

New (Optimization) Perspectives on GANs

Gauthier Gidel

- I. A Variational Inequality Perspective on GANs.

- II. Reducing Noise in GANs with Variance Reduced Methods.

A Variational Inequality Perspective on GANs

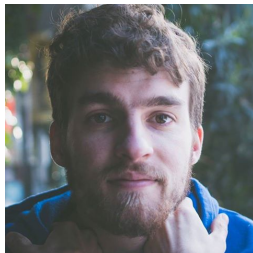
Gauthier Gidel^{*1}, Hugo Berard^{*12}, Gaëtan Vignoud¹,
Pascal Vincent¹², Simon Lacoste-Julien¹

**equal contribution*

¹ Mila, Université de Montréal

² Facebook AI Research (FAIR), Montréal

Hugo
Berard



Gaëtan
Vignoud



Pascal
Vincent



Simon
Lacoste-Julien



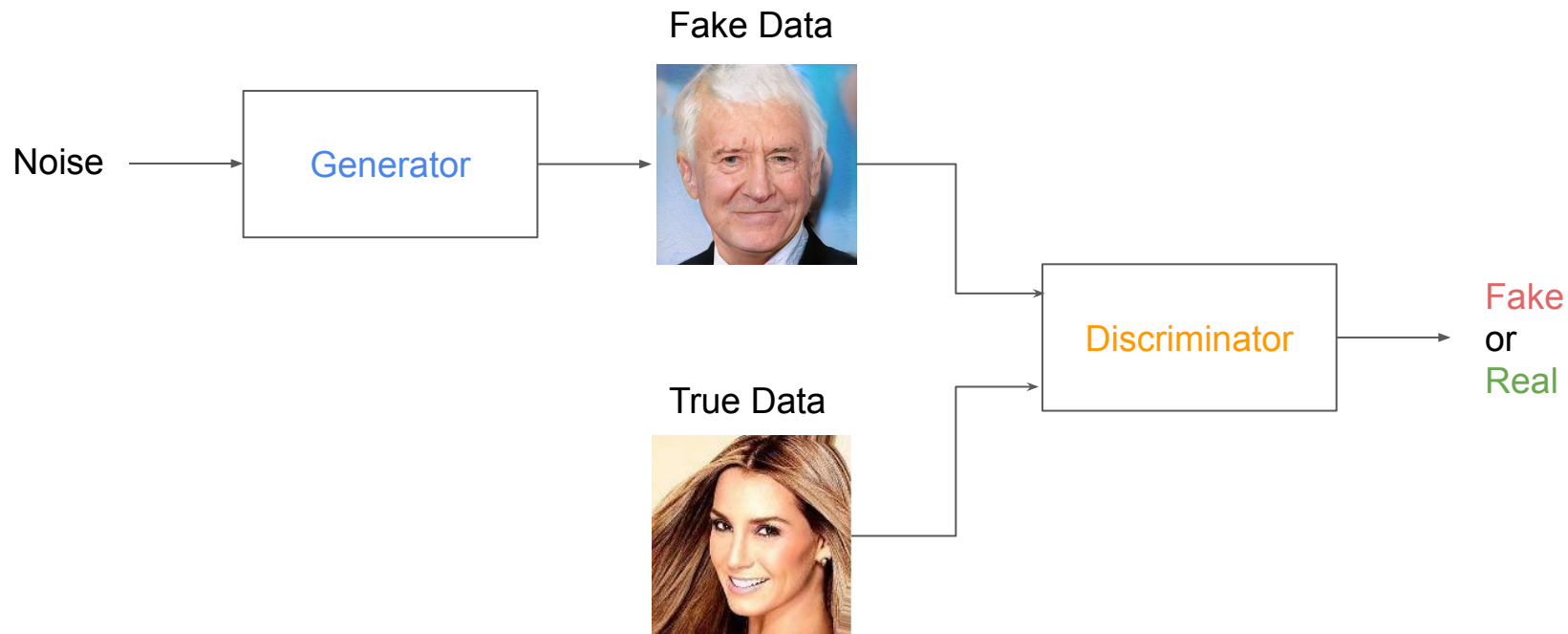
1. Quick Recap on GANs and two-player games.
2. GAN as a Variational Inequality Problem.
3. Optimization of Variational Inequality.
4. Experimental results.
5. Conclusion.

NB: All the citations in this talk are in my arXiv submission.

Quick recap on Generative Adversarial Networks (GANs) (and two-player games)

Generative Adversarial Networks (GANs)

[Goodfellow et al. NIPS 2014]



Generative Adversarial Networks (GANs)

[Goodfellow et al. NIPS 2014]

Discriminator

Generator

$$\min_{\theta} \max_{\phi} \underbrace{\mathbb{E}_{x \sim p_{\mathcal{D}}} [\log(D_{\phi}(x))] + \mathbb{E}_{z \sim p_{\mathcal{Z}}} [\log(1 - D_{\phi}(G_{\theta}(z)))]}_{\text{Total Loss}}$$

If \mathcal{D} is non-parametric: $L(\theta) = \text{JSD}(p_{\mathcal{D}} || p_{\theta})$

Non-saturating GAN: “much stronger gradient in early learning”

Loss of Generator

Loss of Discriminator

$$\min_{\theta} -\mathbb{E}_{z \sim p_{\mathcal{Z}}} [\log(D_{\phi}(G_{\theta}(z)))]$$

$$\max_{\phi} \mathbb{E}_{x \sim p_{\mathcal{D}}} [\log(D_{\phi}(x))] + \mathbb{E}_{z \sim p_{\mathcal{Z}}} [\log(1 - D_{\phi}(G_{\theta}(z)))]$$

Two-player Games

Player 1

Player 2

$$\theta^* \in \arg \min_{\theta \in \Theta} \mathcal{L}^{(\theta)}(\theta, \varphi^*) \quad \text{and} \quad \varphi^* \in \arg \min_{\varphi \in \Phi} \mathcal{L}^{(\varphi)}(\theta^*, \varphi)$$

Zero-sum game if: $\mathcal{L}^{(\theta)} = -\mathcal{L}^{(\varphi)}$ also called *Saddle Point* (SP).

Example: WGAN formulation [Arjovsky et al. 2017]

$$\min_{\theta} \max_{\phi, \|f_{\phi}\|_L \leq 1} \underbrace{\mathbb{E}_{x \sim p_{\mathcal{D}}} [f_{\phi}(x)] - \mathbb{E}_{z \sim p_{\mathcal{Z}}} [f_{\phi}(g_{\theta}(z))]}_{\mathcal{L}^{(\theta)} = -\mathcal{L}^{(\varphi)}}$$

Two-player Games

Player 1

Player 2

$$\theta^* \in \arg \min_{\theta \in \Theta} \mathcal{L}^{(\theta)}(\theta, \varphi^*) \quad \text{and} \quad \varphi^* \in \arg \min_{\varphi \in \Phi} \mathcal{L}^{(\varphi)}(\theta^*, \varphi)$$



- In games we want to **converge** to the Saddle Point.
- Different from **single** objective **minimization** where we want to avoid saddle points.
- ~~Saddle point~~ -> **Zero-sum game (or Minmax)**

Two-player Games

Player 1

Player 2

$$\theta^* \in \arg \min_{\theta \in \Theta} \mathcal{L}^{(\theta)}(\theta, \varphi^*) \quad \text{and} \quad \varphi^* \in \arg \min_{\varphi \in \Phi} \mathcal{L}^{(\varphi)}(\theta^*, \varphi)$$

Non zero-sum game if we **do not** have: $\mathcal{L}^{(\theta)} = -\mathcal{L}^{(\varphi)}$

Example: Non-saturating GAN: [Goodfellow et al. 2014]

Loss of Generator

Loss of Discriminator

$$\min_{\theta} -\mathbb{E}_{z \sim p_Z} [\log(D_{\phi}(G_{\theta}(z)))] \quad \max_{\phi} \mathbb{E}_{x \sim p_{\mathcal{D}}} [\log(D_{\phi}(x))] + \mathbb{E}_{z \sim p_Z} [\log(1 - D_{\phi}(G_{\theta}(z)))]$$

Minmax training is ~~hard~~ different !

Minmax training is ~~hard~~ different !

(You can replace “minmax” with two-player games)

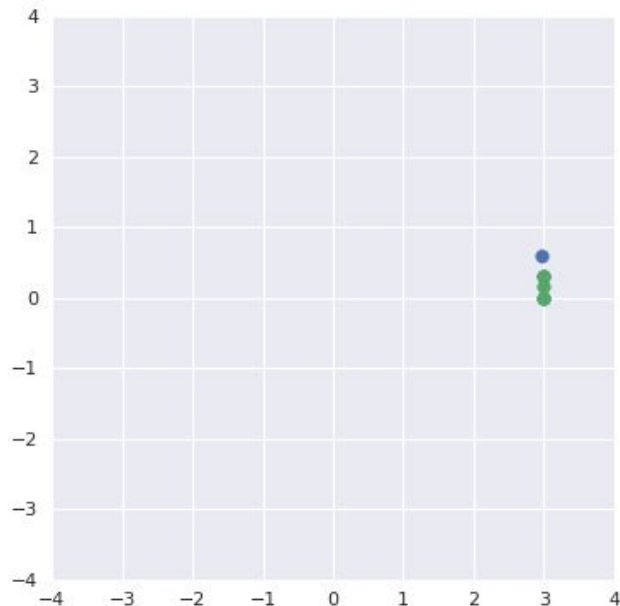
“Minmax Training is Hard ...”

Example: WGAN with **linear discriminator** and **generator**

Gradient vector field: $F(\theta, \phi) = \begin{pmatrix} \nabla_{\theta} L_{\theta}(\theta, \phi) \\ \nabla_{\phi} L_{\phi}(\theta, \phi) \end{pmatrix}$

Bilinear saddle point = Linear in θ and ϕ
⇒ “Cycling behavior” (see right).

$$\min_{\theta} \max_{\phi, \|f_{\phi}\|_L \leq 1} \phi^T \mathbb{E}_{x \sim p_{\mathcal{D}}} [x] - \phi^T \theta \mathbb{E}_{z \sim p_{\mathcal{Z}}} [z]$$



Generative Adversarial Networks as a Variational Inequality Problem (VIP)

GANs as a Variational Inequality

New perspective for GANs:

- Based on **stationary conditions**.
- Relates to vast literature with standard algorithms.

Nash-Equilibrium: $\begin{cases} \theta^* = \arg \min_{\theta} L_{\theta}(\theta, \phi^*) \\ \phi^* = \arg \min_{\phi} L_{\phi}(\theta^*, \phi) \end{cases} \longleftarrow$ No player can improve its cost

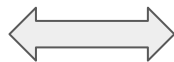
Stationary Conditions: $\begin{cases} \nabla_{\theta} L_{\theta}(\theta^*, \phi^*)^T (\theta - \theta^*) \geq 0 \\ \nabla_{\phi} L_{\phi}(\theta^*, \phi^*)^T (\phi - \phi^*) \geq 0 \end{cases} \quad \forall (\theta, \phi) \in \Theta \times \Phi$

can be **constraint sets**.

GANs as a Variational Inequality

Nash-Equilibrium:

$$\begin{cases} \theta^* = \arg \min_{\theta} L_{\theta}(\theta, \phi^*) \\ \phi^* = \arg \min_{\phi} L_{\phi}(\theta^*, \phi) \end{cases}$$



Stationary Conditions:

$$\begin{cases} \nabla_{\theta} L_{\theta}(\theta^*, \phi^*)^T (\theta - \theta^*) \geq 0 \\ \nabla_{\phi} L_{\phi}(\theta^*, \phi^*)^T (\phi - \phi^*) \geq 0 \end{cases} \quad \forall (\theta, \phi) \in \Theta \times \Phi$$

Same problem but **different** perspective.

Joint Minimization vs. Stationary point

GANs as a Variational Inequality

Stationary Conditions:
$$\begin{cases} \nabla_{\theta} L_{\theta}(\theta^*, \phi^*)^T (\theta - \theta^*) \geq 0 \\ \nabla_{\phi} L_{\phi}(\theta^*, \phi^*)^T (\phi - \phi^*) \geq 0 \end{cases} \quad \forall (\theta, \phi) \in \Theta \times \Phi$$

Can be written as:
$$F(\omega) = \begin{pmatrix} \nabla_{\theta} L_{\theta}(\omega) \\ \nabla_{\phi} L_{\phi}(\omega) \end{pmatrix}$$

$$\omega \stackrel{\uparrow}{=} (\theta, \phi)$$

$$F(\omega^*)^T (\omega - \omega^*) \geq 0 \quad \forall \omega \in \Omega$$

ω^* solves the **Variational Inequality**

GANs as a Variational Inequality

Stationary Conditions: $F(\omega^*)^T (\omega - \omega^*) \geq 0 \quad \forall \omega \in \Omega$

Unconstrained (or optimum in the interior):

$$\|\nabla_{\theta} \mathcal{L}^{(\theta)}(\theta^*, \varphi^*)\| = \|\nabla_{\varphi} \mathcal{L}^{(\varphi)}(\theta^*, \varphi^*)\| = 0.$$

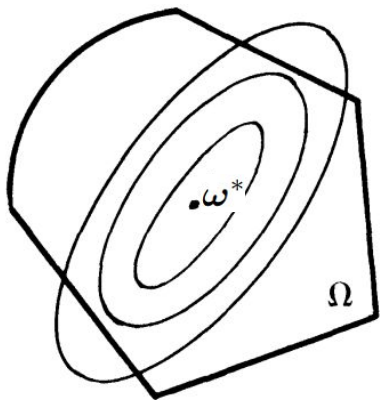


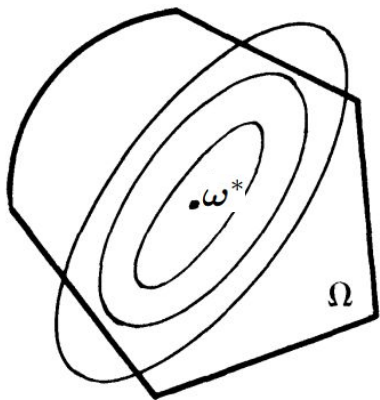
Figure from [Dunn 1979]

GANs as a Variational Inequality

Stationary Conditions: $F(\omega^*)^\top (\omega - \omega^*) \geq 0 \quad \forall \omega \in \Omega$

Unconstrained (or ω^* in the interior):

$$\|\nabla_{\theta} \mathcal{L}^{(\theta)}(\theta^*, \varphi^*)\| = \|\nabla_{\varphi} \mathcal{L}^{(\varphi)}(\theta^*, \varphi^*)\| = 0.$$



Constrained and ω^* on the boundary:

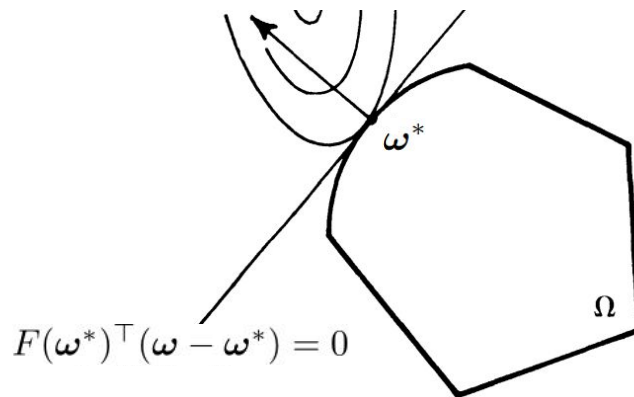


Figure from [Dunn 1979]

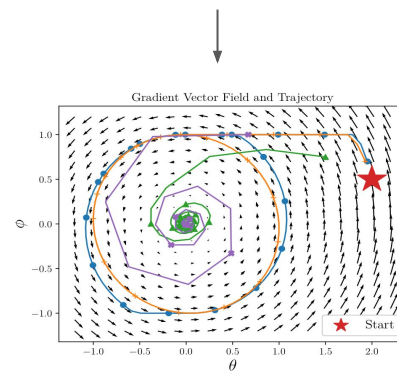
GANs as a Variational Inequality

Takeaways:

- GAN can be formulated as a **Variational Inequality**.
- Encompass most of GANs formulations.
- **Standard algorithms** from Variational Inequality can be used for GANs.
- **Theoretical Guarantees** (for convex and stochastic cost functions).

$$\begin{cases} \theta^* = \arg \min_{\theta} L_{\theta}(\theta, \phi^*) \\ \phi^* = \arg \min_{\phi} L_{\phi}(\theta^*, \phi) \end{cases}$$

$$F(\omega^*)^T (\omega - \omega^*) \geq 0 \quad \forall \omega \in \Omega$$



Techniques to optimize VIP (Batch setting)

Standard Algorithms from Variational Inequality

Method 1: **Averaging**

- **Converge** even for “*cycling behavior*”.
- Easy to implement. (out of the training loop)
- Can be combined with any method.

$$\bar{\omega}_T \stackrel{\text{def}}{=} \frac{\sum_{t=0}^{T-1} \rho_t \omega_t}{S_T}, \quad S_T \stackrel{\text{def}}{=} \sum_{t=0}^{T-1} \rho_t.$$

Averaging schemes can be efficiently implemented in an **online** fashion:

$$\bar{\omega}_t = (1 - \tilde{\rho}_t) \bar{\omega}_{t-1} + \tilde{\rho}_t \omega_t \quad \text{where} \quad 0 \leq \tilde{\rho}_t \leq 1.$$

Standard Algorithms from Variational Inequality

Method 1: **Averaging**

- **Converge** even for “cycling behavior”.
- Easy to implement. (out of the training loop)
- Can be combined with any method.

General Online averaging:

$$\bar{\omega}_t = (1 - \tilde{\rho}_t)\bar{\omega}_{t-1} + \tilde{\rho}_t\omega_t \quad \text{where} \quad 0 \leq \tilde{\rho}_t \leq 1.$$

Example 1: **Uniform** averaging

$$\tilde{\rho}_t = \frac{1}{t}, t \geq 0 : \quad \bar{\omega}_T = \frac{1}{T} \sum_{k=0}^{T-1} \omega_k$$

Example 2:

Exponential moving averaging
(EMA)

$$\tilde{\rho}_t = 1 - \beta < 1, t \geq 0 : \quad \bar{\omega}_T = \frac{1}{1 - \beta} \sum_{k=0}^{T-1} \beta^k \omega_k$$

Standard Algorithms from Variational Inequality

Method 1: **Averaging**

- **Converge** even for “*cycling behavior*”.
- Easy to implement. (out of the training loop)
- Can be combined with any method.

General Online averaging:

$$\bar{\omega}_t = (1 - \tilde{\rho}_t)\bar{\omega}_{t-1} + \tilde{\rho}_t\omega_t \quad \text{where} \quad 0 \leq \tilde{\rho}_t \leq 1.$$

Example 1: **Uniform** averaging

$$\tilde{\rho}_t = \frac{1}{t}, t \geq 0 : \quad \bar{\omega}_T = \frac{1}{T} \sum_{k=0}^{T-1} \omega_k$$

Example 2:

Exponential moving
averaging (EMA)

$$\tilde{\rho}_t = 1 - \beta < 1, t \geq 0 : \quad \bar{\omega}_T = (1 - \beta) \sum_{t=1}^T \beta^{T-t} \omega_t + \beta^T \omega_0$$

Standard Algorithms from Variational Inequality

Method 1: **Averaging**

- **Converge** even for “cycling behavior”.
- Easy to implement. (out of the training loop)
- Can be combined with any method.

General Online averaging: $\bar{\omega}_t = (1 - \tilde{\rho}_t)\bar{\omega}_{t-1} + \tilde{\rho}_t\omega_t$ where $0 \leq \tilde{\rho}_t \leq 1$.

Example 1: Uniform averaging $\tilde{\rho}_t = \frac{1}{t}, t \geq 0 : \bar{\omega}_T = \frac{1}{T} \sum_{k=0}^{T-1} \omega_k$

Example 2:
Exponential moving averaging (EMA) $\tilde{\rho}_t = 1 - \beta < 1, t \geq 0 : \bar{\omega}_T = \frac{1}{1 - \beta} \sum_{k=0}^{T-1} \beta^k \omega_k$

Standard Algorithms from Variational Inequality

Method 1: **Averaging**

- **Converge** even for “cycling behavior”.
- Easy to implement. (out of the training loop)
- Can be combined with any method.

General Online averaging: $\bar{\omega}_t = (1 - \tilde{\rho}_t)\bar{\omega}_{t-1} + \tilde{\rho}_t\omega_t$ where $0 \leq \tilde{\rho}_t \leq 1$.

Example 1: Uniform averaging $\tilde{\rho}_t = \frac{1}{t}, t \geq 0 : \bar{\omega}_T = \frac{1}{T} \sum_{k=0}^{T-1} \omega_k$

Example 2:

Exponential moving averaging (EMA)

$$\tilde{\rho}_t = 1 - \beta < 1, t \geq 0 : \bar{\omega}_T = (1 - \beta) \sum_{t=1}^T \beta^{T-t} \omega_t + \beta^T \omega_0$$

Standard Algorithms from Variational Inequality

Method 1: **Averaging**

Simple Minmax problem: $\min_{\theta \in \mathbb{R}} \max_{\phi \in \mathbb{R}} \theta \cdot \phi \implies (\theta^*, \phi^*) = (0, 0)$.

Simultaneous update: $\begin{cases} \theta_{t+1} = \theta_t - \eta \phi_t \\ \phi_{t+1} = \phi_t + \eta \theta_t \end{cases}$, Alternated update: $\begin{cases} \theta_{t+1} = \theta_t - \eta \phi_t \\ \phi_{t+1} = \phi_t + \eta \theta_{t+1} \end{cases}$

Standard Algorithms from Variational Inequality

Method 1: **Averaging**

Simple Minmax problem: $\min_{\theta \in \mathbb{R}} \max_{\phi \in \mathbb{R}} \theta \cdot \phi \implies (\theta^*, \phi^*) = (0, 0)$.

Simultaneous update:
$$\begin{cases} \theta_{t+1} = \theta_t - \eta \phi_t \\ \phi_{t+1} = \phi_t + \eta \theta_t \end{cases},$$

$$(\bar{\theta}_T, \bar{\phi}_T) := \frac{1}{T} \sum_{k=0}^{T-1} (\theta_k, \phi_k) \rightarrow \infty$$

$$(\theta_T, \phi_T) \rightarrow \infty$$

Alternated update:
$$\begin{cases} \theta_{t+1} = \theta_t - \eta \phi_t \\ \phi_{t+1} = \phi_t + \eta \theta_{t+1} \end{cases}$$

$$0 < m \leq \|\theta_T, \phi_T\| \leq M$$

$$(\bar{\theta}_T, \bar{\phi}_T) \rightarrow (0, 0)$$

Standard Algorithms from Variational Inequality

Method 1: **Averaging**

Simultaneous Vs. **Alternating** more developed in
Negative Momentum for Improved Game Dynamics
Gidel, Askari Hemmat, Pezeshki, Lepriol, Huang, Lacoste-Julien and Mitliagkas

Simultaneous update:
$$\begin{cases} \theta_{t+1} = \theta_t - \eta\phi_t \\ \phi_{t+1} = \phi_t + \eta\theta_t \end{cases},$$

$$(\bar{\theta}_T, \bar{\phi}_T) := \frac{1}{T} \sum_{k=0}^{T-1} (\theta_k, \phi_k) \rightarrow \infty$$

$$(\theta_T, \phi_T) \rightarrow \infty$$

Alternated update:
$$\begin{cases} \theta_{t+1} = \theta_t - \eta\phi_t \\ \phi_{t+1} = \phi_t + \eta\theta_{t+1} \end{cases}$$

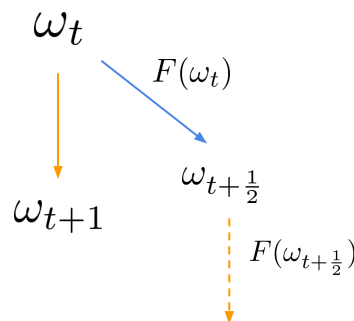
$$0 < m \leq \|\theta_T, \phi_T\| \leq M \quad (\bar{\theta}_T, \bar{\phi}_T) \rightarrow (0, 0)$$

Standard Algorithms from Variational Inequality

Method 2: **Extragradient**

- Step 1: $\omega_{t+\frac{1}{2}} = \omega_t - \gamma_t F(\omega_t)$

- Step 2: $\omega_{t+1} = \omega_t - \gamma_t F(\omega_{t+\frac{1}{2}})$



- **Standard** in the literature.
- Does not require **averaging**.
- *Theoretically and empirically faster.*

Intuition:

1. Game perspective: Look one step in the future and anticipate next move of adversary.
2. Euler's method: Extrapolation is close to an **implicit** method because $\omega_{t+1/2} \approx \omega_{t+1}$

$$\omega_{t+1} - \omega_{t+1/2} = O(\gamma_t^2)$$

Standard Algorithms from Variational Inequality

Method 2: **Extragradient**

Intuition: *Extrapolation is close to an **implicit** method because $\omega_{t+1/2} \approx \omega_{t+1}$*

Implicit step: $\omega_{t+1} = \omega_t - \eta F(\omega_{t+1})$

Unknown:
Require to solve a
non-linear system

Standard Algorithms from Variational Inequality

Method 2: **Extragradient**

Intuition: *Extrapolation is close to an **implicit** method*

$$\min_{\theta \in \mathbb{R}} \max_{\phi \in \mathbb{R}} \theta \cdot \phi \quad \text{and} \quad (\theta^*, \phi^*) = (0, 0).$$

$$\text{Implicit: } \begin{cases} \theta_{t+1} = \theta_t - \eta \phi_{t+1} \\ \phi_{t+1} = \phi_t + \eta \theta_{t+1} \end{cases}, \quad \text{Extrapolation: } \begin{cases} \theta_{t+1} = \theta_t - \eta(\phi_t + \eta \theta_t) \\ \phi_{t+1} = \phi_t + \eta(\theta_t - \eta \phi_t) \end{cases}. \quad (*)$$

Proposition 2. *The squared norm of the iterates $N_t \stackrel{\text{def}}{=} \theta_t^2 + \phi_t^2$, where the update rule of θ_t and ϕ_t are defined in $(*)$, decreases geometrically for any $\eta < 1$ as,*

$$\text{Implicit: } N_{t+1} = (1 - \eta^2 + \eta^4 + \mathcal{O}(\eta^6))N_t, \quad \text{Extrapolation: } N_{t+1} = (1 - \eta^2 + \eta^4)N_t.$$

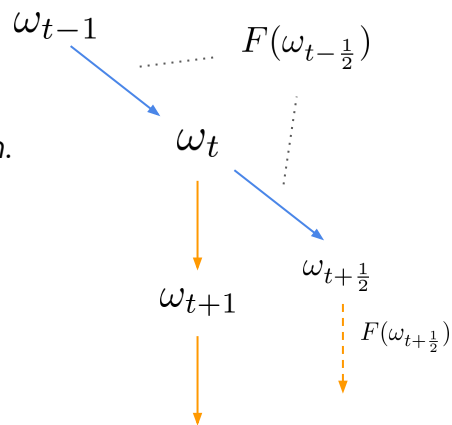
almost the same

Extrapolation from the past: Re-using the gradients

Problem: Extragradient requires to compute **two** gradients at each step.

Solution: **Extrapolation from the past** ← **Re-use** gradient.

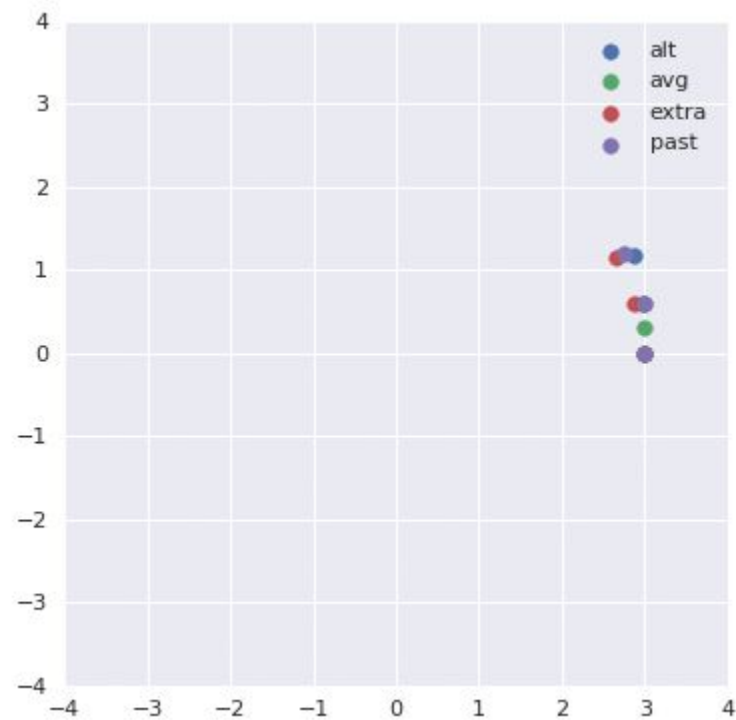
- Step 1: $\omega_{t+\frac{1}{2}} = \omega_t - \gamma_t F(\omega_{t-\frac{1}{2}})$ ← **Re-use** from previous iteration.
- Step 2: $\omega_{t+1} = \omega_t - \gamma_t F(\omega_{t+\frac{1}{2}})$ ← (same as **extragradient**).



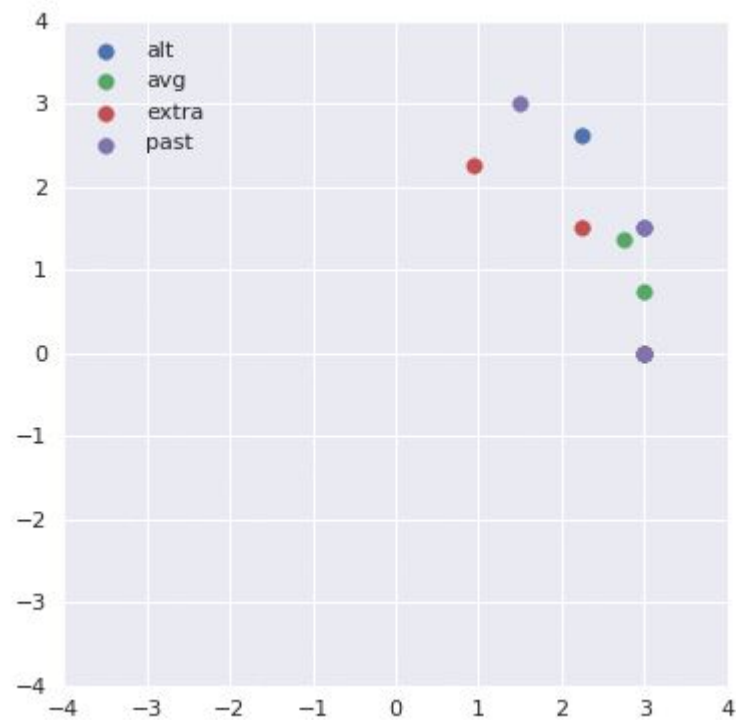
New Method !!!

Related to [Daskalakis et al., 2018]

step-size = 0.2



step-size = 0.5

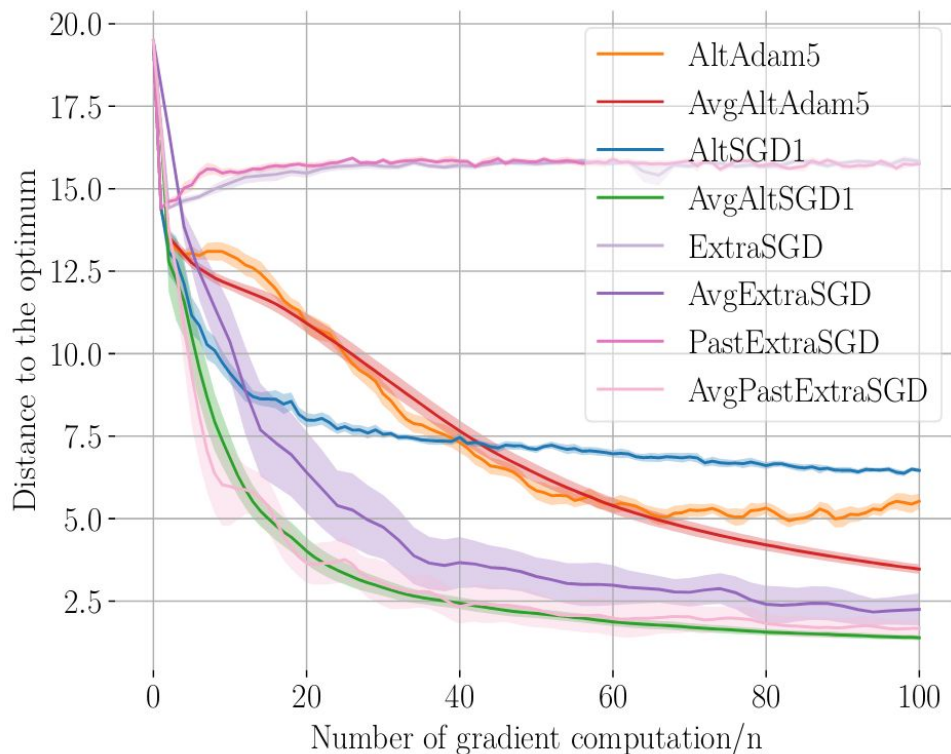


Experimental Results

Experimental Results

Bilinear Stochastic Objective: (with constraints)

$$\frac{1}{n} \sum_{i=1}^n \left(\mathbf{x}^\top \mathbf{M}^{(i)} \mathbf{y} + \mathbf{x}^\top \mathbf{a}^{(i)} + \mathbf{y}^\top \mathbf{b}^{(i)} \right).$$



Algorithm 4 Extra-Adam: proposed Adam with extrapolation step.

input: step-size η , decay rates for moment estimates β_1, β_2 , access to the stochastic gradients $\nabla \ell_t(\cdot)$ and to the projection $P_\Omega[\cdot]$ onto the constraint set Ω , initial parameter ω_0 , averaging scheme $(\rho_t)_{t \geq 1}$
for $t = 0 \dots T - 1$ **do**

Option 1: Standard extrapolation.

Sample new minibatch and compute stochastic gradient: $g_t \leftarrow \nabla \ell_t(\omega_t)$

Option 2: Extrapolation from the past

Load previously saved stochastic gradient: $g_t = \nabla \ell_{t-1/2}(\omega_{t-1/2})$

Update estimate of first moment for extrapolation: $m_{t-1/2} \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$

Update estimate of second moment for extrapolation: $v_{t-1/2} \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$

Correct the bias for the moments: $\hat{m}_{t-1/2} \leftarrow m_{t-1/2} / (1 - \beta_1^{2t-1})$, $\hat{v}_{t-1/2} \leftarrow v_{t-1/2} / (1 - \beta_2^{2t-1})$

Perform *extrapolation* step from iterate at time t : $\omega_{t-1/2} \leftarrow P_\Omega[\omega_t - \eta \frac{m_{t-1/2}}{\sqrt{v_{t-1/2} + \epsilon}}]$

Sample new minibatch and compute stochastic gradient: $g_{t+1/2} \leftarrow \nabla \ell_{t+1/2}(\omega_{t+1/2})$

Update estimate of first moment: $m_t \leftarrow \beta_1 m_{t-1/2} + (1 - \beta_1) g_{t+1/2}$

Update estimate of second moment: $v_t \leftarrow \beta_2 v_{t-1/2} + (1 - \beta_2) g_{t+1/2}^2$

Compute bias corrected for first and second moment: $\hat{m}_t \leftarrow m_t / (1 - \beta_1^{2t})$, $\hat{v}_t \leftarrow v_t / (1 - \beta_2^{2t})$

Perform *update* step from the iterate at time t : $\omega_{t+1} \leftarrow P_\Omega[\omega_t - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}}]$

end for

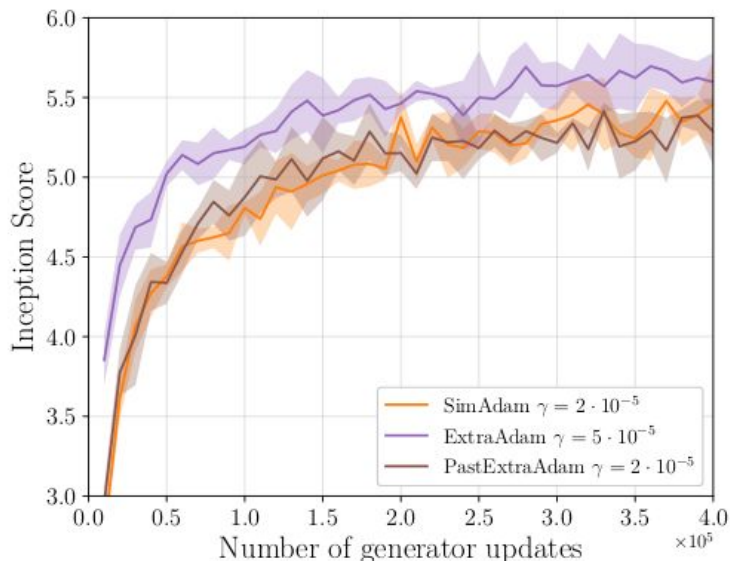
Output: $\omega_{T-1/2}, \omega_T$ or $\bar{\omega}_T = \sum_{t=0}^{T-1} \rho_{t+1} \omega_{t+1/2} / \sum_{t=0}^{T-1} \rho_{t+1}$ (see (8) for online averaging)

Extrapolation
(Adam style)

Update
(Adam style)

Experimental Results: WGAN on CIFAR10

Inception Score vs
nb of generator updates



Inception Score on CIFAR10

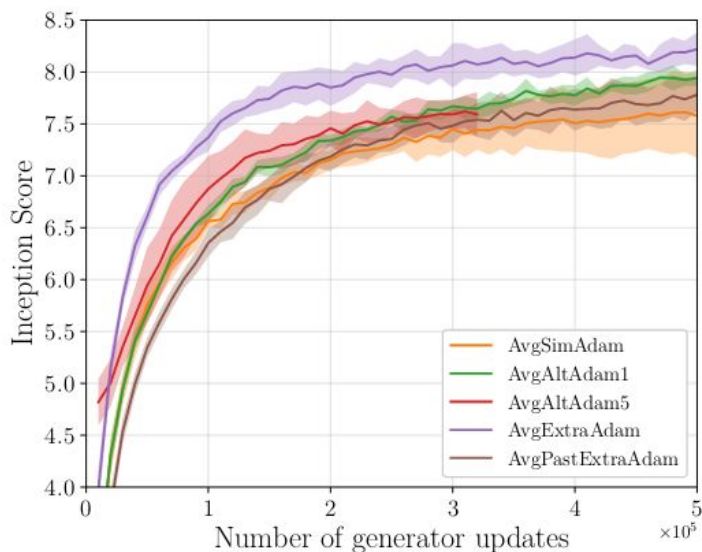
Model	WGAN		
	no averaging	uniform avg	EMA
SimAdam	$6.05 \pm .12$	$5.83 \pm .16$	$6.08 \pm .10$
AltAdam5	$5.45 \pm .08$	$5.72 \pm .06$	$5.49 \pm .05$
ExtraAdam	$6.38 \pm .09$	$6.38 \pm .20$	$6.37 \pm .08$
PastExtraAdam	5.98 ± 0.15	6.07 ± 0.19	6.01 ± 0.11
OptimAdam	5.74 ± 0.10	5.80 ± 0.08	5.78 ± 0.05

↑
Extragradient Methods

↑
Averaging

Experimental Results: WGAN-GP (ResNet) on CIFAR10

Inception Score vs
Number of



Model	WGAN-GP (ResNet)	
Method	no averaging	uniform avg
SimAdam	$7.54 \pm .21$	$7.74 \pm .27$
AltAdam5	$7.20 \pm .06$	$7.67 \pm .15$
ExtraAdam	$7.79 \pm .09$	$8.26 \pm .12$
PastExtraAdam	$7.71 \pm .12$	$7.84 \pm .18$
OptimAdam	$7.80 \pm .07$	$7.99 \pm .12$

↑
Extragradient Methods

↑
Averaging

To sum-up

- GAN can be formulated as a **Variational Inequality**.
- Bring *standard methods* from *optimization literature* to the GAN community.
- **Averaging** helps improve the inception score (further evidence by [Yazici et al. 2018]).
- **Extrapolation** is **faster** and achieve better convergence.
- Introduce **Extrapolation from the past** a **cheaper** version of *extragradient*.
- We can design better algorithm for GANs inspired from Variational Inequality.

Noise in GANs

Reducing Noise in GAN Training with Variance Reduced Extragradient

Tatjana Chavdarova^{*12}, **Gauthier Gidel**^{*1}, François
Fleuret¹², Simon Lacoste-Julien¹

**equal contribution*

¹ Mila, Université de Montréal

² EPFL, IDIAP



Tatjana
Chavdarova



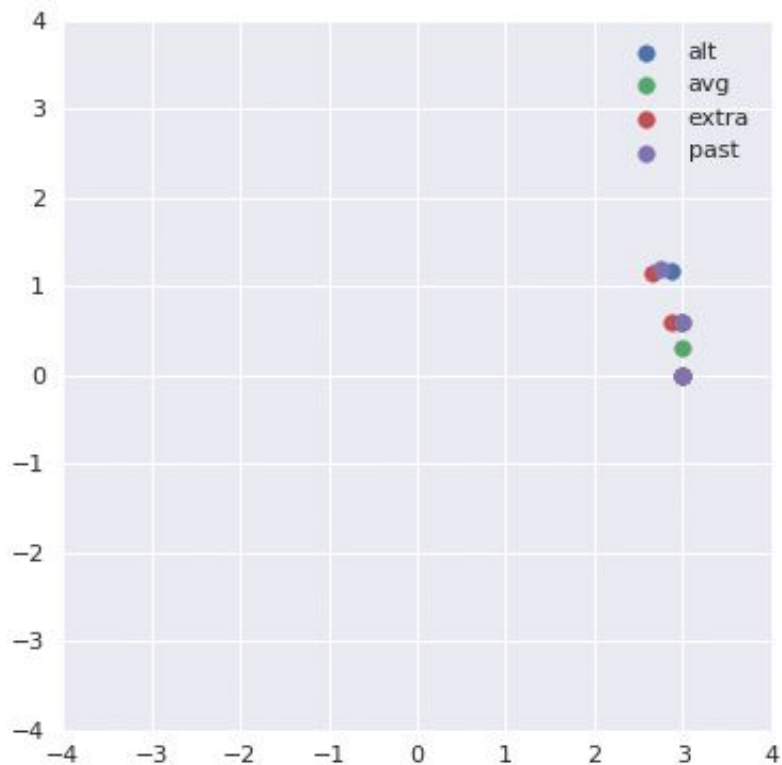
François
Fleuret



Simon
Lacoste-Julien



Reminder: Need for Averaging or/and Extragradients.



Reminder: Need for Averaging or/and Extragradient.



No signal from the average iterate.

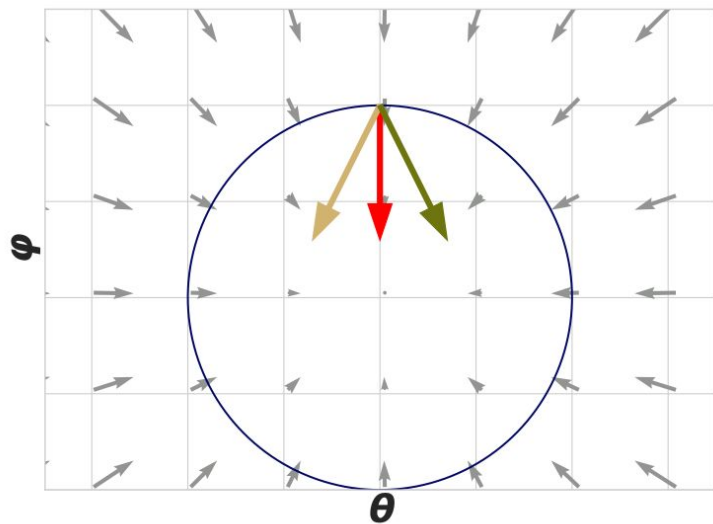
The green sequence **do not stop** at the optimum.

We need **last iterate** convergence.
(Not Convergence of the averaged iterate)

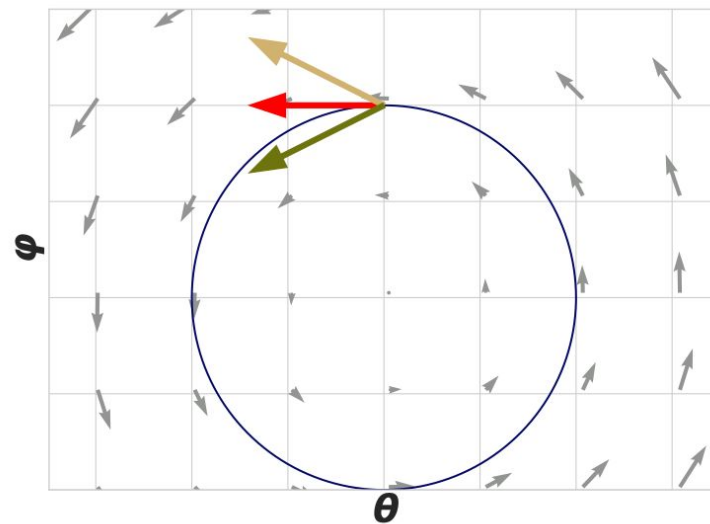
Focus on Extragradient.

Issue: We did not consider **noise**.

Minimization

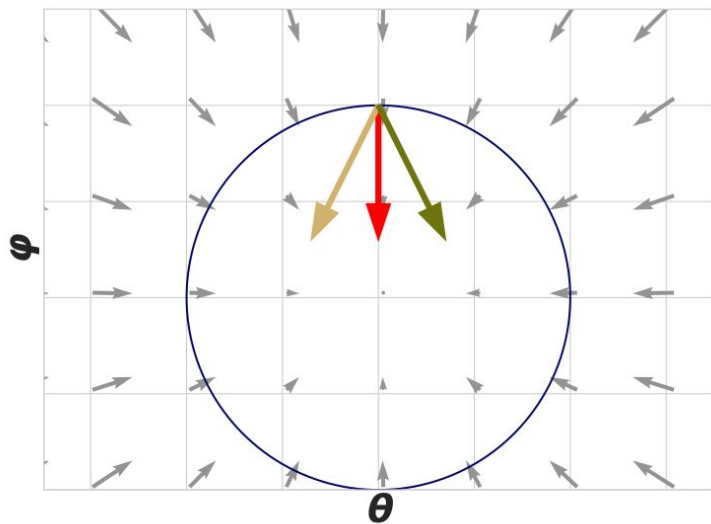


Game

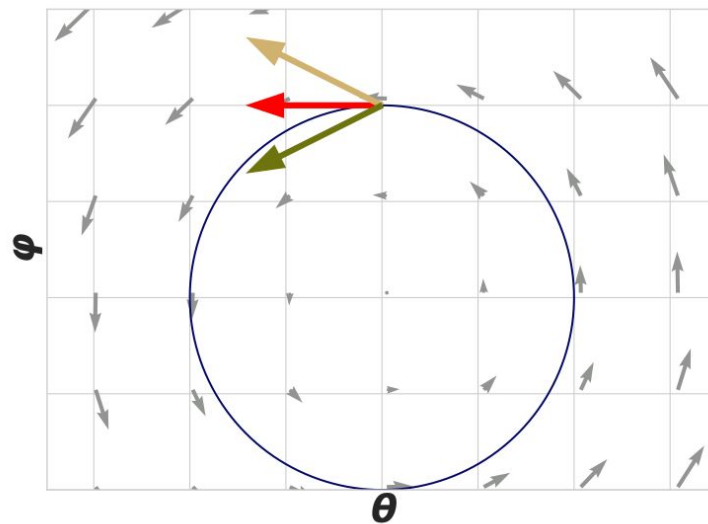


Issue: We did not consider **noise**.


Far from the objective:
“approximately” the right direction



Far from the objective:
Direction with noise can be “bad”.



Standard methods to solve (bilinear) games:

	Gradient method	Extragradient
Batch Method	Diverge to ∞	$\ \omega_t - \omega^*\ \leq (1 - \rho)^t \ \omega_0 - \omega^*\ $
Stochastic Method	 No hope for convergence	????

Noise breaks Extragradient.

Theorem 2 (Noise may induce divergence). *There exists a zero-sum stochastic game such that if $\omega_0 \neq \omega^*$, then for any step-size $\eta > 0$, the iterates (ω_t) computed by the stochastic extragradient method diverge geometrically, i.e., there exists $\rho > 0$, such that $\mathbb{E}[\|\omega_t - \omega^*\|^2] > \|\omega_0 - \omega^*\|^2(1 + \rho)^t$.*

Noise breaks Extragradient.

Theorem 2 (Noise may induce divergence). *There exists a zero-sum stochastic game such that if $\omega_0 \neq \omega^*$, then for any step-size $\eta > 0$, the iterates (ω_t) computed by the stochastic extragradient method diverge geometrically, i.e., there exists $\rho > 0$, such that $\mathbb{E}[\|\omega_t - \omega^*\|^2] > \|\omega_0 - \omega^*\|^2(1 + \rho)^t$.*

Intuition:

$$\min_{\theta \in \mathbb{R}^d} \max_{\varphi \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \theta^\top \mathbf{A}_i \varphi$$

Extragradient Updates:
$$\begin{cases} \theta_{t+1} = \theta_t - \eta \mathbf{A}_I(\varphi_t + \eta \mathbf{A}_J \theta_t) \\ \varphi_{t+1} = \varphi_t + \eta \mathbf{A}_I(\theta_t - \eta \mathbf{A}_J \varphi_t) \end{cases}$$

Noise breaks Extragradient.

Theorem 2 (Noise may induce divergence). *There exists a zero-sum stochastic game such that if $\omega_0 \neq \omega^*$, then for any step-size $\eta > 0$, the iterates (ω_t) computed by the stochastic extragradient method diverge geometrically, i.e., there exists $\rho > 0$, such that $\mathbb{E}[\|\omega_t - \omega^*\|^2] > \|\omega_0 - \omega^*\|^2(1 + \rho)^t$.*

Intuition:

$$\min_{\theta \in \mathbb{R}^d} \max_{\varphi \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \theta^\top A_i \varphi$$

Extragradient Updates:
(Sample i and j)

$$\begin{cases} \theta_{t+1} = \theta_t - \eta A_I(\varphi_t + \eta A_J \theta_t) \\ \varphi_{t+1} = \varphi_t + \eta A_I(\underbrace{\theta_t - \eta A_J \varphi_t}_{\text{Extrapolation part}}) \end{cases}$$

$A_i A_j = 0 \Leftrightarrow$ **No extrapolation**
 \Leftrightarrow **Diverge as GD.**

Extrapolation part

Reducing noise with Variance reduction methods.

- Idea: take advantage of **the finite sum**.
- Finite sum in ML: Expectation of a **finite** number of sample.
- Generator and discriminator losses can be written as:

$$\mathcal{L}^{(\theta)}(\omega) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i^{(\theta)}(\omega), \quad \mathcal{L}^{(\varphi)}(\omega) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i^{(\varphi)}(\omega)$$

SVRG estimate of the gradient.

- Full batch gradient **expensive** but **tractable**.

$$d_i^{(\theta)}(\omega) := \left(\nabla \mathcal{L}_i^{(\theta)}(\omega) - \nabla \mathcal{L}_i^{(\theta)}(\omega^S) \right) + \mu_\theta^S$$

$$d_i^{(\varphi)}(\omega) := \left(\nabla \mathcal{L}_i^{(\varphi)}(\omega) - \nabla \mathcal{L}_i^{(\varphi)}(\omega^S) \right) + \mu_\varphi^S.$$

SVRG estimate of the gradient.

Snapshot network

- Full batch gradient **expensive** but **tractable**.

$$d_i^{(\theta)}(\omega) := \left(\nabla \mathcal{L}_i^{(\theta)}(\omega) - \nabla \mathcal{L}_i^{(\theta)}(\omega^S) \right) + \mu_\theta^S$$

$$d_i^{(\varphi)}(\omega) := \left(\nabla \mathcal{L}_i^{(\varphi)}(\omega) - \nabla \mathcal{L}_i^{(\varphi)}(\omega^S) \right) + \mu_\varphi^S.$$

SVRG estimate of the gradient.

- Full batch gradient **expensive** but **tractable**.

$$d_i^{(\theta)}(\omega) := \left(\nabla \mathcal{L}_i^{(\theta)}(\omega) - \nabla \mathcal{L}_i^{(\theta)}(\omega^S) \right) + \mu_{\theta}^S$$
$$d_i^{(\varphi)}(\omega) := \left(\nabla \mathcal{L}_i^{(\varphi)}(\omega) - \nabla \mathcal{L}_i^{(\varphi)}(\omega^S) \right) + \mu_{\varphi}^S.$$

Snapshot network

Full gradient at the snapshot network

SVRG estimate of the gradient.

- Full batch gradient **expensive** but **tractable**.

$$d_i^{(\theta)}(\omega) := \left(\nabla \mathcal{L}_i^{(\theta)}(\omega) - \nabla \mathcal{L}_i^{(\theta)}(\omega^S) \right) + \mu_{\theta}^S$$
$$d_i^{(\varphi)}(\omega) := \left(\nabla \mathcal{L}_i^{(\varphi)}(\omega) - \nabla \mathcal{L}_i^{(\varphi)}(\omega^S) \right) + \mu_{\varphi}^S.$$

Snapshot network

Full gradient at the snapshot network

- **Unbiased** estimates: $\mathbb{E}[d_i^{(\theta)}(\omega)] = \frac{1}{n} \sum_{i=1}^n \nabla \mathcal{L}_i^{(\theta)}(\omega) = \nabla \mathcal{L}^{(\theta)}(\omega)$

SVRG estimate of the gradient.

- Full batch gradient **expensive** but **tractable**.

$$d_i^{(\theta)}(\omega) := \left(\nabla \mathcal{L}_i^{(\theta)}(\omega) - \nabla \mathcal{L}_i^{(\theta)}(\omega^S) \right) + \mu_{\theta}^S$$
$$d_i^{(\varphi)}(\omega) := \left(\nabla \mathcal{L}_i^{(\varphi)}(\omega) - \nabla \mathcal{L}_i^{(\varphi)}(\omega^S) \right) + \mu_{\varphi}^S.$$

Snapshot network

Full gradient at the snapshot network

- **Unbiased** estimates: $\mathbb{E}[d_i^{(\theta)}(\omega)] = \frac{1}{n} \sum_{i=1}^n \nabla \mathcal{L}_i^{(\theta)}(\omega) = \nabla \mathcal{L}^{(\theta)}(\omega)$
- Compute the snapshot only **once per pass**.

Variance Reduced Extragradient: SVRE

- Combine Extragradient + Variance Reduction for finite sum.

Variance Reduction of Strongly Monotone Games:

Method	Complexity	μ -adaptivity
SVRG	$\ln(1/\epsilon) \times (n + \frac{\bar{L}^2}{\mu^2})$	✗
Acc. SVRG	$\ln(1/\epsilon) \times (n + \sqrt{n} \frac{\bar{L}}{\mu})$	✗
SVRE (This paper)	$\ln(1/\epsilon) \times (n + \frac{\bar{L}}{\mu})$	✓

SVRG and Acc. SVRG are from [Palaniapan and Bach 2016]

Why is this convergence rate not desirable ?

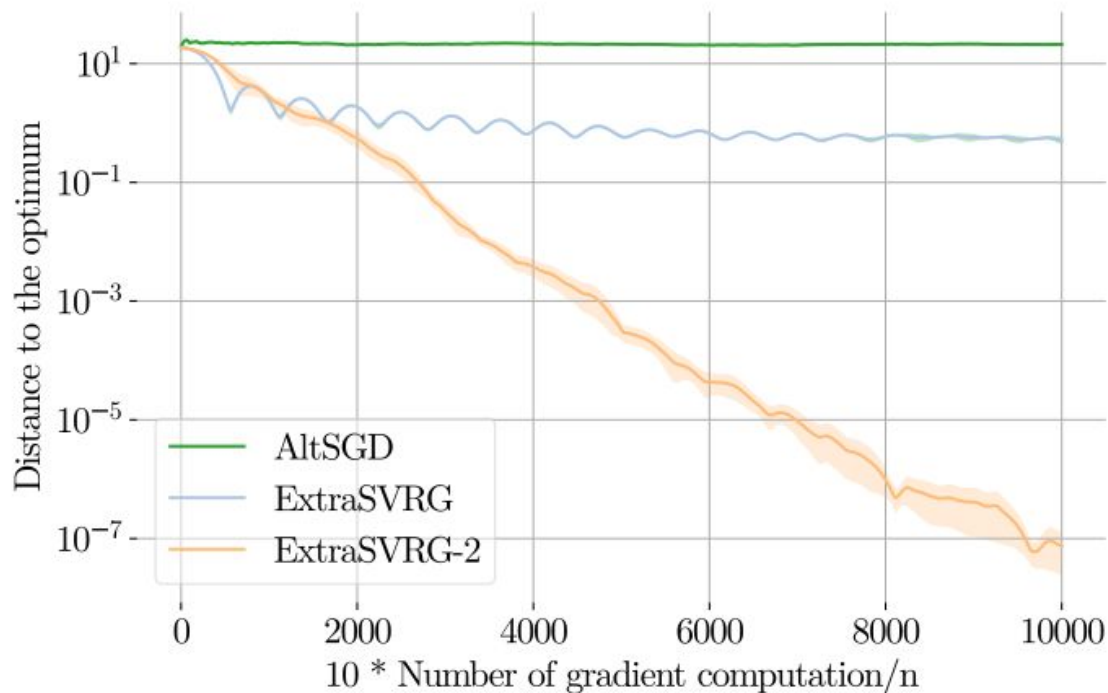
$$\mathbb{E}[Err(\bar{\omega}_T)] \leq O\left(\frac{\sup_{\omega \in \Omega} \|\omega - \omega_0\|^\sigma}{\sqrt{T}}\right) \Rightarrow \text{Does not handle **Unconstrained case.**
No restart possible.}$$

Vs.

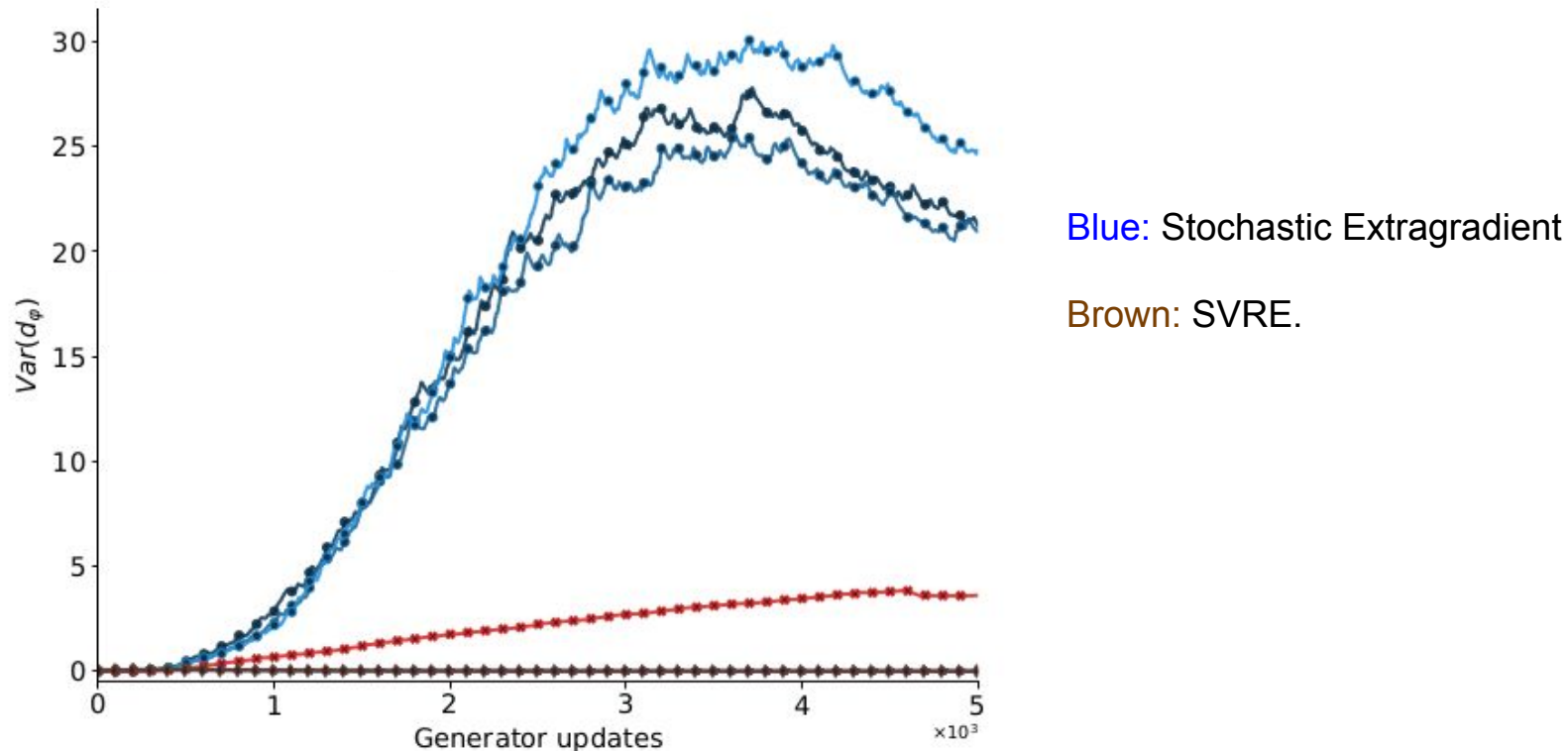
$$\mathbb{E}[Err(\bar{\omega}_T)] \leq O\left(\frac{\|\omega^* - \omega_0\|^\sigma}{\sqrt{T}}\right) \Rightarrow \text{Does handle **Unconstrained case.**
Restart possible.}$$

SVRE on bilinear Game:

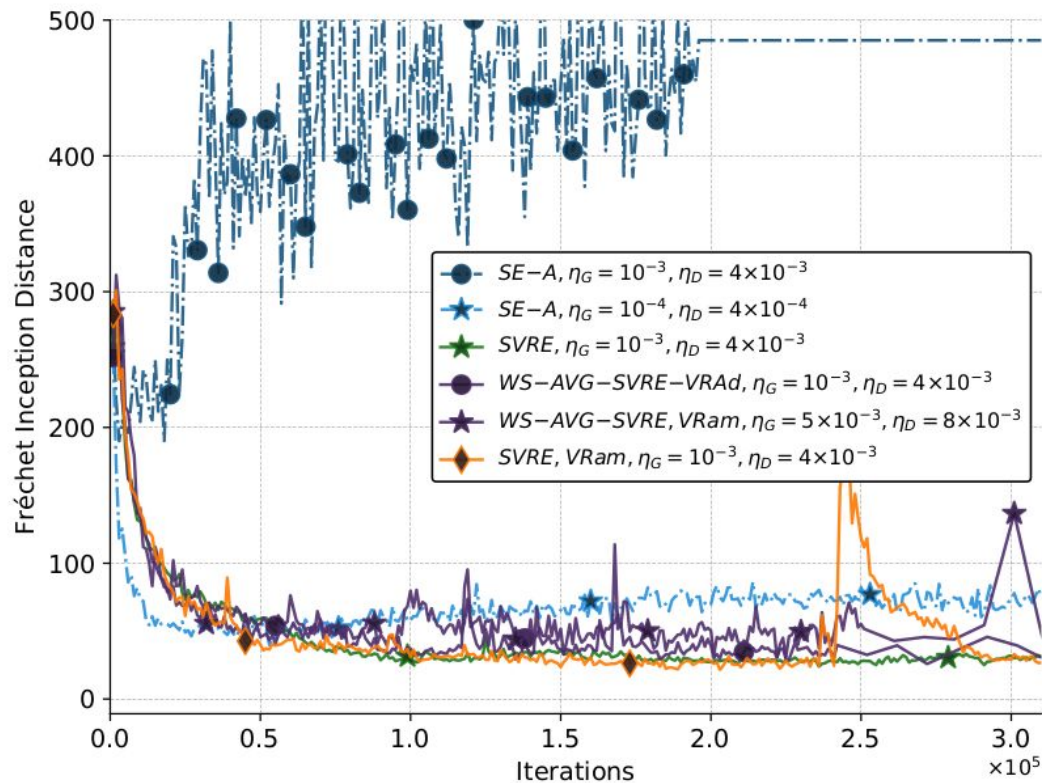
(Exact example where stochastic extragradient breaks)



First point, SVRE effectively reduces the variance:



Second point SVRE allows larger step-sizes: (SVHN)



SE: Stochastic Extragradient.

SVRE: Variance Reduced Extragradient.

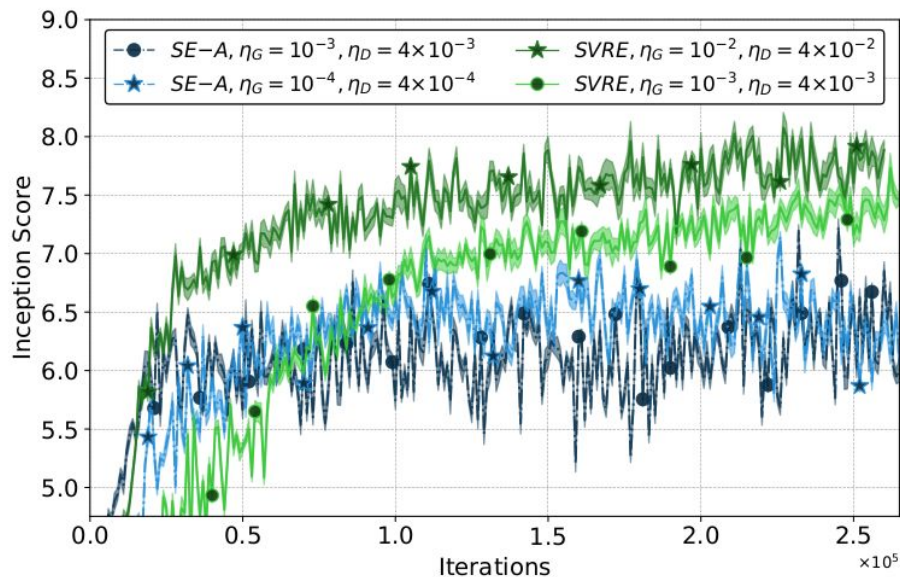
-A: Adam

WS: Warm Start.

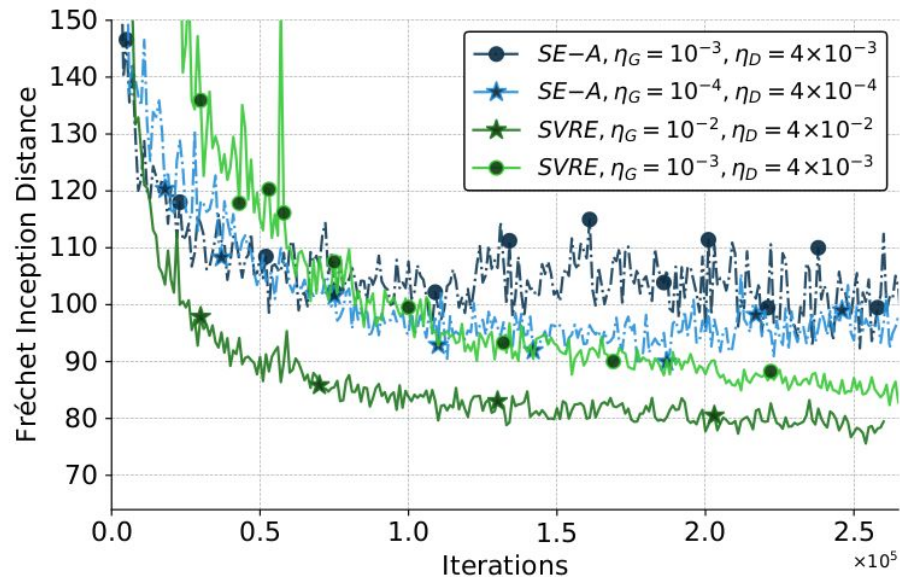
AVG: Average.

-VRAd (VRam): variant of Adam for SVRE.

Second point SVRE allows larger step-sizes: (ImageNet)



(a) IS (higher is better)



(b) FID (lower is better)

To sum-up

- **Noise** may be an issue in GANs.
- Proposed to combine VR + Extragradient to tackle **both** game and noise aspects.
- Unlike in single-objective minimization, we observed that **variance reduction could improve the performance** of deep learning models for GAN training.
- highlights the difference between **game** optimization and **standard minimization**.