# Online Learning in Games

### DRAFT

Prof Volkan Cevher volkan.cevher@epfl.ch

### Lecture 11: Equilibrium computation for nonmonontone problems

Laboratory for Information and Inference Systems (LIONS) École Polytechnique Fédérale de Lausanne (EPFL)

**EE-735** (Spring 2024)

















## License Information for Online Learning in Games Slides

- ▶ This work is released under a <u>Creative Commons License</u> with the following terms:
- Attribution
  - ► The licensor permits others to copy, distribute, display, and perform the work. In return, licensees must give the original authors credit.
- Non-Commercial
  - ► The licensor permits others to copy, distribute, display, and perform the work. In return, licensees may not use the work for commercial purposes unless they get the licensor's permission.
- ▶ Share Alike
  - ► The licensor permits others to distribute derivative works only under a license identical to the one that governs the licensor's work.
- ► Full Text of the License

## Acknowledgements

These slides were originally prepared by Dongyang Fan and Wenjie Xu.

## **Outline**



#### Introduction

Min-Max Optimization has wide applications in machine learning:

- Nonsmooth optimization
- Generative adversarial networks
- Distributionally robust optimization

## Example (Minimax optimization)

$$\min_{x \in \mathbb{R}^{n_x}} \max_{y \in \mathbb{R}^{n_y}} \mathcal{L}(x, y) \coloneqq \phi(x, y) + g(x) - h(y) \tag{1}$$

where  $\phi$  is not necessarily convex in x and concave in y.  $g(\cdot)$  and  $h(\cdot)$  are lower semi-continuous and convex, for example, they are usually chosen as  $l_1$ ,  $l_2$  regularization terms or indicator functions (constrained case).

#### Remarks:

- $\circ Fz = (\nabla_x \phi(x, y), -\nabla_y \phi(x, y)), Az = (\partial g(x), \partial h(y))$
- $\circ\,$  From an operator point of view: the problem can be translated as finding zeros of  ${\bf F}+{\bf A}$
- $\circ$  We will start our analysis from the unconstrained case, i.e.  $A \equiv 0$

## Recap - Extragradient & Monotone Case

## Definition (Operators)

- $lackbox{ An operator } F:\mathbb{R}^d o \mathbb{R}^d ext{ is said to be monotone if } \langle Fz-Fz',z-z' 
  angle \geq 0, orall z,z' \in \mathbb{R}^d$
- F is maximal monotone if there is no monotone operator that properly contains it.
- $lackbox{ An operator } F:\mathbb{R}^d o \mathbb{R}^d ext{ is said to be } L ext{-Lipschitz for } L>0 ext{ if } \|Fz-Fz'\| \leq L\|z-z'\|, \forall z,z' \in \mathbb{R}^d$
- o Recall the extragradient (EG) algorithm from ? ], which is a de facto algorithm for monotone VI problems:

$$\begin{split} &\bar{z}^t = z^t - \gamma F z^t, \\ &z^{t+1} = z^t - \gamma F \bar{z}^t. \end{split} \tag{EG}$$

- $\circ$  EG can be seen as an approximation to the proximal point method  $z^{t+1}=z^t-\gamma F(z^{t+1})$
- o For F monotone and L-Lipschitz ( $\phi$  convex-concave, L-Lipschitz smooth), EG has tight last iterate convergence rate  $\mathcal{O}(\frac{1}{\sqrt{T}})$  in terms of gap function  $^1$ . [? ]. While average iterate enjoys  $\mathcal{O}(\frac{1}{T})$  convergence rate. [? ]

<sup>&</sup>lt;sup>1</sup>Gap function:  $GAP_{Z,F,D}(z) = \max_{z' \in Z \cap \mathcal{B}(z,D)} \langle Fz, z - z' \rangle$ 



### Convergence plot for ExtraGradient – Bilinear Game



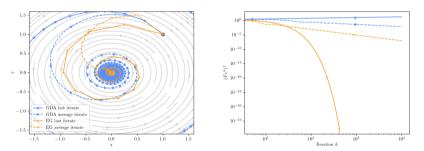


Figure: Left: Convergence plot of EG for bilinear game; Right: decrease of  $||Fz^t||^2$  over number of iterations

For strongly monotone and affine operators (e.g. bilinear games), last iterate converges exponentially fast [?] (compared with  $\mathcal{O}(1/T)$  for average iterate)

## Convergence plot of EG – Forsaken game (non-monontone)

▶ For non-monotone case, we do *not* have guarantees for average iterates (Jensen's ineq not applicable), and even no guarantees for best iterate

$$\begin{aligned} \min_{|x| \leq 3/2} \max_{|y| \leq 3/2} \phi(x,y) &\coloneqq x(y-0.45) + \psi(x) - \psi(y) \\ \text{where } \psi(z) &= \frac{1}{4}z^2 - \frac{1}{2}z^4 + \frac{1}{6}z^6 \end{aligned}$$
 (Forsaken)

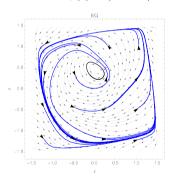


Figure: Last iterate convergence plot (limit cycle marked with dark cicle)

### Go beyond monotone cases?

- ▶ Why Non-monotone? Nonmonotone captures nonconvex-nonconcave minimax by generalizing it.
- ▶ Weak Minty is the (*structured*) type of nonmonotonicity that we will study.
  - it captures the forsaken example we considered in the last slide
  - with weak Minty VI assumption, we could establish nice descent inequality (will appear in the next slide)

## Definition (Variational Inequalities)

- ▶ Minty VI is a sufficient condition for optimality  $:\langle F(z), z z^{\star} \rangle \geq 0$
- ▶ Weak Minty VI:  $\langle F(z), z z^* \rangle \ge \rho ||F(z)||^2$ , for some  $\rho < 0$

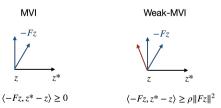


Figure: MVI versus weak MVI

## Why can we have positive result for weak Minty?

- Consider a class of non-monotone problems characterized by Weak MVI
- ▶ Descent inequality for proximal point method ( $z^{t+1} = z^t \gamma F(z^{t+1})$ ):

## **Descent Inequality**

$$||z^{t+1} - z^{\star}||^{2} = ||z^{t} - z^{\star}||^{2} + \gamma^{2} ||F(z^{t+1})||^{2} - 2\langle z^{t} - z^{\star}, \gamma F(z^{t+1}) \rangle$$

$$= ||z^{t} - z^{\star}||^{2} + \gamma^{2} ||F(z^{t+1})||^{2} - 2\langle z^{t+1} - z^{\star}, \gamma F(z^{t+1}) \rangle - 2\langle z^{t+1}, \gamma F(z^{t+1}) \rangle$$

$$= ||z^{t} - z^{\star}||^{2} - \gamma^{2} ||z^{t+1} - z^{t}||^{2} - 2\langle z^{t+1} - z^{\star}, \gamma F(z^{t+1}) \rangle - 2\langle z^{t+1}, \gamma F(z^{t+1}) \rangle$$

$$= ||z^{t} - z^{\star}||^{2} - \gamma^{2} ||z^{t+1} - z^{t}||^{2} - 2\langle z^{t} - z^{t+1}, \gamma F(z^{t+1}) \rangle$$

$$\geq \rho ||F(z^{t+1})||^{2}$$

$$(2)$$

To ensure strict descent between neighboring steps, we need

$$\gamma^2 + 2\gamma \rho > 0 \Rightarrow \rho > -\frac{\gamma}{2} \tag{3}$$

Weak Minty is the largest class for which PPM still has a decrease of distance to optimum every iteration.

## An explicit scheme in unconstrained case: ExtraGradient+

### EG+ is EG with a smaller 2nd step size

? ] introduced EG+, which provably converges to a stationary point for a class of non-monotone problems provided Weak Minty Variational Inequalities (MVI).

$$ar{z}^t = (\mathrm{id} - \gamma F) z^t,$$
 
$$z^{t+1} = z^t - \alpha \gamma F ar{z}^t.$$
 (EG+)

We now show the choices of  $\alpha$ , following the tighter analysis in ? ].

$$||z^{t+1} - z^{\star}||^2 = ||z^t - z^{\star}||^2 + \alpha^2 ||\gamma F \bar{z}^t||^2 - 2\alpha \langle \gamma F \bar{z}^t, z^t - z^{\star} \rangle$$
(4)

Cocoercivity assumption can "convert"  $\langle F \bar{z}^t, z^t - z^\star \rangle$  into  $\|F \bar{z}^t\|^2!$ 

## Cocoercivity of $id - \gamma F$

## Definition (Cocoercivity)

An operator  $F: \mathbb{R}^d \to \mathbb{R}^d$  is said to be  $\beta$ -cocoercive for  $\beta > 0$  if  $\langle Fz - Fz', z - z' \rangle \geq \beta \|Fz - Fz'\|^2$ , for all  $z, z' \in \mathbb{R}^d$ .

#### Claim

Suppose F is L-Lipschitz and  $\gamma \leq 1/L$ . Then, the mapping  $H = \operatorname{id} - \gamma F$  is 1/2-cocoercive

#### Proof.

$$\langle Hz - H\bar{z}, z - \bar{z} \rangle = \langle Hz - H\bar{z}, Hz - H\bar{z} + \gamma Fz - \gamma F\bar{z} \rangle$$

$$= \|Hz - H\bar{z}\|^{2} + \gamma \langle Hz - H\bar{z}, Fz - F\bar{z} \rangle$$

$$= \frac{1}{2} \|Hz - H\bar{z}\|^{2} - \frac{\gamma^{2}}{2} \|F\bar{z} - Fz\|^{2} + \frac{1}{2} \|\bar{z} - z\|^{2}$$

$$= \frac{1}{2} \|Hz - H\bar{z}\|^{2} + \left(\frac{1}{2} - \frac{\gamma^{2}L^{2}}{2}\right) \|\bar{z} - z\|^{2} \ge \frac{1}{2} \|Hz - H\bar{z}\|^{2}$$
(5)

?: 
$$\frac{1}{2}\|\bar{z} - z\|^2 = \frac{1}{2}\|Hz - H\bar{z}\|^2 + \frac{\gamma^2}{2}\|F\bar{z} - Fz\|^2 + \gamma\langle H\bar{z} - Hz, F\bar{z} - Fz\rangle$$

### EG+ requires a small constant step size

We use cocoercivity assumption and weak-Minty VI to bound the  $\langle \gamma F \bar{z}^t, z^t - z^\star \rangle$  term:

$$\langle \gamma F \bar{z}^t, z^t - z^* \rangle = \underbrace{\langle \gamma F \bar{z}^t, z^t - \bar{z}^t \rangle}_{\text{id} - \gamma F \text{ is } 1/2 \text{ cocoercive}} + \underbrace{\langle \gamma F \bar{z}^t, \bar{z}^t - z^* \rangle}_{\text{weak MVI}}$$

$$\geq \frac{1}{2} ||\gamma F \bar{z}^t||^2 + \gamma \rho ||F \bar{z}^t||^2$$
(6)

Thus, the descent inequality becomes

$$||z^{t+1} - z^{\star}||^{2} \le ||z^{t} - z^{\star}||^{2} + \underbrace{\alpha \gamma^{2} \left(\alpha - 1 - 2\rho/\gamma\right)}_{\text{make sure it is smaller than 0}} ||F\bar{z}^{t}||^{2}$$
(7)

Fixing  $\gamma=1/L$ , to ensure a strict descent, we must have  $\alpha<1+2\rho L$ . For a more negative  $\rho$ , i.e. a wider range of class,  $\alpha$  needs to be very small.

## Extend to a more general setting?

- Adding constraints characterized by the A operator: find  $0 \in Tz := Az + Fz$  for Lipschitz continuous operator F and maximally monotone operator A
- $\blacktriangleright$  Allow for larger stepsizes (bigger  $\alpha$ ) and a wider range of problem

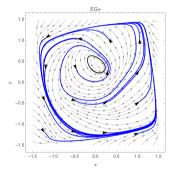


Figure: When the weak MVI constant  $\rho$  does not satisfy algorithmic requirements of (EG+), i.e.  $\rho \in (-1/8L, 0]$ , and (EG+) does not converge to a stationary point but rather the attracting limit cycle (marked in dark circle)

# Algorithm

Compile error

## **Assumptions**

## Assumptions

For T = A + F. Given Assumptions

- ▶ Operator  $A: \mathbb{R}^n \to \mathbb{R}^n$  is maximally monotone
- ightharpoonup Operator  $F:\mathbb{R}^n \to \mathbb{R}^n$  is L-lipschitz continuous
- ▶ Weak Minty condition holds: i.e.  $\exists$  nonempty  $S^* \in \operatorname{zer} T$ , s.t.  $\forall z^* \in S^*$  and  $\rho \in (-1/2L, +\infty)$

$$\langle v, z - z^* \rangle \ge \rho ||v||^2, \quad \forall (z, v) \in \mathbf{gph}(T)$$

#### Generalization of EG+

## Theorem (Convergence to stationarity under weak MVI)

Let  $\lambda_t \in (0,2)$ ,  $\gamma_t \in \left(\lfloor -2\rho \rfloor_+, 1/L \right]$  where  $\lfloor x \rfloor_+ \coloneqq \max\{0,x\}$ ,  $\delta_t \in (-\gamma_t/2,\rho]$ ,  $\liminf_{k \to \infty} \lambda_k (2-\lambda_k) > 0$ , and  $\liminf_{k \to \infty} (\delta_t + \gamma_t/2) > 0$ . Consider the sequences  $(z^t)_{k \in \mathbb{N}}$ ,  $(\bar{z}^t)_{k \in \mathbb{N}}$  generated by  $\left(\operatorname{id} + \gamma_t A \right)^{-1} \left(z^t - \gamma_t F z^t \right)$ . Then for all  $z^* \in \mathcal{S}^*$ , One can establish the following convergence result:

$$\min_{k=0,\dots,m} \frac{1}{\gamma_t^2} ||H\bar{z}^t - Hz^t||^2 \le \frac{1}{\kappa(m+1)} ||z^0 - z^\star||^2$$
 (8)

where  $\kappa = \liminf_{k \to \infty} \lambda_k (2 - \lambda_k) (\delta_k + \gamma_k/2)^2$ .

#### Remark:

- $\circ \ \mathbf{gph}(A) \coloneqq \left\{ (x, y) \in \mathbb{R}^n \times \mathbb{R}^d \mid y \in Ax \right\}$
- $\circ$  LHS denotes the gap between two updates:  $z^{t+1} = z^t + \lambda_t \alpha_t (H\bar{z}^t Hz^t)$

#### **Proof ideas**

Construct a half space that contains the solution set:

$$\mathcal{D}(z) = \left\{ w \mid \langle \bar{v}, \bar{z} - w \rangle \ge \rho ||\bar{v}||^2 \right\} \tag{9}$$

onto which we could iteratively project:  $\Pi_{\mathcal{D}}(z^t)$ 

- ▶ To evaluate  $\bar{v}$ , use  $\bar{v} \in T(\bar{z}^t)$ . As  $H = \mathrm{id} \gamma F$ , the update of  $\bar{z}$  provides a way to evaluate T (show in the next slide)
- ▶ The algorithm update can be casted as KM iteration:

$$z^{t+1} = (1 - \lambda_t) z^t + \lambda_t \mathbf{\Pi}_{\mathcal{D}^t}(z^t)$$
with  $\mathbf{\Pi}_{\mathcal{D}^t}(z^t) = z^t + \alpha_t (H\bar{z}^t - Hz^t)$  (10)

Use best iterate convergence result of KM iterate

### **Proof I – Preparation**

▶ To use weak-MVI, we need to verify  $\left(\bar{z}^t, 1/\gamma_t(Hz^t - H\bar{z}^t)\right) \in \mathbf{gph}(T)$ As  $H = \mathrm{id} - \gamma_t F$  and  $\bar{z}^t = \left(\mathrm{id} + \gamma_t A\right)^{-1} \left(z^t - \gamma_t F z^t\right)$ , we have

$$\bar{z}^t + \gamma_t A \bar{z}^t = H z^t \tag{11}$$

$$\bar{z}^t - \gamma_t F \bar{z}^t + \gamma_t A \bar{z}^t = H z^t - \gamma_t F \bar{z}^t \tag{12}$$

$$\frac{1}{\gamma_t}(Hz^t - H\bar{z}^t) \in A\bar{z}^t + F\bar{z}^t = T\bar{z}^t \tag{13}$$

▶ Verify H is 1/2-cocoercive: refer to slide 11

## **Proof II – Construct** a half-space containing solution set

Following Lectue 3: each step of EG is a projection onto a particular hyperplane. We are aiming at the following construction:

$$\bar{z}^t = \left(\operatorname{id} + \gamma_t A\right)^{-1} \left(z^t - \gamma_t F z^t\right) \tag{14}$$

$$z^{t+1} = (1 - \lambda_t)z^t + \lambda_t \mathbf{\Pi}_{\mathcal{D}_t}(z^t)$$
(15)

where  $\mathcal{D}_t$  is constructed to *contain the solution set*. As  $\frac{1}{\gamma_t}(Hz^t - H\bar{z}^t) \in gph(T)$ , we have

$$\left\langle \frac{1}{\gamma_t} (Hz^t - H\bar{z}^t), \bar{z}^t - z^\star \right\rangle \overset{\text{Weak-MVI}}{\geq} \frac{\rho}{\gamma_t^2} \|Hz^t - H\bar{z}^t\|^2 \overset{\rho \geq \delta_t}{\geq} \frac{\delta_t}{\gamma_t^2} \|Hz^t - H\bar{z}^t\|^2 \tag{16}$$

Thus we can construct  $\mathcal{D}_t$  as

$$\mathcal{D}_{t} \coloneqq \left\{ w \mid \langle Hz^{t} - H\bar{z}^{t}, \bar{z}^{t} - w \rangle \ge \frac{\delta_{t}}{\gamma_{t}} \|Hz^{t} - H\bar{z}^{t}\|^{2} \right\}$$

$$\tag{17}$$

Note due to the 1/2-cocoercivity of H, our chosen step size  $\alpha_t$  is positive and bounded away from 0:

$$\alpha_t = \frac{\delta_t}{\gamma_t} + \frac{\langle \bar{z}^t - z^t, H\bar{z}^t - Hz^t \rangle}{\|H\bar{z}^t - Hz^t\|^2} \ge \frac{1}{2} + \frac{\delta_t}{\gamma_t}$$
(18)

## Proof III - Projection onto the constructed half space

#### Lemma

The projection  $\Pi_{\mathcal{D}}(x) := \arg\min_{z \in \mathcal{D}} \frac{1}{2} \|z - x\|^2$  onto the set  $\mathcal{D} = \{z \mid \langle a, z \rangle \geq b\}$  is given for  $x \notin \mathcal{D}$  as,

$$\Pi_{\mathcal{D}}(x) = x - \frac{\langle a, x \rangle - b}{\|a\|^2} a. \tag{19}$$

The projection onto  $\mathcal{D}_t$  for any  $v \notin \mathcal{D}_t$  is given by

$$\Pi_{\mathcal{D}_t}(v) = v + \frac{\langle \bar{z}^t - v, Hz^t - H\bar{z}^t \rangle - \frac{\delta_t}{\gamma_t} ||Hz^t - H\bar{z}^t||^2}{||Hz^t - H\bar{z}^t||^2} (Hz^t - H\bar{z}^t)$$

For  $v = z^t$ , we have

$$\mathbf{\Pi}_{\mathcal{D}_t}(z^t) = z^t + \alpha_t (H\bar{z}^t - Hz^t) \tag{20}$$

We can thus rewrite an update step as

$$z^{t+1} = (1 - \lambda_t)z^t + \lambda_t \left(z^t - \alpha(Hz^t - H\bar{z}^t)\right)$$
(21)

#### Proof IV - Formulate as KM iteration

### Krasnosel'skii-Mann (KM) iteration

Let  $S:\mathbb{R}^d \to \mathbb{R}^d$  be an operator and  $\lambda>0$ . The KM iteration is given by

$$z^{t+1} = (1 - \lambda)z^t + \lambda Sz^t \tag{KM}$$

#### Remarks

lacktriangle An operator  $S:\mathbb{R}^d o \mathbb{R}^d$  is said to be firmly nonexpansive if

$$||Sz - Sz'||^2 + ||(\mathrm{id} - S)z - (\mathrm{id} - S)z'||^2 \le ||z - z'||^2 \quad \forall z, z' \in \mathbb{R}^d.$$
 (22)

Let C be a nonempty closed convex subset. Then the projector  $\Pi_C$  is firmly nonexpansive.

As  $\mathcal{D}_t$  is closed convex set,  $\Pi_{\mathcal{D}_t}$  is firmly non-expansive

### Proof V - Best iterate result of KM

## Theorem (Best iterate of KM - Recall Lecture 3)

Suppose  $S: \mathbb{R}^d \to \mathbb{R}^d$  is firmly nonexpansive. Consider the sequence  $(z^t)_{t \in \mathbb{N}}$  generated by  $\ref{eq:sequence}$  with  $\lambda \in (0,2)$ . Then, for all  $z^\star \in \operatorname{fix} S$ , it holds that

$$\min_{t \in \{0, \dots, T-1\}} \|Sz^t - z^t\|^2 \le \frac{\|z^0 - z^\star\|^2}{\lambda(2-\lambda)T}.$$
 (23)

Directly apply the theorem, we have

$$\min_{t=0,\dots,m} \|\mathbf{\Pi}_{\mathcal{D}_t} z^t - z^t\|^2 \le \frac{\|z^0 - z^\star\|^2}{\lambda_t (2 - \lambda_t)(m+1)}$$
(24)

$$\|\mathbf{\Pi}_{\mathcal{D}_t} z^t - z^t\|^2 = \alpha_t^2 \|H\bar{z}^t - Hz^t\|^2 \ge \left(\frac{1}{2} + \frac{\delta_t}{\gamma_t}\right)^2 \|H\bar{z}^t - Hz^t\|^2 \tag{25}$$

Choose  $\kappa=\liminf_{t\to\infty}\lambda_t(2-\lambda_t)(\frac{\gamma_t}{2}+\delta_t)^2$ , we directly have

$$\min_{t=0,\dots,m} \frac{1}{\gamma_t^2} ||Hz^t - H\bar{z}^t||^2 \le \frac{1}{\kappa(1+m)} ||z^0 - z^\star||^2$$
 (26)

### Additional results on the last iterate – $\bar{z}^t$

Limit points of  $\bar{z}^t$  belong to zer T.

### Claim 1

Following the same assumptions as the main theorem,  $(\bar{z}^t)_{k\in\mathbb{N}}$  is bounded and its limit points belong to  $\operatorname{zer} T$ ;

#### Proof Ideas.

▶ Since  $\liminf_{t\to\infty} \varepsilon_t \coloneqq \lambda_t (2-\lambda_t)(\frac{\gamma_t}{2}+\delta_t)^2 > 0$ , as

$$\|z^{t+1} - z^{\star}\|^2 \leq \|z^t - z^{\star}\|^2 - \frac{\varepsilon_t}{\gamma_t^2} \|Hz^t - H\bar{z}^t\|^2$$

$$(rac{1}{\gamma_t^2}\|Har{z}^t-Hz^t\|^2)_{k\in\mathbb{N}}$$
 converges to zero.

- ▶  $\|z^t z^*\|^2$  converges, we have  $z^t$  bounded.  $\gamma_t$  is bounded, F and the resolvents  $(\mathrm{id} + \gamma_t A)^{-1}$  are Lipschitz continuous, so is their composition. Thus,  $\bar{z}^t$  is also bounded.
- $\triangleright$   $(\bar{z}^t, 1/\gamma_t(Hz^t H\bar{z}^t)) \in gph(T)$ , that is

$$1/\gamma_t \left( Hz^t - H\bar{z}^t \right) \in T\bar{z}^t$$

As  $\lim_{t\to\infty} 1/\gamma_t \left(Hz^t - H\bar{z}^t\right) = 0$ , we have  $0 \in T\bar{z}^t$ 

### Additional results on the last iterate – $z^t$

Limit points of  $z^t$  stay close to  $\bar{z}^t$ .

### Claim 2

Following the same assumptions as the main theorem, if in addition  $\limsup_{k\to\infty}\gamma_t<1/L$  and  $\mathcal{S}^\star=\operatorname{zer} T$ , then  $(z^t)_{k\in\mathbb{N}},\,(\bar{z}^t)_{k\in\mathbb{N}}$  both converge to some  $z^\star\in\operatorname{zer} T$ .

#### Lemma

If A is L-Lipschitz, then  $T=\operatorname{id}-\eta A$ ,  $\eta\in(0,1/L)$ , is  $(1-\eta L)$ -monotone, and in particular  $\|Tu-Tv\|\geq (1-\eta L)\|u-v\|$  for all  $u,v\in\mathbb{R}^n$ .

• if  $\gamma = \limsup_{k \to \infty} \gamma_t < 1/L$ , following Lemma ??, we have

$$(1 - \gamma L) \|\bar{z}^t - z^t\| \le \|H\bar{z}^t - Hz^t\|$$

Therefore,  $(\|\bar{z}^t - z^t\|)_{k \in \mathbb{N}}$  converges to zero. Hence,  $(z^t)_{k \in \mathbb{N}} [k \in K]$  also converges to  $\bar{z} \in \operatorname{zer} T$ .

## Connection to other existent algorithms

### A constant step size variant

$$\bar{z}^t = (\mathrm{id} + \gamma_k A)^{-1} (z^t - \gamma F z^t), \quad z^{t+1} = z^t - \bar{\alpha}_t (H \bar{z}^t - H z^t). \tag{CEG+}$$

According to the main theorem, convergence is guaranteed when  $\lambda_t \in (0,2)$ , that is when  $\bar{\alpha}_t < 2\alpha_t$ . We could choose constant stepsize  $\bar{\alpha} \in \left(0,1+\frac{2\delta}{\gamma}\right)$ 

$$\bar{\alpha}_t < 1 + \frac{2\delta_t}{\gamma_t} \le \frac{2\delta_t}{\gamma_t} + \frac{2\langle \bar{z}^t - z^t, H\bar{z}^t - Hz^t \rangle}{\|H\bar{z}^t - Hz^t\|^2} = 2\alpha_t \tag{27}$$

- ▶ Larger second stepsize than EG+. As the setting of EG+, in the unconstrained case, i.e.  $A \equiv 0$ , for the smallest  $\rho$  that EG+ allows ( $\rho = -1/8L$ ), EG+ chooses  $\gamma_t = 1/L$  and  $\alpha_t = 1/2$ . While we can select  $\gamma_t = 1/L$  and  $\alpha_t \in (0,3/4)$  here.
- ▶ Connection to FBF. In the monotone case. i.e.  $\rho=0$ , choose  $\gamma_t\in(0,1/L)$  and  $\alpha=1$ , we have  $z^{t+1}=\bar{z}^t+\gamma Fz^t-\gamma F\bar{z}^t$ , which is the forward-backward-forward algorithm proposed by ?

## Adaptive Step Size according to Local Curvature

- lacktriangle A larger stepsize  $\gamma_t$  would guarantee global convergence to a wider class of problem, as  $ho > -rac{\gamma_t}{2}$
- ► Global Lipschitz constant is inherently pessimistic use local curvature instead
- $\circ$  Adaptive in  $\gamma_t$

#### Algorithm Lipschitz constant backtracking

$$\begin{array}{ll} \mbox{Initialize} \ z^t \in \mathbb{R}^n, \ \tau \in (0,1), \ \nu \in (0,1) \\ \mbox{Set initial guess} \ \gamma \ = \ \nu \|JF(z^t)\|^{-1}, \ \mbox{and let} \ G_{\gamma}(z^t) \ \coloneqq \\ \left(\mbox{id} + \gamma A\right)^{-1} \left(z^t - \gamma F z^t\right) \\ \mbox{while} \ \gamma \|F(G_{\gamma}(z^t)) - F z^t\| > \nu \|G_{\gamma}(z^t) - z^t\| \ \mbox{do} \ \gamma \leftarrow \tau \gamma \\ \mbox{Return} \ \gamma_t = \gamma \ \mbox{and} \ \bar{z}^t = G_{\gamma}(z^t) \end{array}$$

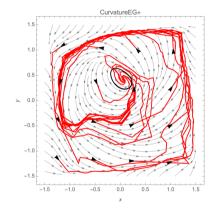


Figure: CurvatureEG+ is able to escape the attractive limit cycle in the Forsaken game

### Overview of the stochastic part

- o In the last lecture,
  - Extra-gradient+ algorithms with adaptive stepsize;
  - Convergence result for deterministic case.
- o This lecture,
  - Extra-gradient+ algorithm in the stochastic setting: Bias-corrected algorithm;
  - Prove the convergence of the bias-corrected algorithm in unconstrained case;
  - Extend the algorithm to constrained case;
  - Extend the algorithm to the primal-dual setting with asynchronous update;
- Thomas Pethick, Olivier Fercoq, Puya Latafat, Panagiotis Patrinos, and Volkan Cevher. Solving Stochastic Weak Minty Variational Inequalities Without Increasing Batch Size. In The Eleventh International Conference on Learning Representations, 2023. URL https://openreview.net/forum?id=ejR4E1jaH9k.

Find 
$$z \in \mathbb{R}^n$$
, such that

$$0 \in Tz \coloneqq Az + Fz,\tag{28}$$

where A maps to a set (captures the constraints) and F maps to a single vector (captures the objective).

Find 
$$z \in \mathbb{R}^n$$
, such that

$$0 \in Tz := Az + Fz, \tag{28}$$

where A maps to a set (captures the constraints) and F maps to a single vector (captures the objective).

o Limitation of the last lecture: Need exact deterministic evaluation of Fz.

Find 
$$z \in \mathbb{R}^n$$
, such that

$$0 \in Tz := Az + Fz, \tag{28}$$

where A maps to a set (captures the constraints) and F maps to a single vector (captures the objective).

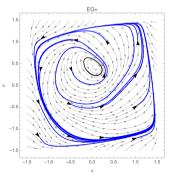
- $\circ$  Limitation of the last lecture: Need exact deterministic evaluation of Fz.
- Such deterministic evaluation can be expensive or even unavailable.
  - ► E.g., the gradient of training a deep neural network.
- o EG+: needs large enough exploration stepsize. Otherwise, may get stuck in a limit cycle.

Find 
$$z \in \mathbb{R}^n$$
, such that

$$0 \in Tz := Az + Fz, \tag{28}$$

where A maps to a set (captures the constraints) and F maps to a single vector (captures the objective).

- $\circ$  Limitation of the last lecture: Need exact deterministic evaluation of Fz.
- o Such deterministic evaluation can be expensive or even unavailable.
  - E.g., the gradient of training a deep neural network.
- o EG+: needs large enough exploration stepsize. Otherwise, may get stuck in a limit cycle.



 $\circ$  Instead, the evaluation of Fz can be **stochastic**.

- $\circ$  Instead, the evaluation of Fz can be **stochastic**.
- $\circ$  We can not directly evaluate the deterministic Fz, but get access to a stochastic oracle  $\hat{F}(z,\xi)$ .

- $\circ$  Instead, the evaluation of Fz can be **stochastic**.
- $\circ$  We can not directly evaluate the deterministic Fz, but get access to a stochastic oracle  $\hat{F}(z,\xi)$ .
- $\circ\,$  E.g., Stochastic Gradient Descent (SGD) for training deep neural network.

  - At step t,  $x^{t+1} = x^t \eta \nabla f_{i_t}(x^t)$ , where  $i_t$  is randomly sampled from [N].
  - $\qquad \qquad \blacksquare \text{ In this example, } z\coloneqq x\text{, } \xi\coloneqq i_t\text{, } \hat{F}(z,\xi)\coloneqq \nabla f_{i_t}(x^t).$

- $\circ$  Instead, the evaluation of Fz can be **stochastic**.
- $\circ$  We can not directly evaluate the deterministic Fz, but get access to a stochastic oracle  $\hat{F}(z,\xi)$ .
- o E.g., Stochastic Gradient Descent (SGD) for training deep neural network.

  - At step t,  $x^{t+1} = x^t \eta \nabla f_{i_t}(x^t)$ , where  $i_t$  is randomly sampled from [N].
  - $\qquad \qquad \blacksquare \text{ In this example, } z\coloneqq x\text{, } \xi\coloneqq i_t\text{, } \hat{F}(z,\xi)\coloneqq \nabla f_{i_t}(x^t).$

# Assumption (Assumptions on $\hat{F}(z,\xi)$ )

For the operator  $\hat{F}(\cdot,\xi):\mathbb{R}^n\to\mathbb{R}^n$ , the following holds.

- 1. Unbiased:  $\mathbb{E}_{\xi}[\hat{F}(z,\xi)] = Fz \quad \forall z \in \mathbb{R}^n$ .
- 2. Bounded variance:  $\mathbb{E}_{\xi}\left[\|\hat{F}(z,\xi) F(z)\|^2\right] \leq \sigma_F^2 \quad \forall z \in \mathbb{R}^n$ .

### Stochastic oracle

- $\circ$  Instead, the evaluation of Fz can be **stochastic**.
- $\circ$  We can not directly evaluate the deterministic Fz, but get access to a stochastic oracle  $\hat{F}(z,\xi)$ .
- o E.g., Stochastic Gradient Descent (SGD) for training deep neural network.

  - At step t,  $x^{t+1} = x^t \eta \nabla f_{i_t}(x^t)$ , where  $i_t$  is randomly sampled from [N].
  - $\qquad \qquad \blacksquare \text{ In this example, } z\coloneqq x\text{, } \xi\coloneqq i_t\text{, } \hat{F}(z,\xi)\coloneqq \nabla f_{i_t}(x^t).$

# Assumption (Assumptions on $\hat{F}(z,\xi)$ )

For the operator  $\hat{F}(\cdot,\xi):\mathbb{R}^n\to\mathbb{R}^n$ , the following holds.

- 1. Unbiased:  $\mathbb{E}_{\xi}[\hat{F}(z,\xi)] = Fz \quad \forall z \in \mathbb{R}^n$ .
- 2. Bounded variance:  $\mathbb{E}_{\xi} \left[ \| \hat{F}(z,\xi) F(z) \|^2 \right] \leq \sigma_F^2 \quad \forall z \in \mathbb{R}^n$ .

### Remarks:

- $\circ$  Assumption 1 holds for SGD when  $i_t$  is uniformly sampled.
- $\circ$  Assumption 2 is common. Easily satisfied when z restricted to a compact set.

- o Stochastic Gradient Descent (SGD).
  - $ightharpoonup \min f(x) \coloneqq \frac{1}{N} \sum_{i=1}^{N} f_i(x)$
  - At step t,  $x^{t+1} = x^t \eta \nabla f_{i_t}(x^t)$ , where  $i_t$  is randomly sampled from [N].
- $\quad \text{O Inbiased: } \mathbb{E}[\nabla f_{i_t}(x^t)] = \tfrac{1}{N} \sum_{i=1}^N \nabla f_i(x) = \nabla f(x^t).$
- $\quad \text{o} \quad \text{Bounded variance: Assume } \mathbb{E}[\|\nabla f_{i_t}(x^t) \nabla f(x^t)\|^2] \leq \sigma_f^2.$

- $\circ x^* \in \arg\min_x f(x).$
- $\circ g_t \coloneqq \nabla f_{i_t}(x^t).$

- $\circ x^* \in \arg\min_x f(x).$
- $\circ \ g_t \coloneqq \nabla f_{i_t}(x^t).$
- o One-step iterate of SGD.

$$||x^{t+1} - x^*||^2$$
  
=  $||x^t - x^* - \eta g_t||^2$ 

- $\circ x^* \in \arg\min_x f(x).$
- $\circ \ g_t \coloneqq \nabla f_{i_t}(x^t).$
- o One-step iterate of SGD.

$$\begin{split} & \|x^{t+1} - x^{\star}\|^2 \\ = & \|x^t - x^{\star} - \eta g_t\|^2 \\ = & \|x^t - \eta g_t\|^2 - 2\eta \langle x^t - x^{\star}, g_t \rangle + \eta^2 \|g_t\|^2 \\ = & \|x^t - x^{\star}\|^2 - 2\eta \langle x^t - x^{\star}, \nabla f(x^t) \rangle - 2\eta \langle x^t - x^{\star}, g_t - \nabla f(x^t) \rangle + \eta^2 \|\nabla f(x^t) + g_t - \nabla f(x^t)\|^2 \end{split}$$

- $\circ x^* \in \arg\min_x f(x).$
- $\circ \ g_t \coloneqq \nabla f_{i_t}(x^t).$
- o One-step iterate of SGD.

## Recall: assumptions

# Assumption (Assumptions on F and A)

1. The operator  $F: \mathbb{R}^n \to \mathbb{R}^n$  is  $L_F$ -Lipschitz with  $L_F \in [0, \infty)$ , i.e.,

$$||Fz - Fz'|| \le L_F ||z - z'|| \quad \forall z, z' \in \mathbb{R}^n.$$

- 2. The operator  $A: \mathbb{R}^n \rightrightarrows \mathbb{R}^n$  is a maximally monotone operator.
- 3. Weak Minty variational inequality (MVI) holds, i.e., there exists a nonempty set  $\mathcal{S}^\star \subseteq \operatorname{zer} T$  such that for all  $z^\star \in \mathcal{S}^\star$  and some  $\rho \in \left(-\frac{1}{2L_F}, \infty\right)$

$$\langle v, z - z^* \rangle \ge \rho \|v\|^2$$
, for all  $(z, v) \in \mathbf{gph}T$ .

# Further assumptions on $\hat{F}(z,\xi)$

## Assumption

- ▶ **Two-point oracle**: The stochastic oracle can be queried for any two points  $z, z' \in \mathbb{R}^n$ ,  $\hat{F}(z, \xi)$ ,  $\hat{F}(z', \xi)$  where  $\xi \sim \mathcal{P}$ .
- **Lipschitz continuity**: The operator  $\hat{F}(\cdot,\xi):\mathbb{R}^n \to \mathbb{R}^n$  is Lipschitz continuous in mean with  $L_{\hat{F}} \in [0,\infty)$ :

$$\mathbb{E}_{\xi} \left[ \left\| \hat{F}(z,\xi) - \hat{F}\left(z',\xi\right) \right\|^{2} \right] \leq L_{\hat{F}}^{2} \left\| z - z' \right\|^{2}$$

for all  $z, z' \in \mathbb{R}^n$ .

# Further assumptions on $\hat{F}(z,\xi)$

## Assumption

- ▶ **Two-point oracle**: The stochastic oracle can be queried for any two points  $z, z' \in \mathbb{R}^n$ ,  $\hat{F}(z, \xi)$ ,  $\hat{F}(z', \xi)$  where  $\xi \sim \mathcal{P}$ .
- ightharpoonup Lipschitz continuity: The operator  $\hat{F}(\cdot,\xi):\mathbb{R}^n o \mathbb{R}^n$  is Lipschitz continuous in mean with  $L_{\hat{F}} \in [0,\infty)$ :

$$\mathbb{E}_{\xi}\left[\left\|\hat{F}(z,\xi) - \hat{F}\left(z',\xi\right)\right\|^{2}\right] \leqslant L_{\hat{F}}^{2}\left\|z - z'\right\|^{2}$$

for all  $z, z' \in \mathbb{R}^n$ .

#### Remarks:

- $\circ$  **Two-point oracle** assumption is reasonable for SGD by choosing the same  $i_t$ .
- o Will see: Lipschitz continuity critical for algorithm design and variance reduction.

### Oblivious extension of EG+?

 $\circ$  Extension of EG+ with  $A \equiv 0$  to the stochastic setting:

$$\bar{z}^{t} = z^{t} - \gamma \frac{1}{B} \sum_{i=1}^{B} \hat{F}\left(z^{t}, \xi_{t,i}\right), \quad z^{t+1} = z^{t} - \alpha_{t} \gamma \frac{1}{B} \sum_{i=1}^{B} \hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t,i}\right). \tag{EG+}$$

### Oblivious extension of EG+?

 $\circ$  Extension of EG+ with  $A\equiv 0$  to the stochastic setting:

$$\bar{z}^t = z^t - \gamma \frac{1}{B} \sum_{i=1}^B \hat{F}\left(z^t, \xi_{t,i}\right), \quad z^{t+1} = z^t - \alpha_t \gamma \frac{1}{B} \sum_{i=1}^B \hat{F}\left(\bar{z}^t, \bar{\xi}_{t,i}\right). \tag{EG+}$$

- $\circ \ \ \mathsf{Let} \ \tilde{\bar{z}}^t \coloneqq z^t \gamma F z^t \ \ \mathsf{and} \ \ \tilde{z}^{t+1} \coloneqq z^t \alpha_t \gamma F \tilde{\bar{z}}^t.$
- $\circ~$  Due to the unbiased assumption of  $\hat{F},$  we have  $\mathbb{E}_{\xi_t}(\bar{z}^t)=\tilde{\bar{z}}^t.$

### Oblivious extension of EG+?

 $\circ\;$  Extension of EG+ with  $A\equiv 0$  to the stochastic setting:

$$\bar{z}^{t} = z^{t} - \gamma \frac{1}{B} \sum_{i=1}^{B} \hat{F}\left(z^{t}, \xi_{t,i}\right), \quad z^{t+1} = z^{t} - \alpha_{t} \gamma \frac{1}{B} \sum_{i=1}^{B} \hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t,i}\right). \tag{EG+}$$

- $\circ \ \ \mathsf{Let} \ \tilde{\bar{z}}^t \coloneqq z^t \gamma F z^t \ \ \mathsf{and} \ \ \tilde{z}^{t+1} \coloneqq z^t \alpha_t \gamma F \tilde{\bar{z}}^t.$
- $\circ$  Due to the unbiased assumption of  $\hat{F}$ , we have  $\mathbb{E}_{\mathcal{E}_t}(ar{z}^t) = \tilde{z}^t$ .

### Observation

When F is linear, we have,

$$\mathbb{E}_{\xi_t,\bar{\xi}_t}(z^{t+1}) = z^t - \alpha_t \gamma \mathbb{E}_{\xi_t} \left( \frac{1}{B} \sum_{i=1}^B \mathbb{E}_{\bar{\xi}_t \mid \xi_t} \hat{F}\left(\bar{z}^t,\bar{\xi}_{t,i}\right) \right) = z^t - \alpha_t \gamma \mathbb{E}_{\xi_t} F \bar{z}^t = z^t - \alpha_t \gamma F(\mathbb{E}_{\xi_t} \bar{z}^t) = \bar{z}^{t+1}$$

$$\text{(Linear } F)$$

**Remark:** The linearity of F leads to unbiased update of z and convergence result when  $A \equiv 0$ .

### What if F is nonlinear or $A \not\equiv 0$ ?

### Observation

In general,

$$\mathbb{E}_{\xi_t,\bar{\xi_t}}(z^{t+1}) = z^t - \alpha_t \gamma \mathbb{E}_{\xi_t} F \bar{z}^t \neq z^t - \alpha_t \gamma F (\mathbb{E}_{\xi_t} \bar{z}^t) = \tilde{z}^{t+1} \tag{Nonlinear } F)$$

 $\circ$  Bias appears with nonlinear F.

### What if F is nonlinear or $A \not\equiv 0$ ?

### Observation

In general,

$$\mathbb{E}_{\xi_t,\bar{\xi_t}}(z^{t+1}) = z^t - \alpha_t \gamma \mathbb{E}_{\xi_t} F \bar{z}^t \neq z^t - \alpha_t \gamma F(\mathbb{E}_{\xi_t} \bar{z}^t) = \tilde{z}^{t+1}$$
 (Nonlinear  $F$ )

- $\circ$  Bias appears with nonlinear F.
- $\circ$  How to reduce the bias in  $z^{t+1}$ ?
  - lacktriangle Ideally, the above equality holds with deterministic extrapolation point  $ar{z}^t$ .
  - ▶ The smaller the variance of  $\bar{z}^t$ , the better.

### What if F is nonlinear or $A \not\equiv 0$ ?

#### Observation

In general,

$$\mathbb{E}_{\xi_t,\bar{\xi}_t}(z^{t+1}) = z^t - \alpha_t \gamma \mathbb{E}_{\xi_t} F \bar{z}^t \neq z^t - \alpha_t \gamma F (\mathbb{E}_{\xi_t} \bar{z}^t) = \bar{z}^{t+1}$$
 (Nonlinear  $F$ )

- $\circ$  Bias appears with nonlinear F.
- $\circ$  How to reduce the bias in  $z^{t+1}$ ?
  - ldeally, the above equality holds with deterministic extrapolation point  $\bar{z}^t$ .
  - ▶ The smaller the variance of  $\bar{z}^t$ , the better.
- Two obvious approaches:
  - Increasing batchsize B. But can be too expensive. [? , Thm. 4.5]
  - lacktriangle Reduce the stepsize  $\gamma$ . But contradicts the requirement of large enough stepsize  $\gamma$  for the convergence in weak MVI. [? ]

## **Bias-corrected Approach**

o Idea: leveraging two-point oracle.

## Algorithm (BC-SEG+) Stochastic algorithm for problem (??) when $A\equiv 0$

- 1: REQUIRE  $z^{-1} = \bar{z}^{-1} = z^0 \in \mathbb{R}^n, \alpha_t \in (0,1), \gamma \in (\lfloor -2\rho \rfloor_+, 1/L_F)$
- 2: for  $t=0,1,\ldots$  until convergence do
- 3: Sample  $\xi_t \sim \mathcal{P}$

$$ar{z}^t = z^t - \gamma \hat{F}\left(z^t, \xi_t
ight) + (1 - lpha_t) \left(ar{z}^{t-1} - z^{t-1} + \gamma \hat{F}\left(z^{t-1}, \xi_t
ight)
ight)$$

- 5: Sample  $ar{\xi}_t \sim \mathcal{P}$
- 6:  $z^{t+1} = z^t \alpha_t \gamma \hat{F}\left(\bar{z}^t, \bar{\xi}_t\right)$
- 7: end for
- 8: RETURN  $z^{t+1}$

## **Bias-corrected Approach**

o Idea: leveraging two-point oracle.

## **Algorithm** (BC-SEG+) Stochastic algorithm for problem (??) when $A\equiv 0$

```
1: REQUIRE z^{-1}=\bar{z}^{-1}=z^0\in\mathbb{R}^n, \alpha_t\in(0,1), \gamma\in(\lfloor-2\rho\rfloor_+,1/L_F)

2: for t=0,1,\ldots until convergence do

3: Sample \xi_t\sim\mathcal{P}

4: \bar{z}^t=z^t-\gamma\hat{F}\left(z^t,\xi_t\right)+(1-\alpha_t)\left(\bar{z}^{t-1}-z^{t-1}+\gamma\hat{F}\left(z^{t-1},\xi_t\right)\right)

5: Sample \bar{\xi}_t\sim\mathcal{P}

6: z^{t+1}=z^t-\alpha_t\gamma\hat{F}\left(\bar{z}^t,\bar{\xi}_t\right)

7: end for

8: RETURN z^{t+1}
```

- Line 4 is the key step.
- o Idea: add a 'correction' term to reduce the variance of  $\bar{z}^t$ . (Similar idea in STORM algorithm [?].)
- $\circ$  Correction term = Real Preceding Extrapolation Point Imagined Stochastic Extrapolation Point with  $\xi_t$

o Approximate error of the stochastic extrapolation point.

$$u^{t} := \bar{z}^{t} - \left(z^{t} - \gamma F\left(z^{t}\right)\right) \tag{29}$$

o Approximate error of the stochastic extrapolation point.

$$u^{t} := \bar{z}^{t} - \left(z^{t} - \gamma F\left(z^{t}\right)\right) \tag{29}$$

 $\circ$  The error  $u^t$  can be decomposed as,

$$\|u^{t}\| = \|\gamma F z^{t} - \gamma \hat{F}\left(z^{t}, \xi_{t}\right) + (1 - \alpha_{t})\left(u^{t-1} - \gamma F z^{t-1} + \gamma \hat{F}\left(z^{t-1}, \xi_{t}\right)\right)\|$$

$$= \|(1 - \alpha_{t})u^{t-1} + \alpha_{t}\gamma\left(F z^{t} - \hat{F}\left(z^{t}, \xi_{t}\right)\right) + (1 - \alpha_{t})\gamma\left(\left(F z^{t} - F z^{t-1}\right) + \left(\hat{F}\left(z^{t-1}, \xi_{t}\right) - \hat{F}\left(z^{t}, \xi_{t}\right)\right)\right)\|$$
(31)

Approximate error of the stochastic extrapolation point.

$$u^{t} := \bar{z}^{t} - \left(z^{t} - \gamma F\left(z^{t}\right)\right) \tag{29}$$

 $\circ$  The error  $u^t$  can be decomposed as,

$$\|u^{t}\| = \|\gamma F z^{t} - \gamma \hat{F}\left(z^{t}, \xi_{t}\right) + (1 - \alpha_{t}) \left(u^{t-1} - \gamma F z^{t-1} + \gamma \hat{F}\left(z^{t-1}, \xi_{t}\right)\right) \|$$

$$= \|(1 - \alpha_{t}) u^{t-1} + \alpha_{t} \gamma \left(F z^{t} - \hat{F}\left(z^{t}, \xi_{t}\right)\right) + (1 - \alpha_{t}) \gamma \left(\left(F z^{t} - F z^{t-1}\right) + \left(\hat{F}\left(z^{t-1}, \xi_{t}\right) - \hat{F}\left(z^{t}, \xi_{t}\right)\right)\right) \|$$

$$\leq (1 - \alpha_{t}) \|u^{t-1}\| + \alpha_{t} \gamma \underbrace{\|F z^{t} - \hat{F}\left(z^{t}, \xi_{t}\right)\|}_{\approx \mathcal{O}(\sigma_{F})} + (1 - \alpha_{t}) \gamma \underbrace{\left(\|F z^{t} - F z^{t-1}\| + \|\hat{F}\left(z^{t-1}, \xi_{t}\right) - \hat{F}\left(z^{t}, \xi_{t}\right)\|\right)}_{\approx \mathcal{O}((L_{F} + L_{\hat{F}})\|z^{t} - z^{t-1}\|)}$$

$$(32)$$

- ightharpoonup Control the second term by  $\alpha_t \to 0$ .
- Control the third term by Lipschitz property.

Approximate error of the stochastic extrapolation point.

$$u^{t} := \bar{z}^{t} - \left(z^{t} - \gamma F\left(z^{t}\right)\right) \tag{29}$$

 $\circ$  The error  $u^t$  can be decomposed as,

$$\|u^{t}\| = \|\gamma F z^{t} - \gamma \hat{F}\left(z^{t}, \xi_{t}\right) + (1 - \alpha_{t})\left(u^{t-1} - \gamma F z^{t-1} + \gamma \hat{F}\left(z^{t-1}, \xi_{t}\right)\right)\|$$

$$= \|(1 - \alpha_{t}) u^{t-1} + \alpha_{t} \gamma\left(F z^{t} - \hat{F}\left(z^{t}, \xi_{t}\right)\right) + (1 - \alpha_{t}) \gamma\left(\left(F z^{t} - F z^{t-1}\right) + \left(\hat{F}\left(z^{t-1}, \xi_{t}\right) - \hat{F}\left(z^{t}, \xi_{t}\right)\right)\right)\|$$

$$\leq (1 - \alpha_{t}) \|u^{t-1}\| + \alpha_{t} \gamma \underbrace{\|F z^{t} - \hat{F}\left(z^{t}, \xi_{t}\right)\|}_{\approx \mathcal{O}(\sigma_{F})} + (1 - \alpha_{t}) \gamma\underbrace{\left(\|F z^{t} - F z^{t-1}\| + \|\hat{F}\left(z^{t-1}, \xi_{t}\right) - \hat{F}\left(z^{t}, \xi_{t}\right)\|\right)}_{\approx \mathcal{O}((L_{F} + L_{\hat{F}})\|z^{t} - z^{t-1}\|)}$$

$$(32)$$

- ightharpoonup Control the second term by  $\alpha_t \to 0$ .
- Control the third term by Lipschitz property.
- $\circ \|u^t\|$  can be controlled.

## Convergence guarantee in unconstrained case

## Theorem (Random Iterate Convergence)

Suppose the three assumptions on  $F,\hat{F},A$  hold. Suppose in addition that  $\gamma\in(\lfloor-2\rho\rfloor_+,1/L_F)$  and  $(\alpha_t)_{t\in[T]}\subset(0,1)$  is a diminishing sequence such that

$$2\gamma L_{\hat{F}} \sqrt{\alpha_0} + \left(1 + \left(\frac{1 + \gamma^2 L_F^2}{1 - \gamma^2 L_F^2} \gamma^2 L_F^2\right) \gamma^2 L_{\hat{F}}^2\right) \alpha_0 \le 1 + \frac{2\rho}{\gamma}$$

Then, the following estimate holds for all  $z^\star \in \mathcal{S}^\star$ 

$$\mathbb{E}\left[\left\|F\left(z^{t_{\star}}\right)\right\|^{2}\right] \leq \frac{\left(1+\eta\gamma^{2}L_{F}^{2}\right)\left\|z^{0}-z^{\star}\right\|^{2}+C\sigma_{F}^{2}\gamma^{2}\sum_{t=0}^{T}\alpha_{t}^{2}}{\mu\sum_{t=0}^{T}\alpha_{t}}$$

where  $C=1+2\eta\left(\left(\gamma^2L_{\hat{F}}^2+1\right)+2\alpha_0\right), \eta=\frac{1}{2}\frac{1+\gamma^2L_F^2}{1-\gamma^2L_F^2}\gamma^2L_F^2+\frac{1}{\gamma L_{\hat{F}}\sqrt{\alpha_0}}, \mu=\gamma^2\left(1-\gamma^2L_F^2\right)/2$  and  $t_\star$  is chosen from  $\{0,1,\ldots,T\}$  according to probability  $\mathcal{P}\left[t_\star=t\right]=\frac{\alpha_t}{\sum_{k=0}^T\alpha_k}$ .

## Convergence guarantee in unconstrained case

## Theorem (Random Iterate Convergence)

Suppose the three assumptions on  $F, \hat{F}, A$  hold. Suppose in addition that  $\gamma \in (\lfloor -2\rho \rfloor_+, 1/L_F)$  and  $(\alpha_t)_{t \in [T]} \subset (0,1)$  is a diminishing sequence such that

$$2\gamma L_{\hat{F}} \sqrt{\alpha_0} + \left(1 + \left(\frac{1 + \gamma^2 L_F^2}{1 - \gamma^2 L_F^2} \gamma^2 L_F^2\right) \gamma^2 L_{\hat{F}}^2\right) \alpha_0 \le 1 + \frac{2\rho}{\gamma}$$

Then, the following estimate holds for all  $z^\star \in \mathcal{S}^\star$ 

$$\mathbb{E}\left[\left\|F\left(z^{t_{\star}}\right)\right\|^{2}\right] \leq \frac{\left(1+\eta\gamma^{2}L_{F}^{2}\right)\left\|z^{0}-z^{\star}\right\|^{2}+C\sigma_{F}^{2}\gamma^{2}\sum_{t=0}^{T}\alpha_{t}^{2}}{\mu\sum_{t=0}^{T}\alpha_{t}}$$

where 
$$C=1+2\eta\left(\left(\gamma^2L_{\hat{F}}^2+1\right)+2\alpha_0\right), \eta=\frac{1}{2}\frac{1+\gamma^2L_F^2}{1-\gamma^2L_F^2}\gamma^2L_F^2+\frac{1}{\gamma L_{\hat{F}}\sqrt{\alpha_0}}, \mu=\gamma^2\left(1-\gamma^2L_F^2\right)/2$$
 and  $t_\star$  is chosen from  $\{0,1,\ldots,T\}$  according to probability  $\mathcal{P}\left[t_\star=t\right]=\frac{\alpha_t}{\sum_{t=0}^T \alpha_t}$ .

 $\circ$  Set  $\alpha_t = \Theta\left(1/\sqrt{t}\right)$ .  $O(\log T/\sqrt{T})$  convergence rate for  $\mathbb{E}[\text{best } ||Fz^t||^2]$  can be achieved.

o Introduce the potential function,

$$\mathcal{U}_{t} := \underbrace{\left\| z^{t} - z^{\star} \right\|^{2}}_{\text{distance to zer}T} + A_{t} \left\| u^{t-1} \right\|^{2} + B_{t} \left\| z^{t} - z^{t-1} \right\|^{2}, \tag{33}$$

where the first term is used in the deterministic case.

o Introduce the potential function,

$$U_{t} := \underbrace{\left\| z^{t} - z^{\star} \right\|^{2}}_{\text{distance to } \mathbf{zer}T} + A_{t} \left\| u^{t-1} \right\|^{2} + B_{t} \left\| z^{t} - z^{t-1} \right\|^{2}, \tag{33}$$

where the first term is used in the deterministic case.

 $\circ$  Goal: show that  $\mathbb{E}[\mathcal{U}_t] \leq \mathbb{E}[\mathcal{U}_{t-1}] - \Omega(\mathbb{E}[\alpha_t || F(z^t)||^2]) - \text{positive coefficient} \cdot \mathbb{E}[|| F(\bar{z}^t)||^2] + \text{small error}$ 

o Introduce the potential function,

$$U_{t} := \underbrace{\left\| z^{t} - z^{\star} \right\|^{2}}_{\text{distance to } \mathbf{zer}T} + A_{t} \left\| u^{t-1} \right\|^{2} + B_{t} \left\| z^{t} - z^{t-1} \right\|^{2}, \tag{33}$$

where the first term is used in the deterministic case.

- $\circ$  Goal: show that  $\mathbb{E}[\mathcal{U}_t] \leq \mathbb{E}[\mathcal{U}_{t-1}] \Omega(\mathbb{E}[\alpha_t \| F(z^t) \|^2]) \text{positive coefficient} \cdot \mathbb{E}[\| F(\bar{z}^t) \|^2] + \text{small error}$
- o Idea: Reduce everything into some terms in the potential function,  $||Fz^t||$  (what we want to bound), or  $F\bar{z}^t$  (which is  $\mathbb{E}_{\bar{\xi}_t}[\hat{F}(\bar{z}^t,\bar{\xi}_t)] = \mathbb{E}_{\bar{\xi}_t}[\frac{z^{t+1}-z^t}{\alpha_t\gamma}]$ ).

o Introduce the potential function,

$$U_{t} := \underbrace{\|z^{t} - z^{\star}\|^{2}}_{\text{distance to } \mathbf{zer}T} + A_{t} \|u^{t-1}\|^{2} + B_{t} \|z^{t} - z^{t-1}\|^{2}, \tag{33}$$

where the first term is used in the deterministic case.

- $\circ$  Goal: show that  $\mathbb{E}[\mathcal{U}_t] \leq \mathbb{E}[\mathcal{U}_{t-1}] \Omega(\mathbb{E}[\alpha_t \| F(z^t) \|^2])$  positive coefficient  $\cdot \mathbb{E}[\| F(\bar{z}^t) \|^2]$  + small error
- o Idea: Reduce everything into some terms in the potential function,  $\|Fz^t\|$  (what we want to bound), or  $F\bar{z}^t$  (which is  $\mathbb{E}_{\bar{\xi}_t}[\hat{F}(\bar{z}^t,\bar{\xi}_t)] = \mathbb{E}_{\bar{\xi}_t}[\frac{z^{t+1}-z^t}{\alpha_t\gamma}]$ ).

## Fenchel-Young Inequality

 $\forall a, b \in \mathbb{R}^n$  and e > 0, we have,

$$2\langle a, b \rangle \le e||a||^2 + \frac{1}{e}||b||^2 \tag{34}$$

$$||a+b||^2 \le (1+e)||a||^2 + (1+\frac{1}{e})||b||^2$$
(35)

o Introduce the potential function,

$$U_{t} := \underbrace{\|z^{t} - z^{\star}\|^{2}}_{\text{distance to } \mathbf{zer}T} + A_{t} \|u^{t-1}\|^{2} + B_{t} \|z^{t} - z^{t-1}\|^{2},$$
(33)

where the first term is used in the deterministic case.

- $\circ$  Goal: show that  $\mathbb{E}[\mathcal{U}_t] \leq \mathbb{E}[\mathcal{U}_{t-1}] \Omega(\mathbb{E}[\alpha_t || F(z^t)||^2])$  positive coefficient  $\cdot \mathbb{E}[||F(\bar{z}^t)||^2] + \text{small error}$
- o Idea: Reduce everything into some terms in the potential function,  $\|Fz^t\|$  (what we want to bound), or  $F\bar{z}^t$  (which is  $\mathbb{E}_{\bar{\xi}_t}[\hat{F}(\bar{z}^t,\bar{\xi}_t)] = \mathbb{E}_{\bar{\xi}_t}[\frac{z^{t+1}-z^t}{\alpha_t\gamma}]$ ).

### Fenchel-Young Inequality

 $\forall a, b \in \mathbb{R}^n$  and e > 0, we have,

$$2\langle a, b \rangle \le e \|a\|^2 + \frac{1}{e} \|b\|^2 \tag{34}$$

$$||a+b||^2 \le (1+e)||a||^2 + (1+\frac{1}{e})||b||^2$$
(35)

### Proof.

$$(1+e)\|a\|^2 + (1+\frac{1}{e})\|b\|^2 - \|a+b\|^2 = e\|a\|^2 + \frac{1}{e}\|b\|^2 - 2\langle a,b\rangle = \|\sqrt{e}a - \frac{1}{\sqrt{e}}b\|^2 \ge 0.$$

o By (??), we can derive

$$\left\|u^{t}\right\|^{2} = \left\|\gamma F\left(z^{t}\right) - \gamma \hat{F}\left(z^{t}, \xi_{t}\right) + (1 - \alpha_{t})\left(\gamma \hat{F}\left(z^{t-1}, \xi_{t}\right) - \gamma F\left(z^{t-1}\right)\right)\right\|^{2} + (1 - \alpha_{t})^{2} \left\|u^{t-1}\right\|^{2} +$$

$$(36)$$

$$2\left(1 - \alpha_{t}\right)\left\langle\bar{z}^{t-1} - z^{t-1} + \gamma F\left(z^{t-1}\right), \underbrace{\gamma\left(F\left(z^{t}\right) - \hat{F}\left(z^{t}, \xi_{t}\right)\right) + (1 - \alpha_{t})\gamma\left(\hat{F}\left(z^{t-1}, \xi_{t}\right) - F\left(z^{t-1}\right)\right)}_{\text{Conditioned on } \mathcal{F}_{t} \text{ (generated by random variables up to } z^{t}\text{), has expectation of zero.}}\right)$$

 $\circ$  Conditioned on  $\mathcal{F}_t$  (generated by random variables up to  $z^t$ )

$$\mathbb{E}\left[\left\|u^{t}\right\|^{2} \mid \mathcal{F}_{t}\right]$$

$$=\gamma^{2}\mathbb{E}\left[\left\|\left(1-\alpha_{t}\right)\left(F\left(z^{t}\right)-\hat{F}\left(z^{t},\xi_{t}\right)+\hat{F}\left(z^{t-1},\xi_{t}\right)-F\left(z^{t-1}\right)\right)+\alpha_{t}\left(F\left(z^{t}\right)-\hat{F}\left(z^{t},\xi_{t}\right)\right)\right\|^{2} \mid \mathcal{F}_{t}\right]$$
(39)

$$+ (1 - \alpha_t)^2 \left\| u^{t-1} \right\|^2 \tag{40}$$

Fenchel-Young 
$$\leq (1 - \alpha_t)^2 \left\| u^{t-1} \right\|^2 + 2 (1 - \alpha_t)^2 \gamma^2 \mathbb{E} \left[ \left\| \hat{F} \left( z^t, \xi_t \right) - \hat{F} \left( z^{t-1}, \xi_t \right) \right\|^2 \mid \mathcal{F}_t \right]$$
 (41)

$$+2\alpha_t^2\gamma^2\mathbb{E}\left[\left\|F\left(z^t\right)-\hat{F}\left(z^t,\xi_t\right)\right\|^2\mid\mathcal{F}_t\right] \tag{42}$$

$$\underset{\text{consided variance}}{\overset{\text{Lipschitz}}{\leq}} (1 - \alpha_t)^2 \left\| u^{t-1} \right\|^2 + 2 (1 - \alpha_t)^2 \gamma^2 L_{\hat{F}}^2 \left\| z^t - z^{t-1} \right\|^2 + 2\alpha_t^2 \gamma^2 \sigma_F^2 \tag{43}$$

o By the step  $z^{t+1} = z^t - \alpha_t \gamma \hat{F}(\bar{z}^t, \bar{\xi}_t)$ , we have

$$\left\|z^{t+1} - z^{\star}\right\|^{2} = \left\|z^{t} - z^{\star}\right\|^{2} - 2\alpha_{t}\gamma\left\langle \hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t}\right), z^{t} - z^{\star}\right\rangle + \alpha_{t}^{2}\gamma^{2}\left\|\hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t}\right)\right\|^{2} \tag{44}$$

o By the step  $z^{t+1} = z^t - \alpha_t \gamma \hat{F}(\bar{z}^t, \bar{\xi}_t)$ , we have

$$\left\|z^{t+1} - z^{\star}\right\|^{2} = \left\|z^{t} - z^{\star}\right\|^{2} - 2\alpha_{t}\gamma\left\langle \hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t}\right), z^{t} - z^{\star}\right\rangle + \alpha_{t}^{2}\gamma^{2}\left\|\hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t}\right)\right\|^{2} \tag{44}$$

 $\circ$  Conditioned on  $\overline{\mathcal{F}}_t$  generated by random variables up to  $\bar{z}^t$ ,

$$\mathbb{E}\left[\left\langle -\gamma \hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t}\right), z^{t} - z^{\star}\right\rangle \mid \overline{\mathcal{F}}_{t}\right] \tag{45}$$

$$\stackrel{\text{unbiased}}{=} -\gamma^{2} \left\langle F\left(\bar{z}^{t}\right), F\left(z^{t}\right) \right\rangle + \gamma \left\langle F\left(\bar{z}^{t}\right), \bar{z}^{t} - z^{t} + \gamma F\left(z^{t}\right) \right\rangle - \gamma \left\langle F\left(\bar{z}^{t}\right), \bar{z}^{t} - z^{\star} \right\rangle \tag{46}$$

o By the step  $z^{t+1} = z^t - \alpha_t \gamma \hat{F}(\bar{z}^t, \bar{\xi}_t)$ , we have

$$\left\|z^{t+1} - z^{\star}\right\|^{2} = \left\|z^{t} - z^{\star}\right\|^{2} - 2\alpha_{t}\gamma\left\langle \hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t}\right), z^{t} - z^{\star}\right\rangle + \alpha_{t}^{2}\gamma^{2}\left\|\hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t}\right)\right\|^{2} \tag{44}$$

 $\circ$  Conditioned on  $\overline{\mathcal{F}}_t$  generated by random variables up to  $\bar{z}^t$ ,

$$\mathbb{E}\left[\left\langle -\gamma \hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t}\right), z^{t} - z^{\star}\right\rangle \mid \overline{\mathcal{F}}_{t}\right] \tag{45}$$

$$\stackrel{\text{unbiased}}{=} -\gamma^{2} \left\langle F\left(\bar{z}^{t}\right), F\left(z^{t}\right) \right\rangle + \gamma \left\langle F\left(\bar{z}^{t}\right), \bar{z}^{t} - z^{t} + \gamma F\left(z^{t}\right) \right\rangle - \gamma \left\langle F\left(\bar{z}^{t}\right), \bar{z}^{t} - z^{\star} \right\rangle \tag{46}$$

Fenchel-Young 
$$\leq \underset{\text{Weak MVI}}{\leq} \gamma^{2} \left( \frac{1}{2} \left\| F\left(\bar{z}^{t}\right) - F\left(z^{t}\right) \right\|^{2} - \frac{1}{2} \left\| F\left(\bar{z}^{t}\right) \right\|^{2} - \frac{1}{2} \left\| F\left(z^{t}\right) \right\|^{2} \right) \tag{47}$$

$$+\frac{\gamma^{2}\varepsilon_{t}}{2}\left\|F\left(\bar{z}^{t}\right)\right\|^{2}+\frac{1}{2\varepsilon_{t}}\left\|\bar{z}^{t}-z^{t}+\gamma F\left(z^{t}\right)\right\|^{2}-\gamma \rho\left\|F\left(\bar{z}^{t}\right)\right\|^{2}$$

$$\tag{48}$$

o By the step  $z^{t+1} = z^t - \alpha_t \gamma \hat{F}(\bar{z}^t, \bar{\xi}_t)$ , we have

$$\left\|z^{t+1} - z^{\star}\right\|^{2} = \left\|z^{t} - z^{\star}\right\|^{2} - 2\alpha_{t}\gamma\left\langle \hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t}\right), z^{t} - z^{\star}\right\rangle + \alpha_{t}^{2}\gamma^{2}\left\|\hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t}\right)\right\|^{2} \tag{44}$$

 $\circ$  Conditioned on  $\overline{\mathcal{F}}_t$  generated by random variables up to  $\bar{z}^t$ ,

$$\mathbb{E}\left[\left\langle -\gamma \hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t}\right), z^{t} - z^{\star}\right\rangle \mid \overline{\mathcal{F}}_{t}\right] \tag{45}$$

$$\stackrel{\text{unbiased}}{=} -\gamma^{2} \left\langle F\left(\bar{z}^{t}\right), F\left(z^{t}\right) \right\rangle + \gamma \left\langle F\left(\bar{z}^{t}\right), \bar{z}^{t} - z^{t} + \gamma F\left(z^{t}\right) \right\rangle - \gamma \left\langle F\left(\bar{z}^{t}\right), \bar{z}^{t} - z^{\star} \right\rangle \tag{46}$$

Fenchel-Young 
$$\stackrel{\leq}{\underset{\text{Weak MVI}}{\leq}} \gamma^{2} \left( \frac{1}{2} \left\| F\left(\bar{z}^{t}\right) - F\left(z^{t}\right) \right\|^{2} - \frac{1}{2} \left\| F\left(\bar{z}^{t}\right) \right\|^{2} - \frac{1}{2} \left\| F\left(z^{t}\right) \right\|^{2} \right)$$
 (47)

$$+\frac{\gamma^{2}\varepsilon_{t}}{2}\left\|F\left(\bar{z}^{t}\right)\right\|^{2}+\frac{1}{2\varepsilon_{t}}\left\|\bar{z}^{t}-z^{t}+\gamma F\left(z^{t}\right)\right\|^{2}-\gamma \rho\left\|F\left(\bar{z}^{t}\right)\right\|^{2}$$
(48)

$$\underset{\text{Fenchel-Young}}{\overset{\text{Lipschitz}}{\leq}} \gamma^{2} L_{F}^{2} \frac{1+b}{2} \left\| u^{t} \right\|^{2} + \frac{1+b^{-1}}{2} \gamma^{4} L_{F}^{2} \left\| F\left(z^{t}\right) \right\|^{2} - \frac{\gamma^{2}}{2} \left\| F\left(\bar{z}^{t}\right) \right\|^{2} \tag{49}$$

$$-\frac{\gamma^{2}}{2}\left\|F\left(z^{t}\right)\right\|^{2}+\frac{\gamma^{2}\varepsilon_{t}}{2}\left\|F\left(\bar{z}^{t}\right)\right\|^{2}+\frac{1}{2\varepsilon_{t}}\left\|u^{t}\right\|^{2}-\gamma\rho\left\|F\left(\bar{z}^{t}\right)\right\|^{2}$$
(50)

o By the step  $z^{t+1} = z^t - \alpha_t \gamma \hat{F}(\bar{z}^t, \bar{\mathcal{E}}_t)$ , we have

$$\|z^{t+1} - z^{\star}\|^{2} = \|z^{t} - z^{\star}\|^{2} - 2\alpha_{t}\gamma \left\langle \hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t}\right), z^{t} - z^{\star}\right\rangle + \alpha_{t}^{2}\gamma^{2} \|\hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t}\right)\|^{2}$$
(44)

 $\circ$  Conditioned on  $\overline{\mathcal{F}}_t$  generated by random variables up to  $\overline{z}^t$ ,

$$\mathbb{E}\left[\left\langle -\gamma \hat{F}\left(\bar{z}^{t}, \bar{\xi}_{t}\right), z^{t} - z^{\star}\right\rangle \mid \overline{\mathcal{F}}_{t}\right] \tag{45}$$

$$\stackrel{\text{unbiased}}{=} -\gamma^{2} \left\langle F\left(\bar{z}^{t}\right), F\left(z^{t}\right) \right\rangle + \gamma \left\langle F\left(\bar{z}^{t}\right), \bar{z}^{t} - z^{t} + \gamma F\left(z^{t}\right) \right\rangle - \gamma \left\langle F\left(\bar{z}^{t}\right), \bar{z}^{t} - z^{\star} \right\rangle \tag{46}$$

Fenchel-Young 
$$\leq \sum_{\text{Weak MVI}} \gamma^2 \left( \frac{1}{2} \left\| F\left(\bar{z}^t\right) - F\left(z^t\right) \right\|^2 - \frac{1}{2} \left\| F\left(\bar{z}^t\right) \right\|^2 - \frac{1}{2} \left\| F\left(z^t\right) \right\|^2 \right)$$

$$+\frac{\gamma^{2}\varepsilon_{t}}{2}\left\|F\left(\bar{z}^{t}\right)\right\|^{2}+\frac{1}{2\varepsilon_{t}}\left\|\bar{z}^{t}-z^{t}+\gamma F\left(z^{t}\right)\right\|^{2}-\gamma \rho\left\|F\left(\bar{z}^{t}\right)\right\|^{2}$$
(48)

$$\underset{\text{enchal Young}}{\text{Lipschitz}} \gamma^{2} L_{F}^{2} \frac{1+b}{2} \left\| u^{t} \right\|^{2} + \frac{1+b^{-1}}{2} \gamma^{4} L_{F}^{2} \left\| F\left(z^{t}\right) \right\|^{2} - \frac{\gamma^{2}}{2} \left\| F\left(\bar{z}^{t}\right) \right\|^{2} \tag{49}$$

$$-\frac{\gamma^{2}}{2}\left\|F\left(z^{t}\right)\right\|^{2}+\frac{\gamma^{2}\varepsilon_{t}}{2}\left\|F\left(\bar{z}^{t}\right)\right\|^{2}+\frac{1}{2\varepsilon_{t}}\left\|u^{t}\right\|^{2}-\gamma\rho\left\|F\left(\bar{z}^{t}\right)\right\|^{2}$$
(50)

$$= \left(\gamma^2 L_F^2 \frac{1+b}{2} + \frac{1}{2\varepsilon_t}\right) \left\|u^t\right\|^2 + \gamma^2 \frac{\gamma^2 L_F^2 \left(1+b^{-1}\right) - 1}{2} \left\|F\left(z^t\right)\right\|^2 + \left(\frac{\gamma^2 \left(\varepsilon_t - 1\right)}{2} - \gamma\rho\right) \left\|F\left(\bar{z}^t\right)\right\|^2.$$

(47)

0

$$\mathbb{E}\left[\left\|\hat{F}\left(\bar{z}^{t},\bar{\xi}_{t}\right)\right\|^{2}\mid\mathcal{F}_{t}\right] = \mathbb{E}\left[\mathbb{E}\left[\left\|\hat{F}\left(\bar{z}^{t},\bar{\xi}_{t}\right) - F(\bar{z}^{t}) + F(\bar{z}^{t})\right\|^{2}\mid\overline{\mathcal{F}}_{t}\right]\mid\mathcal{F}_{t}\right] \stackrel{\mathsf{Unbiased}}{\underset{\mathsf{Bounded variance}}{\leq}} \left\|F\left(\bar{z}^{t}\right)\right\|^{2} + \sigma_{F}^{2}. \tag{52}$$

0

$$\mathbb{E}\left[\left\|\hat{F}\left(\bar{z}^{t},\bar{\xi_{t}}\right)\right\|^{2}\mid\mathcal{F}_{t}\right] = \mathbb{E}\left[\mathbb{E}\left[\left\|\hat{F}\left(\bar{z}^{t},\bar{\xi_{t}}\right) - F(\bar{z}^{t}) + F(\bar{z}^{t})\right\|^{2}\mid\overline{\mathcal{F}}_{t}\right]\mid\mathcal{F}_{t}\right] \stackrel{\mathsf{Unbiased}}{\leq} \left\|F\left(\bar{z}^{t}\right)\right\|^{2} + \sigma_{F}^{2}. \tag{52}$$

 $\circ$  The bound on the term  $\|\hat{F}(\bar{z}^t,\bar{\xi}_t)\|^2$  can further imply,

$$\mathbb{E}\left[\|z^{t+1} - z^t\|^2 \mid \mathcal{F}_t\right] = \alpha_t^2 \gamma^2 \mathbb{E}\left[\|\hat{F}(\bar{z}^t, \bar{\xi}_t)\|^2 \mid \mathcal{F}_t\right] \le \alpha_t^2 \gamma^2 \mathbb{E}\left[\|F\bar{z}^t\|^2 \mid \mathcal{F}_t\right] + \alpha_t^2 \gamma^2 \sigma_F^2 \tag{53}$$

0

$$\mathbb{E}\left[\left\|\hat{F}\left(\bar{z}^{t},\bar{\xi}_{t}\right)\right\|^{2}\mid\mathcal{F}_{t}\right] = \mathbb{E}\left[\mathbb{E}\left[\left\|\hat{F}\left(\bar{z}^{t},\bar{\xi}_{t}\right) - F(\bar{z}^{t}) + F(\bar{z}^{t})\right\|^{2}\mid\overline{\mathcal{F}}_{t}\right]\mid\mathcal{F}_{t}\right] \stackrel{\mathsf{Unbiased}}{\leq} \underset{\mathsf{Bounded variance}}{\left\|F\left(\bar{z}^{t}\right)\right\|^{2} + \sigma_{F}^{2}}.$$
(52)

 $\circ$  The bound on the term  $\|\hat{F}(\bar{z}^t,\bar{\xi}_t)\|^2$  can further imply,

$$\mathbb{E}\left[\|z^{t+1} - z^t\|^2 \mid \mathcal{F}_t\right] = \alpha_t^2 \gamma^2 \mathbb{E}\left[\|\hat{F}(\bar{z}^t, \bar{\xi}_t)\|^2 \mid \mathcal{F}_t\right] \le \alpha_t^2 \gamma^2 \mathbb{E}\left[\|F\bar{z}^t\|^2 \mid \mathcal{F}_t\right] + \alpha_t^2 \gamma^2 \sigma_F^2 \tag{53}$$

Combining the above bounds yields,

$$\mathbb{E}[\|z^{t+1} - z^{\star}\|^{2} + A_{t+1}\|u^{t}\|^{2} + B_{t+1}\|z^{t+1} - z^{t}\|^{2} \mid \mathcal{F}_{t}]$$

$$\leq \|z^{t} - z^{\star}\|^{2} + \left(\underbrace{A_{t+1} + \alpha_{t} \left(\gamma^{2} L_{F}^{2} (1 + b) + \frac{1}{\varepsilon_{t}}\right)}_{X_{1}^{t}}\right) \mathbb{E}[\|u^{t}\|^{2} \mid \mathcal{F}_{t}] - \alpha_{t} \mu \|F(z^{t})\|^{2}$$

$$+ \left(\underbrace{\alpha_t \left(\gamma^2 (\varepsilon_t - 1) - 2\gamma\rho\right) + \alpha_t^2 \gamma^2}_{X_t^t} + B_{t+1} \alpha_t^2 \gamma^2\right) \mathbb{E}[\|F(\bar{z}^t)\|^2 \mid \mathcal{F}_t] + (1 + B_{t+1}) \alpha_t^2 \gamma^2 \sigma_F^2. \tag{54}$$

 $\circ$  Further using the bound on  $||u^t||^2$ .

$$\mathbb{E}\left[\left\|u^{t}\right\|^{2} \mid \mathcal{F}_{t}\right] \leq (1 - \alpha_{t})^{2} \left\|u^{t-1}\right\|^{2} + 2(1 - \alpha_{t})^{2} \gamma^{2} L_{\hat{F}}^{2} \left\|z^{t} - z^{t-1}\right\|^{2} + 2\alpha_{t}^{2} \gamma^{2} \sigma_{F}^{2}$$

$$(55)$$

 $\circ$  Further using the bound on  $\|u^t\|^2$ .

$$\mathbb{E}\left[\left\|u^{t}\right\|^{2} \mid \mathcal{F}_{t}\right] \leq (1 - \alpha_{t})^{2} \left\|u^{t-1}\right\|^{2} + 2(1 - \alpha_{t})^{2} \gamma^{2} L_{\hat{F}}^{2} \left\|z^{t} - z^{t-1}\right\|^{2} + 2\alpha_{t}^{2} \gamma^{2} \sigma_{F}^{2}$$
(55)

We have.

$$\mathbb{E}[\mathcal{U}_{t+1} \mid \mathcal{F}_t] - \mathcal{U}_t \leq -\alpha_t \mu \|F(z^t)\|^2 + \left(X_1^t (1 - \alpha_t)^2 - A_t\right) \|u^{t-1}\|^2$$

$$+ 2\gamma^2 L_{\hat{F}}^2 \left(X_1^t (1 - \alpha_t)^2 - \frac{B_t}{2\gamma^2 L_{\hat{F}}^2}\right) \|z^t - z^{t-1}\|^2 + \left(X_2^t + B_{t+1}\alpha_t^2 \gamma^2\right) \mathbb{E}[\|F(\bar{z}^t)\|^2 \mid \mathcal{F}_t]$$

$$+ \left(B_{t+1}\alpha_t^2 + \alpha_t^2 + 2X_1^t \alpha_t^2\right) \gamma^2 \sigma_F^2.$$
(56)

 $\circ$  Further using the bound on  $\|u^t\|^2$ .

$$\mathbb{E}\left[\left\|u^{t}\right\|^{2} \mid \mathcal{F}_{t}\right] \leq (1 - \alpha_{t})^{2} \left\|u^{t-1}\right\|^{2} + 2(1 - \alpha_{t})^{2} \gamma^{2} L_{\hat{F}}^{2} \left\|z^{t} - z^{t-1}\right\|^{2} + 2\alpha_{t}^{2} \gamma^{2} \sigma_{F}^{2}$$
(55)

o We have.

$$\mathbb{E}[\mathcal{U}_{t+1} \mid \mathcal{F}_{t}] - \mathcal{U}_{t} \leq -\alpha_{t}\mu \|F(z^{t})\|^{2} + \left(X_{1}^{t}(1-\alpha_{t})^{2} - A_{t}\right) \|u^{t-1}\|^{2}$$

$$+ 2\gamma^{2}L_{\hat{F}}^{2} \left(X_{1}^{t}(1-\alpha_{t})^{2} - \frac{B_{t}}{2\gamma^{2}L_{\hat{F}}^{2}}\right) \|z^{t} - z^{t-1}\|^{2} + \left(X_{2}^{t} + B_{t+1}\alpha_{t}^{2}\gamma^{2}\right) \mathbb{E}[\|F(\bar{z}^{t})\|^{2} \mid \mathcal{F}_{t}]$$

$$+ \left(B_{t+1}\alpha_{t}^{2} + \alpha_{t}^{2} + 2X_{1}^{t}\alpha_{t}^{2}\right) \gamma^{2}\sigma_{F}^{2}.$$

$$(56)$$

Having established (??), set  $A_t=A$ ,  $B_t=2A\gamma^2L_{\hat{F}}^2$ , and  $\varepsilon_t=\varepsilon$  to obtain by the law of total expectation that

$$\mathbb{E}[\mathcal{U}_{t+1}] - \mathbb{E}[\mathcal{U}_{t}] \leq -\alpha_{t}\mu\mathbb{E}[\|F(z^{t})\|^{2}] + \left(X_{1}^{t}(1-\alpha_{t})^{2} - A\right)\mathbb{E}[\|u^{t-1}\|^{2}]$$

$$+ 2\gamma^{2}L_{\hat{F}}^{2}\left(X_{1}^{t}(1-\alpha_{t})^{2} - A\right)\mathbb{E}[\|z^{t} - z^{t-1}\|^{2}] + \left(X_{2}^{t} + 2A\gamma^{4}L_{\hat{F}}^{2}\alpha_{t}^{2}\right)\mathbb{E}[\|F(\bar{z}^{t})\|^{2}]$$

$$+ \left(2A\gamma^{2}L_{\hat{F}}^{2} + 1 + 2X_{1}^{t}\right)\alpha_{t}^{2}\gamma^{2}\sigma_{F}^{2}.$$

$$(57)$$

 $\circ$  By carefully choosing the constant A and arepsilon, we can enforce the following negative coefficients,

$$X_1^t(1-\alpha_t)^2-A\leq 0\quad\text{and}\quad X_2^t+2A\gamma^4L_{\hat{F}}^2\alpha_t^2\leq 0.$$

Recall that,

$$X_1^t = \alpha_t \left( \gamma^2 L_F^2(1+b) + \frac{1}{\varepsilon} \right) + A, \quad X_2^t = \alpha_t \left( \gamma^2(\varepsilon - 1) - 2\rho\gamma + \alpha_t \gamma^2 \right).$$

Pick  $A=\frac{1}{2}\left((b+1)\gamma^2L_F^2+\frac{1}{\varepsilon}\right)$  and  $\varepsilon=\gamma L_{\hat{F}}\sqrt{\alpha_0}$ , the negative requirements are satisfied.

 $\circ$  By carefully choosing the constant A and arepsilon, we can enforce the following negative coefficients,

$$X_1^t (1 - \alpha_t)^2 - A \le 0$$
 and  $X_2^t + 2A\gamma^4 L_{\hat{F}}^2 \alpha_t^2 \le 0$ .

Recall that,

$$X_1^t = \alpha_t \left( \gamma^2 L_F^2(1+b) + \frac{1}{\varepsilon} \right) + A, \quad X_2^t = \alpha_t \left( \gamma^2(\varepsilon - 1) - 2\rho\gamma + \alpha_t \gamma^2 \right).$$

Pick  $A=\frac{1}{2}\left((b+1)\gamma^2L_F^2+\frac{1}{\varepsilon}\right)$  and  $\varepsilon=\gamma L_{\hat{F}}\sqrt{\alpha_0}$ , the negative requirements are satisfied.

o Thus we can derive the recursion,

$$\mathbb{E}[\mathcal{U}_{t+1}] - \mathbb{E}[\mathcal{U}_t] \le -\alpha_t \mu \mathbb{E}[\|F(z^t)\|^2] + C\alpha_t^2 \gamma^2 \sigma_F^2.$$

 $\circ$  By carefully choosing the constant A and arepsilon, we can enforce the following negative coefficients,

$$X_1^t (1 - \alpha_t)^2 - A \le 0$$
 and  $X_2^t + 2A\gamma^4 L_{\hat{F}}^2 \alpha_t^2 \le 0$ .

Recall that,

$$X_1^t = \alpha_t \left( \gamma^2 L_F^2(1+b) + \frac{1}{\varepsilon} \right) + A, \quad X_2^t = \alpha_t \left( \gamma^2(\varepsilon - 1) - 2\rho\gamma + \alpha_t \gamma^2 \right).$$

Pick  $A=\frac{1}{2}\left((b+1)\gamma^2L_F^2+\frac{1}{\varepsilon}\right)$  and  $\varepsilon=\gamma L_{\hat{F}}\sqrt{\alpha_0}$ , the negative requirements are satisfied.

o Thus we can derive the recursion,

$$\mathbb{E}[\mathcal{U}_{t+1}] - \mathbb{E}[\mathcal{U}_t] \le -\alpha_t \mu \mathbb{E}[\|F(z^t)\|^2] + C\alpha_t^2 \gamma^2 \sigma_F^2.$$

o Telescoping the above inequality completes the proof.

- o Random iterate convergence is weak.
  - ▶ Hard to check and report the result, since evaluating the deterministic *F* is expensive.
  - ► The result can be volatile.

- o Random iterate convergence is weak.
  - ▶ Hard to check and report the result, since evaluating the deterministic *F* is expensive.
  - The result can be volatile.
- Almost sure convergence is more desirable.
  - No variance to be taken care of.
  - Can just report the final solution after running long enough.

- o Random iterate convergence is weak.
  - $\blacktriangleright$  Hard to check and report the result, since evaluating the deterministic F is expensive.
  - ► The result can be volatile.
- Almost sure convergence is more desirable.
  - No variance to be taken care of.
  - ► Can just report the final solution after running long enough.

## Theorem (Almost Sure Convergence)

Suppose that the three assumptions on  $F, \hat{F}, A$  hold. Suppose  $\gamma \in (\lfloor -2\rho \rfloor_+, 1/L_F)$ ,  $\alpha_t = \frac{1}{t+r}$  for any positive natural number r and

$$(\gamma L_{\hat{F}} + 1)\alpha_t + 2\left(\frac{1+\gamma^2 L_F^2}{1-\gamma^2 L_F^2}\gamma^4 L_F^2 L_{\hat{F}}^2 \alpha_{t+1} + \gamma L_{\hat{F}}\right) (\alpha_{t+1} + 1)\alpha_{t+1} \le 1 + \frac{2\rho}{\gamma}.$$
 (58)

Then, the sequence  $(z^t)_{k\in\mathbb{N}}$  generated by the Alg. ?? converges almost surely to some  $z^\star\in\operatorname{zer} T$ .

- o Random iterate convergence is weak.
  - ▶ Hard to check and report the result, since evaluating the deterministic F is expensive.
  - ► The result can be volatile.
- o Almost sure convergence is more desirable.
  - No variance to be taken care of.
  - ► Can just report the final solution after running long enough.

# Theorem (Almost Sure Convergence)

Suppose that the three assumptions on  $F, \hat{F}, A$  hold. Suppose  $\gamma \in (\lfloor -2\rho \rfloor_+, 1/L_F)$ ,  $\alpha_t = \frac{1}{t+r}$  for any positive natural number r and

$$(\gamma L_{\hat{F}} + 1)\alpha_t + 2\left(\frac{1 + \gamma^2 L_F^2}{1 - \gamma^2 L_F^2} \gamma^4 L_F^2 L_{\hat{F}}^2 \alpha_{t+1} + \gamma L_{\hat{F}}\right) (\alpha_{t+1} + 1)\alpha_{t+1} \le 1 + \frac{2\rho}{\gamma}.$$
(58)

Then, the sequence  $(z^t)_{k\in\mathbb{N}}$  generated by the Alg. ?? converges almost surely to some  $z^\star\in\operatorname{zer} T$ .

#### Remark:

- ▶ Step size  $\alpha_t = \frac{1}{t+r} = \Theta(\frac{1}{t})$ , diminishes faster than  $\Theta(\frac{1}{\sqrt{t}})$ .
- $lackbox{Loses $\tilde{\mathcal{O}}(\frac{1}{\sqrt{T}})$ random iterate convergence rate, but obtain almost sure convergence.}$

#### Constrained case

o Using the prox/resolvent to give the BC-PSEG+ algorithm,

#### Algorithm (BC-PSEG+) Stochastic algorithm for constrained problem.

- 1: REQUIRE  $z^{-1} = z^0 \in \mathbb{R}^n$ ,  $h^{-1} \in \mathbb{R}^n$ ,  $\alpha_t \in (0,1)$ ,  $\gamma \in (\lfloor -2\rho \rfloor_+, 1/L_F)$
- 2: for  $t = 0, 1, \ldots$  until convergence do
- 3: Sample  $\xi_t \sim \mathcal{P}$

4: 
$$h^t = (z^t - \gamma \hat{F}(z^t, \xi_t)) + (1 - \alpha_t) \left( h^{t-1} - (z^{t-1} - \gamma \hat{F}(z^{t-1}, \xi_t)) \right)$$

- 5:  $\bar{z}^t = (\operatorname{id} + \gamma A)^{-1} h_t$
- 6: Sample  $ar{\xi}_t \sim \mathcal{P}$
- 7:  $z^{t+1} = z^t \alpha_t \left( h^t \bar{z}^t + \gamma \hat{F}(\bar{z}^t, \bar{\xi}_t) \right)$
- 8: end for
- 9: Return  $z^{t+1}$

#### Constrained case

o Using the prox/resolvent to give the BC-PSEG+ algorithm,

#### $\textbf{Algorithm} \ (\mathsf{BC}\text{-}\mathsf{PSEG}+) \ \mathsf{Stochastic} \ \mathsf{algorithm} \ \mathsf{for} \ \mathsf{constrained} \ \mathsf{problem}.$

- 1: REQUIRE  $z^{-1} = z^0 \in \mathbb{R}^n$ ,  $h^{-1} \in \mathbb{R}^n$ ,  $\alpha_t \in (0,1)$ ,  $\gamma \in (|-2\rho|_+, 1/L_F)$
- 2: for  $t=0,1,\ldots$  until convergence do
- 3: Sample  $\xi_t \sim \mathcal{P}$

4: 
$$h^t = (z^t - \gamma \hat{F}(z^t, \xi_t)) + (1 - \alpha_t) \left( h^{t-1} - (z^{t-1} - \gamma \hat{F}(z^{t-1}, \xi_t)) \right)$$

- 5:  $\bar{z}^t = (\operatorname{id} + \gamma A)^{-1} h_t$
- 6: Sample  $ar{\xi}_t \sim \mathcal{P}$
- 7:  $z^{t+1} = z^t \alpha_t \left( h^t \bar{z}^t + \gamma \hat{F}(\bar{z}^t, \bar{\xi}_t) \right)$
- 8: end for
- 9: Return  $z^{t+1}$
- o Key idea: apply similar bias-corrected technique to the extrapolation point before the resolvent.
- Properties:
  - Converges to zerT.

# Nonlinear preconditioned primal dual extragradient (NP-PDEG)

$$\min_{x \in \mathbb{R}^n} \max_{y \in \mathbb{R}^r} \quad f(x) + \varphi(x, y) - g(y). \tag{59}$$

where  $\varphi(x,y) := \mathbb{E}_{\xi}[\hat{\varphi}(x,y,\xi)].$ 

# Nonlinear preconditioned primal dual extragradient (NP-PDEG)

$$\min_{x \in \mathbb{R}^n} \max_{y \in \mathbb{R}^r} f(x) + \varphi(x, y) - g(y).$$
 (59)

where  $\varphi(x,y) := \mathbb{E}_{\xi}[\hat{\varphi}(x,y,\xi)].$ 

#### Algorithm Nonlinearly preconditioned primal dual extragradient (NP-PDEG)

- 1: REQUIRE  $z^{-1} = z^0 = (x^0, y^0)$  with  $x^0, x^{-1}, \hat{x}^{-1}, \bar{x}^{-1} \in \mathbb{R}^n, y^0, y^{-1} \in \mathbb{R}^r, \theta \in [0, \infty), \Gamma_1 \succ 0, \Gamma_2 \succ 0$
- 2: **for**  $t = 0, 1, \ldots$  until convergence **do**
- 3:  $\hat{x}^t = x^t \Gamma_1 \nabla_x \hat{\varphi}(z^t, \xi_t) + (1 \alpha_t) (\hat{x}^{t-1} x^{t-1} + \Gamma_1 \nabla_x \hat{\varphi}(x^{t-1}, y^{t-1}, \xi_t)), \xi_t \sim \mathcal{P}$
- 4:  $\bar{\boldsymbol{x}}^t = \operatorname{prox}_f^{\Gamma_1^{-1}} (\hat{x}^t)$
- $5: \quad \hat{y}^t = y^t + \Gamma_2 \nabla_y \hat{\varphi}(\bar{\boldsymbol{x}}^t, y^t, \xi_t') + (1 \alpha_t) \left( \hat{y}^{t-1} y^{t-1} \Gamma_2 \nabla_y \hat{\varphi}(\bar{\boldsymbol{x}}^{t-1}, y^{t-1}, \xi_t') \right), \ \xi_t' \sim \mathcal{P}$
- 6:  $\bar{y}^t = \operatorname{prox}_g^{\Gamma_2^{-1}} \left( \hat{y}^t \right)$
- 7:  $\bar{\xi}_t \sim \mathcal{P}$
- 8:  $x^{t+1} = x^t + \alpha_t \left( \bar{x}^t \hat{x}^t \Gamma_1 \nabla_x \hat{\varphi}(\bar{z}^t, \bar{\xi}_t) \right), \quad y^{t+1} = y^t + \alpha_t \left( \bar{y}^t \hat{y}^t + \Gamma_2 \nabla_y \hat{\varphi}(\bar{z}^t, \bar{\xi}_t) \right)$
- 9: end for
- 10: Return  $z^{t+1} = (x^{t+1}, y^{t+1})$

# Nonlinear preconditioned primal dual extragradient (NP-PDEG)

$$\min_{x \in \mathbb{R}^n} \max_{y \in \mathbb{R}^r} f(x) + \varphi(x, y) - g(y). \tag{59}$$

where  $\varphi(x,y) := \mathbb{E}_{\xi}[\hat{\varphi}(x,y,\xi)].$ 

#### Algorithm Nonlinearly preconditioned primal dual extragradient (NP-PDEG)

- 1: REQUIRE  $z^{-1} = z^0 = (x^0, y^0)$  with  $x^0, x^{-1}, \hat{x}^{-1} \in \mathbb{R}^n, y^0, y^{-1} \in \mathbb{R}^r, \theta \in [0, \infty), \Gamma_1 \succ 0, \Gamma_2 \succ 0$
- 2: for  $t = 0, 1, \ldots$  until convergence do
- 3:  $\hat{x}^t = x^t \Gamma_1 \nabla_x \hat{\varphi}(z^t, \xi_t) + (1 \alpha_t) (\hat{x}^{t-1} x^{t-1} + \Gamma_1 \nabla_x \hat{\varphi}(x^{t-1}, y^{t-1}, \xi_t)), \xi_t \sim \mathcal{P}$
- 4:  $\bar{\boldsymbol{x}}^t = \operatorname{prox}_f^{\Gamma_1^{-1}} (\hat{x}^t)$
- $\text{5:} \quad \hat{y}^t = y^t + \Gamma_2 \nabla_y \hat{\varphi}(\bar{\boldsymbol{x}}^t, y^t, \xi_t') + (1 \alpha_t) \left( \hat{y}^{t-1} y^{t-1} \Gamma_2 \nabla_y \hat{\varphi}(\bar{\boldsymbol{x}}^{t-1}, y^{t-1}, \xi_t') \right), \ \xi_t' \sim \mathcal{P}$
- 6:  $\bar{y}^t = \operatorname{prox}_g^{\Gamma_2^{-1}} \left( \hat{y}^t \right)$
- 7:  $\bar{\mathcal{E}}_t \sim \mathcal{P}$
- 8:  $x^{t+1} = x^t + \alpha_t \left( \bar{x}^t \hat{x}^t \Gamma_1 \nabla_x \hat{\varphi}(\bar{z}^t, \bar{\xi}_t) \right), \quad y^{t+1} = y^t + \alpha_t \left( \bar{y}^t \hat{y}^t + \Gamma_2 \nabla_y \hat{\varphi}(\bar{z}^t, \bar{\xi}_t) \right)$
- 9: end for
- 10: Return  $z^{t+1} = (x^{t+1}, y^{t+1})$
- o Alternatively sampling relaxes the Lipschitz assumption.



## Experiment 1

Example (Unconstrained quadratic game [? , Ex. 5])

Consider,

$$\min_{x \in \mathbb{R}} \max_{y \in \mathbb{R}} \phi(x, y) := axy + \frac{b}{2}x^2 - \frac{b}{2}y^2, \tag{60}$$

where  $a \in \mathbb{R}_+$  and  $b \in \mathbb{R}$ .

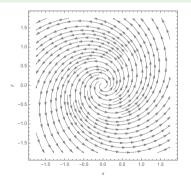


Figure: Unstable dynamics!

### Experiment 2

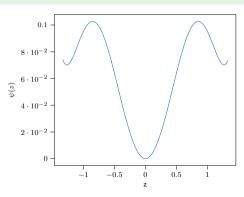
Example (Constrained minimax [? , Ex. 4])

Consider

$$\min_{|x| \leq 4/3} \max_{|y| \leq 4/3} \phi(x,y) := xy + \psi(x) - \psi(y),$$

 $(\mathsf{GlobalForsaken})$ 

where 
$$\psi(z) = \frac{2z^6}{21} - \frac{z^4}{3} + \frac{z^2}{3}$$
.



## Experiment 2

Example (Constrained minimax [? , Ex. 4])

Consider

$$\min_{|x| \leq 4/3} \max_{|y| \leq 4/3} \phi(x,y) := xy + \psi(x) - \psi(y),$$

(GlobalForsaken)

where 
$$\psi(z) = \frac{2z^6}{21} - \frac{z^4}{3} + \frac{z^2}{3}$$
.

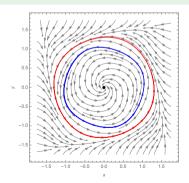


Figure: Limit cycles.

#### **Experimental Results**

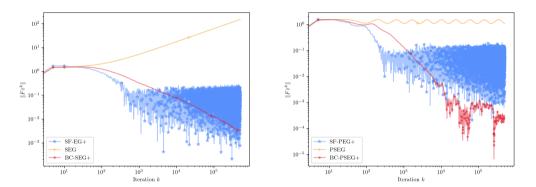


Figure: Comparison of methods in the unconstrained setting and the constrained setting.

#### Remark:

- ► SEG diverges (As expected, since the dynamics is unstable.) and PSEG cycles;
- ▶ Both oblivious extensions (SF-EG+) and (SF-PEG+) only converge to a neighborhood.
- ▶ Only BC-SEG+ and BC-PSEG+ converges properly, with probability 1 as established;

# References I

