

# A Novel Catalyst Scheme for Stochastic Minimax Optimization

Yan Li

Georgia Institute of Technology

Joint work with George Lan

## Minimax Optimization and Proximal Point Method

# Minimax Optimization

## ▷ Problem of interest

$$\min_{x \in X} \left\{ f(x) := \max_{y \in Y} F(x, y) \right\}$$

- Many applications: machine learning (GAN, adversarial training); planning (robust MDPs, Markov games);
- Basic assumptions we will make:
  - 1  $\mu_p$  strongly convex in  $x$ ;  $\mu_d$  strongly concave in  $y$ ;
  - 2  $\nabla F$  is  $L$ -Lipschitz
  - 3 We focus in this talk that  $\mu_d, \mu_p > 0$ , and WLOG  $\mu_d \geq \mu_p$ .

# A Brief Review of Existing Development

## ▷ Many ways to solve this problem

- Variational inequality (VI) based methods (Kotsalis et al., '20; Zhang et al., '23; many others): solve VI associated with the optimality condition

$$\langle \nabla_x F(x^*), x - x^* \rangle - \langle \nabla_y F(y^*), y - y^* \rangle \geq 0$$

- 1 Not optimal in the deterministic setting when  $\mu_p \neq \mu_d$ :  
 $\mathcal{O}(L / \min \{ \mu_p, \mu_d \} \log(1/\epsilon))$

# A Brief Review of Existing Development

## ▷ Many ways to solve this problem

- Variational inequality (VI) based methods (Kotsalis et al., '20; Zhang et al., '23; many others): solve VI associated with the optimality condition

$$\langle \nabla_x F(x^*), x - x^* \rangle - \langle \nabla_y F(y^*), y - y^* \rangle \geq 0$$

- ① Not optimal in the deterministic setting when  $\mu_p \neq \mu_d$ :  
 $\mathcal{O}(L / \min\{\mu_p, \mu_d\} \log(1/\epsilon))$
- Primal-based methods: apply approximate proximal point framework
  - ① (Near) Optimal complexity in the deterministic setting (Lin et al. '20):

$$\tilde{\mathcal{O}}(L / \sqrt{\mu_p \mu_d} \log(1/\epsilon))$$

- ② Stochastic setting: seemingly no easy extension

# A Brief Review of Existing Development

## ▷ Many ways to solve this problem

- Variational inequality (VI) based methods (Kotsalis et al., '20; Zhang et al., '23; many others): solve VI associated with the optimality condition

$$\langle \nabla_x F(x^*), x - x^* \rangle - \langle \nabla_y F(y^*), y - y^* \rangle \geq 0$$

- 1 Not optimal in the deterministic setting when  $\mu_p \neq \mu_d$ :  
 $\mathcal{O}(L / \min\{\mu_p, \mu_d\} \log(1/\epsilon))$
- Primal-based methods: apply approximate proximal point framework
  - 1 (Near) Optimal complexity in the deterministic setting (Lin et al. '20):

$$\tilde{\mathcal{O}}(L / \sqrt{\mu_p \mu_d} \log(1/\epsilon))$$

- 2 Stochastic setting: seemingly no easy extension
- Can we design methods with optimal complexities in both deterministic and stochastic settings?

# Accelerated Proximal Point Method

---

**Algorithm** Accelerated Proximal Point Method

---

**Input:** initial points  $\bar{x}_0 = \tilde{x}_0$

**for**  $k = 1, 2, \dots, K$  **do**

$$\hat{x}_k = \gamma_k \bar{x}_{k-1} + (1 - \gamma_k) x_{k-1}.$$

$$x_k = \min_{x \in X} f(x) + \frac{\beta_k}{2} \|x - \hat{x}_k\|^2$$

$$\bar{x}_k = \frac{1}{\alpha_k \gamma_k + \mu(1 - \gamma_k)} [\beta_k x_k + (\mu_p - \beta_k)(1 - \gamma_k) x_{k-1}].$$

**end for**

**Output:**  $\tilde{x}_K$

---

- Convergence rate: with  $\beta_k = \beta > 0$ ,

$$f(\tilde{x}_K) - f(x^*) = \mathcal{O}\left(\frac{\beta}{K^2} \|x^* - \bar{x}_0\|^2\right)$$

# Accelerated Proximal Point Method

---

**Algorithm** Accelerated Proximal Point Method

---

**Input:** initial points  $\bar{x}_0 = \tilde{x}_0$

**for**  $k = 1, 2, \dots, K$  **do**

$$\hat{x}_k = \gamma_k \bar{x}_{k-1} + (1 - \gamma_k) x_{k-1}.$$

$$x_k = \min_{x \in X} f(x) + \frac{\beta_k}{2} \|x - \hat{x}_k\|^2$$

$$\bar{x}_k = \frac{1}{\alpha_k \gamma_k + \mu(1 - \gamma_k)} [\beta_k x_k + (\mu_p - \beta_k)(1 - \gamma_k) x_{k-1}].$$

**end for**

**Output:**  $\tilde{x}_K$

---

- Convergence rate: with  $\beta_k = \beta > 0$ ,

$$f(\tilde{x}_K) - f(x^*) = \mathcal{O}\left(\frac{\beta}{K^2} \|x^* - \bar{x}_0\|^2\right)$$

- **High-level idea of Catalyst:** do approximate computation of proximal update using simple non-optimal methods



# Accelerated Proximal Point Method

---

## Algorithm Accelerated Proximal Point Method

---

**Input:** initial points  $\bar{x}_0 = \tilde{x}_0$

**for**  $k = 1, 2, \dots, K$  **do**

$$\hat{x}_k = \gamma_k \bar{x}_{k-1} + (1 - \gamma_k) x_{k-1}.$$

$$x_k = \min_{x \in X} f(x) + \frac{\beta_k}{2} \|x - \hat{x}_k\|^2$$

$$\bar{x}_k = \frac{1}{\alpha_k \gamma_k + \mu(1 - \gamma_k)} [\beta_k x_k + (\mu_p - \beta_k)(1 - \gamma_k) x_{k-1}].$$

**end for**

**Output:**  $\tilde{x}_K$

---

- Convergence rate: with  $\beta_k = \beta > 0$ ,

$$f(\tilde{x}_K) - f(x^*) = \mathcal{O}\left(\frac{\beta}{K^2} \|x^* - \bar{x}_0\|^2\right)$$

- **Essential question:** what error condition should we consider to measure “approximate computation”?

## Catalyst for Convex Optimization

# Catalyst for Convex Optimization

▷ **Problem of interest:**  $\min_{x \in X} f(x)$ ,  $f$  is  $\mu$ -strongly-convex and  $L$ -smooth

---

**Algorithm Catalyst( $\mathcal{A}$ ):** catalyst scheme for convex optimization

---

**Input:** initial points  $\bar{x}_0 = \tilde{x}_0$ , to-be-catalyzed method  $\mathcal{A}$ .

**for**  $k = 1, 2, \dots, K$  **do**

$$\hat{x}_k = \gamma_k \bar{x}_{k-1} + (1 - \gamma_k) \tilde{x}_{k-1}.$$

$$(\tilde{x}_k, x_k) = \mathcal{A}(\phi_k, \hat{x}_k)$$

$$\bar{x}_k = \frac{1}{\alpha_k \gamma_k + \mu(1 - \gamma_k)} [\alpha_k x_k + (\mu - \alpha_k)(1 - \gamma_k) \tilde{x}_{k-1}].$$

**end for**

**Output:**  $\tilde{x}_K$

---

# Catalyst for Convex Optimization

▷ **Problem of interest:**  $\min_{x \in X} f(x)$ ,  $f$  is  $\mu$ -strongly-convex and  $L$ -smooth

---

**Algorithm Catalyst( $\mathcal{A}$ ):** catalyst scheme for convex optimization

---

**Input:** initial points  $\bar{x}_0 = \tilde{x}_0$ , to-be-catalyzed method  $\mathcal{A}$ .

**for**  $k = 1, 2, \dots, K$  **do**

$$\hat{x}_k = \gamma_k \bar{x}_{k-1} + (1 - \gamma_k) \tilde{x}_{k-1}.$$

$$(\tilde{x}_k, x_k) = \mathcal{A}(\phi_k, \hat{x}_k)$$

$$\bar{x}_k = \frac{1}{\alpha_k \gamma_k + \mu(1 - \gamma_k)} [\alpha_k x_k + (\mu - \alpha_k)(1 - \gamma_k) \tilde{x}_{k-1}].$$

**end for**

**Output:**  $\tilde{x}_K$

---

- ①  $\phi_k$ : subproblem of proximal step:  $\phi_k(x) = f(x) + \frac{\beta_k}{2} \|x - \hat{x}_k\|^2$
- ②  $\mathcal{A}(\phi_k, \hat{x}_k)$ : minimize  $\phi_k$  using  $\mathcal{A}$  (can be stochastic) initialized from  $\hat{x}_k$
- ③ **Error condition:**

$$\mathbb{E}[\phi_k(\tilde{x}_k) - \phi_k(\tilde{x}) + \frac{\alpha_k}{2} \|\tilde{x} - x_k\|^2] \leq \frac{\epsilon_k}{2} \|\tilde{x} - \hat{x}_k\|^2 + \delta_k, \quad \forall \tilde{x} \in X,$$

- $(\alpha_k, \epsilon_k)$  will depend both on  $\beta_k$  and how long we run  $\mathcal{A}$ . Note that  $\alpha_k$  is required for running Catalyst.

# Catalyst for Convex Optimization

## ▷ Basic recursion

$$\begin{aligned} & \mathbb{E}[f(\tilde{x}_k) - f(x) + \frac{\alpha_k \gamma_k^2 + \gamma_k(1-\gamma_k)\mu}{2} \|x - \bar{x}_k\|^2] \\ & \leq (1 - \gamma_k)[f(\tilde{x}_{k-1}) - f(x)] + \frac{(\beta_k + \varepsilon_k)\gamma_k^2}{2} \|x - \bar{x}_{k-1}\|^2 + \delta_k. \end{aligned}$$

## Lemma

Suppose  $\mu = 0$ . Run  $\text{Catalyst}(\mathcal{A})$  with

$$\gamma_k = \frac{2}{k+1}, \quad \beta_k = \frac{(k+1)L}{k}.$$

In addition, suppose  $\{\alpha_k\}$  is chosen such that there exists  $\{(\varepsilon_k, \delta_k)\}$  certifying error condition with

$$\begin{aligned} \alpha_k &= \beta_k(1 + \varepsilon), \quad \varepsilon_k = \beta_k \varepsilon, \quad \varepsilon \leq 1, \\ \delta_k &\leq \delta. \end{aligned}$$

Then we have

$$\mathbb{E}[f(\tilde{x}_K) - f(x^*)] \leq \frac{4L}{K^2} \|x^* - \bar{x}_0\|^2 + 2K\delta.$$

# Catalyzing SGD

▷ What method to be catalyzed? stochastic gradient descent

---

**Algorithm**  $\text{SGD}(\phi_k; \hat{x}_k)$ : SGD for solving  $\phi_k$ , initialized at  $\hat{x}_k$

---

**Input:** stepsizes  $\{\eta_t\}$ , total number of steps  $n > 0$

**for**  $t = 1, 2, \dots, T$  **do**

    Form  $g_{t-1} = \nabla f(u_{t-1}; \xi_{t-1}) + \beta_k(u - \hat{x}_k)$ .

$u_t = \operatorname{argmin}_{w \in X} \langle g_{t-1}, w \rangle + \frac{1}{2\eta_t} \|w - u_{t-1}\|^2$ .

**end for**

    Compute  $\bar{u}_T =$  proper ergodic mean of  $\{u_t\}$

**Output:**  $(\bar{u}_T, u_T)$ .

---

▷ SGD satisfies **error condition**, with  $(\alpha_k, \varepsilon_k)$  depending on  $\beta_k$  and  $T$ .

## Proposition

Let  $(\tilde{x}_k, x_k)$  be the output of  $\text{SGD}(\phi_k; \hat{x}_k)$ , then

$$\mathbb{E} \left[ \phi_k(\tilde{x}_k) - \phi_k(u) + \frac{\alpha \mu_{\phi_k}}{2} \|\tilde{u} - x_k\|^2 \right] \leq \frac{\varepsilon \mu_{\phi_k}}{2} \|u - \hat{x}_k\|^2 + \delta, \quad \forall u \in X,$$

with  $\varepsilon, \alpha, \delta$  depending on  $T$ , and  $\mu_{\phi_k} = \mu + \beta_k$ .

# Catalyzing SGD

## Theorem

Suppose  $\mu = 0$ . For any  $\epsilon > 0$ , run *Catalyst(SGD)* with parameters

$$K = 4\sqrt{\frac{L\|x^* - \bar{x}_0\|^2}{\epsilon}}, \quad \gamma_k = \frac{2}{k+1}, \quad \beta_k = \frac{(k+1)L}{k}, \quad \alpha_k = \frac{\beta_k}{1-\Lambda_T},$$

where

$$T = 8 + \frac{32\sigma^2 K}{L\epsilon}, \quad \Lambda_T = \frac{90}{(T+9)(T+10)}.$$

At iteration  $k$ , the proximal step is approximately solved by running *SGD*( $\phi_k, \hat{x}_k$ ) for  $T$  steps with stepsize

$$\eta_t = \frac{2}{\beta_k(t+8)}.$$

Then we have  $\mathbb{E}[f(\tilde{x}_K) - f(x^*)] \leq \epsilon$ . The number of calls to *SFO* is bounded by

$$\mathcal{O}\left(\sqrt{\frac{L\|x^* - \bar{x}_0\|^2}{\epsilon}} + \frac{\sigma^2\|x^* - \bar{x}_0\|^2}{\epsilon^2}\right).$$

# Catalyzing SGD

## Theorem

Suppose  $\mu > 0$ . Apply restarting strategy to the Catalyst framework. Then to obtain

$$\mathbb{E}[f(\tilde{x}_K) - f(x^*)] \leq \epsilon.$$

The number of calls to SFO is bounded by

$$\mathcal{O}\left(\sqrt{\frac{L}{\mu}} \log_2\left(\frac{f(x_{(0)}) - f(x^*)}{\epsilon}\right) + \frac{\sigma^2}{\mu\epsilon}\right).$$

- The complexities are optimal in both deterministic ( $\sigma^2 = 0$ ) and stochastic ( $\sigma^2 > 0$ ) settings, for both convex and strongly-convex objectives.



## Catalyst for Minimax Optimization

# Catalyst for Minimax Optimization

Denote  $z = (x, y) \in X \times Y$ .

---

**Algorithm** A minimax catalyst scheme

---

**Input:** initial points  $\bar{z}_0 = \tilde{z}_0 = z_0$

**for**  $k = 1, 2, \dots$ , **do**

$$\hat{x}_k = \gamma_k \bar{x}_{k-1} + (1 - \gamma_k) \tilde{x}_{k-1}.$$

$$(\tilde{z}_k, z_k) = \mathcal{A}(\Phi_k, (\hat{x}_k, y_{k-1}))$$

$$\bar{x}_k = \frac{1}{\alpha_k \gamma_k + \mu_p (1 - \gamma_k)} [\alpha_k x_k + (\mu_p - \alpha_k)(1 - \gamma_k) \tilde{x}_{k-1}].$$

**end for**

---

# Catalyst for Minimax Optimization

Denote  $z = (x, y) \in X \times Y$ .

---

**Algorithm A** minimax catalyst scheme

---

**Input:** initial points  $\bar{z}_0 = \tilde{z}_0 = z_0$

**for**  $k = 1, 2, \dots$ , **do**

$$\hat{x}_k = \gamma_k \bar{x}_{k-1} + (1 - \gamma_k) \tilde{x}_{k-1}.$$

$$(\tilde{z}_k, z_k) = \mathcal{A}(\Phi_k, (\hat{x}_k, y_{k-1}))$$

$$\bar{x}_k = \frac{1}{\alpha_k \gamma_k + \mu_p (1 - \gamma_k)} [\alpha_k x_k + (\mu_p - \alpha_k) (1 - \gamma_k) \tilde{x}_{k-1}].$$

**end for**

---

- The minimax Catalyst scheme looks almost identical to that of convex optimization
- $\mathcal{A}(\Phi_k, (\hat{x}_k, y_{k-1}))$ :  $\min_{x \in X} \max_{y \in Y} \Phi_k(x, y) := F(x, y) + \frac{\beta_k}{2} \|x - \hat{x}_k\|^2$  initialized at  $(\hat{x}_k, y_{k-1})$ .
- **Error condition:**

$$\begin{aligned} & \mathbb{E} \left[ \Phi_k(\tilde{x}_k, \tilde{y}) - \Phi_k(\tilde{x}, \tilde{y}_k) + \frac{\alpha_k}{2} \|\tilde{x} - x_k\|^2 + \frac{\alpha_k}{2} \|\tilde{y} - y_k\|^2 \right] \\ & \leq \mathbb{E} \left[ \frac{\varepsilon'_k}{2} \|\tilde{x} - \hat{x}_k\|^2 + \frac{\varepsilon_k}{2} \|\tilde{y} - y_{k-1}\|^2 \right] + \delta_k \end{aligned}$$

# Catalyst for Minimax Optimization

## ▷ Basic recursion

$$\begin{aligned} & \left(1 - \frac{4\varepsilon_k}{\mu_d}\right) \mathbb{E}[f(\tilde{x}_k) - f(x^*)] + \frac{\alpha_k \gamma_k^2 + \gamma_k(1-\gamma_k)\mu_p}{2} \mathbb{E}[\|x^* - \bar{x}_k\|^2] + \frac{\alpha_k}{2} \mathbb{E}[\|\tilde{y}_k^* - y_k\|^2] \\ \leq & \left(1 - \gamma_k + \frac{4\varepsilon_k}{\mu_d}\right) \mathbb{E}[f(\tilde{x}_{k-1}) - f(x^*)] + \frac{(\beta_k + \varepsilon'_k)\gamma_k^2}{2} \mathbb{E}[\|x^* - \bar{x}_{k-1}\|^2] \\ & + \varepsilon_k \mathbb{E}[\|\tilde{y}_{k-1}^* - y_{k-1}\|^2] + \delta_k. \end{aligned}$$

# Catalyst for Minimax Optimization

## ▷ Basic recursion

$$\begin{aligned} & \left(1 - \frac{4\varepsilon_k}{\mu_d}\right) \mathbb{E}[f(\tilde{x}_k) - f(x^*)] + \frac{\alpha_k \gamma_k^2 + \gamma_k(1-\gamma_k)\mu_p}{2} \mathbb{E}[\|x^* - \bar{x}_k\|^2] + \frac{\alpha_k}{2} \mathbb{E}[\|\tilde{y}_k^* - y_k\|^2] \\ & \leq \left(1 - \gamma_k + \frac{4\varepsilon_k}{\mu_d}\right) \mathbb{E}[f(\tilde{x}_{k-1}) - f(x^*)] + \frac{(\beta_k + \varepsilon'_k)\gamma_k^2}{2} \mathbb{E}[\|x^* - \bar{x}_{k-1}\|^2] \\ & \quad + \varepsilon_k \mathbb{E}[\|\tilde{y}_{k-1}^* - y_{k-1}\|^2] + \delta_k. \end{aligned}$$

### Lemma

Fix total iterations  $K \geq 1$  a priori. Choose

$$\gamma_k = \frac{2}{k+1}, \quad \beta_k = \frac{\mu_d(k+1)}{4(k+2)}.$$

In addition, suppose  $\alpha_k$  is chosen such that there exists  $\varepsilon_k$  certifying error condition with

$$\alpha_k = \beta_k(1 + \varepsilon'), \quad \varepsilon'_k = \beta_k \varepsilon', \quad \varepsilon_k = \beta_k \varepsilon, \quad (3.1)$$

$$\delta_k \leq \delta, \quad (3.2)$$

for some  $\delta > 0$  and  $\varepsilon \leq \min \left\{ \frac{1}{12}, \frac{1}{(K+1)(K+2)}, \frac{\|x^* - \tilde{x}_0\|^2}{2[f(\bar{x}_0) - f(x^*)]} \right\}$ ,  $\varepsilon' \leq 1$ . Then

$$\mathbb{E}[f(\tilde{x}_K) - f(x^*)] \leq \frac{24D_0\mu_d}{K^2} + 64K\delta, \quad D_0 = \|x^* - \tilde{x}_0\|^2 + \|\tilde{y}_0^* - y_0\|^2.$$

# Catalyst for Minimax Optimization

- ▷ **What methods do we catalyze?** The solution of the to-be-catalyzed method needs to certify the error condition.

# Catalyst for Minimax Optimization

▷ **What methods do we catalyze?** The solution of the to-be-catalyzed method needs to certify the error condition.

---

**Algorithm**  $\text{SEG}(H; z_0)$ : extragradient for  $\min_{x \in X} \max_{y \in Y} H(x, y)$

---

**Input:** stepsizes  $\{\eta_t\}$ , total number of steps  $n > 0$ , initial point  $z_0 \in Z$   
**for**  $t = 0, 1, \dots, T-1$  **do**

    Define  $G(z, \xi) = [\nabla_x H(z; \xi); -\nabla_y H(z, \xi)]$ . Sample  $\xi_t, \hat{\xi}_t$ , and update

$$\hat{z}_t = \underset{z \in Z}{\operatorname{argmin}} \eta_t \langle G(z_t, \xi_t), z \rangle + \frac{1}{2} \|z - z_t\|^2;$$

$$z_{t+1} = \underset{z \in Z}{\operatorname{argmin}} \eta_t \left[ \langle G(\hat{z}_t, \hat{\xi}_t), z \rangle + \frac{\mu}{2} \|z - \hat{z}_t\|^2 \right] + \frac{1}{2} \|z - z_t\|^2.$$

**end for**

Construct  $\bar{z}_T =$  proper ergodic mean of  $\{z_t\}$ .

**Output:**  $(\bar{z}_T, z_T)$ .

---

# Catalyst for Minimax Optimization

## Lemma

Suppose

$$L\eta_t \leq 1/2, \quad t \geq 0. \quad (3.3)$$

Then for any  $T \geq 1$ , we have

$$\begin{aligned} & \mathbb{E} \left[ F(\bar{x}_T, y) - F(x, \bar{y}_T) + \frac{\mu\Lambda_T}{2(\Lambda_T - \Lambda_0)} \|z - z_T\|^2 \right] \\ & \leq \frac{\mu\Lambda_0}{2(\Lambda_T - \Lambda_0)} \|z - z_0\|^2 + \frac{8\mu\sigma^2}{\Lambda_T - \Lambda_0} \sum_{t=0}^{T-1} \eta_t^2 \Lambda_t, \end{aligned}$$

where

$$\Lambda_t = \begin{cases} 1, & t = 0; \\ (1 + \mu\eta_{t-1})\Lambda_{t-1}, & t \geq 1. \end{cases}$$



# Catalyst for Minimax Optimization

## Theorem

Suppose  $\mu_p = 0, \mu_d > 0$ . For any  $\epsilon > 0$ , choose total number of iterations

$$K \geq \sqrt{\frac{24[2\mu_d\|x^* - \tilde{x}_0\|^2 + \mu_d\|\tilde{y}_0^* - y_0\|^2]}{\epsilon}}.$$

Choose  $\{(\gamma_k, \beta_k, \alpha_k)\}$  properly. In addition, at the  $k$ -th iteration of the catalyst scheme, run SEG procedure for a total of  $T$  steps with proper stepsizes  $\{\eta_t\}$  and  $T$ . Then we obtain

$$\mathbb{E}[f(\tilde{x}_K) - f(x^*)] \leq \epsilon.$$

The total number of calls to SFO can be bounded by

$$\tilde{\mathcal{O}}\left(\frac{LD_0}{\sqrt{\mu_d\epsilon}} + \frac{\sigma^2 D_0}{\epsilon^2}\right),$$

where  $D_0 = \|x^* - \tilde{x}_0\|^2 + \|\tilde{y}_0^* - y_0\|^2$ .

# Catalyst for Minimax Optimization

## Theorem

Suppose  $\mu_d \geq \mu_p > 0$ . Apply the restarting strategy to the Catalyst scheme. Then to obtain

$$\mathbb{E} [f(\tilde{x}_{(\epsilon)}) - f(x^*)] \leq \epsilon. \quad (3.4)$$

The total number of calls to SFO can be bounded by

$$\tilde{O} \left( \frac{L}{\sqrt{\mu_d \mu_p}} \log_2 \left( \frac{\Delta_0}{\epsilon} \right) + \frac{\sigma^2}{\mu_p \epsilon} \right),$$

$$\text{where } \Delta_0 = f(\tilde{x}_{(0)}) - f(x^*) + \frac{\mu_d}{12} \|\tilde{y}_{(0)}^* - y_{(0)}\|^2.$$

- Optimal complexities in both deterministic and stochastic setting.