

Variational Inequality Problem

$$\text{find } x^* \in Q \subseteq \mathbb{R}^d \text{ such that } \langle F(x^*), x - x^* \rangle \geq 0, \forall x \in Q \quad (\text{VIP-C})$$

- $F : Q \rightarrow \mathbb{R}^d$ is L -Lipschitz operator: $\forall x, y \in Q$

$$\|F(x) - F(y)\| \leq L\|x - y\| \quad (1)$$

- F is monotone: $\forall x, y \in Q$

$$\langle F(x) - F(y), x - y \rangle \geq 0 \quad (2)$$

Variational Inequality Problem: Examples

- Minimization problems:

$$\min_{x \in Q} f(x) \quad (4)$$

If f is convex, then (4) is equivalent to finding a stationary point of f , i.e., it is equivalent to (VIP-C) with

$$F(x) = \nabla f(x)$$

Non-Convergence of GD

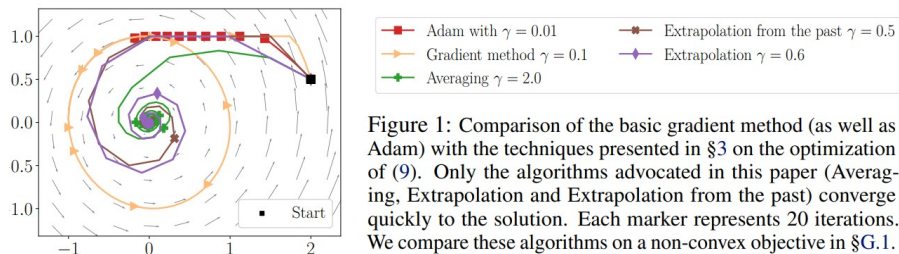


Figure: Behavior of GD on the problem $\min_{u \in \mathbb{R}} \max_{v \in \mathbb{R}} uv$ [Gidel et al., 2019]

9 / 73

Measures of Convergence

- **Restricted gap function:** $\text{Gap}_F(x^K) = \max_{y \in \mathbb{R}^d: \|y - x^*\| \leq R} \langle F(y), x^K - y \rangle$, where $R \sim \|x^0 - x^*\|$ [Nesterov, 2007]
 - ✓ $\text{Gap}_F(x^K)$ can be seen as a natural extension of optimization error for (VIP), when F is monotone
 - ✗ It is unclear how to tightly estimate $\text{Gap}_F(x^K)$ in practice and how to generalize it to non-monotone case
- **Squared norm of the operator:** $\|F(x^K)\|^2$
 - ✗ In general, it provides weaker guarantees than $\text{Gap}_F(x^K)$
 - ✓ $\|F(x^K)\|^2$ is easier to compute than $\text{Gap}_F(x^K)$

In this talk, we focus on the guarantees for $\|F(x^K)\|^2$

When F is monotone and L -Lipschitz the following results are known for EG:

- **Averaged- and best-iterate guarantees:**
 - $\text{Gap}_F(\bar{x}^K) = \mathcal{O}(1/K)$ for $\bar{x}^K = \frac{1}{K+1} \sum_{k=0}^K x^k$ [Nemirovski, 2004, Mokhtari et al., 2019, Hsieh et al., 2019, Monteiro and Svaiter, 2010, Auslender and Teboulle, 2005]
 - $\min_{k=0,1,\dots,K} \|F(x^k)\|^2 = \mathcal{O}(1/K)$ [Solodov and Svaiter, 1999, Ryu et al., 2019]
- **Lower bounds for the last-iterate [Golowich et al., 2020]:**
 - $\text{Gap}_F(x^K) = \Omega(1/\sqrt{K})$
 - $\|F(x^K)\|^2 = \Omega(1/K)$
- **Upper bounds for the last-iterate [Golowich et al., 2020]:** *if additionally the Jacobian $\nabla F(x)$ is Λ -Lipschitz, then*
 - $\text{Gap}_F(x^K) = \mathcal{O}(1/\sqrt{K})$
 - $\|F(x^K)\|^2 = \mathcal{O}(1/K)$

12 / 73

GD Converges Under Star-Cocoercivity

Theorem 1 (Random-iterate convergence of GD)

Let $F : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be ℓ -**star-cocoercive**. Then for all $K \geq 0$ we have

$$\mathbb{E} \|F(\hat{x}^K)\|^2 \leq \frac{\ell \|x^0 - x^*\|^2}{\gamma(K+1)}, \quad (7)$$

where \hat{x}^K is chosen uniformly at random from the set of iterates $\{x^0, x^1, \dots, x^K\}$ produced by GD with $0 < \gamma \leq 1/\ell$.

... and the proof is trivial!

GD Converges Under Star-Cocoercivity

Proof of Theorem 1

Using the update rule of (GD) we derive

$$\begin{aligned} \|x^{k+1} - x^*\|^2 &= \|x^k - \gamma F(x^k) - x^*\|^2 \\ &= \|x^k - x^*\|^2 - 2\gamma \langle x^k - x^*, F(x^k) \rangle + \gamma^2 \|F(x^k)\|^2 \\ &\stackrel{(6)}{\leq} \|x^k - x^*\|^2 - \gamma \left(\frac{2}{\ell} - \gamma \right) \|F(x^k)\|^2. \end{aligned}$$

Rearranging the terms we get

$$\gamma \left(\frac{2}{\ell} - \gamma \right) \|F(x^k)\|^2 \leq \|x^k - x^*\|^2 - \|x^{k+1} - x^*\|^2. \quad (8)$$

It remains to average the above inequalities for $k = 0, 1, \dots, K$.

GD Converges Under Cocoercivity

Theorem 2 (Last-iterate convergence of GD)

Let $F : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be ℓ -cocoercive. Then for all $K \geq 0$ we have

$$\|F(x^K)\|^2 \leq \frac{\ell \|x^0 - x^*\|^2}{\gamma(K+1)}, \quad (9)$$

where x^K is produced by GD with $0 < \gamma \leq 1/\ell$.

The proof is also simple and consist of two steps:

- ① Derivation of $\|F(x^{k+1})\| \leq \|F(x^k)\|$ using ℓ -cocoercivity at x^k and x^{k+1}
- ② Application of the above inequality to the previous result

Key idea: if we manage to show that $F_{\text{EG}, \gamma_1}(x^k)$ is ℓ -cocoercive with some $\ell > 0$ for any monotone and L -Lipschitz F and for a reasonable choice of γ_1 , then we can simply apply the results for GD and we will get the desired last-iterate $\mathcal{O}(1/\kappa)$ convergence rate.

Warm-up: Proximal-Point Operator is Cocoercive

Theorem 3

Let $F : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be monotone and $\gamma > 0$. Then, $F_{\text{PP},\gamma}(x)$ is $2/\gamma$ -cocoercive.

Proof of Theorem 3

In view of Lemma 1, it is enough to prove that $\text{Id} - \gamma F_{\text{PP},\gamma}$ is non-expansive. Consider arbitrary $x, y \in \mathbb{R}^d$ and define \hat{x} and \hat{y} as follows:

$$\hat{x} = x - \gamma F(\hat{x}) = x - \gamma F_{\text{PP},\gamma}(x), \quad \hat{y} = y - \gamma F(\hat{y}) = y - \gamma F_{\text{PP},\gamma}(y).$$

Warm-up: Proximal-Point Operator is Cocoercive

Proof of Theorem 3

Using this notation, we derive

$$\begin{aligned}
 \|\hat{x} - \hat{y}\|^2 &= \|x - y\|^2 - 2\gamma \langle x - y, F(\hat{x}) - F(\hat{y}) \rangle + \gamma^2 \|F(\hat{x}) - F(\hat{y})\|^2 \\
 &= \|x - y\|^2 - 2\gamma \langle \hat{x} + \gamma F(\hat{x}) - \hat{y} - \gamma F(\hat{y}), F(\hat{x}) - F(\hat{y}) \rangle \\
 &\quad + \gamma^2 \|F(\hat{x}) - F(\hat{y})\|^2 \\
 &= \|x - y\|^2 - 2\gamma \langle \hat{x} - \hat{y}, F(\hat{x}) - F(\hat{y}) \rangle - \gamma^2 \|F(\hat{x}) - F(\hat{y})\|^2 \\
 &\stackrel{(2)}{\leq} \|x - y\|^2 - \gamma^2 \|F(\hat{x}) - F(\hat{y})\|^2 \\
 &\leq \|x - y\|^2.
 \end{aligned}$$

That is, $\text{Id} - \gamma F_{\text{PP},\gamma}$ is non-expansive, and, as a result, $F_{\text{PP},\gamma}$ is $2/\gamma$ -cocoercive.

EG is “an Approximation” of PP

$$\text{PP} : x^{k+1} = x^k - \gamma F(x^{k+1}) = x^k - \gamma F_{\text{PP},\gamma}(x^k)$$

$$\text{EG} : x^{k+1} = x^k - \gamma F(x^k - \gamma F(x^k)) = x^k - \gamma F_{\text{EG},\gamma}(x^k)$$

- **Informal explanation:** gradient step $x^k - \gamma F(x^k)$ “approximates” the next point x^{k+1}
- **Formal explanation:** if F is L -Lipschitz and x^{k+1} is obtained via EG, then

$$\begin{aligned} \|F(x^{k+1}) - F(x^k - \gamma F(x^k))\| &\leq L \|x^{k+1} - x^k - \gamma F(x^k)\| \\ &= L\gamma \|F(x^k - \gamma F(x^k)) - F(x^k)\| \\ &\leq L^2\gamma^2 \|F(x^k)\|, \end{aligned}$$

so, for the difference between update directions decreases quadratically in γ

EG and Cocoercivity: What We Obtained

Assume that F is monotone and L -Lipschitz and consider

$$x^{k+1} = x^k - \gamma_2 \underbrace{F(x^k - \gamma_1 F(x^k))}_{F_{EG, \gamma_1}(x^k)} = x^k - \gamma_2 F_{EG, \gamma_1}(x^k)$$

- ✓ If F is linear, i.e., for any $\alpha, \beta \in \mathbb{R}$ and $x, y \in \mathbb{R}^d$ the operator satisfies $F(\alpha x + \beta y) = \alpha F(x) + \beta F(y)$, then operator $F_{EG, \gamma_1}(x)$ with $\gamma_1 \leq 1/L$ is $2/\gamma_1$ -cocoercive $\implies \|F_{EG, \gamma_1}(x^K)\|^2 = \mathcal{O}(1/K)$
- ✓ If $F(x) = Ax + b$ for some $A \in \mathbb{R}^{d \times d}$, $b \in \mathbb{R}^d$, then operator $F_{EG, \gamma_1}(x)$ with $\gamma_1 \leq 1/L$ is $2/\gamma_1$ -cocoercive $\implies \|F_{EG, \gamma_1}(x^K)\|^2 = \mathcal{O}(1/K)$
- ✓✗ If $F(x)$ is not necessarily affine but is star-monotone, i.e., $\langle F(x), x - x^* \rangle \geq 0$ for all $x \in \mathbb{R}^d$, then operator $F_{EG, \gamma_1}(x)$ with $\gamma_1 \leq 1/L$ is $2/\gamma_1$ -star-cocoercive $\implies \min_{k=0,1,\dots,K} \|F_{EG, \gamma_1}(x^k)\|^2 = \mathcal{O}(1/K)$

Proofs are relatively simple and based mainly on Lemmas 1 and 2

PEP for Expansiveness

- Problem (15) is hard to solve since it is infinitely dimensional
 - Let us try to come up with an equivalent finite-dimensional formulation.
- Naive idea №1:** consider the following problem

$$\begin{aligned} \max \quad & \frac{\|x - \gamma_2 x_{F_2} - y + \gamma_2 y_{F_2}\|^2}{\|x - y\|^2} \\ \text{s.t.} \quad & F \text{ is mon. \& } L\text{-Lip.}, x, y \in \mathbb{R}^d, x \neq y, \\ & x_{F_2} = F(x - \gamma_1 x_{F_1}), x_{F_1} = F(x), \\ & y_{F_2} = F(y - \gamma_1 y_{F_1}), y_{F_1} = F(y) \end{aligned} \quad (13)$$

- It is equivalent to (12) but the new problem is finite-dimensional. However, it is still unclear how to check that there exists a monotone and L -Lipschitz operator F such that $F(x) = x_{F_1}, F(y) = y_{F_1}, F(x - \gamma_1 x_{F_1}) = x_{F_2}, F(y - \gamma_1 y_{F_1}) = y_{F_2}$

PEP for Expansiveness

- **Naive idea №2:** consider the following problem

$$\begin{aligned}
 \max \quad & \frac{\|x - \gamma_2 x_{F_2} - y + \gamma_2 y_{F_2}\|^2}{\|x - y\|^2} \\
 \text{s.t.} \quad & \|z_1 - z'_1\|^2 \leq L^2 \|z - z'\|^2, \\
 & \langle z_1 - z'_1, z - z' \rangle \geq 0, \\
 & \text{for each two pairs } (z, z_1), (z', z'_1) \\
 & \text{from } \{(x, x_{F_1}), (y, y_{F_1}), (x - \gamma_1 x_{F_1}, x_{F_2}), (y - \gamma_1 y_{F_1}, y_{F_2})\}
 \end{aligned} \tag{14}$$

- **Bad news:** problem (14) is not equivalent to (13) [Ryu et al., 2020]: feasible set in (14) contains some points that are not feasible for (13), i.e., some feasible points for (14) cannot be interpolated by any monotone and L -Lipschitz operator.

PEP for Expansiveness

- **Good news:** one can circumvent this issue if we focus on a different problem. Let us try to show that for any $\ell > 0$ and any $\gamma_1, \gamma_2 > 0$ there exists a ℓ -cocoercive operator F such that $\text{Id} - \gamma F_{\text{EG}, \gamma_1}$ is not non-expansive.
- In other words, our goal is to show that for all $L, \gamma_1, \gamma_2 > 0$ the quantity

$$\begin{aligned} \rho_{\text{EG}}(\ell, \gamma_1, \gamma_2) = \max \quad & \frac{\|\hat{x} - \hat{y}\|^2}{\|x - y\|^2} \\ \text{s.t.} \quad & F \text{ is } \ell\text{-cocoercive,} \\ & x, y \in \mathbb{R}^d, \ x \neq y, \\ & \hat{x} = x - \gamma_2 F(x - \gamma_1 F(x)), \\ & \hat{y} = y - \gamma_2 F(y - \gamma_1 F(y)) \end{aligned} \tag{15}$$

is bigger than 1, i.e., $\rho_{\text{EG}}(\ell, \gamma_1, \gamma_2) > 1$.

PEP for Expansiveness

- Consider an equivalent finite-dimensional problem:

$$\max \frac{\|x - \gamma_2 x_{F_2} - y + \gamma_2 y_{F_2}\|^2}{\|x - y\|^2} \quad (16)$$

$$\begin{aligned} \text{s.t. } & F \text{ is } \ell\text{-cocoercive, } x, y \in \mathbb{R}^d, x \neq y, \\ & x_{F_2} = F(x - \gamma_1 x_{F_1}), x_{F_1} = F(x), \\ & y_{F_2} = F(y - \gamma_1 y_{F_1}), y_{F_1} = F(y). \end{aligned}$$

- Next, for all $\alpha > 0$ the following equivalence holds:

$$F \text{ is } \ell\text{-cocoercive} \iff (\alpha^{-1}\text{Id}) \circ F \circ (\alpha\text{Id}) \text{ is } \ell\text{-cocoercive.}$$

33 / 73

PEP for Expansiveness

- Proposition 2 from Ryu et al. [2020] implies that (17) is equivalent to

$$\begin{aligned}
 \max \quad & \|x - \gamma_2 x_{F_2} - y + \gamma_2 y_{F_2}\|^2 \\
 \text{s.t.} \quad & x, y, x_{F_1}, y_{F_1}, x_{F_2}, y_{F_2} \in \mathbb{R}^d, \|x - y\|^2 = 1, \\
 & \ell \langle x_{F_1} - x_{F_2}, \gamma_1 x_{F_1} \rangle \geq \|x_{F_1} - x_{F_2}\|^2, \\
 & \ell \langle x_{F_1} - y_{F_1}, x - y \rangle \geq \|x_{F_1} - y_{F_1}\|^2, \\
 & \ell \langle x_{F_1} - y_{F_2}, x - y + \gamma_1 y_{F_1} \rangle \geq \|x_{F_1} - y_{F_2}\|^2, \\
 & \ell \langle x_{F_2} - y_{F_1}, x - \gamma_1 x_{F_1} - y \rangle \geq \|x_{F_2} - y_{F_1}\|^2, \\
 & \ell \langle x_{F_2} - y_{F_2}, x - \gamma_1 x_{F_1} - y + \gamma_1 y_{F_1} \rangle \geq \|x_{F_2} - y_{F_2}\|^2, \\
 & \ell \langle y_{F_1} - y_{F_2}, \gamma_1 y_{F_1} \rangle \geq \|y_{F_1} - y_{F_2}\|^2.
 \end{aligned} \tag{18}$$

The problem is linear in terms of the pairwise inner products of
 $x, y, x_{F_1}, y_{F_1}, x_{F_2}, y_{F_2}$

PEP for Expansiveness

- Therefore, problem (18) is equivalent to the following SDP problem:

$$\begin{aligned} \max \quad & \text{Tr}(\mathbf{M}_0 \mathbf{G}) \\ \text{s.t.} \quad & \mathbf{G} \in \mathbb{S}_+^6, \\ & \text{Tr}(\mathbf{M}_i \mathbf{G}) \geq 0, \quad i = 1, 2, \dots, 6, \\ & \text{Tr}(\mathbf{M}_7 \mathbf{G}) = 1, \end{aligned} \quad (19)$$

where M_0, \dots, M_7 are some symmetric matrices.

PEP for Expansiveness: M_0

$$M_0 = \begin{pmatrix} 1 & -1 & 0 & 0 & -\gamma_2 & \gamma_2 \\ -1 & 1 & 0 & 0 & \gamma_2 & -\gamma_2 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ -\gamma_2 & \gamma_2 & 0 & 0 & \gamma_2^2 & -\gamma_2^2 \\ \gamma_2 & -\gamma_2 & 0 & 0 & -\gamma_2^2 & \gamma_2^2 \end{pmatrix}$$

PEP for Expansiveness: M_1

$$M_1 = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \ell\gamma_1 - 1 & 0 & 1 - \frac{\ell\gamma_1}{2} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 - \frac{\ell\gamma_1}{2} & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

PEP for Expansiveness: M_2

$$M_2 = \begin{pmatrix} 0 & 0 & \frac{\ell}{2} & -\frac{\ell}{2} & 0 & 0 \\ 0 & 0 & -\frac{\ell}{2} & \frac{\ell}{2} & 0 & 0 \\ \frac{\ell}{2} & -\frac{\ell}{2} & -1 & 1 & 0 & 0 \\ -\frac{\ell}{2} & \frac{\ell}{2} & 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

PEP for Expansiveness: M_3

$$M_3 = \begin{pmatrix} 0 & 0 & \frac{\ell}{2} & 0 & 0 & -\frac{\ell}{2} \\ 0 & 0 & -\frac{\ell}{2} & 0 & 0 & \frac{\ell}{2} \\ \frac{\ell}{2} & -\frac{\ell}{2} & -1 & \frac{\ell\gamma_1}{2} & 0 & 1 \\ 0 & 0 & \frac{\ell\gamma_1}{2} & 0 & 0 & -\frac{\ell\gamma_1}{2} \\ 0 & 0 & 0 & 0 & 0 & 0 \\ -\frac{\ell}{2} & \frac{\ell}{2} & 1 & -\frac{\ell\gamma_1}{2} & 0 & -1 \end{pmatrix}$$

PEP for Expansiveness: M_4

$$M_4 = \begin{pmatrix} 0 & 0 & 0 & -\frac{\ell}{2} & \frac{\ell}{2} & 0 \\ 0 & 0 & 0 & \frac{\ell}{2} & -\frac{\ell}{2} & 0 \\ 0 & 0 & 0 & \frac{\ell\gamma_1}{2} & -\frac{\ell\gamma_1}{2} & 0 \\ -\frac{\ell}{2} & \frac{\ell}{2} & \frac{\ell\gamma_1}{2} & -1 & 1 & 0 \\ \frac{\ell}{2} & -\frac{\ell}{2} & -\frac{\ell\gamma_1}{2} & 1 & -1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

PEP for Expansiveness: M_5

$$M_5 = \begin{pmatrix} 0 & 0 & 0 & 0 & \frac{\ell}{2} & -\frac{\ell}{2} \\ 0 & 0 & 0 & 0 & -\frac{\ell}{2} & \frac{\ell}{2} \\ 0 & 0 & 0 & 0 & -\frac{\ell\gamma_1}{2} & \frac{\ell\gamma_1}{2} \\ 0 & 0 & 0 & 0 & \frac{\ell\gamma_1}{2} & -\frac{\ell\gamma_1}{2} \\ \frac{\ell}{2} & -\frac{\ell}{2} & -\frac{\ell\gamma_1}{2} & \frac{\ell\gamma_1}{2} & -1 & 1 \\ -\frac{\ell}{2} & \frac{\ell}{2} & \frac{\ell\gamma_1}{2} & -\frac{\ell\gamma_1}{2} & 1 & -1 \end{pmatrix}$$

PEP for Expansiveness: M_6

$$M_6 = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \ell\gamma_1 - 1 & 0 & 1 - \frac{\ell\gamma_1}{2} \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 - \frac{\ell\gamma_1}{2} & 0 & -1 \end{pmatrix}$$

PEP for Expansiveness: M_7

$$M_7 = \begin{pmatrix} 1 & -1 & 0 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

F_{EG, γ_1} is Not Cocoercive: Numerical Proof

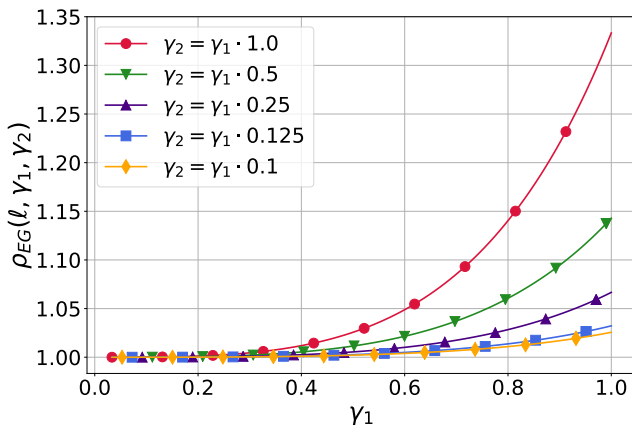


Figure: Numerical estimation of $\rho_{EG}(\ell, \gamma_1, \gamma_2)$ defined in (15) for $\ell = 1$ and different γ_1, γ_2 .

F_{EG,γ_1} is Not Cocoercive?

- ✓ We obtained the answer numerically for different choices of γ_1 and γ_2
- ✗ It is not a rigorous proof: probably, for smaller stepsize F_{EG,γ_1} is cocoercive, but we cannot check it because of the numerical inaccuracies

Analytical example is required

How to Construct Analytical Example?

- Try to solve the problem symbolically after some simplifications of the problem
 - ✓ Ryu et al. [2020] use this trick and obtained quite impressive results that are almost impossible to obtain by hands
 - ✗ Unfortunately, this approach does not always work and it did not help us to get the example
- Try to solve the problem numerically for different parameters γ_1, γ_2 and ℓ to guess the dependencies using visualization
 - ✓ Gu and Yang [2019] successfully applied this technique to derive worst-case examples for PP
 - ✗ It is hard to visualize d -dimensional examples with $d \geq 3$

Low-Dimensional Examples: Trace Heuristic

Instead of solving

$$\begin{aligned} \max \quad & \text{Tr}(M_0 G) \\ \text{s.t.} \quad & G \in \mathbb{S}_+^6, \\ & \text{Tr}(M_i G) \geq 0, \quad i = 1, 2, \dots, 6, \\ & \text{Tr}(M_7 G) = 1, \end{aligned}$$

which gives us 5-dimensional examples of G , we consider another problem [Taylor et al., 2017a]:

$$\begin{aligned} \max \quad & \text{Tr}(G) \\ \text{s.t.} \quad & G \in \mathbb{S}_+^6, \\ & \text{Tr}(M_0 G) \geq 1.0005, \\ & \text{Tr}(M_i G) \geq 0, \quad i = 1, 2, \dots, 6, \\ & \text{Tr}(M_7 G) = 1. \end{aligned}$$

Low-Dimensional Examples: Log-Det Heuristic

To overcome this issue, we consider another problem with so called Log-det heuristic [Fazel et al., 2003]:

$$\begin{aligned}
 \min \quad & \log \det (G + \delta I) \\
 \text{s.t.} \quad & G \in \mathbb{S}_+^6, \\
 & \text{Tr}(M_0 G) \geq 1.0005, \\
 & \text{Tr}(M_i G) \geq 0, \quad i = 1, 2, \dots, 6, \\
 & \text{Tr}(M_7 G) = 1,
 \end{aligned} \tag{20}$$

where $\delta > 0$ is some small positive regularization parameter. For simplicity we used $\gamma_2 = \gamma_1$ in some interval and $\ell = 1$.

- ✓ We obtained solutions of rank 2, i.e., we obtained $x, y, x_{F_1}, y_{F_1}, x_{F_2}, y_{F_2}$ in \mathbb{R}^2
- ✓ We observed that $x = -y$ for all tested values of γ_1

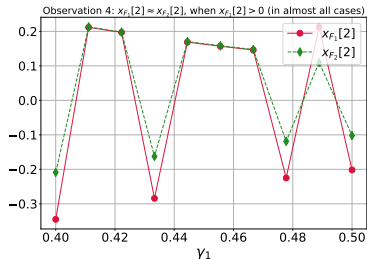
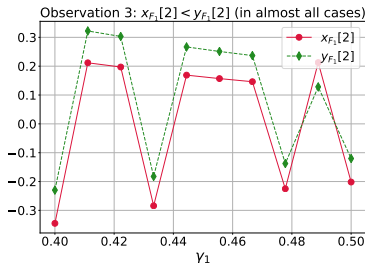
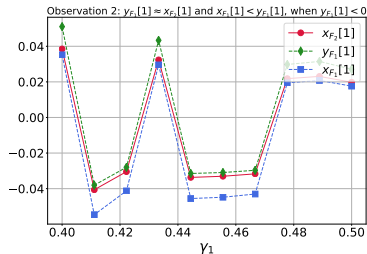
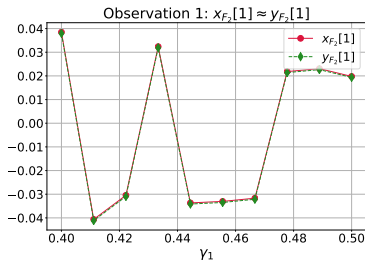
Low-Dimensional Examples: Log-Det Heuristic

- ✗ However, numerical solutions were not consistent enough to guess the right dependencies
- ✓ To overcome this issue, we
 - rotated $x, y, x_{F_1}, y_{F_1}, x_{F_2}, y_{F_2}$ in such a way that $x = (-1/2, 0)^\top$, $y = (1/2, 0)^\top$,
 - plotted the components of $x_{F_1}, y_{F_1}, x_{F_2}, y_{F_2}$ for different γ_1
- ✓ Although the resulting dependencies were not perfect, the obtained plots helped us to sequentially construct the needed example:

$$\begin{aligned}
 x &= \begin{pmatrix} -\frac{1}{2} \\ 0 \end{pmatrix}, & y &= \begin{pmatrix} \frac{1}{2} \\ 0 \end{pmatrix}, & x_{F_1} &= \begin{pmatrix} -\frac{1}{2\gamma_1} \\ \frac{1}{2\gamma_1} \end{pmatrix}, & y_{F_1} &= \begin{pmatrix} -\frac{1-\gamma_1\ell}{2\gamma_1} \\ \frac{1+\gamma_1\ell}{2\gamma_1} \end{pmatrix}, \\
 x_{F_2} &= \begin{pmatrix} -\frac{1-\gamma_1\ell}{2\gamma_1} \\ \frac{1}{2\gamma_1} \end{pmatrix}, & y_{F_2} &= \begin{pmatrix} -\frac{1-\gamma_1\ell}{2\gamma_1} \\ \frac{1-\gamma_1^2\ell^2}{2\gamma_1} \end{pmatrix}
 \end{aligned} \tag{21}$$

- Required several days of playing with plots to get the needed insights

Non-Cocoercivity of F_{EG, γ_1} : Four Observations



Non-Cocoercivity of F_{EG, γ_1} : Four Observations

Mimicking the observed dependencies, we assumed that

$$\begin{aligned} x_{F_2}[1] &= y_{F_2}[1], \\ y_{F_1}[1] &= x_{F_2}[1] \quad \text{and} \quad x_{F_1}[1] < y_{F_1}[1] < 0, \\ 0 &< x_{F_1}[2] < y_{F_1}[2], \\ x_{F_1}[2] &= x_{F_2}[2] \end{aligned}$$

Non-Cocoercivity of F_{EG,γ_1} : Handling Four Observations

After that, we plugged them in the interpolation conditions from (18), and obtained the following inequalities:

$$y_{F_1}[1] \leq (1 - \gamma_1)x_{F_1}[1],$$

$$y_{F_1}[2] \leq \frac{y_{F_2}[2]}{1 - \gamma_1},$$

$$y_{F_1}[2] \leq (1 + \gamma_1)x_{F_2}[2],$$

$$x_{F_2}[2] \leq \frac{y_{F_2}[2]}{1 - \gamma_1^2}$$

55 / 73

Non-Cocoercivity of F_{EG, γ_1} : Making More Assumptions

- Using these dependencies in the remaining interpolation conditions, we derived

$$x_{F_1}[1] + \gamma_1(x_{F_1}[1])^2 + \frac{\gamma_1(y_{F_2}[2])^2}{(1 - \gamma_1^2)^2} \leq 0.$$

- After that, we assumed that

$$y_{F_2}[2] = -x_{F_1}[1](1 - \gamma_1^2).$$

- Together with previous inequality it gives

$$x_{F_1}[1] + 2\gamma_1(x_{F_1}[1])^2 \leq 0.$$

- Next, we chose $x_{F_1}[1] = -1/2\gamma_1$ and put it in all previously derived dependencies.
- Finally, we generalized the example to the case of non-unit ℓ using “physical-dimension” arguments and got (21).

EG and Cocoercivity: Preliminary Conclusions

- Now we know that one cannot apply analysis of GD to prove last-iterate $\mathcal{O}(1/\kappa)$ convergence for EG
- We observed another significant difference between PP and EG: $F_{PP,\gamma}$ is cocoercive and $F_{EG,\gamma}$ is not
- But does it mean that it is impossible to prove $\mathcal{O}(1/\kappa)$ convergence rate for EG in the considered setup (F is monotone and L -Lipschitz)? No, it does not!

Example of the Proof [De Klerk et al., 2017]

Set $f_i = f(\mathbf{x}_i)$ and $\mathbf{g}_i = \nabla f(\mathbf{x}_i)$ for $i \in \{*, 0, 1\}$. Note that $\mathbf{g}_* = \mathbf{0}$. The following five inequalities are now satisfied:

$$\begin{aligned} 1: \quad & f_0 \geq f_1 + \mathbf{g}_1^\top (\mathbf{x}_0 - \mathbf{x}_1) + \frac{1}{2(1-\mu/L)} \left(\frac{1}{L} \|\mathbf{g}_0 - \mathbf{g}_1\|^2 + \mu \|\mathbf{x}_0 - \mathbf{x}_1\|^2 - 2\frac{\mu}{L} (\mathbf{g}_1 - \mathbf{g}_0)^\top (\mathbf{x}_1 - \mathbf{x}_0) \right) \\ 2: \quad & f_* \geq f_0 + \mathbf{g}_0^\top (\mathbf{x}_* - \mathbf{x}_0) + \frac{1}{2(1-\mu/L)} \left(\frac{1}{L} \|\mathbf{g}_* - \mathbf{g}_0\|^2 + \mu \|\mathbf{x}_* - \mathbf{x}_0\|^2 - 2\frac{\mu}{L} (\mathbf{g}_0 - \mathbf{g}_*)^\top (\mathbf{x}_0 - \mathbf{x}_*) \right) \\ 3: \quad & f_* \geq f_1 + \mathbf{g}_1^\top (\mathbf{x}_* - \mathbf{x}_1) + \frac{1}{2(1-\mu/L)} \left(\frac{1}{L} \|\mathbf{g}_* - \mathbf{g}_1\|^2 + \mu \|\mathbf{x}_* - \mathbf{x}_1\|^2 - 2\frac{\mu}{L} (\mathbf{g}_1 - \mathbf{g}_*)^\top (\mathbf{x}_1 - \mathbf{x}_*) \right) \\ 4: \quad & -\mathbf{g}_0^\top \mathbf{g}_1 \geq 0 \\ 5: \quad & \mathbf{g}_1^\top (\mathbf{x}_0 - \mathbf{x}_1) \geq 0. \end{aligned}$$

Indeed, the first three inequalities are the $\mathcal{F}_{\mu,L}$ -interpolability conditions, the fourth inequality is a relaxation of (4), and the fifth inequality is a relaxation of (3).

We aggregate these five inequalities by defining the following positive multipliers,

$$y_1 = \frac{L - \mu}{L + \mu}, \quad y_2 = 2\mu \frac{(L - \mu)}{(L + \mu)^2}, \quad y_3 = \frac{2\mu}{L + \mu}, \quad y_4 = \frac{2}{L + \mu}, \quad y_5 = 1, \quad (9)$$

and adding the five inequalities together after multiplying each one by the corresponding multiplier.

The result is the following inequality (as may be verified directly):

$$f_1 - f_* \leq \left(\frac{L-\mu}{L+\mu} \right)^2 (f_0 - f_*) - \frac{\mu L(L+3\mu)}{2(L+\mu)^2} \left\| \mathbf{x}_0 - \frac{L+\mu}{L+3\mu} \mathbf{x}_1 - \frac{2\mu}{L+3\mu} \mathbf{x}_* - \frac{3L+\mu}{L^2+3\mu L} \mathbf{g}_0 - \frac{L+\mu}{L^2+3\mu L} \mathbf{g}_1 \right\|^2 - \frac{2L\mu^2}{L^2+2L\mu-3\mu^2} \left\| \mathbf{x}_1 - \mathbf{x}_* - \frac{(L-\mu)^2}{2\mu L(L+\mu)} \mathbf{g}_0 - \frac{L+\mu}{2\mu L} \mathbf{g}_1 \right\|^2. \quad (10)$$

Last-Iterate $\mathcal{O}(1/\kappa)$ Rate for EG

- However, guessing the dependencies is not always an easy task: the dependencies on the parameters of the problem like L, γ_1, γ_2 might be quite tricky
- To simplify the process of guessing the proof, we consider a simpler problem:

$$\begin{aligned} \Delta_{\text{EG}}(L, \gamma_1, \gamma_2) = & \max \quad \|F(x^1)\|^2 - \|F(x^0)\|^2 \\ \text{s.t.} \quad & F \text{ is monotone and } L\text{-Lipschitz, } x^0 \in \mathbb{R}^d, \\ & \|x^0 - x^*\|^2 \leq 1, \\ & x^1 = x^0 - \gamma_2 F(x^0 - \gamma_1 F(x^0)) \end{aligned} \quad (24)$$

with $\gamma_1 = \gamma_2 = \gamma$

Last-Iterate $\mathcal{O}(1/\kappa)$ Rate for EG

- In the numerical tests, we observed that $\Delta_{\text{EG}}(L, \gamma_1, \gamma_2) \approx 0$ for all tested pairs of L and γ
- Moreover, the dual variables $\lambda_1, \lambda_2, \lambda_3$ that correspond to the constraints

$$0 \leq \frac{1}{\gamma} \langle F(x^k) - F(x^{k+1}), x^k - x^{k+1} \rangle,$$

$$0 \leq \frac{1}{\gamma} \langle F(x^k - \gamma F(x^k)) - F(x^{k+1}), x^k - \gamma F(x^k) - x^{k+1} \rangle,$$

$$\|F(x^k - \gamma F(x^k)) - F(x^{k+1})\|^2 \leq L^2 \|x^k - \gamma F(x^k) - x^{k+1}\|^2$$

are always close to the constants 2, $1/2$, and $3/2$

- Although λ_2 and λ_3 were sometimes slightly smaller, e.g., sometimes we had $\lambda_2 \approx 3/5$ and $\lambda_3 \approx 13/20$, we simplified these dependencies and simply summed up the corresponding inequalities with weights $\lambda_1 = 2$, $\lambda_2 = 1/2$ and $\lambda_3 = 3/2$ respectively
- After that it was just needed to rearrange the terms and apply Young's inequality to some inner products.

On Our Failures Towards the Proof

The proof that I presented was not the first idea of what we tried to obtain. Here are some of the claims that we tried to prove first:

- We tried to show that $\|F_{\text{EG}, \gamma_1}(x^{k+1})\| \leq \|F_{\text{EG}, \gamma_1}(x^k)\|$ for a reasonable choice of γ_1 and γ_2
 - ✗ Perhaps, surprisingly, but this is not true even for L -cocoercive F : we observed this phenomenon via PEP
- We also tried to show that $\|F(x^{k+1})\| \leq \|F(x^k)\|$ when $\gamma_2 < \gamma_1$
 - ✗ Again, via PEP we observed that this is not true even for L -cocoercive F

The usage of PEP helped us to save a lot of time from trying to prove the claims that do not hold in general!

Other Results in the Paper

- We obtain some “pessimistic” results on Optimistic Gradient method (OG):
 - ✗ for the two popular representations of OG (for the classical one and for the extrapolation from the past) we proved that corresponding operators are not even star-cocoercive
- For Hamiltonian Gradient method (HGM) [Balduzzi et al., 2018] we also
 - ✗ showed non-cocoercivity of corresponding operator when F is monotone and Lipschitz
 - ✓ derived best-iterate $\mathcal{O}(1/\kappa)$ convergence rate in terms of the squared norm of the gradient of the Hamiltonian function $\mathcal{H}(x) = \frac{1}{2}\|F(x)\|^2$ when F and ∇F are Lipschitz-continuous but F is not necessary monotone

A. Auslender and M. Teboulle. Interior projection-like methods for monotone variational inequalities. *Mathematical programming*, 104(1):39–68, 2005.

D. Balduzzi, S. Racaniere, J. Martens, J. Foerster, K. Tuyls, and T. Graepel. The mechanics of n-player differentiable games. In *International Conference on Machine Learning*, pages 354–363. PMLR, 2018.

H. H. Bauschke, P. L. Combettes, et al. *Convex analysis and monotone operator theory in Hilbert spaces*, volume 408. Springer, 2011.

A. Ben-Tal, L. El Ghaoui, and A. Nemirovski. *Robust optimization*. Princeton university press, 2009.

E. De Klerk. *Aspects of semidefinite programming: interior point algorithms and selected applications*, volume 65. Springer Science & Business Media, 2006.

E. De Klerk, F. Glineur, and A. B. Taylor. On the worst-case complexity of the gradient method with exact line search for smooth strongly convex functions. *Optimization Letters*, 11(7):1185–1199, 2017.

- Y. Drori and M. Teboulle. Performance of first-order methods for smooth convex minimization: a novel approach. *Mathematical Programming*, 145(1):451–482, 2014.
- M. Fazel, H. Hindi, and S. P. Boyd. Log-det heuristic for matrix rank minimization with applications to hankel and euclidean distance matrices. In *Proceedings of the 2003 American Control Conference, 2003.*, volume 3, pages 2156–2162. IEEE, 2003.
- G. Gidel, H. Berard, P. Vincent, and S. Lacoste-Julien. A variational inequality perspective on generative adversarial nets. In *ICLR*, 2019.
- N. Golowich, S. Pattathil, C. Daskalakis, and A. Ozdaglar. Last iterate is slower than averaged iterate in smooth convex-concave saddle point problems. In *Conference on Learning Theory*, pages 1758–1784. PMLR, 2020.
- I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.

- I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. *ICLR 2015*, 2015.
- G. Gu and J. Yang. Optimal nonergodic sublinear convergence rate of proximal point algorithm for maximal monotone inclusion problems. *arXiv preprint arXiv:1904.05495*, 2019.
- M. Hast, K. J. Åström, B. Bernhardsson, and S. Boyd. Pid design by convex-concave optimization. In *2013 European Control Conference (ECC)*, pages 4460–4465. IEEE, 2013.
- Y.-G. Hsieh, F. Iutzeler, J. Malick, and P. Mertikopoulos. On the convergence of single-call stochastic extra-gradient methods. *Advances in Neural Information Processing Systems*, 32:6938–6948, 2019.
- D. Kim and J. A. Fessler. Optimized first-order methods for smooth convex minimization. *Mathematical programming*, 159(1):81–107, 2016.
- G. M. Korpelevich. The extragradient method for finding saddle points and other problems. *Matecon*, 12:747–756, 1976.

- L. Lessard, B. Recht, and A. Packard. Analysis and design of optimization algorithms via integral quadratic constraints. *SIAM Journal on Optimization*, 26(1):57–95, 2016.
- A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to adversarial attacks. *ICLR 2018*, 2018.
- B. Martinet. Regularisation d’inequations variationelles par approximations successives. *Revue Francaise d’Informatique et de Recherche Operationelle*, 4: 154–159, 1970.
- A. Mokhtari, A. Ozdaglar, and S. Pattathil. Proximal point approximations achieving a convergence rate of $O(1/k)$ for smooth convex-concave saddle point problems: Optimistic gradient and extra-gradient methods. *arXiv preprint arXiv:1906.01115*, 2019.
- R. D. Monteiro and B. F. Svaiter. On the complexity of the hybrid proximal extragradient method for the iterates and the ergodic mean. *SIAM Journal on Optimization*, 20(6):2755–2787, 2010.

References V

- A. Nemirovski. Prox-method with rate of convergence $O(1/t)$ for variational inequalities with lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 15(1): 229–251, 2004.
- Y. Nesterov. Dual extrapolation and its applications to solving variational inequalities and related problems. *Mathematical Programming*, 109(2):319–344, 2007.
- L. D. Popov. A modification of the arrow-hurwicz method for search of saddle points. *Mathematical notes of the Academy of Sciences of the USSR*, 28(5): 845–848, 1980.
- R. T. Rockafellar. Monotone operators and the proximal point algorithm. *SIAM journal on control and optimization*, 14(5):877–898, 1976.
- E. K. Ryu, K. Yuan, and W. Yin. Ode analysis of stochastic gradient methods with optimism and anchoring for minimax problems. *arXiv preprint arXiv:1905.10899*, 2019.

