Complex Momentum for Optimization in Games

Artificial Intelligence and Statistics (AISTATS) 2022

Jonathan Lorraine^{1,2}, David Acuna^{1,2,3}, Paul Vicol^{1,2}, David Duvenaud^{1,2} March 14th, 2022

University of Toronto¹, Vector Institute², NVIDIA³

• *Games* generalize single objective opt. to have multiple objectives, each with their own parameters.

- *Games* generalize single objective opt. to have multiple objectives, each with their own parameters.
- But, single-objective opt. methods often don't work for games.

- *Games* generalize single objective opt. to have multiple objectives, each with their own parameters.
- But, single-objective opt. methods often don't work for games.
- Gidel et al. [1] show a negative momentum coefficient can help opt. in *adversarial games*.

- *Games* generalize single objective opt. to have multiple objectives, each with their own parameters.
- But, single-objective opt. methods often don't work for games.
- Gidel et al. [1] show a negative momentum coefficient can help opt. in *adversarial games*.
- But, can harm performance if the game is *non-adversarial*, and requires alternating updates costing 2 grad. evaluations per iteration.

1. Converges with 1 grad. eval. per iteration in adversarial games.

- 1. Converges with 1 grad. eval. per iteration in adversarial games.
- 2. Is tractable to analyze with Euler's formula.

- 1. Converges with 1 grad. eval. per iteration in adversarial games.
- 2. Is tractable to analyze with Euler's formula.
- 3. Is trivial to implement.

- 1. Converges with 1 grad. eval. per iteration in adversarial games.
- 2. Is tractable to analyze with Euler's formula.
- 3. Is trivial to implement.
- 4. Only introduces a single optimization hyperparameter.

- 1. Converges with 1 grad. eval. per iteration in adversarial games.
- 2. Is tractable to analyze with Euler's formula.
- 3. Is trivial to implement.
- 4. Only introduces a single optimization hyperparameter.
- 5. Robustly converges in non-adversarial games.

changes in green

```
mass = .8 + .3j
def momentum(step_size, mass):
    ...
    def update(i, g, state):
    x, velocity = state
    velocity = mass * velocity + g
    x=x-jnp.real(step_size(i)*velocity)
    return x, velocity
    ...
```

Background: Minimization with momentum

The standard minimization problem for loss *L* with parameters *θ*:

$$\boldsymbol{\theta}^* \coloneqq \arg\min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) \tag{1}$$

Background: Minimization with momentum

• The standard minimization problem for loss *L* with parameters *θ*:

$$\theta^* := \operatorname{arg\,min}_{\theta} \mathcal{L}(\theta)$$
 (1)

• The loss gradient at parameters θ^j denoted by:

$$\boldsymbol{g}^{j} := \boldsymbol{g}(\boldsymbol{\theta}^{j}) := \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})|_{\boldsymbol{\theta}^{j}}$$
(2)

Background: Minimization with momentum

The standard minimization problem for loss *L* with parameters *θ*:

$$\boldsymbol{\theta}^* := \arg\min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) \tag{1}$$

• The loss gradient at parameters θ^{j} denoted by:

$$\boldsymbol{g}^{j} := \boldsymbol{g}(\boldsymbol{\theta}^{j}) := \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})|_{\boldsymbol{\theta}^{j}}$$
(2)

• We can locally optimize loss using SGD with momentum:

$$\boldsymbol{\mu}^{j+1} = \beta \boldsymbol{\mu}^j - \boldsymbol{g}^j, \qquad \boldsymbol{\theta}^{j+1} = \boldsymbol{\theta}^j + \alpha \boldsymbol{\mu}^{j+1}$$
 (SGDm)

Setting up the notation for our method

• Optimization in games generalizes single-objective minimization:

 $\theta_{A}^{*} := \arg \min_{\theta_{A}} \mathcal{L}_{A}(\theta_{A}, \theta_{B}^{*}(\theta_{A})), \theta_{B}^{*}(\theta_{A}) := \arg \min_{\theta_{B}} \mathcal{L}_{B}(\theta_{A}, \theta_{B})$ (3)

Setting up the notation for our method

• Optimization in games generalizes single-objective minimization:

$$\theta_{A}^{*} := \arg \min_{\theta_{A}} \mathcal{L}_{A}(\theta_{A}, \theta_{B}^{*}(\theta_{A})), \theta_{B}^{*}(\theta_{A}) := \arg \min_{\theta_{B}} \mathcal{L}_{B}(\theta_{A}, \theta_{B})$$
(3)

• Notation: Concatenate the players' params and grads together

$$\boldsymbol{\omega} := [\boldsymbol{\theta}_A, \boldsymbol{\theta}_B], \quad \hat{\boldsymbol{g}}^j :== [\boldsymbol{g}_A^j, \boldsymbol{g}_B^j]$$
(4)

Setting up the notation for our method

• Optimization in games generalizes single-objective minimization:

$$\theta_{A}^{*} := \arg \min_{\theta_{A}} \mathcal{L}_{A}(\theta_{A}, \theta_{B}^{*}(\theta_{A})), \theta_{B}^{*}(\theta_{A}) := \arg \min_{\theta_{B}} \mathcal{L}_{B}(\theta_{A}, \theta_{B})$$
(3)

• Notation: Concatenate the players' params and grads together

$$\boldsymbol{\omega} := [\boldsymbol{\theta}_{A}, \boldsymbol{\theta}_{B}], \quad \hat{\boldsymbol{g}}^{j} := = [\boldsymbol{g}_{A}^{j}, \boldsymbol{g}_{B}^{j}]$$
(4)

6

• Simultaneous SGD is a common gradient-based opt. strategy:

$$\boldsymbol{\theta}_{A}^{j+1} = \boldsymbol{\theta}_{A}^{j} - \alpha \boldsymbol{g}_{A}^{j}, \, \boldsymbol{\theta}_{B}^{j+1} = \boldsymbol{\theta}_{B}^{j} - \alpha \boldsymbol{g}_{B}^{j} \iff \boldsymbol{\omega}^{j+1} = \boldsymbol{\omega}^{j} - \alpha \hat{\boldsymbol{g}}^{j}$$
(SimSGD)

• Simultaneous SGD is a common gradient-based opt. strategy:

$$\boldsymbol{\theta}_{A}^{j+1} = \boldsymbol{\theta}_{A}^{j} - \alpha \boldsymbol{g}_{A}^{j}, \boldsymbol{\theta}_{B}^{j+1} = \boldsymbol{\theta}_{B}^{j} - \alpha \boldsymbol{g}_{B}^{j} \iff \boldsymbol{\omega}^{j+1} = \boldsymbol{\omega}^{j} - \alpha \hat{\boldsymbol{g}}^{j}$$
(SimSGD)

• Or, (simultaneous) SGD with (negative [1]) momentum:

$$\boldsymbol{\mu}^{j+1} = \beta \boldsymbol{\mu}^j - \hat{\boldsymbol{g}}^j, \quad \boldsymbol{\omega}^{j+1} = \boldsymbol{\omega}^j + \alpha \boldsymbol{\mu}^{j+1}$$
 (SimSGDm)

• Or, (simultaneous) SGD with (negative [1]) momentum:

$$\boldsymbol{\mu}^{j+1} = \beta \boldsymbol{\mu}^j - \hat{\boldsymbol{g}}^j, \quad \boldsymbol{\omega}^{j+1} = \boldsymbol{\omega}^j + \alpha \boldsymbol{\mu}^{j+1}$$
 (SimSGDm)

We make β ∈ C instead of R. Then μ^{j+1} ∈ C, but need
 R-valued parameter update. We simply take the real part:

$$\boldsymbol{\mu}^{j+1} = \beta \boldsymbol{\mu}^j - \hat{\boldsymbol{g}}^j, \quad \boldsymbol{\omega}^{j+1} = \boldsymbol{\omega}^j + \Re(\alpha \boldsymbol{\mu}^{j+1}) \qquad (\mathsf{SimCM})$$



CM connects Classic & Negative Momentum

• By combining Euler's Formula and $\mu^{j+1} = -\sum_{k=0}^{k=j} \beta^k \hat{g}^{j-k}$:

$$\Re(\boldsymbol{\mu}^{j+1}) = -\sum_{k=0}^{k=j} \cos(k \arg(\beta)) |\beta|^k \hat{\boldsymbol{g}}^{j-k}$$
(5)

CM's $j = 50^{th}$ update's dependence on past grads



Comparing optimization algorithms



Algorithm 1 Complex Adam variant without momentum bias-correction

- 1: $\beta_1 \in \mathbb{C}, \beta_2 \in [0,1)$
- 2: $\alpha \in \mathbb{R}^+, \epsilon \in \mathbb{R}^+$
- 3: for i = 1 ... N do
- 4: $\mu^{j+1} = \beta_1 \mu^j g^j$
- 5: $\mathbf{v}^{j+1} = \beta_2 \mathbf{v}^j + (1 \beta_2) (\mathbf{g}^j)^2$
- 6: $\hat{\boldsymbol{v}}^{j+1} = \frac{\boldsymbol{v}^{j+1}}{1 (\beta_2)^j}$
- 7: $\omega^{j+1} = \omega^j + \alpha \frac{\Re(\mu^j)}{\sqrt{\hat{\mu}^{j+1}} + \epsilon}$ return ω^N

CIFAR-10 IS

Discriminator β_1	Max	Mean
0, BigGAN's [3]	9.10	8.93
$.8 \exp(i\pi/8)$, ours	9.25(+.15)	8.97(+.04)
.8, ablation	9.05(<mark>05</mark>)	7.19(-1.7)

 Using an *almost-positive* complex momentum – i.e., arg(β) is near, but not 0 – allows convergence in any setup with 1 grad. eval.

- Using an *almost-positive* complex momentum i.e., arg(β) is near, but not 0 – allows convergence in any setup with 1 grad. eval.
- Almost-positive momentum approaches classical momentum, gaining similar acceleration properties in *cooperative games*.

- Using an *almost-positive* complex momentum i.e., arg(β) is near, but not 0 – allows convergence in any setup with 1 grad. eval.
- Almost-positive momentum approaches classical momentum, gaining similar acceleration properties in *cooperative games*.
- A default of $\arg(\beta) = \pi/8$ did well across our experiments.

- Using an *almost-positive* complex momentum i.e., arg(β) is near, but not 0 – allows convergence in any setup with 1 grad. eval.
- Almost-positive momentum approaches classical momentum, gaining similar acceleration properties in *cooperative games*.
- A default of $\arg(\beta) = \pi/8$ did well across our experiments.
- We were able to extend our method to more sophisticated optimizers, training BigGAN to better inception scores.

Jonathan Lorraine



David Acuna



Paul Vicol



David Duvenaud



References

 Gauthier Gidel, Reyhane Askari Hemmat, Mohammad Pezeshki, Rémi Le Priol, Gabriel Huang, Simon Lacoste-Julien, and Ioannis Mitliagkas. Negative momentum for improved game dynamics. In *AISTATS*, 2019.

[2] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, and Skye Wanderman-Milne. JAX: composable transformations of Python+NumPy programs, 2018. URL http://github.com/google/jax.

[3] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. In International Conference on Learning Representations, 2018.



- Script for associated talk is here.
- Animation of our method is here.

Comparing optimization algorithms



Max adversarialness γ_{max}

Future Directions



- How to best extend to Adam?
- What assumptions about game eigenstructure are useful for realistic problems like GANs?