

Program:

- ⊙ The notion of a min-max problem
- ⊙ Game-theoretic interpretation
- ⊙ Some applications

⊙ Convex-concave case, saddle points, minimax thm & duality, duality gap.

▶ Variational inequality viewpoint

▶ "Baseline": gradient descent-ascent (GDA), aka. Saddle-Point Mirror Descent).

- Stochastic Oracles & ~~Non-Euclidean structure.~~

▶ Extragradient method (aka. Mirror Prox):

- Motivation & Key Idea

- $O(\frac{1}{\epsilon})$ convergence result

~~▶ Alternative approach: Nesterov's smoothing~~

▶ Sampling for bilinear problems.

⊙ Convex-strongly-concave case

⊙ Nonconvex-[strongly]-concave case:

?

⊙ Algorithm for finding 1st-order stationary points

⊙ Convergence rates

let $X \subseteq \mathbb{R}^d$ be a convex body (set with nonempty interior)

Def.: $f: X \rightarrow \mathbb{R}$ is (closed) convex if $\forall x \in X$ one has:

$$f(y) \geq f(x) + \langle f'(x), y - x \rangle$$

for some vector $f'(x)$ called subgradient of $f(\cdot)$ at x , and $\forall y \in X$

Def.: Subdifferential: $\partial f(x) = \{ \text{all such } f'(x) \}$

⊙ If f is differentiable at x , then $\partial f(x) = \{ \nabla f(x) \}$

⊙ If $f(x) = \max_{y \in Y} F(x, y)$ where $F(\cdot, y)$ is convex and

differentiable for each $y \in Y$, then $\partial f(x)$ is the convex hull of the

active^o gradients: $\partial f(x) = \text{Conv} \left(\{ \nabla_x F(x, y^*) \mid y^* \in \text{Argmax}_{y \in Y} F(x, y) \} \right)$

In particular, $\nabla f(x) = \nabla_x F(x, y^*(x))$ if $y^*(x)$ is a unique maximizer.

Black-Box model of convex optimization:

$$\min_{x \in X} f(x)$$

- ⊙ $X \subseteq E_x$ is a known convex body \Rightarrow can project on it
- ⊙ $f: X \rightarrow \mathbb{R}$ is an unknown function from certain class, available via a 1st-order oracle that returns $f'(x)$ at any $x \in X$.

Problem classes:

▶ Convex and B -Lipschitz: $\|f'(x)\| \leq B$.

▶ Convex and L -smooth (differentiable with L -Lipschitz gradient):

$$0 \leq f(y) - f(x) - \langle \nabla f(x), y - x \rangle \leq \frac{L}{2} \|y - x\|^2$$

▶ L -smooth and λ -strongly convex: ($\lambda > 0$):

$$\frac{\lambda}{2} \|y - x\|^2 \leq f(y) - f(x) - \langle \nabla f(x), y - x \rangle \leq \frac{L}{2} \|y - x\|^2$$

[Black-Box model of convex optimization:]

$$\min_{x \in X} f(x)$$

- ⊙ $X \subseteq E_x$ is a known convex body \Rightarrow can project on it
- ⊙ $f: X \rightarrow \mathbb{R}$ is an unknown function from certain class, available via a 1st-order oracle that returns $f'(x)$ at any $x \in X$.

Accuracy Measures $f(\hat{x}) - f^*$ where $f^* := \min_{x \in X} f(x)$ - objective error.

First-order algorithms construct a sequence of points $(\hat{x}_1, \dots, \hat{x}_T)$ such that $f(\hat{x}_T) - f^* \leq \epsilon(T)$ for $\epsilon(T)$ depending on the problem parameters.

#1: $\epsilon(T) = O(1) \frac{BD}{\sqrt{T}}$ for B -Lipschitz & convex and $\text{diam}(X) = D$,
with (projected) subgradient method.

#2: $\epsilon(T) = O(1) \frac{LD^2}{T}$ for L -smooth & convex, with (projected) grad. descent.

Minimax Problems

$$\min_{x \in X} \max_{y \in Y} f(x, y)$$

⊙ Game-theoretic interpretation: zero-sum game between Min and Max. $f(x, y)$ is the payoff of Max and the loss of Min when they choose a pair of strategies (x, y) .

Who goes first, Min or Max?

⊙ Primal & Dual problems:

$$(P) \min_{x \in X} \left\{ \varphi(x) := \max_{y \in Y} f(x, y) \right\} \geq \max_{y \in Y} \left\{ \psi(y) := \min_{x \in X} f(x, y) \right\} (D)$$

⊙ Weak duality: $(P) \geq (D)$ "Choosing the second is better."

Proof: Let $x^* \in \text{Argmin } \varphi(x)$, $y^* \in \text{Argmax } \psi(y)$. Then
 $\varphi(x^*) = \max_{y \in Y} f(x^*, y) \geq f(x^*, y^*) \geq \min_{x \in X} f(x, y^*) = \psi(y^*)$.

Examples

① Robust system design:

Ex: Adversarially robust training

$$\min_{x \in X} \max_{\|y - \bar{y}\| \leq \delta} f(x, y)$$

\uparrow control parameter \uparrow nominal input \nwarrow cost function

② Linear Regression in l_p norms:

Bilinear SPP;
 $Y =$ the unit dual norm ball ($\frac{1}{p} + \frac{1}{q} = 1$)

$$\min_{x \in X} \frac{\|Ax - b\|_p}{\varphi(x)}$$

$$= \min_{x \in X} \max_{\|y\|_q \leq 1} \langle y, Ax - b \rangle$$

③ Minimization of a maximum of m functions:

$$\min_{x \in X} \max_{i \in \{1, \dots, m\}} g_i(x) = \min_{x \in X} \max_{y \in \Delta_m} \langle y, G(x) \rangle$$

where Δ_m is the standard simplex; $G(x) = \begin{bmatrix} g_1(x) \\ \vdots \\ g_m(x) \end{bmatrix}$.

Objective is linear (\Rightarrow concave) in y ;

⊙ smooth if all $g_i(x)$ are smooth, and convex if $g_i(x)$ are so.

④ Finite-sum minimization:

$$\min_{x \in X} \sum_{i=1}^n \ell(x, z_i) \quad \text{where } \ell(x, z) \geq 0$$

Then $\sum_{i=1}^n \ell(x, z_i) = \max_{y \in [0, 1]^n} \langle y, L_n(x) \rangle$ where $L_n(x) = \begin{bmatrix} \ell(x, z_1) \\ \vdots \\ \ell(x, z_n) \end{bmatrix}$.

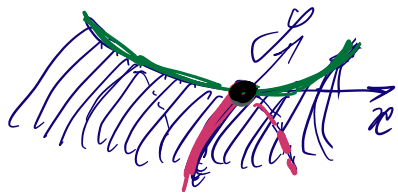
Convex-Concave Saddle-Point Problems

$$\min_{x \in X} \max_{y \in Y} f(x, y)$$

where $X \subseteq E_x$ and $Y \subseteq E_y$ are convex sets in their spaces;
 $f(\cdot, y)$ is convex $\forall y \in Y$; $f(x, \cdot)$ is concave $\forall x \in X$.

⊙ Strong duality (a.k.a. Minimax Theorem; M. Sion 1958)

If f is convex-concave, and X or Y is compact, then $\varphi^* = \psi^*$
and there exists a saddle point — a point (x^*, y^*) for which



$$\varphi^* = \varphi(x^*) = f(x^*, y^*) = \psi(y^*) = \psi^*$$

$$f(x^*, y) \leq f(x^*, y^*) \leq f(x, y^*) \quad \forall x \in X, \forall y \in Y$$

Also works with non-compact sets, but strong convexity or concavity

Convex - Concave Saddle-Point Problems:

$$\min_{x \in X} \max_{y \in Y} f(x, y)$$

$X \subseteq E_x$ and $Y \subseteq E_y$ are convex sets in their spaces;
 $f(\cdot, y)$ is convex $\forall y \in Y$; $f(x, \cdot)$ is concave $\forall x \in X$.

Task: Find a saddle point (x^*, y^*) :

$$\varphi^* = \varphi(x^*) = f(x^*, y^*) = \psi(y^*) = \psi^*$$

Duality Gap: How good is a candidate solution $(\tilde{x}, \tilde{y}) \in X \times Y$:

$$\text{Gap}(\tilde{x}, \tilde{y}) := \varphi(\tilde{x}) - \psi(\tilde{y}) = \underbrace{\varphi(\tilde{x}) - \varphi^*}_{\text{Primal gap} \geq 0} + \underbrace{\psi^* - \psi(\tilde{y})}_{\text{Dual gap} \geq 0}$$

- ⊙ Thus, $\text{Gap}(\tilde{x}, \tilde{y}) \leq \varepsilon$ guarantees that \tilde{x} is ε -suboptimal for (P), and \tilde{y} is ε -suboptimal for (D).
- ⊙ Typically we might upper-bound or even compute it - certificate.

Black-Box model for CCSPP:

$$\min_{x \in X} \max_{y \in Y} f(x, y)$$

⊙ X and Y known (can project)

⊙ $f: X \times Y \rightarrow \mathbb{R}$ is available via 1st-order oracle:

$$(x, y) \mapsto \left[\underset{\uparrow}{f'_x(x, y)}; -\underset{\uparrow}{f'_y(x, y)} \right]$$

subgradient of $f(\cdot, y)$ at x supergradient of $f(x, \cdot)$ at y

Goal: In T queries, find $(\hat{x}_T, \hat{y}_T) \in X \times Y$ such that $\text{Gap}(\hat{x}_T, \hat{y}_T) \leq \epsilon_T$.

Certificate:

$$\begin{aligned} \text{Gap}(\hat{x}, \hat{y}) &= \rho(\hat{x}) - \psi(\hat{y}) \leq f(\hat{x}, \hat{y}) - \min_{x \in X} f(x, \hat{y}) + \max_{y \in Y} f(\hat{x}, y) - f(\hat{x}, \hat{y}) \\ &\leq \sup_{x \in X} \langle f'_x(\hat{x}, \hat{y}), \hat{x} - x \rangle + \sup_{y \in Y} \langle f'_y(\hat{x}, \hat{y}), \hat{y} - y \rangle \end{aligned}$$

The RHS can be bounded whenever we have linear maximization oracles (LMO) for X, Y , which is a very weak assumption.

CCSPP as a monotone variational inequality (MVI)

$$\boxed{\min_{x \in X} \max_{y \in Y} f(x, y)} \quad \begin{array}{l} z = (x, y) \\ Z = X \times Y \end{array} \quad \mapsto \quad \begin{array}{l} F(z) = [f'_x(x, y); -f'_y(x, y)] \\ \uparrow \\ \text{vector field} \end{array}$$

Optimality condition: $z^* = (x^*, y^*)$ is a SP iff $\langle F(z^*), z - z^* \rangle \geq 0 \quad \forall z \in Z$.

Variational Inequality: for a vector field $F(\cdot)$ on Z , find z^* such that

$$\boxed{\langle F(z^*), z - z^* \rangle \geq 0 \quad \forall z \in Z}$$

⊙ For any SPP (not necessarily a convex-concave one), a solution to the corresponding (strong) VI is a first-order Nash equilibrium ("SPP for the locally linearized objective")

⊙ For CCSPPs, the operator $F(\cdot)$ is monotone: one has that

$$\langle F(z') - F(z), z' - z \rangle \geq 0 \quad \forall z', z \in Z.$$

(Compare with the case of a subgradient of a convex function.)

CCSPP as a monotone variational inequality (MVI)

$$\boxed{\min_{x \in X} \max_{y \in Y} f(x, y)} \quad \begin{array}{l} z = (x, y) \\ \mathbb{Z} = X \times Y \end{array} \mapsto \begin{array}{l} F(z) = [f'_x(x, y); -f'_y(x, y)] \\ \uparrow \\ \text{vector field} \end{array}$$

Optimality condition: $z^* = (x^*, y^*)$ is a SP iff $\langle F(z^*), z - z^* \rangle \geq 0 \quad \forall z \in \mathbb{Z}$.

Variational Inequality: for a vector field $F(\cdot)$ on \mathbb{Z} , find z^* such that

$$\langle F(z^*), z - z^* \rangle \geq 0 \quad \forall z \in \mathbb{Z}$$

⊙ Recall that $\text{Gap}(z) = \varphi(x) - \psi(y) \leq \sup_{z \in \mathbb{Z}} \langle F(z), \hat{z} - z \rangle$.

⊙ ε -approximate VI solution; i.e. $\hat{z} \in \mathbb{Z}$ such that

$$\langle F(\hat{z}), \hat{z} - z \rangle \leq \varepsilon \quad \forall z \in \mathbb{Z}$$

is also an ε -approximate saddle point, i.e. $\text{Gap}(\hat{z}) \leq \varepsilon$.

⊙ Hence, we can focus on solving MVIs (approximately).

CCSPP with a smooth objective (MVI's with continuous operator)

- ⊙ The common feature in all examples is that a non-smooth (primal) minimization problem $\min_{x \in X} \{ \varphi(x) = \max_{y \in Y} f(x, y) \}$ translates to a smooth CCSPP, provided that we have a max-type representation for $\varphi(x)$.
- ⊙ Recall that non-smooth (Lipshitz) convex minimization can be done in $O(\frac{1}{\sqrt{T}})$ in the black-box model - oracle $x \mapsto \varphi'(x)$ - via SGM. But here we "go out of the black box" in terms of φ : new oracle is $(x, y) \mapsto [\nabla_x f(x, y); -\nabla_y f(x, y)]$. Hopefully, we can do faster?
- ⊙ For smooth convex minimization, projected gradient descent converges as $O(\frac{1}{\sqrt{T}})$, and Nesterov as $O(\frac{1}{T})$ which is optimal.
- ▶ Can we generalize some of these results to MVI's w/ continuous operators?
- ⊙ Yes, we can attain $\text{Gap}(\bar{x}_T, \bar{y}_T) = O(\frac{1}{T})$ - and this is optimal already for bilinear CCSPPs.

Gradient Descent-Ascent (GDA) & the Loopy Problem

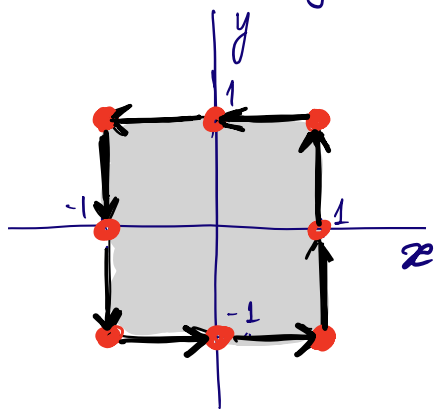
(Projected) GDA is a naive adaptation of (projected) gradient descent

Problem	Oracle	Method
$\min_{x \in X} g(x)$ C^1, convex	$x \mapsto \nabla g(x)$	$x_{t+1} = \Pi_X \left[x_t - \frac{1}{L} \nabla g(x_t) \right]$ (PGD)
$\min_{x \in X} \max_{y \in Y} f(x, y)$ f is C^1 and convex-concave	$z = (x, y) \rightarrow F(z)$ where $F(z) = \begin{bmatrix} \nabla_x f(x, y) \\ -\nabla_y f(x, y) \end{bmatrix}$	$z_{t+1} = \Pi_Z \left[z_t - \frac{1}{L} F(z_t) \right]$ $\Leftrightarrow z$ $\begin{bmatrix} x_{t+1} \\ y_{t+1} \end{bmatrix} = \begin{bmatrix} \Pi_X \left[x_t - \frac{1}{L} \nabla_x f(x, y) \right] \\ \Pi_Y \left[y_t + \frac{1}{L} \nabla_y f(x, y) \right] \end{bmatrix}$ (PGDA)

Unfortunately, vanilla PGDA can "go in circles" even for bilinear SPPs

$$\min_{x \in [-1, 1]} \max_{y \in [-1, 1]} xy$$

Unique SP: $(x^*, y^*) = (0, 0)$



$$F(x, y) = \begin{bmatrix} y \\ -x \end{bmatrix}$$

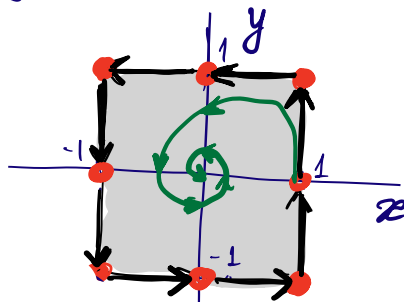
$$\begin{bmatrix} x_{t+1} \\ y_{t+1} \end{bmatrix} = \begin{bmatrix} \Pi_{[-1, 1]} [x_t - y_t] \\ \Pi_{[-1, 1]} [y_t + x_t] \end{bmatrix}$$

Ergodic Convergence of GDA

⊙ We can enforce PGDA convergence via overlapping & smaller steps.

$$\hat{z}_T = \frac{1}{T} \sum_{t=1}^T z_t$$

where $z_{t+1} = z_t - \gamma F(z_t)$



Theorem

Assume $\max\{\text{diam}(X), \text{diam}(Y)\} \leq D$,

and $\sup_{z \in Z} \|F(z)\| \leq B$.

Then PGDA with $\gamma = \frac{D}{B\sqrt{T}}$ ensures

$$\text{Gap}(\hat{z}_T) \leq \frac{BD}{\sqrt{T}}$$

⊙ More precisely, if $\text{diam}(X) \leq D_x$, $\text{diam}(Y) \leq D_y$, and F is (B_x, B_y) -Lipschitz, we can attain $\text{Gap}(\hat{z}_T) \leq \frac{B_x D_x + B_y D_y}{\sqrt{T}}$

⊙ In the C^1 case with $D_x f(z_x) = D_y f(z_y) = 0$ for some z_x, z_y :

$$\text{Gap}(\hat{z}_T) \leq \frac{L_{xy} D_x D_y + L_{xx} D_x^2 + L_{yy} D_y^2}{\sqrt{T}}$$

Note that $L_{xx} = L_{yy} = 0$ in the bilinear case.

GDA with a stochastic oracle

$$\left. \begin{aligned} \|F(x', y) - F(x, y)\| &\leq L_{xx} \|x' - x\|; \\ \|F(x', y) - F(x, y)\| &\leq L_{xy} \|y' - y\|; \\ \|F(x, y') - F(x, y)\| &\leq L_{yy} \|y' - y\|. \end{aligned} \right\} \text{for all } x', x \in X; y', y \in Y.$$

$$\text{Gap}(\bar{z}_T) \lesssim \frac{L_{xy} D_x D_y + L_{xx} D_x^2 + L_{yy} D_y^2}{\sqrt{T}}, \quad \text{with } L_{xx} = L_{yy} = 0 \text{ in the bilinear case}$$

⊙ Stochastic oracle: $\tilde{F}(z) = F(z) + \zeta(z)$ where $\zeta(z) = \frac{\zeta_x(z)}{\zeta_y(z)}$

- satisfies $\mathbb{E}[\zeta(z)] = 0$ and $\mathbb{E} \begin{bmatrix} \|\zeta_x(z)\|^2 \\ \|\zeta_y(z)\|^2 \end{bmatrix} \leq \begin{bmatrix} \sigma_x^2 \\ \sigma_y^2 \end{bmatrix}$ for all $z \in Z$

$$\text{Gap}(\bar{z}_T) \lesssim [\dots] + \frac{\sigma_x D_x + \sigma_y D_y}{\sqrt{T}}$$

where $[\dots]$ is the "deterministic" error (as if $\sigma_x = \sigma_y = 0$).

Extragradient method.

⊙ f is convex-concave and L -smooth (i.e. F is L -Lipschitz).

Then
$$\text{Gap}(\bar{z}_T) \leq \frac{LD^2}{T} + \frac{\sigma_x D_x + \sigma_y D_y}{\sqrt{T}}.$$

Key Ideas:

① Consider the following "conceptual method" (implicit update):

$$z_{t+1} = \Pi_Z \left[z_t - \frac{c}{L} F(z_{t+1}) \right] \quad \text{with some } c \geq 0 \text{ (to be chosen later).}$$

Instead of mimicking PGD, we mimic the proximal point method (PPM)

⊙ It turns out that the analysis of PPM generalizes directly, and this "conceptual method" converges as $\frac{LD^2}{T}$ when $c < 1$.

How to implement the conceptual method?

Extragradient method: basic version

[Understood since 1970s:]

Conceptual update $z_{t+1} = \Pi_Z [z_t - \frac{c}{L} F(z_{t+1})]$ is a fixed point of the operator $P_{z_t}(z) = \Pi_Z [z - \frac{c}{L} F(z)]$. (By definition!)

Lemma. Assume $F(\cdot)$ is L -Lipschitz, and $c < 1$. Then $P_{z_t}(\cdot)$ is a contraction.

Proof: $\|P_{z_t}(z') - P_{z_t}(z)\| = \|\Pi_Z [\bar{z} - \frac{c}{L} F(z')] - \Pi_Z [\bar{z} - \frac{c}{L} F(z)]\|$
[by the non-expansiveness of projections] $\leq \|\bar{z} - \frac{c}{L} F(z') - \bar{z} + \frac{c}{L} F(z)\| \leq c \|z' - z\|$

Hence, given $\tilde{z}_t \approx z_t$ we can approximate z_{t+1} with K iterations of the form:

$$\tilde{z}_t \mapsto z_t^{(1)} = P_{z_t}(z_t) \mapsto z_t^{(2)} = P_{z_t}(z_t^{(1)}) \mapsto \dots \tilde{z}_{t+1} := z_t^{(K)} = P_{z_t}(z_t^{(K-1)})$$

with a linear (i.e. exponentially fast) convergence towards z_{t+1} .

This is "K-lockhead extragradient" update, and it gives the desired convergence rate.

Extragradient method with $K=2$

⊙ Note that $K=1$ gives PGDA. A. Nemirovski discovered that $K=2$ suffices!

Theorem: (Nemirovski, 2003)

Let $c \leq \frac{1}{\sqrt{2}}$. Then the following method converges as $O(\frac{1}{T})$:

Input: z_t

$$z_{t+\frac{1}{2}} = \Pi_Z \left[z_t - \frac{c}{L} F(z_t) \right]$$

Return: $z_{t+1} = \Pi_Z \left[z_t - \frac{c}{L} F(z_{t+\frac{1}{2}}) \right]$

(Convergence is for z_T , i.e. in terms of the last iterate.)

The proof is non-trivial, and heavily utilizes monotonicity and the algebraic structure of Bregman divergences.

BONUS: Weak / Strong MVIDs, Certificates.

Monotone VIs, weak and strong solutions.

Strong solution: $\langle F(z^*), z - z^* \rangle \geq 0 \quad \forall z \in Z.$

Monotonicity: $\langle F(z) - F(z^*), z - z^* \rangle \geq 0$

$\Rightarrow z^*$ is also a weak solution: $\langle F(z), z - z^* \rangle \geq 0 \quad \forall z \in Z.$

⊙ Moreover, in a monotone VI with a continuous operator, any weak solution is also a strong solution (Minty's Lemma).

Proof: let $x_s = (1-s)x^* + sx$ for $s \in [0,1]$, where x^* is a weak solution.

then $0 \leq \langle F(x_s), x_s - x^* \rangle = s \langle F(x_s), x - x^* \rangle$, so $\langle F(x_s), x - x^* \rangle \geq 0$.
Take $s \rightarrow 0$. \square

⊙ Hence, for CCSPs with C_1 -smooth objectives, SP_s are weak VI solutions.

Certificates

① Certificate: let $\begin{bmatrix} \hat{x}_T \\ \hat{y}_T \end{bmatrix} = \sum_{t=1}^T w_t \begin{bmatrix} x_t \\ y_t \end{bmatrix}$ where x_1, \dots, x_T are the iterates, (w_1, \dots, w_T) are the convex weights.

$$\begin{aligned} \text{Gap}(\hat{x}_T, \hat{y}_T) &= \varphi(\hat{x}_T) - \psi(\hat{y}_T) \leq \sum_{t=1}^T w_t \left[\varphi(x_t) - f(x_t, y_t) + f(x_t, y_t) - \psi(y_t) \right] \\ &\leq \sum_{t=1}^T w_t \left[\langle f'_y(x_t, y_t), y_t^* - y_t \rangle + \langle f'_x(x_t, y_t), x_t - x_t^* \rangle \right] \end{aligned}$$

Can upper-bound the RHS by maximizing linear functions on X and Y .