if for any $z_1, z_2 \in \mathbb{R}^{m+n}$

$$\langle F(z_1) - F(z_2), z_1 - z_2 \rangle \geq 0$$

**Remark:** If $h : \mathbb{R}^n \to \mathbb{R}$ is convex, then $\nabla h : \mathbb{R}^n \to \mathbb{R}^n$ is monotone

# Outline

Motivation

Background

Algorithms
   ▷ Gradient Descent Ascent (GDA)
   ▷ Proximal Point Algorithm (PPA)
   ▷ Optimistic Gradient Descent Ascent (OGDA)
   ▷ Extragradient Method (EG)

Last Iterate

# GDA (I)

## Algorithm: Gradient Descent Ascent

- **Initialization:** $x_0 \in \mathbb{R}^m, y_0 \in \mathbb{R}^n$ and step size $\eta > 0$
- **Iteration:**

  $\min_x f(x,y)$

  $$x_{k+1} = x_k - \eta \nabla_x f(x_k, y_k) \qquad \text{Gradient Descent}$$

  $$y_{k+1} = y_k + \eta \nabla_y f(x_k, y_k) \qquad \text{Gradient Ascent}$$

  $\max_y f(x,y)$

# GDA (I)

## Algorithm: Gradient Descent Ascent

- **Initialization:** $\boldsymbol{x}_0 \in \mathbb{R}^m, \boldsymbol{y}_0 \in \mathbb{R}^n$ and step size $\eta > 0$
- **Iteration:**

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \eta \nabla_{\boldsymbol{x}} f(\boldsymbol{x}_k, \boldsymbol{y}_k) \qquad \text{Gradient Descent}$$

$$\boldsymbol{y}_{k+1} = \boldsymbol{y}_k + \eta \nabla_{\boldsymbol{y}} f(\boldsymbol{x}_k, \boldsymbol{y}_k) \qquad \text{Gradient Ascent}$$

- Even for the simplest case, GDA diverges

- Consider the following bilinear problem

$$\min_{\boldsymbol{x} \in \mathbb{R}^d} \max_{\boldsymbol{y} \in \mathbb{R}^d} f(\boldsymbol{x}, \boldsymbol{y}) = \boldsymbol{x}^\mathsf{T} \boldsymbol{y}$$

- The GDA updates for this problem

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \eta \boldsymbol{y}_k$$

$$\boldsymbol{y}_{k+1} = \boldsymbol{y}_k + \eta \boldsymbol{x}_k$$

# GDA (II)

▶ The GDA updates for this problem

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \eta \boldsymbol{y}_k \quad \text{and} \quad \boldsymbol{y}_{k+1} = \boldsymbol{y}_k + \eta \boldsymbol{x}_k$$

▶ At the $k$-th of GDA, we have

$$\|\boldsymbol{x}_{k+1}\|^2 + \|\boldsymbol{y}_{k+1}\|^2 = (1 + \eta^2)(\|\boldsymbol{x}_k\|^2 + \|\boldsymbol{y}_k\|^2)$$

▶ GDA diverges because $1 + \eta^2 > 1$   $(1+\eta^2)^{(k)}(\|x_0\|^2 + \|y_0\|^2)$

$\uparrow \quad \Rightarrow \text{ diverge}$



● Saddle $(0, 0)$    ● Initial $(10, 10)$

# Outline

Motivation

Background

Algorithms
- ▷ Gradient Descent Ascent (GDA)
- ▷ Proximal Point Algorithm (PPA)
- ▷ Optimistic Gradient Descent Ascent (OGDA)
- ▷ Extragradient Method (EG)

Last Iterate

# PPA (I)

## Algorithm: Proximal Point Algorithm

▶ **Initialization:** $\boldsymbol{x}_0 \in \mathbb{R}^m, \boldsymbol{y}_0 \in \mathbb{R}^n$ and step size $\eta > 0$

▶ **Iteration:** The pair $(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1})$ is the unique solution to

$$\min_{\boldsymbol{x} \in \mathbb{R}^m} \max_{\boldsymbol{y} \in \mathbb{R}^n} \left\{ f(\boldsymbol{x}, \boldsymbol{y}) + \frac{1}{2\eta} \|\boldsymbol{x} - \boldsymbol{x}_k\|^2 - \frac{1}{2\eta} \|\boldsymbol{y} - \boldsymbol{y}_k\|^2 \right\} \; = g(x,y)$$

strongly convex to x

strongly concave to y

Optimality condition:

$$\begin{cases} \nabla_x g(x_{k+1}, y_{k+1}) = 0 \\ \\ \nabla_y g(x_{k+1}, y_{k+1}) = 0 \end{cases} \longrightarrow$$

$$\nabla g(x,y) = \nabla_x f(x_{k+1}, y_{k+1}) + \frac{1}{\eta}(x_{k+1} - x_k)$$

$$\iff x_{k+1} = x_k - \eta \cdot \nabla_x f(x_{k+1}, y_{k+1})$$

# PPA (I)

**Algorithm: Proximal Point Algorithm**

▸ **Initialization:** $x_0 \in \mathbb{R}^m, y_0 \in \mathbb{R}^n$ and step size $\eta > 0$

▸ **Iteration:** The pair $(x_{k+1}, y_{k+1})$ is the unique solution to

$$\min_{x \in \mathbb{R}^m} \max_{y \in \mathbb{R}^n} \left\{ f(x, y) + \frac{1}{2\eta} \|x - x_k\|^2 - \frac{1}{2\eta} \|y - y_k\|^2 \right\}$$

**Remark:** Iterative steps of PPA can be written as

$$x_{k+1} = x_k - \eta \nabla_x f(x_{k+1}, y_{k+1})$$

$$y_{k+1} = y_k + \eta \nabla_y f(x_{k+1}, y_{k+1})$$

Different from GDA steps

$$x_{k+1} = x_k - \eta \nabla_x f(x_k, y_k)$$

$$y_{k+1} = y_k + \eta \nabla_y f(x_k, y_k)$$

# PPA (II)

▶ PPA for $f(\boldsymbol{x}, \boldsymbol{y}) = \boldsymbol{x}^\mathsf{T}\boldsymbol{y}$

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

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

▶ At the $k$-th iteration of PPA, we have

$$\|\boldsymbol{x}_{k+1}\|^2 + \|\boldsymbol{y}_{k+1}\|^2 = \frac{1}{1+\eta^2}(\|\boldsymbol{x}_k\|^2 + \|\boldsymbol{y}_k\|^2)$$

# PPA (II)

- ▶ PPA for $f(\mathbf{x}, \mathbf{y}) = \mathbf{x}^{\mathsf{T}} \mathbf{y}$

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

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

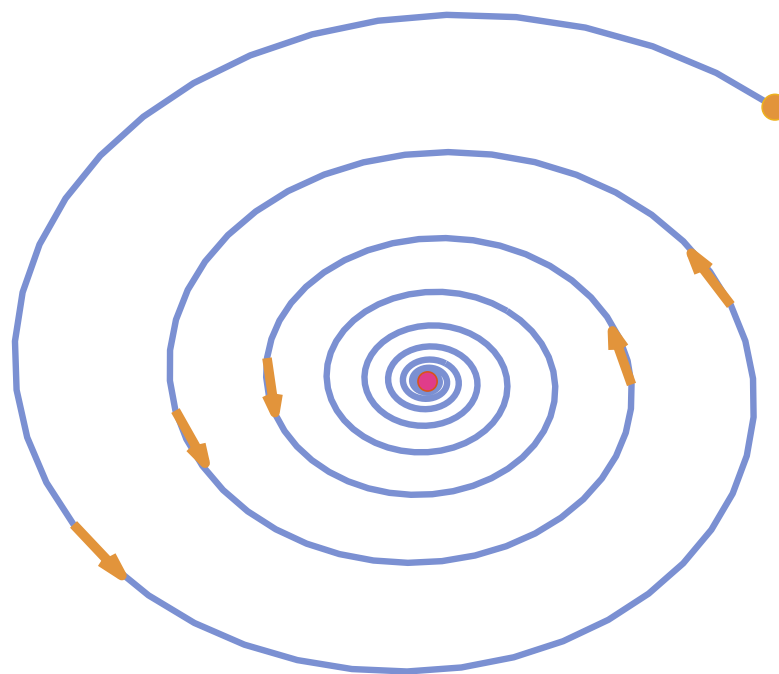- ▶ At the $k$-th iteration of PPA, we have

$$\|\mathbf{x}_{k+1}\|^2 + \|\mathbf{y}_{k+1}\|^2 = \frac{1}{1+\eta^2}(\|\mathbf{x}_k\|^2 + \|\mathbf{y}_k\|^2)$$

- ▶ True iterative steps

$$\mathbf{x}_{k+1} = \frac{\mathbf{x}_k - \eta \mathbf{y}_k}{1+\eta^2}$$

$$\mathbf{y}_{k+1} = \frac{\mathbf{y}_k + \eta \mathbf{x}_k}{1+\eta^2}$$

- ▶ PPA converges to saddle point

● Saddle $(0,0)$  ● Initial $(10,10)$

# PPA (III)

- Let iterates $(\boldsymbol{x}_k, \boldsymbol{y}_k)$ be generated by PPA with step size $\eta$
- Define the averaged iterates $(\bar{\boldsymbol{x}}_k, \bar{\boldsymbol{y}}_k)$ as

$$\bar{\boldsymbol{x}}_k := \frac{1}{k}\sum_{i=1}^{k} \boldsymbol{x}_i \quad \text{and} \quad \bar{\boldsymbol{y}}_k := \frac{1}{k}\sum_{i=1}^{k} \boldsymbol{y}_i$$

## Theorem: Convergence of Averaged Iterates

- If $f$ is convex-concave and $L$-smooth
- Then, we have

$$\max_{\boldsymbol{y}:(\bar{\boldsymbol{x}}_k, \boldsymbol{y}) \in \mathcal{S}} f(\bar{\boldsymbol{x}}_k, \boldsymbol{y}) - \min_{\boldsymbol{x}:(\boldsymbol{x}, \bar{\boldsymbol{y}}_k) \in \mathcal{S}} f(\boldsymbol{x}, \bar{\boldsymbol{y}}_k) \leq \frac{D}{\eta k}$$

**Remark:** PPA involves operator inversion and is not easy to implement

**Require:** Efficient algorithms that behave like PPA!

# Outline

Motivation

Background

Algorithms
  ▷ Gradient Descent Ascent (GDA)
  ▷ Proximal Point Algorithm (PPA)
  ▷ Optimistic Gradient Descent Ascent (OGDA)
  ▷ Extragradient Method (EG)

Last Iterate

# OGDA (I)

## Algorithm: Optimistic Gradient Descent Ascent

- **Initialization:** $\boldsymbol{x}_0 \in \mathbb{R}^m, \boldsymbol{y}_0 \in \mathbb{R}^n$ and step size $\eta > 0$
- **Iteration:**

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

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

$$x_{k+1} = x_k - \eta \nabla_x f(x_{k+1}, y_{k+1}) + \eta \underbrace{\left[ \nabla_x f(x_{k+1}, y_{k+1}) - \nabla_x f(x_k, y_k) \right] - \eta \left[ \nabla_x f(x_k, y_k) - \nabla_x f(x_{k-1}, y_{k-1}) \right]}_{\text{gradient difference}}$$

$$\Downarrow$$

$$\nabla_x f(x_{k+1}, y_{k+1}) - \nabla_x f(x_k, y_k) \approx \nabla_x f(x_k, y_k) - \nabla_x f(x_{k-1}, y_{k-1})$$

# OGDA (I)

**Algorithm: Optimistic Gradient Descent Ascent**

- **Initialization:** $\boldsymbol{x}_0 \in \mathbb{R}^m, \boldsymbol{y}_0 \in \mathbb{R}^n$ and step size $\eta > 0$
- **Iteration:**

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

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

**Remark:** OGDA can be seen as PPA with error term

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \eta \nabla_{\boldsymbol{x}} f(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1}) + \eta \varepsilon_{\boldsymbol{x},k}$$

$$\boldsymbol{y}_{k+1} = \boldsymbol{y}_k + \eta \nabla_{\boldsymbol{y}} f(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1}) - \eta \varepsilon_{\boldsymbol{y},k}$$

Approximate using linear extrapolation of the previous gradients

$$\nabla f(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1}) \overset{\approx}{\approx} \nabla f(\boldsymbol{x}_k, \boldsymbol{y}_k) + [\nabla f(\boldsymbol{x}_k, \boldsymbol{y}_k) - \nabla f(\boldsymbol{x}_{k-1}, \boldsymbol{y}_{k-1})]$$

*Approximate*

# OGDA (II)

## Algorithm: Optimistic Gradient Descent Ascent

- **Initialization:** $\boldsymbol{x}_0 \in \mathbb{R}^m, \boldsymbol{y}_0 \in \mathbb{R}^n$ and step size $\eta > 0$
- **Iteration:**

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

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

- Consider $f(\boldsymbol{x}, \boldsymbol{y}) = \boldsymbol{x}^{\mathsf{T}} \boldsymbol{y}$

- Convergence paths are similar

- OGDA approximates PPA



Converge ↑

— PPA    —· OGDA ✓

# OGDA (III)

- Let iterates $(\boldsymbol{x}_k, \boldsymbol{y}_k)$ be generated by OGDA with step size $\eta$
- Define the averaged iterates $(\bar{\boldsymbol{x}}_k, \bar{\boldsymbol{y}}_k)$ as

$$\bar{\boldsymbol{x}}_k := \frac{1}{k}\sum\nolimits_{i=1}^{k} \boldsymbol{x}_i \quad \text{and} \quad \bar{\boldsymbol{y}}_k := \frac{1}{k}\sum\nolimits_{i=1}^{k} \boldsymbol{y}_i$$

## Theorem: Convergence of Averaged Iterates

- If $f$ is convex-concave and $L$-smooth
- Then, we have

$$\max_{\boldsymbol{y}:(\bar{\boldsymbol{x}}_k,\boldsymbol{y})\in\mathcal{S}} f(\bar{\boldsymbol{x}}_k, \boldsymbol{y}) - \min_{\boldsymbol{x}:(\boldsymbol{x},\bar{\boldsymbol{y}}_k)\in\mathcal{S}} f(\boldsymbol{x}, \bar{\boldsymbol{y}}_k) \leq \frac{5D}{\eta k}$$

**Remark:**

- OGDA is an implementable version of PPA
- OGDA enjoys similar convergence guarantee $\mathcal{O}(1/k)$

# Outline

Motivation

Background

Algorithms

    ▷ Gradient Descent Ascent (GDA)

    ▷ Proximal Point Algorithm (PPA)

    ▷ Optimistic Gradient Descent Ascent (OGDA)

    ▷ Extragradient Method (EG)

Last Iterate

# EG (I)

## Algorithm: Extragradient Method

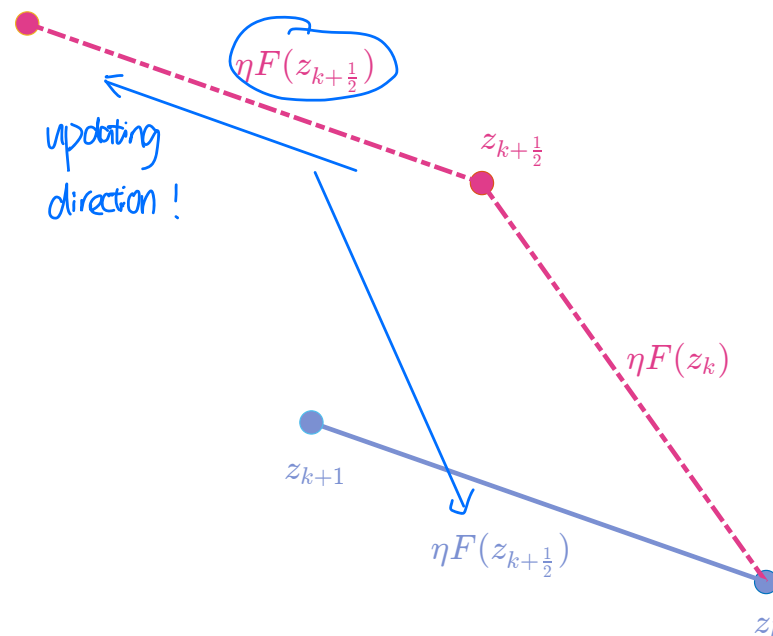▶ **Initialization:** $x_0 \in \mathbb{R}^m$, $y_0 \in \mathbb{R}^n$ and step size $\eta > 0$

▶ **Iteration:**

$$[x,y]^\top \qquad [\nabla_x f(x,y), \nabla_y f(x,y)]^\top$$

$$z_{k+\frac{1}{2}} = z_k - \eta F(z_k)$$

$$z_{k+1} = z_k - \eta F(z_{k+\frac{1}{2}})$$

▶ Define vector $z := [x; y]$

▶ Define the operator $F$ as

$$F(z) := [\nabla_x f(x, y); -\nabla_y f(x, y)]$$

▶ EG utilizes the gradient of midpoint to update

$\eta F(z_{k+\frac{1}{2}})$

updating direction !

$z_{k+\frac{1}{2}}$

$\eta F(z_k)$

$z_{k+1}$

$\eta F(z_{k+\frac{1}{2}})$

$z_k$

# EG (II)

## Algorithm: Extragradient Method

▶ **Initialization:** $x_0 \in \mathbb{R}^m, y_0 \in \mathbb{R}^n$ and step size $\eta > 0$

▶ **Iteration:**

$$z_{k+\frac{1}{2}} = z_k - \eta F(z_k) \qquad z_{k-\frac{1}{2}} = z_{k-1} - \eta F(z_{k-1})$$

$$z_{k+1} = z_k - \eta F(z_{k+\frac{1}{2}}) \qquad z_{k+\frac{1}{2}} = z_k - \eta F(z_{k-\frac{1}{2}})$$

▶ Updates can be written as

$$z_{k+\frac{1}{2}} = z_{k-\frac{1}{2}} - \eta F(z_{k-\frac{1}{2}}) - \eta[F(z_k) - F(z_{k-1})]$$

$$z_{k+\frac{1}{2}} = z_k - \eta F(z_k)$$
$$= (z_{k-1} - \eta F(z_{k-\frac{1}{2}})) - \eta F(z_k)$$
$$= \left[ z_{k-\frac{1}{2}} + \eta F(z_{k-1}) - \eta F(z_{k-\frac{1}{2}}) \right] - \eta F(z_k)$$
$$= z_{k-\frac{1}{2}} - \eta F(z_{k-\frac{1}{2}}) - \eta \left[ F(z_k) - F(z_{k-1}) \right]$$
$$= z_{k-\frac{1}{2}} - \eta F(z_{k+\frac{1}{2}}) + \eta \left[ (F(z_{k+\frac{1}{2}}) - F(z_{k-\frac{1}{2}})) - (F(z_k) - F(z_{k-1})) \right]$$

Approximate

# EG (II)

**Algorithm: Extragradient Method**

▶ **Initialization:** $\boldsymbol{x}_0 \in \mathbb{R}^m, \boldsymbol{y}_0 \in \mathbb{R}^n$ and step size $\eta > 0$

▶ **Iteration:**

$$\boldsymbol{z}_{k+\frac{1}{2}} = \boldsymbol{z}_k - \eta F(\boldsymbol{z}_k)$$

$$\boldsymbol{z}_{k+1} = \boldsymbol{z}_k - \eta F(\boldsymbol{z}_{k+\frac{1}{2}})$$

▶ Updates can be written as

$$\boldsymbol{z}_{k+\frac{1}{2}} = \boldsymbol{z}_{k-\frac{1}{2}} - \eta F(\boldsymbol{z}_{k-\frac{1}{2}}) - \eta[F(\boldsymbol{z}_k) - F(\boldsymbol{z}_{k-1})]$$

▶ When the variations are close to each other, i.e.,

$$F(\boldsymbol{z}_k) - F(\boldsymbol{z}_{k-1}) \approx F(\boldsymbol{z}_{k+\frac{1}{2}}) - F(\boldsymbol{z}_{k-\frac{1}{2}}) \qquad \text{EG}$$

EG method approximates PPA $\qquad \downarrow$ Approximate

$$\boldsymbol{z}_{k+\frac{1}{2}} \approx \boldsymbol{z}_{k-\frac{1}{2}} - \eta F(\boldsymbol{z}_{k+\frac{1}{2}})$$

# EG (III)

## Algorithm: Extragradient Method

▶ **Initialization:** $\boldsymbol{x}_0 \in \mathbb{R}^m, \boldsymbol{y}_0 \in \mathbb{R}^n$ and step size $\eta > 0$

▶ **Iteration:**

$$\boldsymbol{z}_{k+\frac{1}{2}} = \boldsymbol{z}_k - \eta F(\boldsymbol{z}_k)$$

$$\boldsymbol{z}_{k+1} = \boldsymbol{z}_k - \eta F(\boldsymbol{z}_{k+\frac{1}{2}})$$

▶ Consider $f(\boldsymbol{x}, \boldsymbol{y}) = \boldsymbol{x}^\mathsf{T} \boldsymbol{y}$

▶ Convergence paths are similar

▶ EG approximates PPA



— PPA    — · EG

# EG (IV)

- Let iterates $(\boldsymbol{x}_k, \boldsymbol{y}_k)$ be generated by EG with step size $\eta$
- Define the averaged iterates $(\bar{\boldsymbol{x}}_k, \bar{\boldsymbol{y}}_k)$ as

$$\bar{\boldsymbol{x}}_k := \frac{1}{k}\sum_{i=1}^{k} \boldsymbol{x}_i \quad \text{and} \quad \bar{\boldsymbol{y}}_k := \frac{1}{k}\sum_{i=1}^{k} \boldsymbol{y}_i$$

## Theorem: Convergence of Averaged Iterates

- If $f$ is convex-concave and $L$-smooth
- Then, we have

$$\max_{\boldsymbol{y}:(\bar{\boldsymbol{x}}_k, \boldsymbol{y})\in\mathcal{S}} f(\bar{\boldsymbol{x}}_k, \boldsymbol{y}) - \min_{\boldsymbol{x}:(\boldsymbol{x}, \bar{\boldsymbol{y}}_k)\in\mathcal{S}} f(\boldsymbol{x}, \bar{\boldsymbol{y}}_k) \leq \frac{16D}{\eta k}$$

**Remark:**

- EG is an implementable version of PPA
- EG enjoys similar convergence guarantee $\mathcal{O}(1/k)$

# Last Iterate Convergence

The averaged iterate is not always what we want!

- Imagine that we are seeking for a sparse solution $x^*$
- Assume $\bar{x} := (x_1 + x_2 + x_3)/3$ reaches $\varepsilon$-accuracy

$$x_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \qquad x_2 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \qquad x_3 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \qquad \bar{x} = \begin{bmatrix} 2/3 \\ 1 \\ 1/3 \end{bmatrix}$$

*sparse*           *non-sparse*

# Last Iterate Convergence

> The averaged iterate is not always what we want!

- ▶ Imagine that we are seeking for a sparse solution $\boldsymbol{x}^*$
- ▶ Assume $\bar{\boldsymbol{x}} := (\boldsymbol{x}_1 + \boldsymbol{x}_2 + \boldsymbol{x}_3)/3$ reaches $\varepsilon$-accuracy

$$
\boldsymbol{x}_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \qquad
\boldsymbol{x}_2 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \qquad
\boldsymbol{x}_3 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \qquad
\bar{\boldsymbol{x}} = \begin{bmatrix} 2/3 \\ 1 \\ 1/3 \end{bmatrix}
$$

## Theorem: Last Iterate Convergence

- ▶ Let iterates $(\boldsymbol{x}_k, \boldsymbol{y}_k)$ be generated by EG/PPA
- ▶ If $f$ is convex-concave and $L$-smooth
- ▶ Then, we have
$$
\max_{\boldsymbol{y}:(\boldsymbol{x}^k, \boldsymbol{y}) \in \mathcal{S}} f(\boldsymbol{x}^k, \boldsymbol{y}) - \min_{\boldsymbol{x}:(\boldsymbol{x}, \boldsymbol{y}^k) \in \mathcal{S}} f(\boldsymbol{x}, \boldsymbol{y}^k) = \Theta\left(\frac{1}{\sqrt{k}}\right)
$$

**Remark:** Slower than the averaged iterate results $\mathcal{O}(1/k)$

# Conclusion

Motivation

Background
- ▷ Minmax Problems
- ▷ Convergence Measure

Algorithms
- ▷ Gradient Descent Ascent (GDA)
- ▷ Proximal Point Algorithm (PPA)
- ▷ Optimistic Gradient Descent Ascent (OGDA)
- ▷ Extragradient Method (EG)

Last Iterate