

---

# Negative Momentum for Improved Game Dynamics

---

Gauthier Gidel\* Reyhane Askari Hemmat\* Mohammad Pezeshki Gabriel Huang  
Rémi Lepriol Simon Lacoste-Julien† Ioannis Mitliagkas  
Mila & DIRO, Université de Montréal

## Abstract

Games generalize the single-objective optimization paradigm by introducing different objective functions for different players. Differentiable games often proceed by simultaneous or alternating gradient updates. In machine learning, games are gaining new importance through formulations like generative adversarial networks (GANs) and actor-critic systems. However, compared to single-objective optimization, game dynamics are more complex and less understood. In this paper, we analyze gradient-based methods with momentum on simple games. Next, we show empirically that alternating gradient updates with a negative momentum term achieves convergence on the notoriously difficult to train saturating GANs.

## 1 Introduction and Background

Recent advances in machine learning are largely driven by the success of gradient-based optimization methods for the training process. Games generalize the standard optimization framework by introducing different objective functions for different optimizing agents, known as *players*. For example generative adversarial networks (GANs) [Goodfellow et al., 2014] use a two-player game formulation. We are commonly interested in finding a local *Nash equilibrium*: a set of parameters from which no player can (locally and unilaterally) improve its objective function. Games with differentiable objectives often proceed by simultaneous or alternating gradient steps on the players’ objectives.

In this work we are interested in studying the effect of two particular algorithmic choices: (i) the choice between simultaneous and alternating updates, and (ii) the choice of step size and momentum value. We summarize our main contributions as follows: We show that the alternating gradient method with negative momentum is the only setting within our study parameters (Tab. 1) that converges on a bilinear smooth game. Using a zero or positive momentum value, or doing simultaneous updates in such games fails to converge. We also show in §3 that, for general dynamics, when the eigenvalues of the Jacobian have a large imaginary part, negative momentum can improve the local convergence properties of the gradient method. Finally, we confirm the benefits of negative momentum for training GANs with the notoriously ill-behaved saturating loss on CIFAR10 and FASHION-MNIST datasets.

**Related work** A lot of work has been done in the context of understanding momentum and its variants but all are restricted to minimization problems [Polyak, 1964, Qian, 1999, Nesterov, 2013, Sutskever et al., 2013]. Balduzzi et al. [2018] develop new methods to understand the dynamics of general games and propose Symplectic Gradient Adjustment which helps in discovering stable fixed points in general games. Mescheder et al. [2017] show how presence of eigenvalues with zero real-part and large imaginary-part result in oscillatory behavior. Nagarajan and Kolter [2017] analyze the local stability of GANs and show that during training of a GAN, the eigenvalues of the Jacobian of the vector field are pushed away from one along the real axis.

---

\*Equal contribution, correspondence to <gauthier.gidel,reyhane.askari.hemmat>umontreal.ca

†CIFAR fellow

Method	$\beta$	Bounded	Converges
Simultaneous	$\beta \in \mathbb{R}$	✗	✗
Alternated	$>0$	✗	✗
	$0$	✓	✗
	$<0$	✓	✓

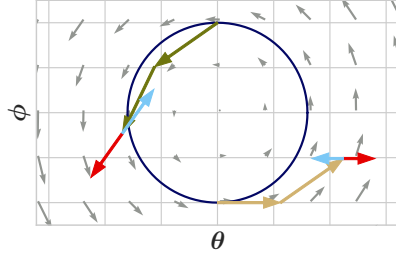


Figure 1: Effect of Gradient Methods on an unconstrained bilinear example:  $\min_{\theta} \max_{\varphi} \theta^{\top} A \varphi$ .

**Game theory formulation of GANs** From a game theory point of view, GAN training is a differentiable two-player game: the discriminator  $D_{\varphi}$  aims at minimizing its cost function  $\mathcal{L}^{(\varphi)}$  and the generator  $G_{\theta}$  aims at minimizing its own cost function  $\mathcal{L}^{(\theta)}$ ,

$$\theta^* \in \arg \min_{\theta \in \Theta} \mathcal{L}^{(\theta)}(\theta, \varphi^*) \quad \text{and} \quad \varphi^* \in \arg \min_{\varphi \in \Phi} \mathcal{L}^{(\varphi)}(\theta^*, \varphi). \quad (1)$$

State  $(\varphi^*, \theta^*)$ . In order to analyze the dynamics of gradient-based methods near a Nash Equilibrium, we look at the *gradient vector field* and its associated *Jacobian*  $\nabla v(\varphi, \theta)$ ,

$$v(\varphi, \theta) := \begin{bmatrix} \nabla_{\varphi} \mathcal{L}^{(\varphi)}(\varphi, \theta) \\ \nabla_{\theta} \mathcal{L}^{(\theta)}(\varphi, \theta) \end{bmatrix} \quad \nabla v(\varphi, \theta) := \begin{bmatrix} \nabla_{\varphi}^2 \mathcal{L}^{(\varphi)}(\varphi, \theta) & \nabla_{\varphi} \nabla_{\theta} \mathcal{L}^{(\varphi)}(\varphi, \theta) \\ \nabla_{\varphi} \nabla_{\theta} \mathcal{L}^{(\theta)}(\varphi, \theta)^T & \nabla_{\theta}^2 \mathcal{L}^{(\theta)}(\varphi, \theta) \end{bmatrix}. \quad (2)$$

**Simultaneous Gradient Method.** Let us consider the dynamics of the Simultaneous Gradient Method. It is defined as the repeated application of the following operator,

$$F_{\eta}(\varphi, \theta) := [\varphi \quad \theta]^{\top} - \eta v(\varphi, \theta), \quad (\varphi, \theta) \in \mathbb{R}^m, \quad (3)$$

where  $\eta$  is the learning rate. Now, for brevity we write the joint parameters  $\omega := (\varphi, \theta) \in \mathbb{R}^m$ . For  $t \in \mathbb{N}$ , let  $\omega_t = (\varphi_t, \theta_t)$  be the  $t^{\text{th}}$  point of the sequence computed by the gradient method. If the gradient method converges, its limit point  $\omega^* = (\varphi^*, \theta^*)$  is a *fixed point* of  $F_{\eta}$ . In addition, if  $\nabla v(\omega^*)$  is positive-definite, then  $\omega^*$  is a local Nash equilibrium.

## 2 Tuning the Step Size

Under certain conditions, linear convergence is guaranteed in a neighborhood around a fixed point.

**Theorem 1** (Prop. 4.4.1 Bertsekas [1999]). *If the spectral radius  $\rho_{\max} := \rho(\nabla F_{\eta}(\omega^*)) < 1$ , then, for  $\omega_0$  in a neighborhood of  $\omega^*$ , the distance of  $\omega_t$  to the stationary point  $\omega^*$  converges at a linear rate of  $\mathcal{O}((\rho_{\max} + \epsilon)^t)$ ,  $\forall \epsilon > 0$ .*

If the eigenvalues of  $\nabla v(\omega^*)$  all have a positive real-part, then for small enough  $\eta$ , the eigenvalues of  $\nabla F_{\eta}(\omega^*)$  are inside a convergence circle of radius  $\rho_{\max} < 1$ , as illustrated in Fig. 2. Then Thm. 1 guarantees the existence of an optimal step-size  $\eta_{\text{best}}$  which yields a non-trivial convergence rate  $\rho_{\max} < 1$ . Thm. 2 gives analytic bounds on the optimal step size  $\eta_{\text{best}}$ , and lower-bounds the best convergence rate  $\rho_{\max}(\eta_{\text{best}})$  we can expect.

**Theorem 2.** *If the eigenvalues of  $\nabla v(\omega^*)$  all have a positive real-part, then, the best step-size  $\eta_{\text{best}}$ , which minimizes the spectral radius  $\rho_{\max}(\eta)$  of  $\nabla F_{\eta}(\varphi^*, \theta^*)$ , is the solution of a (convex) quadratic by parts problem, and satisfies,*

$$\max_{1 \leq k \leq m} \sin(\psi_k)^2 \leq \rho_{\max}(\eta_{\text{best}})^2 \leq 1 - \Re(1/\lambda_1)\delta \quad \text{and} \quad \Re(1/\lambda_1) \leq \eta_{\text{best}} \leq 2\Re(1/\lambda_1), \quad (4)$$

where  $\delta := \min_{1 \leq k \leq m} |\lambda_k|^2 (2\Re(1/\lambda_k) - \Re(1/\lambda_1))$  and  $(\lambda_k = r_k e^{i\psi_k})_{1 \leq k \leq m} = \text{Sp}(\nabla v(\varphi^*, \theta^*))$  are sorted such that  $0 < \Re(1/\lambda_1) \leq \dots \leq \Re(1/\lambda_m)$ . Particularly, when  $\eta_{\text{best}} = \Re(1/\lambda_1)$  we are in the case of the top plot of Fig.2 and  $\rho_{\max}(\eta_{\text{best}})^2 = \sin(\psi_1)^2$ .

When  $\nabla v$  is positive-definite, the  $\eta_{\text{best}}$  is achieved because of either one or several limiting eigenvalues, we illustrate and interpret these two cases in Fig. 2. In (4) we can see that if the Jacobian of  $v$  has an almost purely imaginary eigenvalue then the convergence rate of the gradient method can be arbitrarily close to 1. Our goal is to use momentum to wrangle game dynamics into convergence.

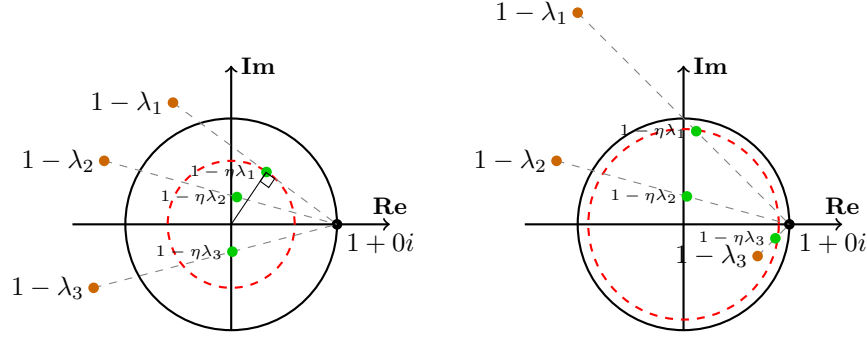


Figure 2: Trajectories of  $1 - \eta\lambda_i$  for growing step sizes, and the optimal step-size for  $\lambda_i \in \text{Sp}(\nabla v(\phi^*, \theta^*))$ . The unit circle is drawn in **black**, and the **red** dashed circle has radius equal to the largest eigenvalue  $\mu_{\max}$ , which is directly related to the convergence rate. Therefore, smaller red circles mean better convergence rates. **Left:** The red circle is limited by the tangent trajectory line  $1 - \eta\lambda_1$ , which means the best convergence rate is limited only by the eigenvalue which will pass furthest from the origin as  $\eta$  grows, i.e.,  $\lambda_1 = \arg \min \frac{|\Re(\lambda_i)|}{|\lambda_i|^2}$ . **Right:** The red circle is cut (not tangent) by the trajectories at points  $1 - \eta\lambda_1$  and  $1 - \eta\lambda_3$ . The  $\eta$  is optimal because any increase in  $\eta$  will push the eigenvalue  $\lambda_1$  out of the red circle, while any decrease in step-size will retract the eigenvalue  $\lambda_3$  out of the red circle, which will lower the convergence rate in any case.

### 3 Negative Momentum

As shown in (4), the presence of eigenvalues with large imaginary parts can restrict us to using small step sizes and lead to slow convergence rates. In order to improve convergence, we add a *negative* momentum term into the update rule. The new momentum term leads to a modification of the *parameter update operator*  $F_\eta(\omega)$  of (3). It requires to augment the state  $\omega_t$  with the previous iterate  $\omega_{t-1}$  (similar to Zhang and Mitiagkas [2017]) to form a compound state  $(\omega_t, \omega_{t-1}) := (\varphi_t, \theta_t, \varphi_{t-1}, \theta_{t-1}) \in \mathbb{R}^{2m}$ . The update rule (3) turns into the following,

$$F_{\eta,\beta}(\omega_t, \omega_{t-1}) = (\omega_{t+1}, \omega_t) \quad \text{where} \quad \omega_{t+1} := \omega_t - \eta v(\omega_t) + \beta(\omega_t - \omega_{t-1}), \quad (5)$$

in which  $\beta \in \mathbb{R}$  is the momentum parameter. Therefore, the Jacobian of  $F_{\eta,\beta}$  has the following form,

$$\nabla F_{\eta,\beta}(\omega_t, \omega_{t-1}) = \begin{bmatrix} \mathbf{I}_n & \mathbf{0}_n \\ \mathbf{I}_n & \mathbf{0}_n \end{bmatrix} - \eta \begin{bmatrix} \nabla v(\omega_t) & \mathbf{0}_n \\ \mathbf{0}_n & \mathbf{0}_n \end{bmatrix} + \beta \begin{bmatrix} \mathbf{I}_n & -\mathbf{I}_n \\ \mathbf{0}_n & \mathbf{0}_n \end{bmatrix} \quad (6)$$

In the following theorem, we provide an explicit equation for the eigenvalues of the Jacobian of  $F_{\eta,\beta}$ .

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

$$\mu_\pm(\beta, \eta, \lambda) := (1 - \eta\lambda + \beta) \frac{1 \pm \Delta^{\frac{1}{2}}}{2}, \quad \lambda \in \text{Sp}(\nabla v(\omega^*)), \quad (7)$$

where  $\Delta := 1 - \frac{4\beta}{(1 - \eta\lambda + \beta)^2}$  and  $\Delta^{\frac{1}{2}}$  is the complex square root of  $\Delta$  with positive real part<sup>3</sup>.

When  $\beta$  is small enough,  $\Delta$  is a complex number close to 1. Consequently,  $\mu_+$  is close to the original eigenvalue for gradient dynamics  $1 - \eta\lambda$ , and  $\mu_-$ , the eigenvalue introduced by the state augmentation, is close to 0. Since our goal is to minimize the largest magnitude of the eigenvalues of  $F_{\eta,\beta}$  computed in Thm. 3, we want to understand the effect of  $\beta$  on these eigenvalues with potential large magnitude. Let  $\lambda \in \text{Sp}(\nabla v(\omega^*))$ , we define the squared magnitude  $\rho_{\lambda,\eta}(\beta)$  we want to optimize,

$$\rho_{\lambda,\eta}(\beta) := \max \left\{ |\mu_+(\beta, \eta, \lambda)|^2, |\mu_-(\beta, \eta, \lambda)|^2 \right\}. \quad (8)$$

We study the local behavior of  $\rho_{\lambda,\eta}$  for small values of  $\beta$ . The following theorem shows that a well suited  $\beta$  decreases  $\rho_{\lambda,\eta}$ , which corresponds to faster convergence.

**Theorem 4.** *For any  $\lambda \in \text{Sp}(\nabla v(\omega^*))$  s.t.  $\Re(\lambda) > 0$ ,*

$$\rho'_{\lambda,\eta}(0) > 0 \Leftrightarrow \eta \in I(\lambda) := \left( \frac{|\lambda| - \Im(\lambda)}{|\lambda| \Re(\lambda)}, \frac{|\lambda| + \Im(\lambda)}{|\lambda| \Re(\lambda)} \right). \quad (9)$$

Particularly,  $\rho'_{\lambda,\Re(1/\lambda)}(0) = 2\Re(\lambda)\Re(1/\lambda) > 0$  and  $|\text{Arg}(\lambda)| \geq \frac{\pi}{4} \Rightarrow (\Re(1/\lambda), 2\Re(1/\lambda)) \subset I(\lambda)$ .

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

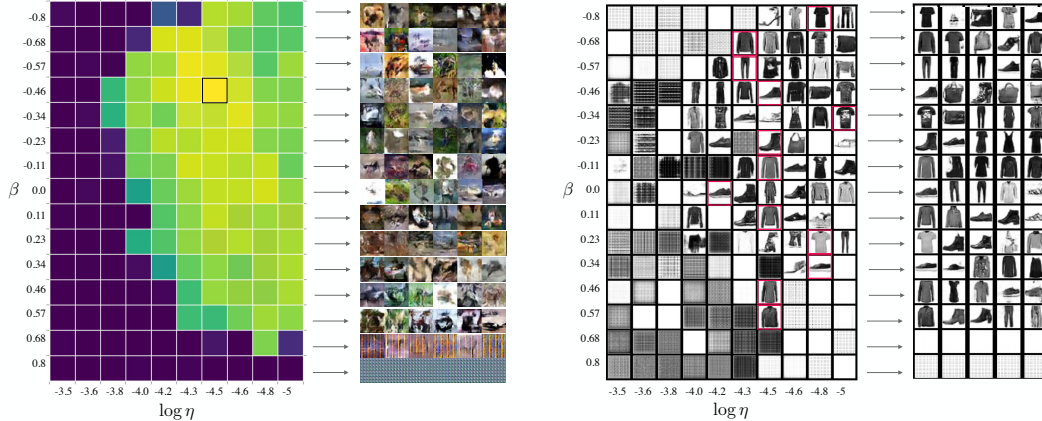


Figure 3: Comparison between negative and positive momentum on GANs with saturating loss on CIFAR-10 (left) and on Fashion MNIST (right) using a residual network. For each dataset, a grid of different values of momentum ( $\beta$ ) and step sizes ( $\eta$ ) is provided which describes the discriminator’s settings while a constant momentum of 0.5 and step size of  $10^{-4}$  is used for the generator. Each cell in CIFAR-10 (or Fashion MNIST) grid contains a single configuration in which its color (or its content) indicates the inception score (or a single sample) of the model. For CIFAR-10 experiments, yellow is higher while blue is the lower inception score. Along each row, the best configuration is chosen and more samples from that configuration are presented on the right side of each grid.

As we have seen previously in Fig. 2 and Thm. 2, there are only few eigenvalues which slow down the convergence. Thm. 4 is a local result showing that a small negative momentum can improve the magnitude of the limiting eigenvalues in the following cases: when there is only one limiting eigenvalue  $\lambda_1$  (since in that case the optimal step-size is  $\eta_{best} = \Re(1/\lambda_1) \in I(\lambda_1)$ ) or when there are several limiting eigenvalues  $\lambda_1, \dots, \lambda_k$  and the intersection  $I(\lambda_1) \cap \dots \cap I(\lambda_k)$  is not empty. We point out that we do not provide any guarantees on whether this intersection is empty or not but note that if the absolute value of the argument of  $\lambda_1$  is larger than  $\pi/4$  then by (4), our theorem provides that the optimal step-size  $\eta_{best}$  belongs to  $I(\lambda_1)$ . Nevertheless, we numerically optimized (8) with respect to  $\beta$  and  $\eta$  and found that for any non-imaginary fixed eigenvalue  $\lambda$ , the optimal  $\beta$  is negative and the associated optimal step size is larger than  $\Re(1/\lambda)$ .

## 4 Experiments and Discussion

We use negative momentum in a GAN setup on CIFAR-10 [Krizhevsky and Hinton, 2009] and Fashion-MNIST [Xiao et al., 2017] with *saturating loss* and alternating steps. Fig. 3 shows the results. We use residual networks for both the generator and the discriminator with no batch-normalization. Each residual block is made of two  $3 \times 3$  convolution layers with *ReLU* activation function. Up-sampling and down-sampling layers are respectively used in the generator and discriminator. We observe that using a negative value generally result in samples with higher quality and inception scores. We use negative momentum only on the discriminator because intuitively, negative momentum slows down the learning process of the discriminator and allows for better gradient flow to the generator.

## 5 Conclusion

We theoretically show that alternating updates with negative momentum is the only method within our study parameters (Tab.1) that converges in bilinear smooth games. Next, we study the effect of using negative values of momentum in a GAN setup both theoretically and experimentally and show that negative momentum can improve the convergence rate of gradient-based methods by shifting the eigenvalues of the Jacobian appropriately into a smaller convergence disk. We test out theory on CIFAR-10 and fashion MNIST datasets. We believe that our results explain a decreasing trend in momentum values reported in GAN literature in the past few years. Some of the most successful papers use zero momentum [Arjovsky et al., 2017, Gulrajani et al., 2017] for architectures that would otherwise call for high momentum values in a non-adversarial setting.

## References

- M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein generative adversarial networks. In *ICML*, 2017.
- D. Balduzzi, S. Racaniere, J. Martens, J. Foerster, K. Tuyls, and T. Graepel. The mechanics of n-player differentiable games. In *ICML*, 2018.
- D. P. Bertsekas. *Nonlinear programming*. Athena scientific Belmont, 1999.
- I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014.
- I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville. Improved training of wasserstein GANs. In *NIPS*, 2017.
- A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- L. Mescheder, S. Nowozin, and A. Geiger. The numerics of GANs. In *NIPS*, 2017.
- V. Nagarajan and J. Z. Kolter. Gradient descent GAN optimization is locally stable. In *NIPS*, 2017.
- Y. Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013.
- B. T. Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 1964.
- N. Qian. On the momentum term in gradient descent learning algorithms. *Neural networks*, 1999.
- I. Sutskever, J. Martens, G. Dahl, and G. Hinton. On the importance of initialization and momentum in deep learning. In *ICML*, 2013.
- H. Xiao, K. Rasul, and R. Vollgraf. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- J. Zhang and I. Mitliagkas. Yellowfin and the art of momentum tuning. *arXiv preprint arXiv:1706.03471*, 2017.