

L4. NC-NC: interaction dominance, PPM and EGM

Road map: - SC-SC of L_S under Inter. Dom. conditions

- GDA on L_S = damped PPM on L
- PPM computed approximately by EGM.

[GLWM 2023] Grimmer, Lu, Worah, Mirrokni (2023). The landscape of the proximal point method for nonconvex–nonconcave minimax optimization. *Mathematical Programming*, 201(1), 373-407.

[HLG 2024] Hajizadeh, Lu, Grimmer (2024). On the linear convergence of extragradient methods for nonconvex–nonconcave minimax problems. *INFORMS Journal on Optimization*, 6(1), pp.19-31.

$$\min_x \max_y L(x, y) \quad L: \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$$

L is C^2 , l -smooth, ρ -WC-WC, $\frac{1}{s} > \rho$

Recall the saddle envelope

$$L_S(x, y) = \arg \min_u \max_v L(u, v) + \frac{1}{2s} \|u - x\|^2 - \frac{1}{2s} \|v - y\|^2$$

$$z = \begin{bmatrix} x \\ y \end{bmatrix}, \quad z_+ = \begin{bmatrix} x_+ \\ y_+ \end{bmatrix} = \text{Prox}_{s, L}(z)$$

$$\partial = \begin{bmatrix} \nabla_x \\ -\nabla_y \end{bmatrix}, \quad \partial^2 = \begin{bmatrix} \nabla_{xx}^2 & \nabla_{xy}^2 \\ -\nabla_{yx}^2 & -\nabla_{yy}^2 \end{bmatrix}$$

$$\begin{cases} \partial L_S(z) = \frac{1}{s} (z - z_+) = \partial L(z_+) \\ \partial^2 L_S(z) = \frac{1}{s} \left(I - (I + s \partial^2 L(z_+))^{-1} \right) \end{cases}$$

The Hessian expression :

$$\begin{bmatrix} \nabla_{xx}^2 L_S & \nabla_{xy}^2 L_S \\ -\nabla_{yx}^2 L_S & -\nabla_{yy}^2 L_S \end{bmatrix}_z = \frac{1}{s} \left(I - \left(I + s \begin{bmatrix} \nabla_{xx}^2 L & \nabla_{xy}^2 L \\ -\nabla_{yx}^2 L & -\nabla_{yy}^2 L \end{bmatrix}_{z_+} \right)^{-1} \right)$$

by inverse of block matrix,

$$\begin{aligned} \nabla_{xx}^2 L_s(z) &= \frac{1}{s} \left(I - \left(I + s \left[\nabla_{xx}^2 L + \nabla_{xy}^2 L \left(\frac{1}{s} I - \nabla_{yy}^2 L \right)^{-1} \nabla_{yx}^2 L \right] \right) \right)^{-1} \\ -\nabla_{yy}^2 L_s(z) &= \frac{1}{s} \left(I - \left(I + s \left[\nabla_{yy}^2 L + \nabla_{yx}^2 L \left(\frac{1}{s} I + \nabla_{xx}^2 L \right)^{-1} \nabla_{xy}^2 L \right] \right) \right)^{-1} \end{aligned}$$

Consider

$$[\text{matrix}] = \nabla_{xx}^2 L + \nabla_{xy}^2 L \underbrace{\left(\frac{1}{s} I - \nabla_{yy}^2 L \right)^{-1} \nabla_{yx}^2 L}_{\geq (\frac{1}{s} - \rho) I} \geq \nabla_{xx}^2 L$$

Thus if $\nabla_{xx}^2 L(z) \succeq 0$, then $[\text{matrix}] \succeq 0$, and then $\nabla_{xx}^2 L_s(z) \succeq 0$. But it is also possible for $[\text{matrix}] \succeq 0$ even if $\nabla_{xx}^2 L$ N.C., if $\nabla_{xy}^2 L$ is large.

(A6) Interaction dominance

Given $s > 0$, $\exists \alpha(s) > 0$ s.t.

$$\left[\nabla_{xx}^2 L + \nabla_{xy}^2 L \left(\frac{1}{s} I - \nabla_{yy}^2 L \right)^{-1} \nabla_{yx}^2 L \right](z) \succeq \alpha I, \quad \forall z \in \mathbb{R}^n \times \mathbb{R}^n$$

α -Inter. Dom. in x

$$\left[-\nabla_{yy}^2 L + \nabla_{yx}^2 L \left(\frac{1}{s} I + \nabla_{xx}^2 L \right)^{-1} \nabla_{xy}^2 L \right](z) \succeq \alpha I, \quad \forall z$$

α -Inter. Dom. in y

Eq. L has α -Inter. Dom. in x if

$$\frac{\nabla_{xy}^2 L(z) \nabla_{yx}^2 L(z)}{\frac{1}{s} + \rho} \succeq -\nabla_{xx}^2 L(z) + \alpha I, \quad \forall z$$

pf. $(\frac{1}{s} - \nabla_{yy}^2 L) \succeq (\frac{1}{s} + \rho) I$ by that $\nabla_y L$ is ρ -Lip. $\#$

Prop (SC-SC of L_s under Inter. Dom.)

Suppose L has α -Inter. Dom. in x , then

$$\frac{1}{s} I \geq \nabla_{xx}^2 L_s(z) \geq \frac{1}{s + 1/\alpha} I, \quad \forall z.$$

If L has α -Inter. Dom. in y , then

$$\frac{1}{s} I \geq -\nabla_{yy}^2 L_s(z) \geq \frac{1}{s + 1/\alpha} I, \quad \forall z.$$

pf. $\nabla_{xx}^2 L_s = \frac{1}{s} (I - (I + s [\text{matrix}])^{-1})$

$[\text{matrix}] \geq \alpha I$ by Inter. Dom. $\Rightarrow \frac{1}{s + 1/\alpha}$ lower bound

For the upper bound, use that $(I + s [\text{matrix}])^{-1} \geq 0$. $\#$

Now we have that L_s is $\bar{\mu}$ -S.C-SC, $\bar{\mu} = \frac{1}{s + 1/\alpha}$,

under (A1)(A2)(A3). L_s also share stationary pts with L .

Idea: GDA applied to L_s will converge exponentially fast

to the (unique) stationary pt z^* of L_s , by the

\Downarrow
the unique stationary pt of L
(saddle pt)

theory in L2.

We still need to show the $\bar{\mu}$ -smoothness of L_s .

Q. How to compute GDA of L_s ?

Key observation: GDA of L_s = damped PPM of L .

damped PPM $z_k = \begin{bmatrix} x_k \\ y_k \end{bmatrix}, \quad 0 < \lambda \leq 1$

$$\begin{cases} z_t \leftarrow \text{prox}_{s, L}(z_k) \\ z_{k+1} \leftarrow (1-\lambda) z_k + \lambda z_t \end{cases}$$

It's vanilla PPM when $\lambda = 1$.

Prop damped PPM of L with $0 < \lambda \leq 1$ is the same as
GDA of L_s with step size $\bar{s} = \lambda s$.

pf. GDA of L_s :

$$\begin{aligned} z_{k+1} &\leftarrow z_k - \bar{s} \underbrace{\nabla L_s(z_k)}_{\frac{1}{s} \nabla (z_k - z_*)} \\ &= z_k - \lambda (z_k - z_*). \quad \# \end{aligned}$$

This means that damped PPM can converge exponentially
fast to the saddle pt, as long as s, λ properly chosen.

$$s < 1/\rho, \quad 0 < \bar{s} < \frac{2\bar{\kappa}}{\bar{\kappa}^2}.$$

• L -smoothness of saddle envelope

let's revisit the smoothness (regularity) of Moreau envelope.

$$f_s(x) = \min_u f(u) + \frac{1}{2s} \|u - x\|^2, \quad f \text{ } L\text{-smooth, } L\text{-W.C.}$$

$$x_* = \arg\min_u, \quad \exists! \text{ when } \frac{1}{s} > \rho.$$

$$\nabla f_s(x) = \frac{1}{s}(x - x_*) = \nabla f(x_*)$$

We want a Lip-constant of ∇f_s , is $x_*(x)$ Lipschitz?

$$\nabla^2 f_s(x) = \frac{1}{s} [I - (I + s \nabla^2 f(x_*))^{-1}]$$

$$0 < (I + s \nabla^2 f(x_*))^{-1} \leq \frac{1}{1 - \rho s} I$$

$$\Rightarrow -\frac{1}{\rho - s} I \leq \nabla^2 f_s(x) < \frac{1}{s} I$$

Thus ∇f_s has a Lip-constant $\bar{L} = \max \left\{ \frac{1}{s}, \frac{1}{\rho - s} \right\}$

Q when f "smoother" f ? $\bar{L} \leq L$?

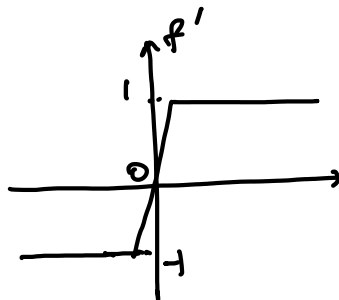
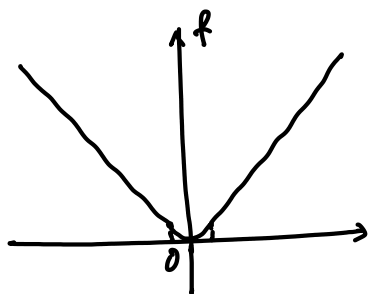
$l \geq \rho$. Eg. $l = 3\rho$, $\frac{1}{s} = 2\rho$, $\frac{1}{\frac{1}{\rho} - s} = \frac{\rho}{1 - \rho s} = 2\rho$,

$\Rightarrow \bar{l} = 2\rho < 3\rho$.

Intuitively, \bar{l} only "see" the interplay btw s and ρ ,
and not the l of f .

Eg. f is close to 1×1 , l is large, but f is O.W.C.

$\Rightarrow \rho > 0$ can be any small positive number



Prop L_s is \bar{l} -smooth with $\bar{l} = \max \left\{ \frac{1}{s}, \frac{1}{|s - \frac{1}{\rho}|} \right\}$.

ref. Prop 2.9 of [GLWM 2023]

- Approximate PPM by extra gradient method (EGM)

The solving of z_+ involves inner-loop,

Recall that $z_k - z_+ = s \partial L(z_+)$

$$z_+ = z_k - s \underbrace{\partial L(z_+)}_{\partial L(z_k) \text{ unknown}}$$

EG: approx z_+ by \hat{z} that can be computed explicitly

formally, $\|z_k - z_+\| = O(s)$,

$$\Rightarrow \|\partial L(z_+) - \partial L(z_k)\| \leq l \|z_k - z_+\| = O(s)$$

$$\left. \begin{array}{l} \text{Then, let } \hat{z} = z_k - s \partial L(z_k) \\ z_+ = z_k - s \partial L(z_+) \end{array} \right\} \Rightarrow \|\hat{z} - z_+\| = O(s^2)$$

idea: z_k is an $O(s)$ approx. to z_* , yet the vanilla gradient step will give z' an $O(s^2)$ approx of z_* .

GD (or GDA) will just use z' , but the extra grad will use z' to evaluate gradient again

$$\left. \begin{aligned} \hat{z} &= z_k - s \partial L(z') \\ z_+ &= z_k - s \partial L(z_+) \end{aligned} \right\} \Rightarrow \|\hat{z} - z_+\| = O(s^3)$$

Then \hat{z} is an $O(s^3)$ approx of z_* .

In many convergence analysis, by choosing small (but finite) step size s , then effect of $O(s^3)$ error can be bounded and then EGM enjoys the same convergence rate as (damped) PPM.

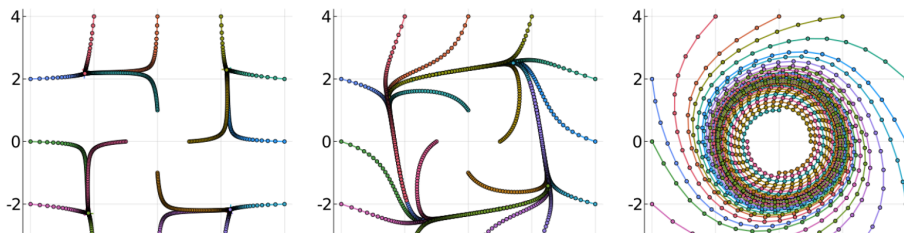
ref: [HLG2024]

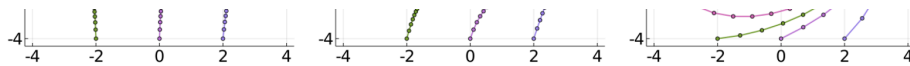
EGM of damped PPM:

$$\left\{ \begin{aligned} z' &\leftarrow z_k - s \partial L(z_k) \\ \hat{z} &\leftarrow z_k - s \partial L(z') \\ z_{k+1} &\leftarrow (1-\lambda)z_k + \lambda \hat{z} \end{aligned} \right. \quad \hat{z} \approx z_*$$

$$\begin{aligned} z_{k+1} &= (1-\lambda)z_k + \lambda(z_k - s \partial L(z')) \\ &= z_k - \lambda s \partial L(z') \end{aligned}$$

$$\min_x \max_y L(x, y) = f(x) + Axy - f(y), \quad \text{with } f(x) = (x-1)(x+1)(x-3)(x+3),$$





(e) EGM w/ $A=1$

(f) EGM w/ $A=10$

(g) EGM w/ $A=100$

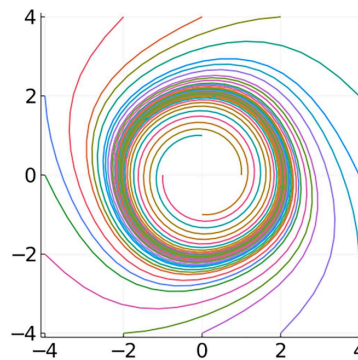
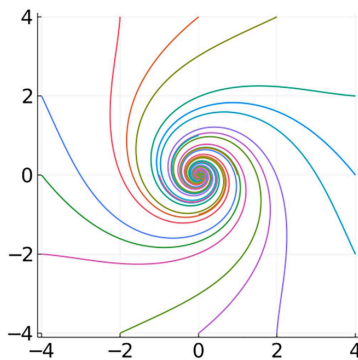
$\lambda=0.01$

(a) Convergence of damped EGM

$\lambda=1$

(b) Cycling of vanilla EGM

$A=100$



Q. What if Inter. Dom. only holds in one of x or y ?

$\left\{ \begin{array}{l} \alpha\text{-Inter. Dom. in both } x \text{ and } y \Rightarrow \text{SC-SC of } L_s \\ \alpha\text{-Inter. Dom. only in } x \text{ or } y \Rightarrow \text{NC-SC of } L_s \end{array} \right.$

The convergence speed degenerate from $O(\log \frac{1}{\epsilon})$ to $O(\frac{1}{\epsilon^2})$
as for the GDA convergence rate for NC-PL type. (L5)

