

Strict Non-Optimality in Minimax Optimization and Its Application to GAN

Donghwan Kim

KAIST

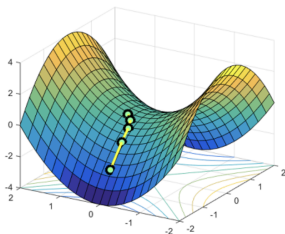
ICCOPT 2025, USC

July 24, 2025

- 1 Background
- 2 Strict Non-Optimality in Minimax Optimization
- 3 Two-Timescale Methods Escape Strict Non-Optimal Points
- 4 Application: GAN

Gradient descent (GD) and strict saddle

- GD: $\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \nabla f(\mathbf{x}_k)$
- Strict saddle point: its Hessian $\nabla^2 f(\mathbf{x})$ has a strictly negative eigenvalue



- By stable manifold theorem, for any strict saddle $\tilde{\mathbf{x}}$, GD satisfies¹

$$P(\lim_k \mathbf{x}_k = \tilde{\mathbf{x}}) = 0.$$

If GD converges, then it is almost surely not a strict saddle point.

¹Lee, Simchowitz, Jordan and Recht, Gradient descent only converges to minimizers, COLT, 2016.

Minimax and gradient descent ascent (GDA)

- Minimax optimization:

$$\min_{\mathbf{x}} \max_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$$

- GDA:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k)$$

$$\mathbf{y}_{k+1} = \mathbf{y}_k + \eta \nabla_{\mathbf{y}} f(\mathbf{x}_k, \mathbf{y}_k)$$

- Let $\mathbf{z} := (\mathbf{x}, \mathbf{y})$ and $\mathbf{F} := (\nabla_{\mathbf{x}} f, -\nabla_{\mathbf{y}} f)$:

$$\mathbf{z}_{k+1} = \mathbf{z}_k - \eta \mathbf{F}(\mathbf{z}_k)$$

Standard optimality in minimax optimization

- Nash equilibrium:

$$f(\mathbf{x}_*, \mathbf{y}) \leq f(\mathbf{x}_*, \mathbf{y}_*) \leq f(\mathbf{x}, \mathbf{y}_*), \quad \forall \mathbf{x}, \mathbf{y}.$$

- Local Nash equilibrium:

$$f(\mathbf{x}_*, \mathbf{y}) \leq f(\mathbf{x}_*, \mathbf{y}_*) \leq f(\mathbf{x}, \mathbf{y}_*), \quad \forall \mathbf{x}, \mathbf{y} \text{ in a neighborhood of } (\mathbf{x}_*, \mathbf{y}_*).$$

- Second-order necessary condition:

$$\nabla f(\mathbf{z}_*) = \mathbf{0}, \quad \nabla_{xx}^2 f(\mathbf{z}_*) \succeq \mathbf{0}, \quad \text{and} \quad \nabla_{yy}^2 f(\mathbf{z}_*) \preceq \mathbf{0}$$

- **Strict non-Nash point:** (analogous to strict saddle)
 at least one of $\nabla_{xx}^2 f(\mathbf{z})$ or $-\nabla_{yy}^2 f(\mathbf{z})$ has a strictly negative eigenvalue
 (Are we done?)

Standard optimality in minimax optimization

- Nash equilibrium:

$$f(\mathbf{x}_*, \mathbf{y}) \leq f(\mathbf{x}_*, \mathbf{y}_*) \leq f(\mathbf{x}, \mathbf{y}_*), \quad \forall \mathbf{x}, \mathbf{y}.$$

- Local Nash equilibrium:

$$f(\mathbf{x}_*, \mathbf{y}) \leq f(\mathbf{x}_*, \mathbf{y}_*) \leq f(\mathbf{x}, \mathbf{y}_*), \quad \forall \mathbf{x}, \mathbf{y} \text{ in a neighborhood of } (\mathbf{x}_*, \mathbf{y}_*).$$

- Second-order necessary condition:

$$\nabla f(\mathbf{z}_*) = \mathbf{0}, \quad \nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{z}_*) \succeq \mathbf{0}, \quad \text{and} \quad \nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{z}_*) \preceq \mathbf{0}$$

- **Strict non-Nash point:** (analogous to strict saddle)
 at least one of $\nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{z})$ or $-\nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{z})$ has a strictly negative eigenvalue
 (Are we done?)

Standard optimality in minimax optimization

- Nash equilibrium:

$$f(\mathbf{x}_*, \mathbf{y}) \leq f(\mathbf{x}_*, \mathbf{y}_*) \leq f(\mathbf{x}, \mathbf{y}_*), \quad \forall \mathbf{x}, \mathbf{y}.$$

- Local Nash equilibrium:

$$f(\mathbf{x}_*, \mathbf{y}) \leq f(\mathbf{x}_*, \mathbf{y}_*) \leq f(\mathbf{x}, \mathbf{y}_*), \quad \forall \mathbf{x}, \mathbf{y} \text{ in a neighborhood of } (\mathbf{x}_*, \mathbf{y}_*).$$

- Second-order necessary condition:

$$\nabla f(\mathbf{z}_*) = \mathbf{0}, \quad \nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{z}_*) \succeq \mathbf{0}, \quad \text{and} \quad \nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{z}_*) \preceq \mathbf{0}$$

- **Strict non-Nash point:** (analogous to strict saddle)
 at least one of $\nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{z})$ or $-\nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{z})$ has a strictly negative eigenvalue
 (Are we done?)

GDA can converge to non-Nash points

- Stability of a dynamic $\mathbf{x}_{k+1} = \mathbf{w}(\mathbf{x}_k)$ is determined by its Jacobian $D\mathbf{w}$.
- GD: $\mathbf{x}_{k+1} = \mathbf{w}_1(\mathbf{x}_k) = \mathbf{x}_k - \eta \nabla f(\mathbf{x}_k)$, $D\mathbf{w}_1 = \mathbf{I} - \eta \nabla^2 f$
 GDA: $\mathbf{z}_{k+1} = \mathbf{w}_2(\mathbf{z}_k) = \mathbf{z}_k - \eta \mathbf{F}(\mathbf{z}_k)$, $D\mathbf{w}_2 = \mathbf{I} - \eta D\mathbf{F}$
- GD almost surely escapes strict saddle, while
 GDA has no such guarantee², since

$$D\mathbf{F} = \begin{bmatrix} \nabla_{xx}^2 f & \nabla_{xy}^2 f \\ -\nabla_{yx}^2 f & -\nabla_{yy}^2 f \end{bmatrix}$$

has no direct connection to $\nabla_{xx}^2 f$ and $\nabla_{yy}^2 f$ in general.

²Daskalakis and Panageas, The limit points of (optimistic) gradient descent in min-max optimization, NeurIPS, 2018.

GDA can converge to non-Nash points

- Stability of a dynamic $\mathbf{x}_{k+1} = \mathbf{w}(\mathbf{x}_k)$ is determined by its Jacobian $D\mathbf{w}$.
- GD: $\mathbf{x}_{k+1} = \mathbf{w}_1(\mathbf{x}_k) = \mathbf{x}_k - \eta \nabla f(\mathbf{x}_k)$, $D\mathbf{w}_1 = \mathbf{I} - \eta \nabla^2 f$
 GDA: $\mathbf{z}_{k+1} = \mathbf{w}_2(\mathbf{z}_k) = \mathbf{z}_k - \eta \mathbf{F}(\mathbf{z}_k)$, $D\mathbf{w}_2 = \mathbf{I} - \eta D\mathbf{F}$
- GD almost surely escapes strict saddle, while
 GDA has no such guarantee², since

$$D\mathbf{F} = \begin{bmatrix} \nabla_{xx}^2 f & \nabla_{xy}^2 f \\ -\nabla_{yx}^2 f & -\nabla_{yy}^2 f \end{bmatrix}$$

has no direct connection to $\nabla_{xx}^2 f$ and $\nabla_{yy}^2 f$ in general.

²Daskalakis and Panageas, The limit points of (optimistic) gradient descent in min-max optimization, NeurIPS, 2018.

GDA can converge to non-Nash points

- Stability of a dynamic $\mathbf{x}_{k+1} = \mathbf{w}(\mathbf{x}_k)$ is determined by its Jacobian $D\mathbf{w}$.
- GD: $\mathbf{x}_{k+1} = \mathbf{w}_1(\mathbf{x}_k) = \mathbf{x}_k - \eta \nabla f(\mathbf{x}_k)$, $D\mathbf{w}_1 = \mathbf{I} - \eta \nabla^2 f$
 GDA: $\mathbf{z}_{k+1} = \mathbf{w}_2(\mathbf{z}_k) = \mathbf{z}_k - \eta \mathbf{F}(\mathbf{z}_k)$, $D\mathbf{w}_2 = \mathbf{I} - \eta D\mathbf{F}$
- GD almost surely escapes strict saddle, while
 GDA has no such guarantee², since

$$D\mathbf{F} = \begin{bmatrix} \nabla_{xx}^2 f & \nabla_{xy}^2 f \\ -\nabla_{yx}^2 f & -\nabla_{yy}^2 f \end{bmatrix}$$

has no direct connection to $\nabla_{xx}^2 f$ and $\nabla_{yy}^2 f$ in general.

²Daskalakis and Panageas, The limit points of (optimistic) gradient descent in min-max optimization, NeurIPS, 2018.

- 1 Background
- 2 Strict Non-Optimality in Minimax Optimization
- 3 Two-Timescale Methods Escape Strict Non-Optimal Points
- 4 Application: GAN

Stackelberg equilibrium

- Nash may not even exist.³⁴ Any broader alternative?
- **Stackelberg equilibrium:**

$$f(x_*, y) \leq f(x_*, y_*) \leq \max_{y' \in \mathcal{Y}} f(x, y'), \quad \forall (x, y) \in \mathcal{X} \times \mathcal{Y}$$

Definition 1 (Jin-Netrapalli-Jordan, ICML, '20)

A point (x^*, y^*) is said to be a **local minimax point** if there exists $\delta_0 > 0$ and a function h satisfying $h(\delta) \rightarrow 0$ as $\delta \rightarrow 0$ such that, for any $\delta \in (0, \delta_0]$ and any (x, y) satisfying $\|x - x_*\| \leq \delta$ and $\|y - y_*\| \leq \delta$, we have

$$f(x_*, y) \leq f(x_*, y_*) \leq \max_{y' : \|y' - y_*\| \leq h(\delta)} f(x, y').$$

³Jin, Netrapalli and Jordan, What is local optimality in nonconvex-nonconcave minimax optimization?, ICML, 2020.

⁴Farnia and Ozdaglar, Do GANs always have Nash equilibria?, ICML, 2020.

Stackelberg equilibrium

- Nash may not even exist.³⁴ Any broader alternative?
- **Stackelberg** equilibrium:

$$f(x_*, y) \leq f(x_*, y_*) \leq \max_{y' \in \mathcal{Y}} f(x, y'), \quad \forall (x, y) \in \mathcal{X} \times \mathcal{Y}$$

Definition 1 (Jin-Netrapalli-Jordan, ICML, '20)

A point (x^*, y^*) is said to be a **local minimax point** if there exists $\delta_0 > 0$ and a function h satisfying $h(\delta) \rightarrow 0$ as $\delta \rightarrow 0$ such that, for any $\delta \in (0, \delta_0]$ and any (x, y) satisfying $\|x - x_*\| \leq \delta$ and $\|y - y_*\| \leq \delta$, we have

$$f(x_*, y) \leq f(x_*, y_*) \leq \max_{y' : \|y' - y_*\| \leq h(\delta)} f(x, y').$$

³Jin, Netrapalli and Jordan, What is local optimality in nonconvex-nonconcave minimax optimization?, ICML, 2020.

⁴Farnia and Ozdaglar, Do GANs always have Nash equilibria?, ICML, 2020.

Stackelberg equilibrium

- Nash may not even exist.³⁴ Any broader alternative?
- **Stackelberg** equilibrium:

$$f(x_*, y) \leq f(x_*, y_*) \leq \max_{y' \in \mathcal{Y}} f(x, y'), \quad \forall (x, y) \in \mathcal{X} \times \mathcal{Y}$$

Definition 1 (Jin-Netrapalli-Jordan, ICML, '20)

A point (x^*, y^*) is said to be a **local minimax point** if there exists $\delta_0 > 0$ and a function h satisfying $h(\delta) \rightarrow 0$ as $\delta \rightarrow 0$ such that, for any $\delta \in (0, \delta_0]$ and any (x, y) satisfying $\|x - x_*\| \leq \delta$ and $\|y - y_*\| \leq \delta$, we have

$$f(x_*, y) \leq f(x_*, y_*) \leq \max_{y' : \|y' - y_*\| \leq h(\delta)} f(x, y').$$

³Jin, Netrapalli and Jordan, What is local optimality in nonconvex-nonconcave minimax optimization?, ICML, 2020.

⁴Farnia and Ozdaglar, Do GANs always have Nash equilibria?, ICML, 2020.

Loose second-order necessary condition

- Second-order necessary condition: (Jin-Netrapalli-Jordan, ICML, '20)

$$\nabla f(z_*) = \mathbf{0}, \quad \nabla_{yy}^2 f(z_*) \preceq \mathbf{0},$$

and if $\nabla_{yy}^2 f(z_*) \prec \mathbf{0}$, then in addition,

$$S(z_*) := \underbrace{[\nabla_{xx}^2 f - \nabla_{xy}^2 f (\nabla_{yy}^2 f)^{-1} \nabla_{yx}^2 f]}_{\text{Schur complement of } D\mathbf{F}}(z_*) \succeq \mathbf{0}$$

(Loose when $\nabla_{yy}^2 f(z_*)$ is not invertible...)

Restricted Schur Complement

- By similarity transform, we may assume that $\nabla_{\mathbf{y}\mathbf{y}}^2 f$ is diagonal such that

$$D\mathbf{F} = \left[\begin{array}{c|ccc} \nabla_{\mathbf{x}\mathbf{x}}^2 f & & & \nabla_{\mathbf{x}\mathbf{y}}^2 f & \\ \hline & \beta_1 & & & \\ & & \ddots & & \\ -\nabla_{\mathbf{y}\mathbf{x}}^2 f & & & \beta_r & \\ \hline & & & & 0 \\ & & & & \vdots \\ & & & & & \ddots \\ & & & & & & 0 \end{array} \right]$$

- Let $\mathbf{\Gamma}$ be the submatrix in the shaded part above.
- Let \mathbf{U} be a matrix whose columns form an orthonormal basis of $\mathcal{R}(\mathbf{\Gamma})^\top$.

Definition 2 (Restricted Schur Complement, Chae-Kim-K., ICLR, '24)

$$\mathbf{S}_{\text{res}} = \mathbf{U}^\top \underbrace{(\nabla_{\mathbf{x}\mathbf{x}}^2 f - \nabla_{\mathbf{x}\mathbf{y}}^2 f (\nabla_{\mathbf{y}\mathbf{y}}^2 f)^\dagger \nabla_{\mathbf{y}\mathbf{x}}^2 f)}_{(\text{Generalized Schur complement of } D\mathbf{F}) =: \mathbf{S}} \mathbf{U}$$

Chae-Kim-K., Two-timescale extragradient for finding local minimax points, ICLR, 2024.

Improved necessary condition and strict non-minimax point

Proposition 1 (Chae-Kim-K., ICLR, '24)

$S_{\text{res}} \succeq \mathbf{0}$ if and only if $\mathbf{v}^\top S \mathbf{v} \geq 0$ for any $\mathbf{v} \in \mathcal{R}(U)$,
or equivalently, for any $(\nabla_{\mathbf{y}\mathbf{x}}^2 f) \mathbf{v} \in \mathcal{R}(\nabla_{\mathbf{y}\mathbf{y}}^2 f)$.

Proposition 2 (Second-order necessary, Chae-Kim-K., ICLR, '24)

Any local minimax point satisfies $\nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{z}) \preceq \mathbf{0}$ and if $h(\delta)$ satisfies $\limsup_{\delta \rightarrow 0+} \frac{h(\delta)}{\delta} < \infty$,^a then $S_{\text{res}}(\mathbf{z}) \succeq \mathbf{0}$.

^aMa, Yao, Ye and Zhang, Calm local optimality for nonconvex-nonconcave minimax problems, Set-Valued and Variational Analysis, 2025

Definition 3 (Chae-Kim-K., ICLR, '24)

A stationary point \mathbf{z} is said to be a **strict non-minimax point** if at least one of $S_{\text{res}}(\mathbf{z})$ or $-\nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{z})$ has a strictly negative eigenvalue.

Improved necessary condition and strict non-minimax point

Proposition 1 (Chae-Kim-K., ICLR, '24)

$S_{\text{res}} \succeq \mathbf{0}$ if and only if $\mathbf{v}^\top S \mathbf{v} \geq 0$ for any $\mathbf{v} \in \mathcal{R}(U)$,
or equivalently, for any $(\nabla_{\mathbf{y}\mathbf{x}}^2 f) \mathbf{v} \in \mathcal{R}(\nabla_{\mathbf{y}\mathbf{y}}^2 f)$.

Proposition 2 (Second-order necessary, Chae-Kim-K., ICLR, '24)

Any local minimax point satisfies $\nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{z}) \preceq \mathbf{0}$ and if $h(\delta)$ satisfies $\limsup_{\delta \rightarrow 0+} \frac{h(\delta)}{\delta} < \infty$,^a then $S_{\text{res}}(\mathbf{z}) \succeq \mathbf{0}$.

^aMa, Yao, Ye and Zhang, Calm local optimality for nonconvex-nonconcave minimax problems, Set-Valued and Variational Analysis, 2025

Definition 3 (Chae-Kim-K., ICLR, '24)

A stationary point \mathbf{z} is said to be a **strict non-minimax point** if at least one of $S_{\text{res}}(\mathbf{z})$ or $-\nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{z})$ has a strictly negative eigenvalue.

Improved necessary condition and strict non-minimax point

Proposition 1 (Chae-Kim-K., ICLR, '24)

$S_{\text{res}} \succeq \mathbf{0}$ if and only if $\mathbf{v}^\top S \mathbf{v} \geq 0$ for any $\mathbf{v} \in \mathcal{R}(U)$,
or equivalently, for any $(\nabla_{\mathbf{y}\mathbf{x}}^2 f) \mathbf{v} \in \mathcal{R}(\nabla_{\mathbf{y}\mathbf{y}}^2 f)$.

Proposition 2 (Second-order necessary, Chae-Kim-K., ICLR, '24)

Any local minimax point satisfies $\nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{z}) \preceq \mathbf{0}$ and if $h(\delta)$ satisfies $\limsup_{\delta \rightarrow 0+} \frac{h(\delta)}{\delta} < \infty$,^a then $S_{\text{res}}(\mathbf{z}) \succeq \mathbf{0}$.

^aMa, Yao, Ye and Zhang, Calm local optimality for nonconvex-nonconcave minimax problems, Set-Valued and Variational Analysis, 2025

Definition 3 (Chae-Kim-K., ICLR, '24)

A stationary point \mathbf{z} is said to be a **strict non-minimax point** if at least one of $S_{\text{res}}(\mathbf{z})$ or $-\nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{z})$ has a strictly negative eigenvalue.

- 1 Background
- 2 Strict Non-Optimality in Minimax Optimization
- 3 Two-Timescale Methods Escape Strict Non-Optimal Points
- 4 Application: GAN

Two-timescale GDA

- Two-timescale τ -GDA⁵ for $\tau > 1$:

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - \frac{\eta}{\tau} \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k), \\ \mathbf{y}_{k+1} &= \mathbf{y}_k + \eta \nabla_{\mathbf{y}} f(\mathbf{x}_k, \mathbf{y}_k) \end{aligned}$$

- In a compact form:

$$\mathbf{z}_{k+1} = \mathbf{z}_k - \eta \Lambda_{\tau} \mathbf{F}(\mathbf{z}_k), \quad \text{where} \quad \Lambda_{\tau} = \begin{bmatrix} \frac{1}{\tau} \mathbf{I} & \\ & \mathbf{I} \end{bmatrix}$$

- Two-timescaled Jacobian:

$$\Lambda_{\tau} D\mathbf{F} = \begin{bmatrix} \frac{1}{\tau} \nabla_{\mathbf{x}\mathbf{x}}^2 f & \frac{1}{\tau} \nabla_{\mathbf{x}\mathbf{y}}^2 f \\ -\nabla_{\mathbf{y}\mathbf{x}}^2 f & -\nabla_{\mathbf{y}\mathbf{y}}^2 f \end{bmatrix}$$

(But why two-timescale?)

⁵Heusel, Ramsauer, Unterthiner, Nessler and Hochreiter, GANs trained by a two time-scale update rule converge to a local Nash equilibrium, NeurIPS, 2017.

Two-timescale GDA

- Two-timescale τ -GDA⁵ for $\tau > 1$:

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - \frac{\eta}{\tau} \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k), \\ \mathbf{y}_{k+1} &= \mathbf{y}_k + \eta \nabla_{\mathbf{y}} f(\mathbf{x}_k, \mathbf{y}_k) \end{aligned}$$

- In a compact form:

$$\mathbf{z}_{k+1} = \mathbf{z}_k - \eta \Lambda_{\tau} \mathbf{F}(\mathbf{z}_k), \quad \text{where} \quad \Lambda_{\tau} = \begin{bmatrix} \frac{1}{\tau} \mathbf{I} & \\ & \mathbf{I} \end{bmatrix}$$

- Two-timescaled Jacobian:

$$\Lambda_{\tau} D\mathbf{F} = \begin{bmatrix} \frac{1}{\tau} \nabla_{\mathbf{x}\mathbf{x}}^2 f & \frac{1}{\tau} \nabla_{\mathbf{x}\mathbf{y}}^2 f \\ -\nabla_{\mathbf{y}\mathbf{x}}^2 f & -\nabla_{\mathbf{y}\mathbf{y}}^2 f \end{bmatrix}$$

(But why two-timescale?)

⁵Heusel, Ramsauer, Unterthiner, Nessler and Hochreiter, GANs trained by a two time-scale update rule converge to a local Nash equilibrium, NeurIPS, 2017.

Spectrum of two-timescaled Jacobian

Lemma 1 (Jin-Netrapalli-Jordan, ICML, '20, Lemma 40)

If $\nabla_{yy}^2 f$ is invertible, the $d_1 + d_2$ complex eigenvalues $\{\lambda_j\}$ of $\Lambda_\tau D\mathbf{F}$ have one of the following asymptotics as $\epsilon = \frac{1}{\tau} \rightarrow 0+$:

$$|\lambda_j - \epsilon \mu_j| = o(\epsilon),$$

$$|\lambda_j - \nu_j| = o(1),$$

where $\{\mu_j\}$ and $\{\nu_j\}$ are the eigenvalues of \mathbf{S} and $-\nabla_{yy}^2 f$, respectively.

(Invertibility of $\nabla_{yy}^2 f$ is too restrictive.)

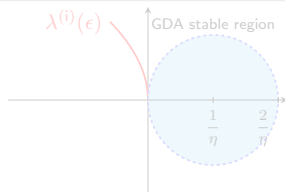
Spectrum of two-timescaled Jacobian (cont'd)

Theorem 1 (Chae-Kim-K., ICLR '24)

If at least one of S_{res} and $\nabla_{yy}^2 f$ is invertible, the $d_1 + d_2$ complex eigenvalues $\{\lambda_j\}$ of $\Lambda_\tau DF$ have one of the following asymptotics as $\epsilon = \frac{1}{\tau} \rightarrow 0+$:

- (i) $|\lambda_j \pm i\sqrt{\epsilon}\sigma_j| = o(\sqrt{\epsilon})$,
- (ii) $|\lambda_j - \epsilon\mu_j| = o(\epsilon)$,
- (iii) $|\lambda_j - \nu_j| = o(1)$,

where $\{\sigma_j\}$ are the singular values of Γ , $\{\mu_j\}$ are the eigenvalues of S_{res} , and $\{\nu_j\}$ are the nonzero eigenvalues of $-\nabla_{yy}^2 f$.



(type (i) makes GDA unstable)

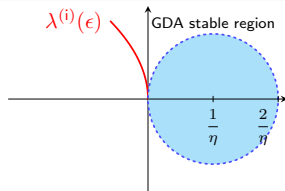
Spectrum of two-timescaled Jacobian (cont'd)

Theorem 1 (Chae-Kim-K., ICLR '24)

If at least one of S_{res} and $\nabla_{yy}^2 f$ is invertible, the $d_1 + d_2$ complex eigenvalues $\{\lambda_j\}$ of $\Lambda_\tau DF$ have one of the following asymptotics as $\epsilon = \frac{1}{\tau} \rightarrow 0+$:

- (i) $|\lambda_j \pm i\sqrt{\epsilon}\sigma_j| = o(\sqrt{\epsilon})$,
- (ii) $|\lambda_j - \epsilon\mu_j| = o(\epsilon)$,
- (iii) $|\lambda_j - \nu_j| = o(1)$,

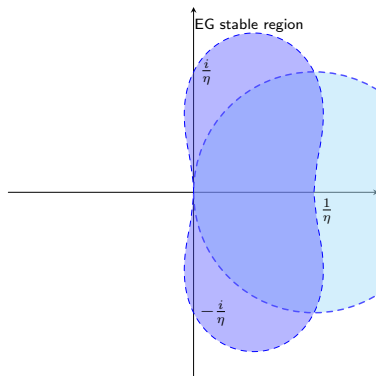
where $\{\sigma_j\}$ are the singular values of Γ , $\{\mu_j\}$ are the eigenvalues of S_{res} , and $\{\nu_j\}$ are the nonzero eigenvalues of $-\nabla_{yy}^2 f$.



(type (i) makes GDA unstable)

Extragradient

- Extragradient (EG): $z_{k+1} = w(z_k) = z_k - \eta F(z_k - \eta F(z_k))$

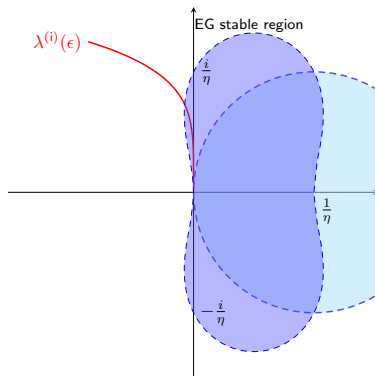


- Kyuwon Kim⁶, Wed 10:30am - 11:45am
(Alternating update + Chambolle-Pock extrapolation helps!)

⁶Kim-K., Double-step alternating extragradient with increasing timescale separation for finding local minimax points: Provable improvements, ICML, 2024

Extragradient

- τ -EG: $z_{k+1} = w(z_k) = z_k - \eta \Lambda_{\tau} F(z_k - \eta \Lambda_{\tau} F(z_k))$

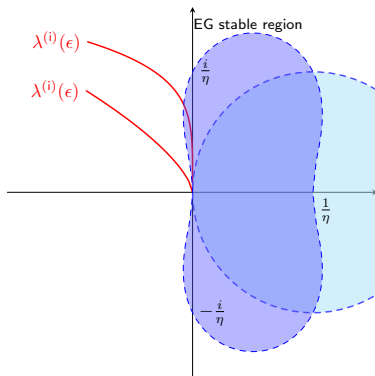


- Kyuwon Kim⁶, Wed 10:30am - 11:45am
(Alternating update + Chambolle-Pock extrapolation helps!)

⁶Kim-K., Double-step alternating extragradient with increasing timescale separation for finding local minimax points: Provable improvements, ICML, 2024

Extragradient

- τ -EG: $z_{k+1} = w(z_k) = z_k - \eta \Lambda_{\tau} F(z_k - \eta \Lambda_{\tau} F(z_k))$



- Kyuwon Kim⁶, Wed 10:30am - 11:45am
(Alternating update + Chambolle-Pock extrapolation helps!)

⁶Kim-K., Double-step alternating extragradient with increasing timescale separation for finding local minimax points: Provable improvements, ICML, 2024

Two-timescale method avoids strict non-minimax points

- By the stable manifold theorem, if at least one of \mathcal{S}_{res} and $\nabla_{yy}^2 f$ is invertible, for any **strict non-minimax** point \tilde{z} , τ -EG satisfies
(Chae-Kim-K., ICLR '24)

$$P(\lim_k z_k = \tilde{z}) = 0,$$

for sufficiently large τ

If τ -EG (and τ -GDA) converges to a point, then it is almost surely not a **strict non-minimax** point.

(But τ -GDA also escapes some optimal points)

- 1 Background
- 2 Strict Non-Optimality in Minimax Optimization
- 3 Two-Timescale Methods Escape Strict Non-Optimal Points
- 4 Application: GAN

Generative models and GAN

- \mathbb{P}_{true} and \mathbb{P}_{θ} : True and generated data distributions
- Generative models aim to train θ so that $\mathbb{P}_{\text{true}} \approx \mathbb{P}_{\theta}$.
- GAN: minimizes the distance between them via minimax optimization

Wasserstein GAN

- Wasserstein GAN⁷: minimizes the Wasserstein-1 distance

$$\begin{aligned}\min_{\theta} W(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) &= \inf_{\gamma \in \Pi(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta})} \mathbb{E}_{(x,z) \sim \gamma} [\|x - z\|] \\ &= \min_{\theta} \left(\max_{f: \|f\|_L \leq 1} \mathbb{E}_{x \sim \mathbb{P}_{\text{true}}} [f(x)] - \mathbb{E}_{z \sim \mathbb{P}_{\theta}} [f(z)] \right).\end{aligned}$$

- Since this is informative even when \mathbb{P}_{true} and \mathbb{P}_{θ} have disjoint supports, it is considered well-suited for gradient-based training.
- Yet enforcing the constraint required heuristics, and optimization remained difficult, leading the community to favor diffusion models for their stable training.
- Avoidance in constrained problem?: projected GD may converge to a strict saddle even when there is only a single linear constraint.⁸

⁷Arjovsky, Chintala and Bottou, Wasserstein generative adversarial networks, ICML, 2017

⁸Nouiehed, Lee and Razaviyayn, Convergence to second-order stationarity for constrained non-convex optimization, arXiv, 2018.

Wasserstein GAN

- Wasserstein GAN⁷: minimizes the Wasserstein-1 distance

$$\begin{aligned}\min_{\theta} W(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) &= \inf_{\gamma \in \Pi(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta})} \mathbb{E}_{(x,z) \sim \gamma} [\|x - z\|] \\ &= \min_{\theta} \left(\max_{f: \|f\|_L \leq 1} \mathbb{E}_{x \sim \mathbb{P}_{\text{true}}} [f(x)] - \mathbb{E}_{z \sim \mathbb{P}_{\theta}} [f(z)] \right).\end{aligned}$$

- Since this is informative even when \mathbb{P}_{true} and \mathbb{P}_{θ} have disjoint supports, it is considered well-suited for gradient-based training.
- Yet enforcing the constraint required heuristics, and optimization remained difficult, leading the community to favor diffusion models for their stable training.
- Avoidance in constrained problem?: projected GD may converge to a strict saddle even when there is only a single linear constraint.⁸

⁷Arjovsky, Chintala and Bottou, Wasserstein generative adversarial networks, ICML, 2017

⁸Nouiehed, Lee and Razaviyayn, Convergence to second-order stationarity for constrained non-convex optimization, arXiv, 2018.

Wasserstein GAN

- Wasserstein GAN⁷: minimizes the Wasserstein-1 distance

$$\begin{aligned}\min_{\theta} W(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) &= \inf_{\gamma \in \Pi(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta})} \mathbb{E}_{(x,z) \sim \gamma} [\|x - z\|] \\ &= \min_{\theta} \left(\max_{f: \|f\|_L \leq 1} \mathbb{E}_{x \sim \mathbb{P}_{\text{true}}} [f(x)] - \mathbb{E}_{z \sim \mathbb{P}_{\theta}} [f(z)] \right).\end{aligned}$$

- Since this is informative even when \mathbb{P}_{true} and \mathbb{P}_{θ} have disjoint supports, it is considered well-suited for gradient-based training.
- Yet enforcing the constraint required heuristics, and optimization remained difficult, leading the community to favor diffusion models for their stable training.
- Avoidance in constrained problem?: projected GD may converge to a strict saddle even when there is only a single linear constraint.⁸

⁷Arjovsky, Chintala and Bottou, Wasserstein generative adversarial networks, ICML, 2017

⁸Nouiehed, Lee and Razaviyayn, Convergence to second-order stationarity for constrained non-convex optimization, arXiv, 2018.

Wasserstein GAN

- Wasserstein GAN⁷: minimizes the Wasserstein-1 distance

$$\begin{aligned}\min_{\theta} W(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) &= \inf_{\gamma \in \Pi(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta})} \mathbb{E}_{(x,z) \sim \gamma} [\|x - z\|] \\ &= \min_{\theta} \left(\max_{f: \|f\|_L \leq 1} \mathbb{E}_{x \sim \mathbb{P}_{\text{true}}} [f(x)] - \mathbb{E}_{z \sim \mathbb{P}_{\theta}} [f(z)] \right).\end{aligned}$$

- Since this is informative even when \mathbb{P}_{true} and \mathbb{P}_{θ} have disjoint supports, it is considered well-suited for gradient-based training.
- Yet enforcing the constraint required heuristics, and optimization remained difficult, leading the community to favor diffusion models for their stable training.
- Avoidance in constrained problem?: projected GD may converge to a strict saddle even when there is only a single linear constraint.⁸

⁷Arjovsky, Chintala and Bottou, Wasserstein generative adversarial networks, ICML, 2017

⁸Nouiehed, Lee and Razaviyayn, Convergence to second-order stationarity for constrained non-convex optimization, arXiv, 2018.

Jensen-Shannon GAN

- Original GAN⁹: minimizes the Jensen-Shannon (JS) divergence:

$$\text{JS}(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) := \frac{1}{2} \text{KL} \left(\mathbb{P}_{\text{true}} \left\| \frac{\mathbb{P}_{\text{true}} + \mathbb{P}_{\theta}}{2} \right\| \right) + \frac{1}{2} \text{KL} \left(\mathbb{P}_{\theta} \left\| \frac{\mathbb{P}_{\text{true}} + \mathbb{P}_{\theta}}{2} \right\| \right),$$

which is

$$\min_{\theta} (2\text{JS}(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) - 2 \log 2) = \min_{\theta} \left(\max_f -\mathbb{E}_{x \sim \mathbb{P}_{\text{true}}} [l(f(x))] - \mathbb{E}_{z \sim \mathbb{P}_{\theta}} [l(-f(z))] \right)$$

where $l(t) = \log(1 + \exp(-t))$ is the logistic loss.

- This is unconstrained, but suffers from a vanishing gradient issue, since for each θ , near-optimal f stays in the tail of the logistic loss.
- Can we design a loss (or distance) that mitigates vanishing gradients without imposing constraints?

⁹Goodfellow et al., Generative adversarial nets, NeurIPS, 2014.

Jensen-Shannon GAN

- Original GAN⁹: minimizes the Jensen-Shannon (JS) divergence:

$$\text{JS}(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) := \frac{1}{2} \text{KL} \left(\mathbb{P}_{\text{true}} \left\| \frac{\mathbb{P}_{\text{true}} + \mathbb{P}_{\theta}}{2} \right\| \right) + \frac{1}{2} \text{KL} \left(\mathbb{P}_{\theta} \left\| \frac{\mathbb{P}_{\text{true}} + \mathbb{P}_{\theta}}{2} \right\| \right),$$

which is

$$\min_{\theta} (2\text{JS}(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) - 2 \log 2) = \min_{\theta} \left(\max_f -\mathbb{E}_{x \sim \mathbb{P}_{\text{true}}} [l(f(x))] - \mathbb{E}_{z \sim \mathbb{P}_{\theta}} [l(-f(z))] \right)$$

where $l(t) = \log(1 + \exp(-t))$ is the logistic loss.

- This is unconstrained, but suffers from a vanishing gradient issue, since for each θ , near-optimal f stays in the tail of the logistic loss.
- Can we design a loss (or distance) that mitigates vanishing gradients without imposing constraints?

⁹Goodfellow et al., Generative adversarial nets, NeurIPS, 2014.

Jensen-Shannon GAN

- Original GAN⁹: minimizes the Jensen-Shannon (JS) divergence:

$$\text{JS}(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) := \frac{1}{2} \text{KL} \left(\mathbb{P}_{\text{true}} \left\| \frac{\mathbb{P}_{\text{true}} + \mathbb{P}_{\theta}}{2} \right\| \right) + \frac{1}{2} \text{KL} \left(\mathbb{P}_{\theta} \left\| \frac{\mathbb{P}_{\text{true}} + \mathbb{P}_{\theta}}{2} \right\| \right),$$

which is

$$\min_{\theta} (2\text{JS}(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) - 2 \log 2) = \min_{\theta} \left(\max_f -\mathbb{E}_{x \sim \mathbb{P}_{\text{true}}} [l(f(x))] - \mathbb{E}_{z \sim \mathbb{P}_{\theta}} [l(-f(z))] \right)$$

where $l(t) = \log(1 + \exp(-t))$ is the logistic loss.

- This is unconstrained, but suffers from a vanishing gradient issue, since for each θ , near-optimal f stays in the tail of the logistic loss.
- Can we design a loss (or distance) that mitigates vanishing gradients without imposing constraints?

⁹Goodfellow et al., Generative adversarial nets, NeurIPS, 2014.

Zero-Infinity GAN

- Zero-Infinity distance:

$$\text{ZI}(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) = \begin{cases} 0, & \mathbb{P}_{\text{true}} = \mathbb{P}_{\theta}, \\ \infty, & \text{otherwise.} \end{cases}$$

- Zero-Infinity (ZI) GAN:¹⁰ (= WGAN w/o constraint)

$$\min_{\theta} \text{ZI}(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) = \min_{\theta} \left(\max_f \mathbb{E}_{x \sim \mathbb{P}_{\text{true}}} [f(x)] - \mathbb{E}_{z \sim \mathbb{P}_{\theta}} [f(z)] \right)$$

- We argue that this least informative distance yields the simplest minimax loss, one that is potentially solvable by gradient methods. (But... really?)

¹⁰Lee-K., Zero-Infinity GAN and implicit bias of extragradient (anchored at zero), 2025

Zero-Infinity GAN

- Zero-Infinity distance:

$$\text{ZI}(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) = \begin{cases} 0, & \mathbb{P}_{\text{true}} = \mathbb{P}_{\theta}, \\ \infty, & \text{otherwise.} \end{cases}$$

- Zero-Infinity (ZI) GAN:¹⁰ (= WGAN w/o constraint)

$$\min_{\theta} \text{ZI}(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) = \min_{\theta} \left(\max_f \mathbb{E}_{x \sim \mathbb{P}_{\text{true}}} [f(x)] - \mathbb{E}_{z \sim \mathbb{P}_{\theta}} [f(z)] \right)$$

- We argue that this least informative distance yields the simplest minimax loss, one that is potentially solvable by gradient methods. (But... really?)

¹⁰Lee-K., Zero-Infinity GAN and implicit bias of extragradient (anchored at zero), 2025

Zero-Infinity GAN

- Zero-Infinity distance:

$$\text{ZI}(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) = \begin{cases} 0, & \mathbb{P}_{\text{true}} = \mathbb{P}_{\theta}, \\ \infty, & \text{otherwise.} \end{cases}$$

- Zero-Infinity (ZI) GAN:¹⁰ (= WGAN w/o constraint)

$$\min_{\theta} \text{ZI}(\mathbb{P}_{\text{true}}, \mathbb{P}_{\theta}) = \min_{\theta} \left(\max_f \mathbb{E}_{x \sim \mathbb{P}_{\text{true}}} [f(x)] - \mathbb{E}_{z \sim \mathbb{P}_{\theta}} [f(z)] \right)$$

- We argue that this least informative distance yields the simplest minimax loss, one that is potentially solvable by gradient methods. (But... really?)

¹⁰Lee-K., Zero-Infinity GAN and implicit bias of extragradient (anchored at zero), 2025

Toy Example: Dirac GAN

- Consider the following simple setting:¹¹

True and generated data: 0 and θ in \mathbb{R} (Linear generator)

True and generated data distributions: $\mathbb{P}_{\text{true}} = \delta_0$ and $\mathbb{P}_{\theta} = \delta_{\theta}$

Linear discriminator: $f(x) = wx$

- Dirac GAN:

$$\min_{\theta} \min_{w \in \mathcal{W}} -l(0) - l(-w\theta)$$

W: $l(t) = -t$ and $\mathcal{W} = \{|w| \leq 1\}$

JS: $l(t) = \log(1 + \exp(-t))$ and $\mathcal{W} = \mathbb{R}$

ZI: $l(t) = -t$ and $\mathcal{W} = \mathbb{R}$ (Unconstrained bilinear)

¹¹Mescheder, Geiger and Nowozin. Which training methods for GANs do actually converge?, ICML, 2018.

Toy Example: Dirac GAN

- Consider the following simple setting:¹¹

True and generated data: 0 and θ in \mathbb{R} (Linear generator)

True and generated data distributions: $\mathbb{P}_{\text{true}} = \delta_0$ and $\mathbb{P}_{\theta} = \delta_{\theta}$

Linear discriminator: $f(x) = wx$

- Dirac GAN:

$$\min_{\theta} \min_{w \in \mathcal{W}} -l(0) - l(-w\theta)$$

W: $l(t) = -t$ and $\mathcal{W} = \{|w| \leq 1\}$

JS: $l(t) = \log(1 + \exp(-t))$ and $\mathcal{W} = \mathbb{R}$

ZI: $l(t) = -t$ and $\mathcal{W} = \mathbb{R}$ (Unconstrained bilinear)

¹¹Mescheder, Geiger and Nowozin. Which training methods for GANs do actually converge?, ICML, 2018.

Toy Example: Dirac GAN (cont'd)

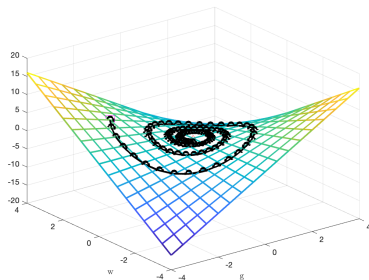
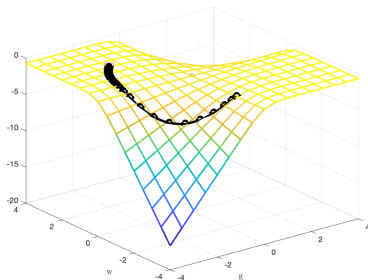


Figure: EG trajectories for Dirac GAN: (L) JS divergence and (R) ZI distance

Strict Non-Optimality in Zero-Infinity GAN

- We investigated the loss landscape of the ZIGAN for a linear generator and a two-layer **neural network** discriminator f . (Lee-K., 2025)
- There are strict non-minimax points in ZIGAN, which two-timescale methods (τ -GDA and τ -EG) can almost surely escape.
- τ -GDA vs. τ -EG?: the latter locally converges to global solutions, while the former does not.

Strict Non-Optimality in Zero-Infinity GAN

- We investigated the loss landscape of the ZIGAN for a linear generator and a two-layer **neural network** discriminator f . (Lee-K., 2025)
- There are strict non-minimax points in ZIGAN, which two-timescale methods (τ -GDA and τ -EG) can almost surely escape.
- τ -GDA vs. τ -EG?: the latter locally converges to global solutions, while the former does not.

Strict Non-Optimality in Zero-Infinity GAN

- We investigated the loss landscape of the ZIGAN for a linear generator and a two-layer **neural network** discriminator f . (Lee-K., 2025)
- There are strict non-minimax points in ZIGAN, which two-timescale methods (τ -GDA and τ -EG) can almost surely escape.
- τ -GDA vs. τ -EG?: the latter locally converges to global solutions, while the former does not.

Conclusion

1. We defined a new strict non-optimality in minimax optimization, named strict non-minimax points, which the two-timescale gradient methods can almost surely escape.
2. We introduced the Zero-Infinity (ZI) GAN that neither requires Lipschitz constraint nor suffers from gradient vanishing.
3. We showed that the ZIGAN has strict non-minimax points.