



### Negative Momentum for Improved Game Dynamics

Gauthier Gidel\*, Reyhane Askari Hemmat\*, Mohammad Pezeshki, Gabriel Huang, Remi Lepriol, Simon Lacoste-Julien, Ioannis Mitliagkas \*equal contribution

### Simple Min-max smooth game:

Gradient dynamic:



 $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \boldsymbol{\phi}_t$  $\phi_{t+1} = \phi_t + \eta \theta_t$ 





# Simple Min-max smooth game:





Gauthier Gidel, Workshop on learning and strategic behavior, August 22, 2018



# $\min_{\boldsymbol{\theta} \in \mathbb{R}} \max_{\boldsymbol{\phi} \in \mathbb{R}} \boldsymbol{\theta} \cdot \boldsymbol{\phi}$



# Simple Min-max smooth game:

#### Gradient dynamic:



# $\min_{\boldsymbol{\theta} \in \mathbb{R}} \max_{\boldsymbol{\phi} \in \mathbb{R}} \boldsymbol{\theta} \cdot \boldsymbol{\phi}$



Gauthier Gidel,

Workshop on learning and strategic behavior, August 22, 2018



Method	update rule	Convergence (iterates)
Simultaneous Gradient	$oldsymbol{ heta}_{t+1} = oldsymbol{ heta}_t - \eta oldsymbol{\phi}_t \ oldsymbol{\phi}_{t+1} = oldsymbol{\phi}_t + \eta oldsymbol{ heta}_t$	×





Method	update rule	Convergence (iterates)
Simultaneous Gradient	$oldsymbol{ heta}_{t+1} = oldsymbol{ heta}_t - \eta oldsymbol{\phi}_t \ oldsymbol{\phi}_{t+1} = oldsymbol{\phi}_t + \eta oldsymbol{ heta}_t$	×
Alternated Gradient	$oldsymbol{ heta}_{t+1} = oldsymbol{ heta}_t - \eta_1 oldsymbol{\phi}_t \ oldsymbol{\phi}_{t+1} = oldsymbol{\phi}_t + \eta_2 oldsymbol{ heta}_{t+1}$	×





Method	update rule	Convergence (iterates)
Simultaneous Gradient	$oldsymbol{ heta}_{t+1} = oldsymbol{ heta}_t - \eta oldsymbol{\phi}_t \ oldsymbol{\phi}_{t+1} = oldsymbol{\phi}_t + \eta oldsymbol{ heta}_t$	×
Alternated Gradient	$oldsymbol{ heta}_{t+1} = oldsymbol{ heta}_t - \eta_1 oldsymbol{\phi}_t \ oldsymbol{\phi}_{t+1} = oldsymbol{\phi}_t + \eta_2 oldsymbol{ heta}_{t+1}$	×
Alternated Gradient + negative momentum	$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_1 \boldsymbol{\phi}_t + \beta_1 (\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1}) \\ \boldsymbol{\phi}_{t+1} = \boldsymbol{\phi}_t + \eta_2 \boldsymbol{\theta}_{t+1} + \beta_2 (\boldsymbol{\phi}_t - \boldsymbol{\phi}_{t-1})$	$\checkmark$

Gauthier Gidel, Workshop on learning and strategic behavior, August 22, 2018

ila



Method	update rule	Convergence (iterates)
Simultaneous Gradient	$ \begin{aligned} \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \eta M \boldsymbol{\phi}_t \\ \boldsymbol{\phi}_{t+1} &= \boldsymbol{\phi}_t + \eta M^\top \boldsymbol{\theta}_t \end{aligned} $	×
Simultaneous Gradient + negative momentum	$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta M \boldsymbol{\phi}_t + \beta (\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1}) \\ \boldsymbol{\phi}_{t+1} = \boldsymbol{\phi}_t + \eta M^\top \boldsymbol{\theta}_t + \beta (\boldsymbol{\phi}_t - \boldsymbol{\phi}_{t-1})$	$\bigstar$ (but less worse)
Alternated Gradient	$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_1 M \boldsymbol{\phi}_t \\ \boldsymbol{\phi}_{t+1} = \boldsymbol{\phi}_t + \eta_2 M^\top \boldsymbol{\theta}_{t+1}$	×
Alternated Gradient + negative momentum	$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_1 M \boldsymbol{\phi}_t + \beta_1 (\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1}) \\ \boldsymbol{\phi}_{t+1} = \boldsymbol{\phi}_t + \eta_2 M^\top \boldsymbol{\theta}_{t+1} + \beta_2 (\boldsymbol{\phi}_t - \boldsymbol{\phi}_{t-1})$	✓







Gauthier Gidel,

Workshop on learning and strategic behavior, August 22, 2018

Universi

de Montréal

Two players aim to minimize their respective cost functions:

$$oldsymbol{ heta}^* \in rgmin_{oldsymbol{ heta}\inoldsymbol{ heta}} \mathcal{L}^{(oldsymbol{ heta})}(oldsymbol{ heta},oldsymbol{\phi}^*) \quad ext{and} \quad oldsymbol{\phi}^* \in rgmin_{oldsymbol{ heta}\inoldsymbol{\phi}} \mathcal{L}^{(oldsymbol{\phi})}(oldsymbol{ heta}^*,oldsymbol{\phi})$$





Two players aim to minimize their respective cost functions:

$$oldsymbol{ heta}^* \in rgmin_{oldsymbol{ heta}\inoldsymbol{ heta}} \mathcal{L}^{(oldsymbol{ heta})}(oldsymbol{ heta},oldsymbol{\phi}^*) \quad ext{and} \quad oldsymbol{\phi}^* \in rgmin_{oldsymbol{\phi}\inoldsymbol{\phi}} \mathcal{L}^{(oldsymbol{\phi})}(oldsymbol{ heta}^*,oldsymbol{\phi})$$

Examples:

• Simple class of zero-sum games: (  $\mathcal{L}^{(\theta)} = -\mathcal{L}^{(\phi)}$ )

 $\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \max_{\boldsymbol{\phi} \in \mathbb{R}^p} \alpha \|\boldsymbol{\theta}\|_2^2 + (1-\alpha)\boldsymbol{\theta}^\top \boldsymbol{M} \boldsymbol{\phi} - \alpha \|\boldsymbol{\phi}\|_2^2, \quad \alpha \in [0,1], \ \boldsymbol{M} \in \mathbb{R}^{d \times p}$ 





Two players aim to minimize their respective cost functions:

$$oldsymbol{ heta}^* \in rgmin_{oldsymbol{ heta}\inoldsymbol{ heta}} \mathcal{L}^{(oldsymbol{ heta})}(oldsymbol{ heta},oldsymbol{\phi}^*) \quad ext{and} \quad oldsymbol{\phi}^* \in rgmin_{oldsymbol{ heta}\inoldsymbol{\phi}} \mathcal{L}^{(oldsymbol{\phi})}(oldsymbol{ heta}^*,oldsymbol{\phi})$$

Examples:

• Simple class of zero-sum games: (  $\mathcal{L}^{(\theta)} = -\mathcal{L}^{(\phi)}$ )

 $\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \max_{\boldsymbol{\phi} \in \mathbb{R}^p} \alpha \|\boldsymbol{\theta}\|_2^2 + (1-\alpha)\boldsymbol{\theta}^\top \boldsymbol{M} \boldsymbol{\phi} - \alpha \|\boldsymbol{\phi}\|_2^2, \quad \alpha \in [0,1], \ \boldsymbol{M} \in \mathbb{R}^{d \times p}$ 

• Generative Adversarial Networks:

$$\mathcal{L}^{(\theta)}(\theta,\phi) \stackrel{\text{def}}{=} - \underset{\boldsymbol{x}' \sim q_{\theta}}{\mathbb{E}} \log f_{\phi}(\boldsymbol{x}') \quad \text{and} \quad \mathcal{L}^{(\phi)}(\theta,\phi) \stackrel{\text{def}}{=} - \underset{\boldsymbol{x} \sim p}{\mathbb{E}} \log f_{\phi}(\boldsymbol{x}) - \underset{\boldsymbol{x}' \sim q_{\theta}}{\mathbb{E}} \log(1 - f_{\phi}(\boldsymbol{x}'))$$

(non-saturating GAN from Goodfellow et al. 2014)



Two players aim to minimize their respective cost functions:

$$\boldsymbol{\theta}^* \in \mathop{\mathrm{arg\,min}}_{\boldsymbol{\theta} \in \boldsymbol{\theta}} \mathcal{L}^{(\boldsymbol{\theta})}(\boldsymbol{\theta}, \boldsymbol{\phi}^*) \quad \text{and} \quad \boldsymbol{\phi}^* \in \mathop{\mathrm{arg\,min}}_{\boldsymbol{\phi} \in \boldsymbol{\phi}} \mathcal{L}^{(\boldsymbol{\phi})}(\boldsymbol{\theta}^*, \boldsymbol{\phi})$$

Dynamics of gradient based method depends on the gradient vector fields:

$$oldsymbol{v}(oldsymbol{\phi},oldsymbol{ heta}) \coloneqq egin{bmatrix} 
abla_{oldsymbol{\phi}}\mathcal{L}^{(oldsymbol{\phi})}(oldsymbol{\phi},oldsymbol{ heta}) \ 
abla_{oldsymbol{ heta}}\mathcal{L}^{(oldsymbol{ heta})}(oldsymbol{\phi},oldsymbol{ heta}) \end{bmatrix}$$



Two players aim to minimize their respective cost functions:

$$oldsymbol{ heta}^* \in rgmin_{oldsymbol{ heta}\inoldsymbol{ heta}} \mathcal{L}^{(oldsymbol{ heta})}(oldsymbol{ heta},oldsymbol{\phi}^*) \quad ext{and} \quad oldsymbol{\phi}^* \in rgmin_{oldsymbol{\phi}\inoldsymbol{\phi}} \mathcal{L}^{(oldsymbol{\phi})}(oldsymbol{ heta}^*,oldsymbol{\phi})$$

Dynamics of gradient based method depends on the gradient vector fields:

$$oldsymbol{v}(oldsymbol{\phi},oldsymbol{ heta}) \coloneqq egin{bmatrix} 
abla_{oldsymbol{\phi}} \mathcal{L}^{(oldsymbol{\phi})}(oldsymbol{\phi},oldsymbol{ heta}) \ 
abla_{oldsymbol{ heta}} \mathcal{L}^{(oldsymbol{ heta})}(oldsymbol{\phi},oldsymbol{ heta}) \end{bmatrix}$$

And its associated Jacobian,

$$\nabla \boldsymbol{v}(\boldsymbol{\phi}, \boldsymbol{\theta}) \coloneqq \begin{bmatrix} \nabla_{\boldsymbol{\phi}}^{2} \mathcal{L}^{(\boldsymbol{\phi})}(\boldsymbol{\phi}, \boldsymbol{\theta}) & \nabla_{\boldsymbol{\phi}} \nabla_{\boldsymbol{\theta}} \mathcal{L}^{(\boldsymbol{\phi})}(\boldsymbol{\phi}, \boldsymbol{\theta}) \\ \nabla_{\boldsymbol{\phi}} \nabla_{\boldsymbol{\theta}} \mathcal{L}^{(\boldsymbol{\theta})}(\boldsymbol{\phi}, \boldsymbol{\theta})^{T} & \nabla_{\boldsymbol{\theta}}^{2} \mathcal{L}^{(\boldsymbol{\theta})}(\boldsymbol{\phi}, \boldsymbol{\theta}) \end{bmatrix}$$

Mila



#### Fixed point dynamics

Gradient method is defined as the repetition of the operator:

$$F_{\eta}(\boldsymbol{\phi}, \boldsymbol{\theta}) = \begin{bmatrix} \boldsymbol{\phi} & \boldsymbol{\theta} \end{bmatrix}^{\top} - \eta \ \boldsymbol{v}(\boldsymbol{\phi}, \boldsymbol{\theta}).$$

Thus, the sequence computed is

$$(\boldsymbol{\phi}_t, \boldsymbol{\theta}_t) = \underbrace{F_{\eta} \circ \ldots \circ F_{\eta}}_{t}(\boldsymbol{\phi}_0, \boldsymbol{\theta}_0) = F_{\eta}^{(t)}(\boldsymbol{\phi}_0, \boldsymbol{\theta}_0)$$





#### Fixed point dynamics

Gradient method is defined as the repetition of the operator:

$$F_{\eta}(\boldsymbol{\phi}, \boldsymbol{\theta}) = \begin{bmatrix} \boldsymbol{\phi} & \boldsymbol{\theta} \end{bmatrix}^{\top} - \eta \ \boldsymbol{v}(\boldsymbol{\phi}, \boldsymbol{\theta}).$$

Thus, the sequence computed is

$$(\boldsymbol{\phi}_t, \boldsymbol{\theta}_t) = \underbrace{F_{\eta} \circ \ldots \circ F_{\eta}}_{t}(\boldsymbol{\phi}_0, \boldsymbol{\theta}_0) = F_{\eta}^{(t)}(\boldsymbol{\phi}_0, \boldsymbol{\theta}_0)$$

We aim to converge to a Nash Equilibrium:





#### Tuning the step size

Jacobian of our fixed point operator:

$$\nabla F_{\eta}(\boldsymbol{\theta}^*, \boldsymbol{\phi}^*) = I_n - \eta \nabla \boldsymbol{v}(\boldsymbol{\theta}^*, \boldsymbol{\phi}^*)$$

- To have fixed point we need  $\, 
  abla v(\phi^*, oldsymbol{ heta}^*)$ to be definite positive.
- Thus, small enough step-size  $\implies$  Eigenvalues in the unit disk.
- Want to find optimal step-size.





#### Fixed point dynamics

Jacobian of our fixed point operator:  $\nabla F_{\eta}(\boldsymbol{\theta}^*, \boldsymbol{\phi}^*) = I_n - \eta \nabla \boldsymbol{v}(\boldsymbol{\theta}^*, \boldsymbol{\phi}^*)$ 

**Theorem 1** (Proposition 4.4.1 Bertsekas (1999)). Let  $\mu_{max}$  be the largest eigenvalue (in magnitude) of  $\nabla F_{\eta}(\phi^*, \theta^*)$ . If  $\rho_{max} := |\mu_{max}| < 1$  then, for  $(\phi_0, \theta_0)$  in a neighborhood of  $(\phi^*, \theta^*)$ , the sequence  $(\phi_t, \theta_t)$  converges to the local Nash equilibrium at a rate of  $\mathcal{O}((\rho_{max} + \epsilon)^t), \forall \epsilon > 0$ .

- Local convergence.
- Stationary point may not be a Nash equilibrium. (See Adolphs et al. 2018)
- But any Nash equilibrium is an stationary point.
- In this talk: local results on stationary points.





Mila



Recall Polyak's momentum :

$$x_{t+1} = x_t - \eta v(x_t) + \beta (x_t - x_{t-1}), \quad x_t = (\theta_t, \phi_t)$$

Fixed point operator requires a state augmentation : (because need previous iterates)

$$F_{\eta,\beta}(\boldsymbol{x}_t, \boldsymbol{x}_{t-1}) := \begin{bmatrix} \boldsymbol{I}_n & \boldsymbol{0}_n \\ \boldsymbol{I}_n & \boldsymbol{0}_n \end{bmatrix} \begin{bmatrix} \boldsymbol{x}_t \\ \boldsymbol{x}_{t-1} \end{bmatrix} - \eta \begin{bmatrix} \boldsymbol{v}(\boldsymbol{x}_t) \\ \boldsymbol{0}_n \end{bmatrix} + \beta \begin{bmatrix} \boldsymbol{I}_n & -\boldsymbol{I}_n \\ \boldsymbol{0}_n & \boldsymbol{0}_n \end{bmatrix} \begin{bmatrix} \boldsymbol{x}_t \\ \boldsymbol{x}_{t-1} \end{bmatrix}$$



**Theorem 3.** The eigenvalues of  $\nabla F_{\eta,\beta}(\phi^*, \theta^*)$  are

$$\mu_{\pm}(\beta,\eta,\lambda) := (1 - \eta\lambda + \beta) \frac{1 \pm \Delta^{\frac{1}{2}}}{2},\tag{9}$$

where  $\Delta := 1 - \frac{4\beta}{(1-\eta\lambda+\beta)^2}$ ,  $\lambda \in Sp(\nabla v(\phi^*, \theta^*))$  and  $\Delta^{\frac{1}{2}}$  is the complex square root of  $\Delta$  with positive real part<sup>3</sup>. Moreover we have the following Taylor approximation,

$$\frac{\mu_{+}(\beta,\eta,\lambda) = 1 - \eta\lambda - \beta \frac{\eta\lambda}{1 - \eta\lambda} + O(\beta^{2}) \quad and \quad \mu_{-}(\beta,\eta,\lambda) = \frac{\beta}{1 - \eta\lambda} + O(\beta^{2})$$
(10)

<sup>3</sup> If  $\Delta$  is a negative real number we set  $\Delta^{\frac{1}{2}} := i\sqrt{-\Delta}$ 



- Fixed momentum.
   (- 0.25)
- Step-size is **not** fixed.

ila

• Helps when the eigenvalue has large imaginary part.





#### What happens in practice?



Fashion MNIST:



Gauthier Gidel,

Workshop on learning and strategic behavior, August 22, 2018



#### What happen in practice?

CIFAR-10:







To sum up:

- Negative momentum seems to improve the behaviour of the "bad" eigenvalues.
- If small enough seems to always help.
- It also allows larger step-size.





#### Thank you !

If you are interested in that topic:

• NIPS Workshop : Smooth Games Optimization and Machine Learning

Co-organized with:

Simon Lacoste-Julien · Ioannis Mitliagkas · Vasilis Syrgkanis · Eva Tardos

· Leon Bottou · Sebastian Nowozin

Soon : Call for contributions !!!



