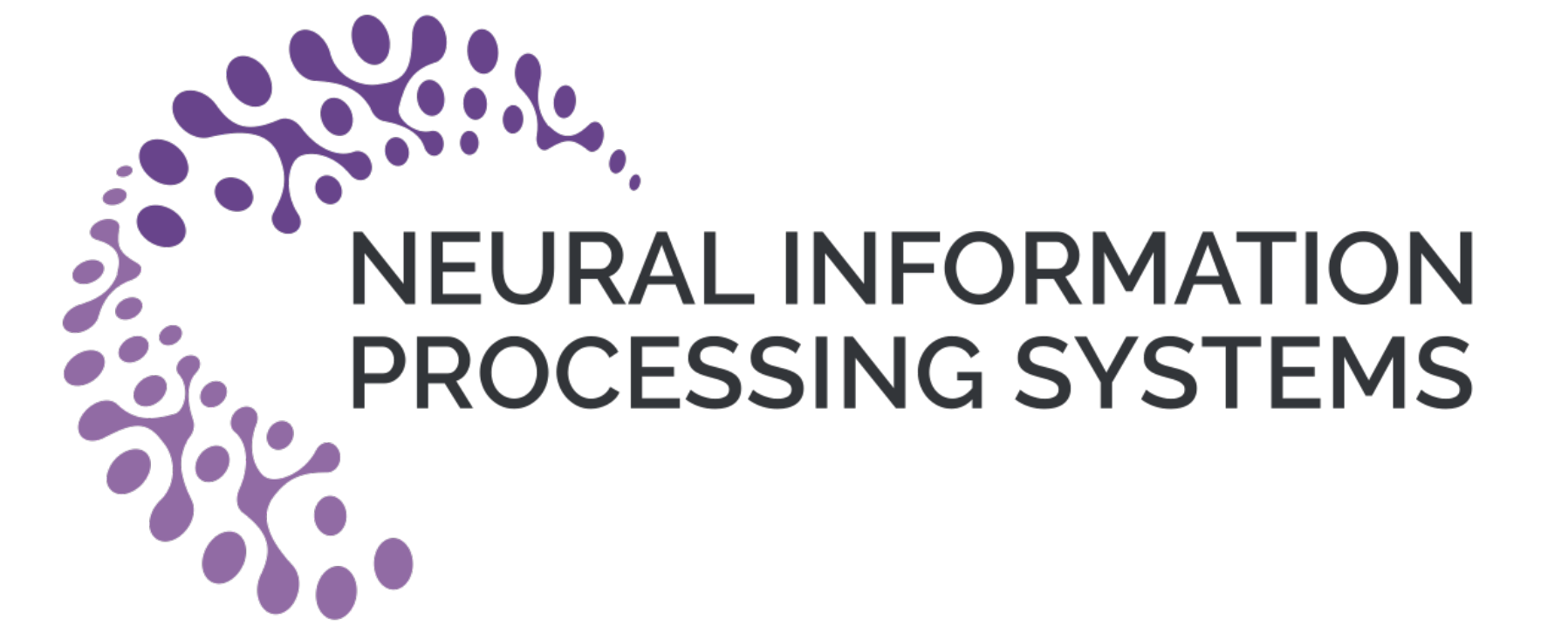


Explore Aggressively, Update Conservatively: Stochastic Extragradient Methods with Variable Stepsize Scaling

Yu-Guan Hsieh, Franck Iutzeler, Jérôme Malick, Panayotis Mertikopoulos (Univ. Grenoble Alpes)



Saddle-point Optimization

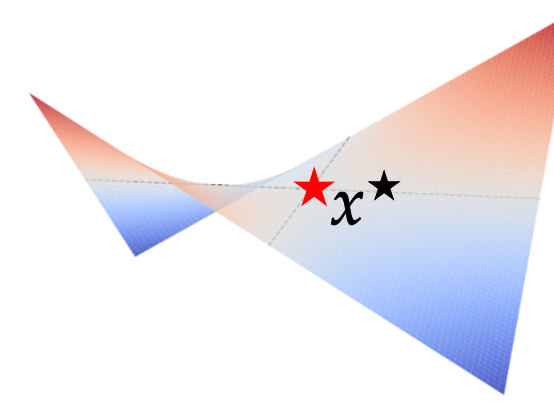
Find $x^* = (\theta^*, \phi^*)$ such that

$$\forall \theta \in \mathbb{R}^{d_1}, \forall \phi \in \mathbb{R}^{d_2}, \mathcal{L}(\theta^*, \phi) \leq \mathcal{L}(\theta^*, \phi^*) \leq \mathcal{L}(\theta, \phi^*)$$

Associated vector field:

$$V(\theta, \phi) = (\nabla_{\theta} \mathcal{L}(\theta, \phi), -\nabla_{\phi} \mathcal{L}(\theta, \phi))$$

First order optimality condition: $V(x^*) = 0$



Applications. • Generative adversarial networks • Adversarial training • Self-play • Robust optimization

Extragradient and its Failure

From gradient to extragradient

$$\mathcal{L} : (\theta, \phi) \in \mathbb{R} \times \mathbb{R} \mapsto \theta \cdot \phi; \quad V(\theta, \phi) = (\phi, -\theta); \quad x^* = (0, 0)$$

Algorithms

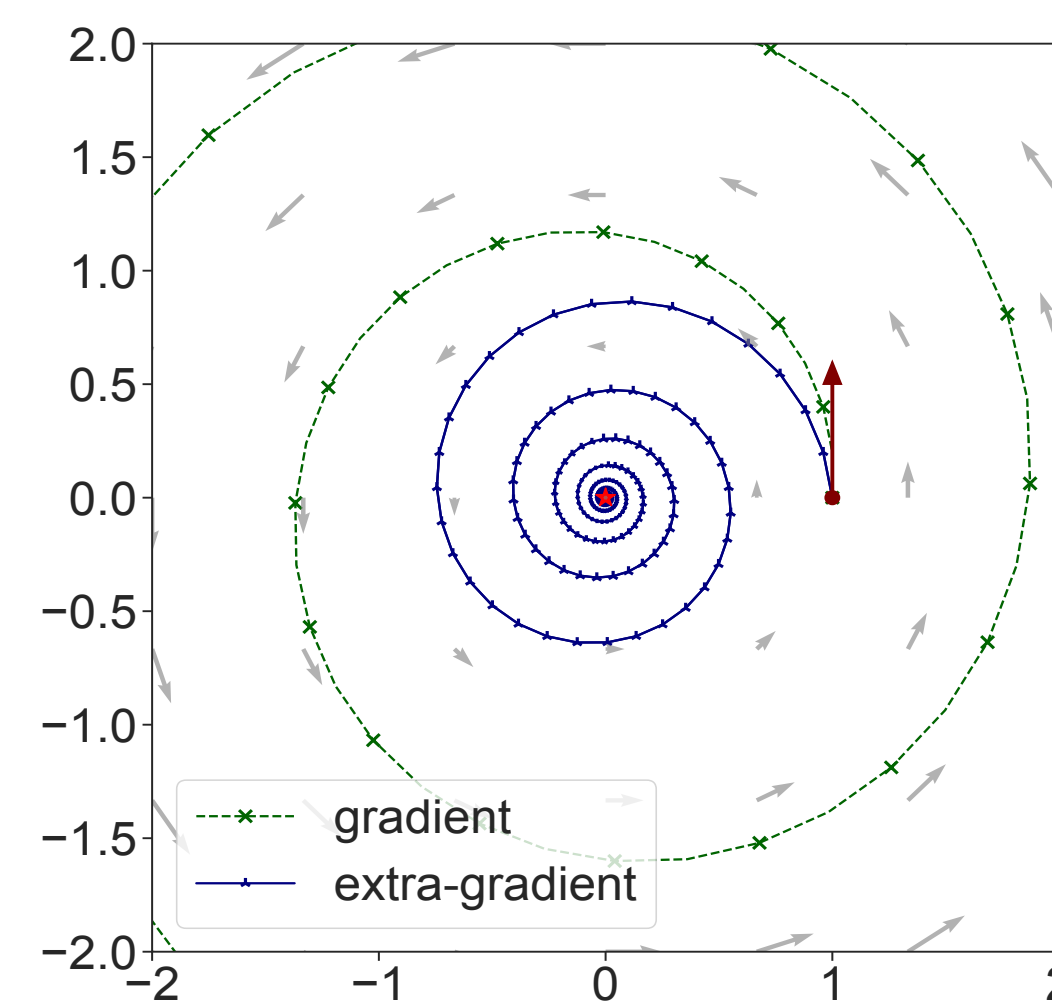
1. Gradient method:

$$X_{t+1} = X_t - \gamma_t V(X_t)$$

2. Extragradient (EG):

$$X_{t+\frac{1}{2}} = X_t - \gamma_t V(X_t)$$

$$X_{t+1} = X_t - \gamma_t V(X_{t+\frac{1}{2}})$$

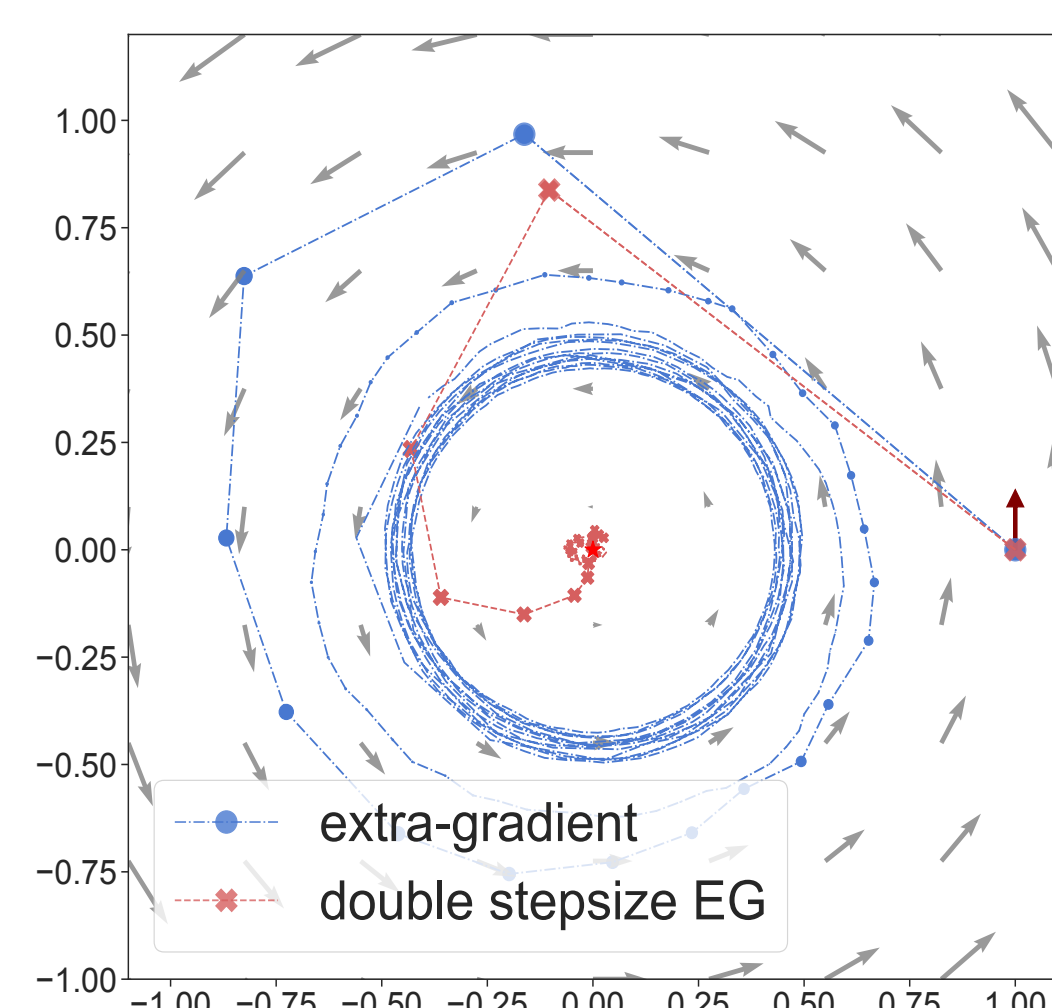


The non-convergence of stochastic EG

$$\hat{V}_t = (\phi_t + \xi_t, -\theta_t); \quad \mathbb{E}[\xi_t] = 0, \quad \mathbb{E}[\xi_t^2] \geq \sigma^2 > 0$$

Proposition. Whatever stepsize is used, running EG with oracle feedback \hat{V}_t leads to $\liminf_{t \rightarrow \infty} \mathbb{E}[\theta_t^2 + \phi_t^2] > 0$.

N.B. EG is known to converge ergodically in $\mathcal{O}(1/\sqrt{t})$ in all stochastic monotone problems. However, in this work, we are interested in its **last-iterate** convergence.



A Remedy with Double Stepsize

DSEG

$$X_{t+\frac{1}{2}} = X_t - \gamma_t \hat{V}_t, \quad X_{t+1} = X_t - \eta_t \hat{V}_{t+\frac{1}{2}}$$

Explore aggressively, update conservatively: $\eta_t \leq \gamma_t$

Assumptions

On the **operator**

1. β -Lipschitz continuity (L): $\|V(x) - V(x')\| \leq \beta \|x - x'\|$
2. Variational stability (VS): $\langle V(x), x - x^* \rangle \geq 0$
3. Error bound (EB): $\exists \tau > 0, \forall x, \|V(x)\| \geq \tau \text{dist}(x, \mathcal{X}^*)$

solution set



On the **noise** ($\forall s \in \mathbb{N}/2, \hat{V}_s = V(X_s) + Z_s$)

1. Unbiasedness: $\mathbb{E}[Z_s | \mathcal{F}_s] = 0$
2. Variance control: $\mathbb{E}[\|Z_s\|^2 | \mathcal{F}_s] \leq (\sigma + \kappa \|X_s - x^*\|)^2, \forall x^* \in \mathcal{X}^*$

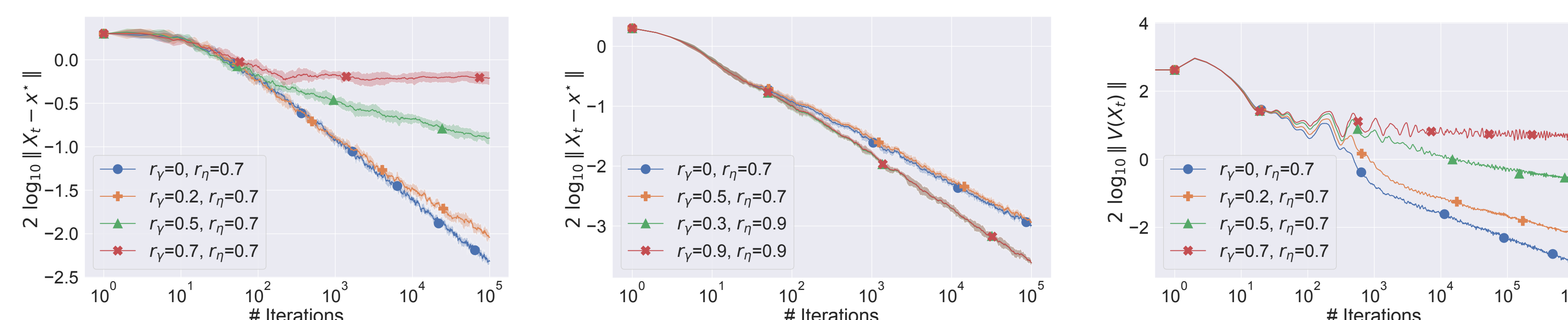
A descent lemma

Assume (L). Let $C_t = 4\gamma_t^2\eta_t\beta + 2\gamma_t^3\eta_t\beta^2 + 4\eta_t^2 + 16\gamma_t^2\eta_t^2\kappa^2$. Then

$$\begin{aligned} \mathbb{E}[\|X_{t+1} - x^*\|^2 | \mathcal{F}_t] &\leq \underbrace{(1 + C_t \kappa^2)}_{\rightarrow 1} \|X_t - x^*\|^2 \\ &\quad - \underbrace{2\eta_t \mathbb{E}[\langle V(X_{t+\frac{1}{2}}), X_{t+\frac{1}{2}} - x^* \rangle | \mathcal{F}_t]}_{\leq 0 \text{ (VS)}} \\ &\quad - \underbrace{\gamma_t \eta_t (1 - \gamma_t^2 \beta^2 - 8\gamma_t \eta_t \kappa^2) \|V(X_t)\|^2}_{< 0 \text{ possible to use (EB)}} + \underbrace{C_t \sigma^2}_{\geq 0} \end{aligned}$$

The decrement term is in $\Theta(\gamma_t \eta_t)$ while the noise term is in $\Theta(\eta_t^2)$

Numerical Illustrations



Left Bilinear zero-sum game
Middle Strongly monotone (biquadratic)
Right Linear quadratic Gaussian WGAN

$$x \sim \mathcal{N}(0, \Sigma)$$

$$G(z) = Yz, D(x) = x^T W x$$



Read the paper

Convergence Result

Asymptotic convergence

Stepsize condition (SC). $\sum_t \gamma_t \eta_t = \infty, \sum_t \eta_t^2 < \infty, \sum_t \gamma_t^2 \eta_t < \infty$

Theorem. Assume (L) + (VS). DSEG with (SC) and $\gamma_t \leq c/\beta$ for some $c < 1$ **converges** to a solution x^* **almost surely**.

Convergence rate

Consider $\gamma_t = \gamma/(t+b)^{r_\gamma}$ and $\eta_t = \eta/(t+b)^{r_\eta}$

$r_\gamma = r_\eta = 1$	Strongly Monotone	$\mathcal{O}(1/t)$
$r_\gamma = 1/3, r_\eta = 2/3$	(EB) + (VS)	$\mathcal{O}(1/t^{1/3})$
$r_\gamma = r_\eta = 0$ ($\eta < \gamma$)	(EB) + (VS) + ($\sigma = 0$)	$\mathcal{O}(e^{-\rho t})$
$r_\gamma = 0, r_\eta = 1$	Affine Monotone	$\mathcal{O}(1/t)$

Local convergence

New!

- (a) Around solution x^* : (L) + (VS) + q -th moment control for Z_t
- (b) $\text{Jac}_V(x^*)$ is defined and invertible \Rightarrow local (EB)

Theorem. Fix tolerance level $\delta \in (0, 1)$. DSEG with close enough initialization and suitable stepsizes guarantees:

- If (a), then $\mathbb{P}(X_t \rightarrow x^*) \geq 1 - \delta$
- If (a)+(b), then there exists event E such that $\mathbb{P}(E) \geq 1 - \delta$ and

$$\mathbb{E}[\|X_t - x^*\|^2 | E] = \mathcal{O}(1/t^{1/3})$$

