

On Fundamental Proof Structures in First-Order Optimization

Baptiste Goujaud



62nd IEEE Conference on Decision and Control (CDC 2023) – Singapore – December, 14th 2023

Collaborators



Pr. Aymeric Dieuleveut



Pr. Adrien Taylor



Goal and Roadmap

Goal of today: Improving our in-depth understanding of how proofs for first order optimization work.

Goal and Roadmap

Goal of today: Improving our in-depth understanding of how proofs for first order optimization work.

Roadmap:

- ◇ From explicit to implicit characterization of function classes.
- ◇ From explicit to implicit characterization of algorithms.
- ◇ Insights from the dual formulation of the problem.
- ◇ Natural performance metrics.
- ◇ Extensions of a certificate to larger set of class and algorithms.
- ◇ Lyapunov approaches.

References

- B.G., D. Scieur, A. Dieuleveut, A. Taylor, F. Pedregosa (2022). Super-acceleration with cyclical step-sizes.
- B.G., A. Taylor, A. Dieuleveut (2022). Quadratic minimization: from conjugate gradient to an adaptive Heavy-ball method with Polyak step-sizes.
- Y. Drori, M. Teboulle (2014). Performance of first-order methods for smooth convex minimization: a novel approach.
- A. Taylor, J. Hendrickx, F. Glineur (2017). Smooth strongly convex interpolation and exact worst-case performance of first-order methods.
- L. Lessard, B. Recht, A. Packard (2014). Analysis and Design of Optimization Algorithms via Integral Quadratic Constraints.
- Y. Drori and A. Taylor (2020). Efficient first-order methods for convex minimization: a constructive approach.
- B.G., A. Taylor, A. Dieuleveut (2022). Optimal first-order methods for convex functions with a quadratic upper bound.
- A. Taylor, B. Van Scoy, L. Lessard (2018). Lyapunov functions for first-order methods: Tight automated convergence guarantees

References

- C. Park, E.K. Ryu (2021). Optimal First-Order Algorithms as a Function of Inequalities
- N. Bouselmi, J.M. Hendrickx, F. Glineur (2023). Interpolation Conditions for Linear Operators and applications to Performance Estimation Problems
- A. Taylor, J. Hendrickx, F. Glineur (2017). Performance Estimation Toolbox (PESTO): automated worst-case analysis of first-order optimization methods.
- B.G., C. Moucer, F. Glineur, J. Hendrickx, A. Taylor, A. Dieuleveut (2022). "PEPit: computer-assisted worst-case analyses of first-order optimization methods in Python."

First order optimization

Optimization: We aim to solve

$$f_{\star} \triangleq \min_{x \in \mathbb{R}^d} f(x)$$

(OPT)

First order optimization

Optimization: We aim to solve

$$f_* \triangleq \min_{x \in \mathbb{R}^d} f(x) \quad (\text{OPT})$$

First order oracle:

- ◇ Allowed: request function values and gradients of f at some points.

Examples:

- Gradient descent: $x_{t+1} = x_t - \gamma_t \nabla f(x_t)$.
- Heavy-ball: $x_{t+1} = x_t - \gamma_t \nabla f(x_t) + \beta_t (x_t - x_{t-1})$.

- ◇ Forbidden: request hessian and higher order derivatives.

Examples:

- Newton method: $x_{t+1} = x_t - \nabla^2 f(x_t)^{-1} \nabla f(x_t)$.

First order optimization

Optimization: We aim to solve

$$f_* \triangleq \min_{x \in \mathbb{R}^d} f(x) \quad (\text{OPT})$$

First order oracle:

- ◇ Allowed: request function values and gradients of f at some points.

Examples:

- Gradient descent: $x_{t+1} = x_t - \gamma_t \nabla f(x_t)$.
- Heavy-ball: $x_{t+1} = x_t - \gamma_t \nabla f(x_t) + \beta_t(x_t - x_{t-1})$.

- ◇ Forbidden: request hessian and higher order derivatives.

Examples:

- Newton method: $x_{t+1} = x_t - \nabla^2 f(x_t)^{-1} \nabla f(x_t)$.

Worst-case study:

- ◇ Fixed first order algorithm \mathcal{A} .
- ◇ Fixed class (set) of functions \mathcal{F} .
- ◇ Question: what *a priori* guarantees after N iterations?

Examples:

$$\forall f \in \mathcal{F}, \quad \|x_N - x_*\|^2 \leq \rho_N \|x_0 - x_*\|^2$$

$$\forall f \in \mathcal{F}, \quad f(x_N) - f_* \leq \rho_N \|x_0 - x_*\|^2$$

$$\forall f \in \mathcal{F}, \quad \|\nabla f(x_N)\|^2 \leq \rho_N (f(x_0) - f_*)$$

Class example: strongly convex quadratic functions

Quadratic functions: $\mathcal{F} = \left\{ f \mid f(x) \triangleq \frac{1}{2}(x - x_*)^T H(x - x_*) + f_* \mid H \text{ is PSD.} \right\}$

Class example: strongly convex quadratic functions

Quadratic functions: $\mathcal{F} = \left\{ f \mid f(x) \triangleq \frac{1}{2}(x - x_*)^T H(x - x_*) + f_* \mid H \text{ is PSD.} \right\}$

First order methods: $x_N = x_{N-1} - \sum_{i=0}^{N-1} h_{N,i} \nabla f(x_i)$

Class example: strongly convex quadratic functions

Quadratic functions: $\mathcal{F} = \left\{ f \mid f(x) \triangleq \frac{1}{2}(x - x_*)^T H(x - x_*) + f_* \mid H \text{ is PSD.} \right\}$

First order methods: $x_N = x_{N-1} - \sum_{i=0}^{N-1} h_{N,i} \nabla f(x_i)$



Link with polynomials: $x_N - x_* = P_N(H)(x_0 - x_*)$.

Class example: strongly convex quadratic functions

Quadratic functions: $\mathcal{F} = \left\{ f \mid f(x) \triangleq \frac{1}{2}(x - x_*)^T H(x - x_*) + f_* \mid H \text{ is PSD.} \right\}$

First order methods: $x_N = x_{N-1} - \sum_{i=0}^{N-1} h_{N,i} \nabla f(x_i)$



Link with polynomials: $x_N - x_* = P_N(H)(x_0 - x_*)$.

$Sp(H) \subset [\mu, L]$

\hookrightarrow Momentum.

Class example: strongly convex quadratic functions

Quadratic functions: $\mathcal{F} = \left\{ f \mid f(x) \triangleq \frac{1}{2}(x - x_*)^T H(x - x_*) + f_* \mid H \text{ is PSD.} \right\}$

First order methods: $x_N = x_{N-1} - \sum_{i=0}^{N-1} h_{N,i} \nabla f(x_i)$



Link with polynomials: $x_N - x_* = P_N(H)(x_0 - x_*)$.

$$Sp(H) \subset [\mu, L]$$

\hookrightarrow Momentum.

$$Sp(H) \subset [\mu_1, L_1] \cup [\mu_2, L_2]$$

\hookrightarrow Cycling step-size.

Class example: strongly convex quadratic functions

Quadratic functions: $\mathcal{F} = \left\{ f \mid f(x) \triangleq \frac{1}{2}(x - x_*)^T H(x - x_*) + f_* \mid H \text{ is PSD.} \right\}$

First order methods: $x_N = x_{N-1} - \sum_{i=0}^{N-1} h_{N,i} \nabla f(x_i)$



Link with polynomials: $x_N - x_* = P_N(H)(x_0 - x_*)$.

$Sp(H) \subset [\mu, L]$
 \hookrightarrow Momentum.

$Sp(H) \subset [\mu_1, L_1] \cup [\mu_2, L_2]$
 \hookrightarrow Cycling step-size.

$Sp(H)$
 \hookrightarrow Conjugate gradient.

Implicit definition of classes

Explicit definition of a class: parametrization.

Example:

$$\diamond \text{ Quadratic functions } \mathcal{F} = \left\{ f \mid f(x) \triangleq \frac{1}{2}(x - x_*)^T H(x - x_*) + f_* \mid H \text{ is PSD.} \right\}$$

Implicit definition of classes

Explicit definition of a class: parametrization.

Example:

$$\diamond \text{ Quadratic functions } \mathcal{F} = \left\{ f \mid f(x) \triangleq \frac{1}{2}(x - x_*)^T H(x - x_*) + f_* \mid H \text{ is PSD.} \right\}$$

Implicit definition of a class: constraints.

Examples:

$$\diamond \text{ Convex functions } \mathcal{F} = \{f \mid \forall x, y \ f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle\}.$$

$$\diamond \text{ Smooth functions } \mathcal{F} = \left\{ f \mid \forall x, y \ f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{L}{2} \|x - y\|^2 \right\}.$$

How to study them? PEP!

Implicit definition of classes

Explicit definition of a class: parametrization.

Example:

$$\diamond \text{ Quadratic functions } \mathcal{F} = \left\{ f \mid f(x) \triangleq \frac{1}{2}(x - x_*)^T H(x - x_*) + f_* \mid H \text{ is PSD.} \right\}$$

Implicit definition of a class: constraints.

Examples:

$$\diamond \text{ Convex functions } \mathcal{F} = \{f \mid \forall x, y \ f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle\}.$$

$$\diamond \text{ Smooth functions } \mathcal{F} = \left\{ f \mid \forall x, y \ f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{L}{2} \|x - y\|^2 \right\}.$$

How to study them? PEP!

$$\tau = \max_{f, x_0, x_1, x_*} \frac{\text{Perf}(x_1)}{\text{Init}(x_0)}$$

$$\text{s.t. } f \in \mathcal{F}$$

Functional class

Implicit definition of classes

Explicit definition of a class: parametrization.

Example:

$$\diamond \text{ Quadratic functions } \mathcal{F} = \left\{ f \mid f(x) \triangleq \frac{1}{2}(x - x_*)^T H(x - x_*) + f_* \mid H \text{ is PSD.} \right\}$$

Implicit definition of a class: constraints.

Examples:

$$\diamond \text{ Convex functions } \mathcal{F} = \{ f \mid \forall x, y \ f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle \}.$$

$$\diamond \text{ Smooth functions } \mathcal{F} = \left\{ f \mid \forall x, y \ f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{L}{2} \|x - y\|^2 \right\}.$$

How to study them? PEP!

$$\tau = \max_{f, x_0, x_1, x_*} \frac{\text{Perf}(x_1)}{\text{Init}(x_0)}$$

$$\text{s.t. } f \in \mathcal{F}$$

$$x_1 = \mathcal{A}(x_0, \nabla f(x_0))$$

$$\nabla f(x_*) = 0$$

Functional class

Algorithm

Optimality of x_*

Derivation of a PEP for convergence rate of a gradient step

Derivation of a PEP for convergence rate of a gradient step

Toy example: What is the smallest τ such that:

$$\|x_1 - x_*\|^2 \leq \tau \|x_0 - x_*\|^2,$$

for all

- ◇ L -smooth and μ -strongly convex function f (notation $f \in \mathcal{F}_{\mu,L}$),
- ◇ x_0 , and x_1 generated by gradient step $x_1 = x_0 - \gamma \nabla f(x_0)$,
- ◇ $x_* = \underset{x}{\operatorname{argmin}} f(x)$?

Derivation of a PEP for convergence rate of a gradient step

Toy example: What is the smallest τ such that:

$$\|x_1 - x_*\|^2 \leq \tau \|x_0 - x_*\|^2,$$

for all

- ◇ L -smooth and μ -strongly convex function f (notation $f \in \mathcal{F}_{\mu,L}$),
- ◇ x_0 , and x_1 generated by gradient step $x_1 = x_0 - \gamma \nabla f(x_0)$,
- ◇ $x_* = \underset{x}{\operatorname{argmin}} f(x)$?

Final goal: optimize τ (as a function of the step-size, γ)

Derivation of a PEP for convergence rate of a gradient step

Toy example: What is the smallest τ such that:

$$\|x_1 - x_*\|^2 \leq \tau \|x_0 - x_*\|^2,$$

for all

- ◇ L -smooth and μ -strongly convex function f (notation $f \in \mathcal{F}_{\mu,L}$),
- ◇ x_0 , and x_1 generated by gradient step $x_1 = x_0 - \gamma \nabla f(x_0)$,
- ◇ $x_* = \underset{x}{\operatorname{argmin}} f(x)$?

Final goal: optimize τ (as a function of the step-size, γ)

First: let's compute τ !

Derivation of a PEP for convergence rate of a gradient step

Toy example: What is the smallest τ such that:

$$\|x_1 - x_\star\|^2 \leq \tau \|x_0 - x_\star\|^2,$$

for all

- ◇ L -smooth and μ -strongly convex function f (notation $f \in \mathcal{F}_{\mu,L}$),
- ◇ x_0 , and x_1 generated by gradient step $x_1 = x_0 - \gamma \nabla f(x_0)$,
- ◇ $x_\star = \underset{x}{\operatorname{argmin}} f(x)$?

Final goal: optimize τ (as a function of the step-size, γ)

First: let's compute τ !

$$\tau = \max_{f, x_0, x_1, x_\star} \frac{\|x_1 - x_\star\|^2}{\|x_0 - x_\star\|^2}$$

$$\text{s.t. } f \in \mathcal{F}_{\mu,L}$$

$$x_1 = x_0 - \gamma \nabla f(x_0)$$

$$\nabla f(x_\star) = 0$$

Functional class

Algorithm

Optimality of x_\star

From infinite to finite dimensional problem

From infinite to finite dimensional problem

Performance estimation problem:

$$\begin{aligned} & \max_{f, x_0, x_1, x_*} \frac{\|x_1 - x_*\|^2}{\|x_0 - x_*\|^2} \\ \text{subject to} & \quad f \text{ is } L\text{-smooth and } \mu\text{-strongly convex,} \\ & \quad x_1 = x_0 - \gamma \nabla f(x_0) \\ & \quad \nabla f(x_*) = 0. \end{aligned}$$

From infinite to finite dimensional problem

Performance estimation problem:

$$\max_{f, x_0, x_1, x_*} \frac{\|x_1 - x_*\|^2}{\|x_0 - x_*\|^2}$$

subject to f is L -smooth and μ -strongly convex,

$$x_1 = x_0 - \gamma \nabla f(x_0)$$

$$\nabla f(x_*) = 0.$$

\Leftrightarrow Variables: f, x_0, x_1, x_* .

From infinite to finite dimensional problem

Performance estimation problem:

$$\max_{f, x_0, x_1, x_*} \frac{\|x_1 - x_*\|^2}{\|x_0 - x_*\|^2}$$

subject to f is L -smooth and μ -strongly convex,

$$x_1 = x_0 - \gamma \nabla f(x_0)$$

$$\nabla f(x_*) = 0.$$

\Leftrightarrow Variables: f, x_0, x_1, x_* .

Sampled version: f is only used at x_0 and x_* .

From infinite to finite dimensional problem

Performance estimation problem:

$$\begin{aligned} \max_{f, x_0, x_1, x_*} \quad & \frac{\|x_1 - x_*\|^2}{\|x_0 - x_*\|^2} \\ \text{subject to} \quad & f \text{ is } L\text{-smooth and } \mu\text{-strongly convex,} \\ & x_1 = x_0 - \gamma \nabla f(x_0) \\ & \nabla f(x_*) = 0. \end{aligned}$$

\hookrightarrow Variables: f, x_0, x_1, x_* .

Sampled version: f is only used at x_0 and x_* .

$$\begin{aligned} \max_{\substack{x_0, x_* \\ g_0, g_* \\ f_0, f_*}} \quad & \frac{\|x_0 - \gamma g_0 - x_*\|^2}{\|x_0 - x_*\|^2} \\ \text{subject to} \quad & f \in \mathcal{F}_{\mu, L} \text{ such that } \begin{cases} f_i = f(x_i) & i = 0, * \\ g_i = \nabla f(x_i) & i = 0, * \end{cases} \\ & g_* = 0. \end{aligned}$$

From infinite to finite dimensional problem

Performance estimation problem:

$$\begin{aligned} & \max_{f, x_0, x_1, x_*} \frac{\|x_1 - x_*\|^2}{\|x_0 - x_*\|^2} \\ \text{subject to} \quad & f \text{ is } L\text{-smooth and } \mu\text{-strongly convex,} \\ & x_1 = x_0 - \gamma \nabla f(x_0) \\ & \nabla f(x_*) = 0. \end{aligned}$$

\hookrightarrow Variables: f, x_0, x_1, x_* .

Sampled version: f is only used at x_0 and x_* .

$$\begin{aligned} & \max_{\substack{x_0, x_* \\ g_0, g_* \\ f_0, f_*}} \frac{\|x_0 - \gamma g_0 - x_*\|^2}{\|x_0 - x_*\|^2} \\ \text{subject to} \quad & f \in \mathcal{F}_{\mu, L} \text{ such that } \begin{cases} f_i = f(x_i) & i = 0, * \\ g_i = \nabla f(x_i) & i = 0, * \end{cases} \\ & g_* = 0. \end{aligned}$$

\hookrightarrow Variables: $x_0, x_*, g_0, g_*, f_0, f_*$.

From infinite to finite dimensional problem

Performance estimation problem:

$$\max_{f, x_0, x_1, x_*} \frac{\|x_1 - x_*\|^2}{\|x_0 - x_*\|^2}$$

subject to f is L -smooth and μ -strongly convex,

$$x_1 = x_0 - \gamma \nabla f(x_0)$$

$$\nabla f(x_*) = 0.$$

\hookrightarrow Variables: f, x_0, x_1, x_* .

Sampled version: f is only used at x_0 and x_* .

$$\max_{\substack{x_0, x_* \\ g_0, g_* \\ f_0, f_*}} \frac{\|x_0 - \gamma g_0 - x_*\|^2}{\|x_0 - x_*\|^2}$$

subject to $\exists f \in \mathcal{F}_{\mu, L}$ such that $\begin{cases} f_i = f(x_i) & i = 0, * \\ g_i = \nabla f(x_i) & i = 0, * \end{cases}$

$$g_* = 0.$$

\hookrightarrow Variables: $x_0, x_*, g_0, g_*, f_0, f_*$.

Interpolation conditions:

$\exists f \in \mathcal{F}_{\mu, L}$ such that $\forall i, f(x_i) = f_i$, and $\nabla f(x_i) = g_i$ iff $\forall i, j$,

$$f_i \geq f_j + \langle g_j, x_i - x_j \rangle + \frac{1}{2L} \|g_i - g_j\|^2 \\ + \frac{\mu}{2(1-\mu/L)} \|x_i - x_j - \frac{1}{L}(g_i - g_j)\|^2.$$

From infinite to finite dimensional problem

Performance estimation problem:

$$\max_{f, x_0, x_1, x_*} \frac{\|x_1 - x_*\|^2}{\|x_0 - x_*\|^2}$$

subject to f is L -smooth and μ -strongly convex,

$$x_1 = x_0 - \gamma \nabla f(x_0)$$

$$\nabla f(x_*) = 0.$$

\hookrightarrow Variables: f, x_0, x_1, x_* .

Sampled version: f is only used at x_0 and x_* .

$$\max_{\substack{x_0, x_* \\ g_0, g_* \\ f_0, f_*}} \frac{\|x_0 - \gamma g_0 - x_*\|^2}{\|x_0 - x_*\|^2}$$

subject to $\exists f \in \mathcal{F}_{\mu, L}$ such that $\begin{cases} f_i = f(x_i) & i = 0, * \\ g_i = \nabla f(x_i) & i = 0, * \end{cases}$

$$g_* = 0.$$

\hookrightarrow Variables: $x_0, x_*, g_0, g_*, f_0, f_*$.

Interpolation conditions:

$\exists f \in \mathcal{F}_{\mu, L}$ such that $\forall i, f(x_i) = f_i$, and $\nabla f(x_i) = g_i$ iff $\forall i, j$,

$$f_i \geq f_j + \langle g_j, x_i - x_j \rangle + \frac{1}{2L} \|g_i - g_j\|^2 + \frac{\mu}{2(1-\mu/L)} \|x_i - x_j - \frac{1}{L}(g_i - g_j)\|^2.$$

Reformulation:

$$\max_{\substack{x_0, x_* \\ g_0, g_* \\ f_0, f_*}} \frac{\|x_0 - \gamma g_0 - x_*\|^2}{\|x_0 - x_*\|^2}$$

subject to $f_* \geq f_0 + \langle g_0, x_* - x_0 \rangle + \frac{1}{2L} \|g_* - g_0\|^2 + \frac{\mu}{2(1-\mu/L)} \|x_* - x_0 - \frac{1}{L}(g_* - g_0)\|^2$
 $f_0 \geq f_* + \langle g_*, x_0 - x_* \rangle + \frac{1}{2L} \|g_0 - g_*\|^2 + \frac{\mu}{2(1-\mu/L)} \|x_0 - x_* - \frac{1}{L}(g_0 - g_*)\|^2$
 $g_* = 0.$

Normalization invariance

Non convex objective:

$$\max_{\substack{x_0, x_* \\ g_0, g_* \\ f_0, f_*}} \frac{\|x_0 - \gamma g_0 - x_*\|^2}{\|x_0 - x_*\|^2}$$

subject to

$$f_* \geq f_0 + \langle g_0, x_* - x_0 \rangle + \frac{1}{2L} \|g_* - g_0\|^2$$
$$+ \frac{\mu}{2(1-\mu/L)} \left\| x_* - x_0 - \frac{1}{L} (g_* - g_0) \right\|^2$$
$$f_0 \geq f_* + \langle g_*, x_0 - x_* \rangle + \frac{1}{2L} \|g_0 - g_*\|^2$$
$$+ \frac{\mu}{2(1-\mu/L)} \left\| x_0 - x_* - \frac{1}{L} (g_0 - g_*) \right\|^2$$
$$g_* = 0.$$

Normalization invariance

Non convex objective:

$$\max_{\substack{x_0, x_* \\ g_0, g_* \\ f_0, f_*}} \frac{\|x_0 - \gamma g_0 - x_*\|^2}{\|x_0 - x_*\|^2}$$

subject to

$$f_* \geq f_0 + \langle g_0, x_* - x_0 \rangle + \frac{1}{2L} \|g_* - g_0\|^2$$
$$+ \frac{\mu}{2(1-\mu/L)} \|x_* - x_0 - \frac{1}{L}(g_* - g_0)\|^2$$
$$f_0 \geq f_* + \langle g_*, x_0 - x_* \rangle + \frac{1}{2L} \|g_0 - g_*\|^2$$
$$+ \frac{\mu}{2(1-\mu/L)} \|x_0 - x_* - \frac{1}{L}(g_0 - g_*)\|^2$$
$$g_* = 0.$$

Homogeneity

$$\max_{\substack{x_0, x_* \\ g_0, g_* \\ f_0, f_*}} \|x_0 - \gamma g_0 - x_*\|^2$$

subject to

$$f_* \geq f_0 + \langle g_0, x_* - x_0 \rangle + \frac{1}{2L} \|g_* - g_0\|^2$$
$$+ \frac{\mu}{2(1-\mu/L)} \|x_* - x_0 - \frac{1}{L}(g_* - g_0)\|^2$$
$$f_0 \geq f_* + \langle g_*, x_0 - x_* \rangle + \frac{1}{2L} \|g_0 - g_*\|^2$$
$$+ \frac{\mu}{2(1-\mu/L)} \|x_0 - x_* - \frac{1}{L}(g_0 - g_*)\|^2$$
$$g_* = 0$$
$$\|x_0 - x_*\|^2 \leq 1.$$

Semidefinite lifting

Quadratic constraints

$$\max_{\substack{x_0, x_* \\ g_0, g_* \\ f_0, f_*}} \|x_0 - \gamma g_0 - x_*\|^2$$

$$\begin{aligned} \text{subject to } f_* &\geq f_0 + \langle g_0, x_* - x_0 \rangle + \frac{1}{2L} \|g_* - g_0\|^2 \\ &\quad + \frac{\mu}{2(1-\mu/L)} \|x_* - x_0 - \frac{1}{L}(g_* - g_0)\|^2 \\ f_0 &\geq f_* + \langle g_*, x_0 - x_* \rangle + \frac{1}{2L} \|g_0 - g_*\|^2 \\ &\quad + \frac{\mu}{2(1-\mu/L)} \|x_0 - x_* - \frac{1}{L}(g_0 - g_*)\|^2 \\ g_* &= 0 \\ \|x_0 - x_*\|^2 &\leq 1. \end{aligned}$$

Semidefinite lifting

Quadratic constraints

$$\max_{\substack{x_0, x_* \\ g_0, g_* \\ f_0, f_*}} \|x_0 - \gamma g_0 - x_*\|^2$$

$$\begin{aligned} \text{subject to} \quad & f_* \geq f_0 + \langle g_0, x_* - x_0 \rangle + \frac{1}{2L} \|g_* - g_0\|^2 \\ & \quad + \frac{\mu}{2(1-\mu/L)} \|x_* - x_0 - \frac{1}{L}(g_* - g_0)\|^2 \\ & f_0 \geq f_* + \langle g_*, x_0 - x_* \rangle + \frac{1}{2L} \|g_0 - g_*\|^2 \\ & \quad + \frac{\mu}{2(1-\mu/L)} \|x_0 - x_* - \frac{1}{L}(g_0 - g_*)\|^2 \\ & g_* = 0 \\ & \|x_0 - x_*\|^2 \leq 1. \end{aligned}$$

Homogeneity and Gram matrix

- ◇ Dependency on x and g in the objective and constraints all are quadratics.
↪ Introduce the new variables $G \succeq 0$ and F

$$G = \begin{bmatrix} \|x_0 - x_*\|^2 & \langle g_0, x_0 - x_* \rangle \\ \langle g_0, x_0 - x_* \rangle & \|g_0\|^2 \end{bmatrix}, \quad F = f_0 - f_*,$$

Semidefinite lifting

Quadratic constraints

$$\max_{\substack{x_0, x_* \\ g_0, g_* \\ f_0, f_*}} \|x_0 - \gamma g_0 - x_*\|^2$$

$$\begin{aligned} \text{subject to } f_* &\geq f_0 + \langle g_0, x_* - x_0 \rangle + \frac{1}{2L} \|g_* - g_0\|^2 \\ &\quad + \frac{\mu}{2(1-\mu/L)} \|x_* - x_0 - \frac{1}{L}(g_* - g_0)\|^2 \\ f_0 &\geq f_* + \langle g_*, x_0 - x_* \rangle + \frac{1}{2L} \|g_0 - g_*\|^2 \\ &\quad + \frac{\mu}{2(1-\mu/L)} \|x_0 - x_* - \frac{1}{L}(g_0 - g_*)\|^2 \\ g_* &= 0 \\ \|x_0 - x_*\|^2 &\leq 1. \end{aligned}$$

Homogeneity and Gram matrix

- ◇ Dependency on x and g in the objective and constraints all are quadratics.
↪ Introduce the new variables $G \succeq 0$ and F

$$G = \begin{bmatrix} \|x_0 - x_*\|^2 & \langle g_0, x_0 - x_* \rangle \\ \langle g_0, x_0 - x_* \rangle & \|g_0\|^2 \end{bmatrix}, \quad F = f_0 - f_*,$$

Semi-definite program (SDP) formulation

$$\begin{aligned} \max_{G, F} \quad & G_{1,1} + \gamma^2 G_{2,2} - 2\gamma G_{1,2} \\ \text{subject to } \quad & F + \frac{L\mu}{2(L-\mu)} G_{1,1} + \frac{1}{2(L-\mu)} G_{2,2} - \frac{L}{L-\mu} G_{1,2} \leq 0 \\ & -F + \frac{L\mu}{2(L-\mu)} G_{1,1} + \frac{1}{2(L-\mu)} G_{2,2} - \frac{\mu}{L-\mu} G_{1,2} \leq 0 \\ & G_{1,1} \leq 1 \\ & G \succeq 0, \end{aligned}$$

Semidefinite lifting

Quadratic constraints

$$\max_{\substack{x_0, x_* \\ g_0, g_* \\ f_0, f_*}} \|x_0 - \gamma g_0 - x_*\|^2$$

$$\begin{aligned} \text{subject to} \quad & f_* \geq f_0 + \langle g_0, x_* - x_0 \rangle + \frac{1}{2L} \|g_* - g_0\|^2 \\ & \quad + \frac{\mu}{2(1-\mu/L)} \|x_* - x_0 - \frac{1}{L}(g_* - g_0)\|^2 \\ & f_0 \geq f_* + \langle g_*, x_0 - x_* \rangle + \frac{1}{2L} \|g_0 - g_*\|^2 \\ & \quad + \frac{\mu}{2(1-\mu/L)} \|x_0 - x_* - \frac{1}{L}(g_0 - g_*)\|^2 \\ & g_* = 0 \\ & \|x_0 - x_*\|^2 \leq 1. \end{aligned}$$

Homogeneity and Gram matrix

- ◇ Dependency on x and g in the objective and constraints all are quadratics.
↪ Introduce the new variables $G \succcurlyeq 0$ and F

$$G = \begin{bmatrix} \|x_0 - x_*\|^2 & \langle g_0, x_0 - x_* \rangle \\ \langle g_0, x_0 - x_* \rangle & \|g_0\|^2 \end{bmatrix}, \quad F = f_0 - f_*,$$

Semi-definite program (SDP) formulation

$$\begin{aligned} \max_{G, F} \quad & G_{1,1} + \gamma^2 G_{2,2} - 2\gamma G_{1,2} \\ \text{subject to} \quad & F + \frac{L\mu}{2(L-\mu)} G_{1,1} + \frac{1}{2(L-\mu)} G_{2,2} - \frac{L}{L-\mu} G_{1,2} \leq 0 \\ & -F + \frac{L\mu}{2(L-\mu)} G_{1,1} + \frac{1}{2(L-\mu)} G_{2,2} - \frac{\mu}{L-\mu} G_{1,2} \leq 0 \\ & G_{1,1} \leq 1 \\ & G \succcurlyeq 0, \end{aligned}$$

Convex problem!!!

Other usages from PEPit

- 1. Unconstrained convex minimization
 - 1.1. Gradient descent
 - 1.2. Subgradient method
 - 1.3. Subgradient method under restricted secant inequality and error bound
 - 1.4. Gradient descent with exact line search
 - 1.5. Conjugate gradient
 - 1.6. Heavy Ball momentum
 - 1.7. Accelerated gradient for convex objective
 - 1.8. Accelerated gradient for strongly convex objective
 - 1.9. Optimized gradient
 - 1.10. Optimized gradient for gradient
 - 1.11. Robust momentum
 - 1.12. Triple momentum
 - 1.13. Information theoretic exact method
 - 1.14. Proximal point
 - 1.15. Accelerated proximal point
 - 1.16. Inexact gradient descent
 - 1.17. Inexact gradient descent with exact line search
 - 1.18. Inexact accelerated gradient
 - 1.19. Epsilon-subgradient method
 - 1.20. Gradient descent for quadratically upper bounded convex objective
 - 1.21. Gradient descent with decreasing step sizes for quadratically upper bounded convex objective
 - 1.22. Conjugate gradient for quadratically upper bounded convex objective
 - 1.23. Heavy Ball momentum for quadratically upper bounded convex objective
- 2. Composite convex minimization
 - 2.1. Proximal gradient
 - 2.2. Accelerated proximal gradient
 - 2.3. Bregman proximal point
 - 2.4. Douglas Rachford splitting
 - 2.5. Douglas Rachford splitting contraction
 - 2.6. Accelerated Douglas Rachford splitting
 - 2.7. Frank Wolfe
 - 2.8. Improved interior method
 - 2.9. No Lips in function value
 - 2.10. No Lips in Bregman divergence
 - 2.11. Three operator splitting
- 3. Non-convex optimization
 - 3.1. Gradient Descent
 - 3.2. No Lips 1
 - 3.3. No Lips 2
- 4. Stochastic and randomized convex minimization
 - 4.1. Stochastic gradient descent
 - 4.2. Stochastic gradient descent in overparametrized setting
 - 4.3. SAGA
 - 4.4. Point SAGA
 - 4.5. Randomized coordinate descent for smooth strongly convex functions
 - 4.6. Randomized coordinate descent for smooth convex functions
- 5. Monotone inclusions and variational inequalities
 - 5.1. Proximal point
 - 5.2. Accelerated proximal point
 - 5.3. Optimal Strongly-monotone Proximal Point
 - 5.4. Douglas Rachford Splitting
 - 5.5. Three operator splitting
 - 5.6. Optimistic gradient
 - 5.7. Past extragradient
- 6. Fixed point
 - 6.1. Halpern iteration
 - 6.2. Optimal Contractive Halpern iteration
 - 6.3. Krasnoselskii-Mann with constant step-sizes
 - 6.4. Krasnoselskii-Mann with increasing step-sizes
- 7. Potential functions
 - 7.1. Gradient descent Lyapunov 1
 - 7.2. Gradient descent Lyapunov 2
 - 7.3. Accelerated gradient method
- 8. Inexact proximal methods
 - 8.1. Accelerated inexact forward backward
 - 8.2. Partially inexact Douglas Rachford splitting
 - 8.3. Relatively inexact proximal point
- 9. Adaptive methods
 - 9.1. Polyak steps in distance to optimum
 - 9.2. Polyak steps in function value
- 10. Low dimensional worst-cases scenarios
 - 10.1. Inexact gradient
 - 10.2. Non-convex gradient descent
 - 10.3. Optimized gradient
 - 10.4. Frank Wolfe
 - 10.5. Proximal point
 - 10.6. Halpern iteration
 - 10.7. Alternate projections
 - 10.8. Averaged projections
 - 10.9. Dykstra
- 11. Continuous-time models
 - 11.1. Gradient flow for strongly convex functions
 - 11.2. Gradient flow for convex functions
 - 11.3. Accelerated gradient flow for strongly convex functions
 - 11.4. Accelerated gradient flow for convex functions
- 12. Tutorials
 - 12.1. Contraction rate of gradient descent

Key ingredients

- ◇ Algorithm: first order update;
- ◇ Class of functions: interpolation constraints expressible linearly in F and G ;
- ◇ Performance metrics: expressible linearly in F and G .

Avoiding semidefinite programming modeling steps?

Performance Estimation Toolbox (**PESTO**) available on

PerformanceEstimation/Performance-Estimation-Toolbox



Python version (**PEPit**) available on PyPI.

PerformanceEstimation/PEPit

- ◇ Contain more than 50 examples.
- ◇ Contain several classical classes of functions.



Learning to PEP.

PerformanceEstimation/Learning-Performance-Estimation

- ◇ Contains theoretical exercises as well as code exercises.

From explicit to implicit algorithms

- ◇ So far algorithms were defined through an explicit update rule.
- ◇ Some update could be characterized by an inequality constraint.
- ◇ From now, we consider the general form:

$$\left| \begin{array}{l} \text{maximize} \\ F, G \succeq 0 \end{array} \right. \quad \langle F, v_P \rangle + \langle G, M_P \rangle$$
$$\left| \begin{array}{l} \text{subject to} \end{array} \right. \left\{ \begin{array}{l} \langle F, v_I \rangle + \langle G, M_I \rangle \leq R^2 \\ \forall k, \langle F, v_{\mathcal{F}}^{(k)} \rangle + \langle G, M_{\mathcal{F}}^{(k)} \rangle \leq 0 \\ \forall l, \langle F, v_{\mathcal{A}}^{(l)} \rangle + \langle G, M_{\mathcal{A}}^{(l)} \rangle \leq 0 \end{array} \right. \quad \text{(PEP-primal)}$$

From primal to dual problem

◇ Primal problem is

$$\left\{ \begin{array}{l} \text{maximize} \\ F, G \succeq 0 \end{array} \right. \langle F, v_P \rangle + \langle G, M_P \rangle$$
$$\left\{ \begin{array}{l} \text{subject to} \\ \forall k, \langle F, v_{\mathcal{F}}^{(k)} \rangle + \langle G, M_{\mathcal{F}}^{(k)} \rangle \leq 0 \quad : \lambda_{\mathcal{F}}^{(k)} \\ \forall l, \langle F, v_{\mathcal{A}}^{(l)} \rangle + \langle G, M_{\mathcal{A}}^{(l)} \rangle \leq 0 \quad : \lambda_{\mathcal{A}}^{(l)} \end{array} \right. \leq R^2 : \tau$$

From primal to dual problem

◇ Primal problem is

$$\begin{array}{l} \text{maximize} \\ F, G \succeq 0 \end{array} \left\{ \begin{array}{l} \langle F, v_P \rangle + \langle G, M_P \rangle \\ \text{subject to} \left\{ \begin{array}{l} \langle F, v_I \rangle + \langle G, M_I \rangle \leq R^2 : \tau \\ \forall k, \langle F, v_{\mathcal{F}}^{(k)} \rangle + \langle G, M_{\mathcal{F}}^{(k)} \rangle \leq 0 : \lambda_{\mathcal{F}}^{(k)} \\ \forall l, \langle F, v_{\mathcal{A}}^{(l)} \rangle + \langle G, M_{\mathcal{A}}^{(l)} \rangle \leq 0 : \lambda_{\mathcal{A}}^{(l)} \end{array} \right. \end{array} \right.$$

◇ Lagrangian is

$$\begin{aligned} \mathcal{L} \triangleq & \langle F, v_P \rangle + \langle G, M_P \rangle - \tau [\langle F, v_I \rangle + \langle G, M_I \rangle - R^2] \\ & - \sum_k \lambda_{\mathcal{F}}^{(k)} [\langle F, v_{\mathcal{F}}^{(k)} \rangle + \langle G, M_{\mathcal{F}}^{(k)} \rangle] \\ & - \sum_l \lambda_{\mathcal{A}}^{(l)} [\langle F, v_{\mathcal{A}}^{(l)} \rangle + \langle G, M_{\mathcal{A}}^{(l)} \rangle] \end{aligned}$$

From primal to dual problem

◇ Primal problem is

$$\begin{array}{l} \text{maximize} \\ F, G \succeq 0 \end{array} \left\{ \begin{array}{l} \langle F, v_P \rangle + \langle G, M_P \rangle \\ \text{subject to} \end{array} \right. \left\{ \begin{array}{l} \langle F, v_I \rangle + \langle G, M_I \rangle \leq R^2 : \tau \\ \forall k, \langle F, v_{\mathcal{F}}^{(k)} \rangle + \langle G, M_{\mathcal{F}}^{(k)} \rangle \leq 0 : \lambda_{\mathcal{F}}^{(k)} \\ \forall l, \langle F, v_{\mathcal{A}}^{(l)} \rangle + \langle G, M_{\mathcal{A}}^{(l)} \rangle \leq 0 : \lambda_{\mathcal{A}}^{(l)} \end{array} \right.$$

◇ Lagrangian is

$$\begin{aligned} \mathcal{L} &\triangleq \langle F, v_P \rangle + \langle G, M_P \rangle - \tau [\langle F, v_I \rangle + \langle G, M_I \rangle - R^2] \\ &\quad - \sum_k \lambda_{\mathcal{F}}^{(k)} [\langle F, v_{\mathcal{F}}^{(k)} \rangle + \langle G, M_{\mathcal{F}}^{(k)} \rangle] \\ &\quad - \sum_l \lambda_{\mathcal{A}}^{(l)} [\langle F, v_{\mathcal{A}}^{(l)} \rangle + \langle G, M_{\mathcal{A}}^{(l)} \rangle] \\ &= \tau R^2 + \left\langle F, v_P - \tau v_I - \sum_k \lambda_{\mathcal{F}}^{(k)} v_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} v_{\mathcal{A}}^{(l)} \right\rangle \\ &\quad + \left\langle G, M_P - \tau M_I - \sum_k \lambda_{\mathcal{F}}^{(k)} M_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} M_{\mathcal{A}}^{(l)} \right\rangle \end{aligned}$$

From primal to dual problem

◇ Primal problem is

$$\begin{array}{l} \text{maximize} \\ F, G \succeq 0 \end{array} \left\{ \begin{array}{l} \langle F, v_P \rangle + \langle G, M_P \rangle \\ \langle F, v_I \rangle + \langle G, M_I \rangle \leq R^2 : \tau \\ \forall k, \langle F, v_{\mathcal{F}}^{(k)} \rangle + \langle G, M_{\mathcal{F}}^{(k)} \rangle \leq 0 : \lambda_{\mathcal{F}}^{(k)} \\ \forall l, \langle F, v_{\mathcal{A}}^{(l)} \rangle + \langle G, M_{\mathcal{A}}^{(l)} \rangle \leq 0 : \lambda_{\mathcal{A}}^{(l)} \end{array} \right.$$

◇ Lagrangian is

$$\begin{aligned} \mathcal{L} = & \tau R^2 + \left\langle F, v_P - \tau v_I - \sum_k \lambda_{\mathcal{F}}^{(k)} v_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} v_{\mathcal{A}}^{(l)} \right\rangle \\ & + \left\langle G, M_P - \tau M_I - \sum_k \lambda_{\mathcal{F}}^{(k)} M_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} M_{\mathcal{A}}^{(l)} \right\rangle \end{aligned}$$

From primal to dual problem

◇ Primal problem is

$$\left\{ \begin{array}{l} \text{maximize} \\ F, G \succeq 0 \end{array} \right. \left\{ \begin{array}{l} \langle F, v_P \rangle + \langle G, M_P \rangle \\ \text{subject to} \end{array} \right. \left\{ \begin{array}{l} \langle F, v_I \rangle + \langle G, M_I \rangle \leq R^2 : \tau \\ \forall k, \langle F, v_{\mathcal{F}}^{(k)} \rangle + \langle G, M_{\mathcal{F}}^{(k)} \rangle \leq 0 : \lambda_{\mathcal{F}}^{(k)} \\ \forall l, \langle F, v_{\mathcal{A}}^{(l)} \rangle + \langle G, M_{\mathcal{A}}^{(l)} \rangle \leq 0 : \lambda_{\mathcal{A}}^{(l)} \end{array} \right.$$

◇ Lagrangian is

$$\begin{aligned} \mathcal{L} = & \tau R^2 + \left\langle F, v_P - \tau v_I - \sum_k \lambda_{\mathcal{F}}^{(k)} v_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} v_{\mathcal{A}}^{(l)} \right\rangle \\ & + \left\langle G, M_P - \tau M_I - \sum_k \lambda_{\mathcal{F}}^{(k)} M_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} M_{\mathcal{A}}^{(l)} \right\rangle \end{aligned}$$

◇ Dual problem is

$$\left\{ \begin{array}{l} \text{minimize} \\ \tau, \lambda_{\mathcal{F}}^{(k)}, \lambda_{\mathcal{A}}^{(l)} \geq 0 \end{array} \right. \tau R^2$$

$$\text{s.t.} \left\{ \begin{array}{l} v_P - \tau v_I - \sum_k \lambda_{\mathcal{F}}^{(k)} v_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} v_{\mathcal{A}}^{(l)} = 0 \\ M_P - \tau M_I - \sum_k \lambda_{\mathcal{F}}^{(k)} M_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} M_{\mathcal{A}}^{(l)} \preceq 0 \end{array} \right.$$

(PEP-dual)

From primal to dual problem

◇ Primal problem is

$$\begin{cases} \text{maximize} & \langle F, v_P \rangle + \langle G, M_P \rangle \\ & F, G \succeq 0 \\ \text{subject to} & \begin{cases} \langle F, v_I \rangle + \langle G, M_I \rangle \leq R^2 : \tau \\ \forall k, \langle F, v_{\mathcal{F}}^{(k)} \rangle + \langle G, M_{\mathcal{F}}^{(k)} \rangle \leq 0 : \lambda_{\mathcal{F}}^{(k)} \\ \forall l, \langle F, v_{\mathcal{A}}^{(l)} \rangle + \langle G, M_{\mathcal{A}}^{(l)} \rangle \leq 0 : \lambda_{\mathcal{A}}^{(l)} \end{cases} \end{cases}$$

◇ Lagrangian is

$$\begin{aligned} \mathcal{L} = & \tau R^2 + \left\langle F, v_P - \tau v_I - \sum_k \lambda_{\mathcal{F}}^{(k)} v_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} v_{\mathcal{A}}^{(l)} \right\rangle \\ & + \left\langle G, M_P - \tau M_I - \sum_k \lambda_{\mathcal{F}}^{(k)} M_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} M_{\mathcal{A}}^{(l)} \right\rangle \end{aligned}$$

◇ Dual problem is

$$\begin{cases} \text{minimize} & \tau R^2 \\ & \tau, \lambda_{\mathcal{F}}^{(k)}, \lambda_{\mathcal{A}}^{(l)} \geq 0 \\ \text{s.t.} & \begin{cases} v_P - \tau v_I - \sum_k \lambda_{\mathcal{F}}^{(k)} v_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} v_{\mathcal{A}}^{(l)} = 0 \\ M_P - \tau M_I - \sum_k \lambda_{\mathcal{F}}^{(k)} M_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} M_{\mathcal{A}}^{(l)} \preceq 0 \end{cases} \end{cases} \quad (\text{PEP-dual})$$

◇ Any *primal* feasible point leads to a *lower* bound on the worst-case convergence rate. The *primal* variables provide adversarial objective.

From primal to dual problem

- Primal problem is

$$\begin{cases} \text{maximize} & \langle F, v_P \rangle + \langle G, M_P \rangle \\ & F, G \succeq 0 \\ \text{subject to} & \begin{cases} \langle F, v_I \rangle + \langle G, M_I \rangle \leq R^2 : \tau \\ \forall k, \langle F, v_{\mathcal{F}}^{(k)} \rangle + \langle G, M_{\mathcal{F}}^{(k)} \rangle \leq 0 : \lambda_{\mathcal{F}}^{(k)} \\ \forall l, \langle F, v_{\mathcal{A}}^{(l)} \rangle + \langle G, M_{\mathcal{A}}^{(l)} \rangle \leq 0 : \lambda_{\mathcal{A}}^{(l)} \end{cases} \end{cases}$$

- Lagrangian is

$$\begin{aligned} \mathcal{L} = & \tau R^2 + \left\langle F, v_P - \tau v_I - \sum_k \lambda_{\mathcal{F}}^{(k)} v_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} v_{\mathcal{A}}^{(l)} \right\rangle \\ & + \left\langle G, M_P - \tau M_I - \sum_k \lambda_{\mathcal{F}}^{(k)} M_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} M_{\mathcal{A}}^{(l)} \right\rangle \end{aligned}$$

- Dual problem is

$$\begin{cases} \text{minimize} & \tau R^2 \\ & \tau, \lambda_{\mathcal{F}}^{(k)}, \lambda_{\mathcal{A}}^{(l)} \geq 0 \\ \text{s.t.} & \begin{cases} v_P - \tau v_I - \sum_k \lambda_{\mathcal{F}}^{(k)} v_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} v_{\mathcal{A}}^{(l)} = 0 \\ M_P - \tau M_I - \sum_k \lambda_{\mathcal{F}}^{(k)} M_{\mathcal{F}}^{(k)} - \sum_l \lambda_{\mathcal{A}}^{(l)} M_{\mathcal{A}}^{(l)} \preceq 0 \end{cases} \end{cases} \quad (\text{PEP-dual})$$

- Any *primal* feasible point leads to a *lower* bound on the worst-case convergence rate. The *primal* variables provide adversarial objective.
- Any *dual* feasible point leads to an *upper* bound on the worst-case convergence rate. The *dual* variables provide the proof of convergence.

Deriving a proof of convergence

◇ Lagrangian is

$$\begin{aligned} \mathcal{L} = & \tau R^2 + \left\langle F, v_P - \tau v_I - \sum_k \lambda_{\mathcal{F}}^{(k)} v_{\mathcal{F}}^{(k)} - \sum_I \lambda_{\mathcal{A}}^{(I)} v_{\mathcal{A}}^{(I)} \right\rangle \\ & + \left\langle G, M_P - \tau M_I - \sum_k \lambda_{\mathcal{F}}^{(k)} M_{\mathcal{F}}^{(k)} - \sum_I \lambda_{\mathcal{A}}^{(I)} M_{\mathcal{A}}^{(I)} \right\rangle \end{aligned}$$

Deriving a proof of convergence

◇ Lagrangian is

$$\mathcal{L} = \tau R^2 + \left\langle F, v_P - \tau v_I - \sum_k \lambda_{\mathcal{F}}^{(k)} v_{\mathcal{F}}^{(k)} - \sum_I \lambda_{\mathcal{A}}^{(I)} v_{\mathcal{A}}^{(I)} \right\rangle \\ + \left\langle G, M_P - \tau M_I - \sum_k \lambda_{\mathcal{F}}^{(k)} M_{\mathcal{F}}^{(k)} - \sum_I \lambda_{\mathcal{A}}^{(I)} M_{\mathcal{A}}^{(I)} \right\rangle$$

◇ Dual problem is

$$\left| \begin{array}{l} \text{minimize } \tau R^2 \\ \tau, \lambda_{\mathcal{F}}^{(k)}, \lambda_{\mathcal{A}}^{(I)} \geq 0 \\ \text{s.t. } \begin{cases} v_P - \tau v_I - \sum_k \lambda_{\mathcal{F}}^{(k)} v_{\mathcal{F}}^{(k)} - \sum_I \lambda_{\mathcal{A}}^{(I)} v_{\mathcal{A}}^{(I)} = 0 \\ M_P - \tau M_I - \sum_k \lambda_{\mathcal{F}}^{(k)} M_{\mathcal{F}}^{(k)} - \sum_I \lambda_{\mathcal{A}}^{(I)} M_{\mathcal{A}}^{(I)} \preceq 0 \end{cases} \end{array} \right. \quad \text{(PEP-dual)}$$

◇ Lagrangian \leq Dual:

Deriving a proof of convergence

◇ Lagrangian is

$$\mathcal{L} = \tau R^2 + \left\langle F, v_P - \tau v_I - \sum_k \lambda_{\mathcal{F}}^{(k)} v_{\mathcal{F}}^{(k)} - \sum_I \lambda_{\mathcal{A}}^{(I)} v_{\mathcal{A}}^{(I)} \right\rangle + \left\langle G, M_P - \tau M_I - \sum_k \lambda_{\mathcal{F}}^{(k)} M_{\mathcal{F}}^{(k)} - \sum_I \lambda_{\mathcal{A}}^{(I)} M_{\mathcal{A}}^{(I)} \right\rangle$$

◇ Dual problem is

$$\begin{cases} \text{minimize } \tau R^2 \\ \tau, \lambda_{\mathcal{F}}^{(k)}, \lambda_{\mathcal{A}}^{(I)} \geq 0 \\ \text{s.t. } \begin{cases} v_P - \tau v_I - \sum_k \lambda_{\mathcal{F}}^{(k)} v_{\mathcal{F}}^{(k)} - \sum_I \lambda_{\mathcal{A}}^{(I)} v_{\mathcal{A}}^{(I)} = 0 \\ M_P - \tau M_I - \sum_k \lambda_{\mathcal{F}}^{(k)} M_{\mathcal{F}}^{(k)} - \sum_I \lambda_{\mathcal{A}}^{(I)} M_{\mathcal{A}}^{(I)} \leq 0 \end{cases} \end{cases} \quad (\text{PEP-dual})$$

◇ Lagrangian \leq Dual: for any feasible primal F, G and feasible dual $\tau, (\lambda_{\mathcal{F}}^{(k)})_k, (\lambda_{\mathcal{A}}^{(I)})_I$,

$$\underbrace{\langle F, v_P \rangle + \langle G, M_P \rangle}_{\text{Performance metric}} - \tau \underbrace{[\langle F, v_I \rangle + \langle G, M_I \rangle]}_{\text{Initialization}} \leq \sum_k \lambda_{\mathcal{F}}^{(k)} \underbrace{[\langle F, v_{\mathcal{F}}^{(k)} \rangle + \langle G, M_{\mathcal{F}}^{(k)} \rangle]}_{\text{Class constraint}} + \sum_I \lambda_{\mathcal{A}}^{(I)} \underbrace{[\langle F, v_{\mathcal{A}}^{(I)} \rangle + \langle G, M_{\mathcal{A}}^{(I)} \rangle]}_{\text{Algorithm constraint}} \leq 0.$$

Deriving a proof of convergence

◇ Lagrangian is

$$\mathcal{L} = \tau R^2 + \left\langle F, v_P - \tau v_I - \sum_k \lambda_{\mathcal{F}}^{(k)} v_{\mathcal{F}}^{(k)} - \sum_I \lambda_{\mathcal{A}}^{(I)} v_{\mathcal{A}}^{(I)} \right\rangle + \left\langle G, M_P - \tau M_I - \sum_k \lambda_{\mathcal{F}}^{(k)} M_{\mathcal{F}}^{(k)} - \sum_I \lambda_{\mathcal{A}}^{(I)} M_{\mathcal{A}}^{(I)} \right\rangle$$

◇ Dual problem is

$$\begin{cases} \text{minimize } \tau R^2 \\ \tau, \lambda_{\mathcal{F}}^{(k)}, \lambda_{\mathcal{A}}^{(I)} \geq 0 \\ \text{s.t. } \begin{cases} v_P - \tau v_I - \sum_k \lambda_{\mathcal{F}}^{(k)} v_{\mathcal{F}}^{(k)} - \sum_I \lambda_{\mathcal{A}}^{(I)} v_{\mathcal{A}}^{(I)} = 0 \\ M_P - \tau M_I - \sum_k \lambda_{\mathcal{F}}^{(k)} M_{\mathcal{F}}^{(k)} - \sum_I \lambda_{\mathcal{A}}^{(I)} M_{\mathcal{A}}^{(I)} \leq 0 \end{cases} \end{cases} \quad (\text{PEP-dual})$$

◇ Lagrangian \leq Dual: for any feasible primal F, G and feasible dual $\tau, (\lambda_{\mathcal{F}}^{(k)})_k, (\lambda_{\mathcal{A}}^{(I)})_I$,

$$\underbrace{\langle F, v_P \rangle + \langle G, M_P \rangle}_{\text{Performance metric}} - \tau \underbrace{[\langle F, v_I \rangle + \langle G, M_I \rangle]}_{\text{Initialization}} \leq \sum_k \lambda_{\mathcal{F}}^{(k)} \underbrace{[\langle F, v_{\mathcal{F}}^{(k)} \rangle + \langle G, M_{\mathcal{F}}^{(k)} \rangle]}_{\text{Class constraint}} + \sum_I \lambda_{\mathcal{A}}^{(I)} \underbrace{[\langle F, v_{\mathcal{A}}^{(I)} \rangle + \langle G, M_{\mathcal{A}}^{(I)} \rangle]}_{\text{Algorithm constraint}} \leq 0.$$

◇ Performance metric - τ initialization $\leq \sum_i \lambda_i$ Constraint $_i \leq 0$.

Some insights

- ◇ The proof of worst-case guarantee simply consists in a positively linear combination of the constraints.

Some insights

- ◇ The proof of worst-case guarantee simply consists in a positively linear combination of the constraints.
- ◇ Interpolation constraints are the right constraints to be considered when studying a class.

Some insights

- ◇ The proof of worst-case guarantee simply consists in a positively linear combination of the constraints.
- ◇ Interpolation constraints are the right constraints to be considered when studying a class.
- ◇ Dual zeros:

Some insights

- ◇ The proof of worst-case guarantee simply consists in a positively linear combination of the constraints.
- ◇ Interpolation constraints are the right constraints to be considered when studying a class.
- ◇ Dual zeros:
 - Each $\lambda = 0$ witnesses the non use of the corresponding inequality.

Some insights

- ◇ The proof of worst-case guarantee simply consists in a positively linear combination of the constraints.
- ◇ Interpolation constraints are the right constraints to be considered when studying a class.
- ◇ Dual zeros:
 - Each $\lambda = 0$ witnesses the non use of the corresponding inequality.
 - By strong duality, the proof is tight for a function. We therefore recover the slackness condition: the corresponding inequality must be an equality.

Some insights

- ◇ The proof of worst-case guarantee simply consists in a positively linear combination of the constraints.
- ◇ Interpolation constraints are the right constraints to be considered when studying a class.
- ◇ Dual zeros:
 - Each $\lambda = 0$ witnesses the non use of the corresponding inequality.
 - By strong duality, the proof is tight for a function. We therefore recover the slackness condition: the corresponding inequality must be an equality.
- ◇ Obtaining a guarantee over the minimum of performance metrics is equivalent to obtaining one on a certain average of those.

Natural performance metric

- ◇ $\frac{1}{\tau}$ Performance metric - initialization $\leq \sum_i \lambda_i \text{Constraint}_i \leq 0$ might be loose.
There sometimes remains a residual.

Natural performance metric

- ◇ $\frac{1}{\tau}$ Performance metric - initialization $\leq \sum_i \lambda_i \text{Constraint}_i \leq 0$ might be loose. There sometimes remains a residual.
- ◇ We can tighten the proof by considering

$$V_t = \underbrace{[\langle F, v_I \rangle + \langle G, M_I \rangle]}_{\text{Initialization}} + \sum_{j \mid \text{only involves values observed before step } t} \lambda^{(j)} \underbrace{[\langle F, v^{(j)} \rangle + \langle G, M^{(j)} \rangle]}_{\text{Constraint}}. \quad (1)$$

Natural performance metric

- ◇ $\frac{1}{\tau}$ Performance metric - initialization $\leq \sum_i \lambda_i \text{Constraint}_i \leq 0$ might be loose. There sometimes remains a residual.
- ◇ We can tighten the proof by considering

$$V_t = \underbrace{[\langle F, v_I \rangle + \langle G, M_I \rangle]}_{\text{Initialization}} + \sum_{j \mid \text{only involves values observed before step } t} \lambda^{(j)} \underbrace{[\langle F, v^{(j)} \rangle + \langle G, M^{(j)} \rangle]}_{\text{Constraint}}. \quad (1)$$

- ◇ Example on NAG:

$$\lambda_{t+1} = \frac{1}{2} + \sqrt{\frac{1}{4} + \lambda_t^2}, \quad y_t = x_t + \frac{\lambda_t - 1}{\lambda_{t+1}}(x_t - x_{t-1}), \quad x_{t+1} = y_t - \frac{1}{L} \nabla f(y_t).$$

$$V_t = \lambda_t^2 (f_t - f_*) + \frac{L}{2} \|\lambda_t(x_t - x_*) + (1 - \lambda_t)(x_{t-1} - x_*)\|^2 \\ + \frac{1}{2L} \sum_{s=1}^{t-1} [\lambda_{s+1}^2 \|\nabla f(x_{s+1})\|^2 + \lambda_{s+1} \|\nabla f(y_s)\| \lambda_s^2 \|\nabla f(y_s) - \nabla f(x_s)\|^2]$$

Extensions of a certificate

- ◇ Fewer class constraints allows extension to larger class,

Extensions of a certificate

- ◇ Fewer class constraints allows extension to larger class,
- ◇ Fewer algorithm constraints allows extension to new algorithms,

Extensions of a certificate

- ◇ Fewer class constraints allows extension to larger class,
- ◇ Fewer algorithm constraints allows extension to new algorithms,
- ◇ Fewer class constraints allows extension to new algorithms.

Extensions of a certificate

- ◇ Fewer class constraints allows extension to larger class,
- ◇ Fewer algorithm constraints allows extension to new algorithms,
- ◇ Fewer class constraints allows extension to new algorithms.
Example of Backtracking line-search [C. Park et al., 2021]:

Extensions of a certificate

- ◇ Fewer class constraints allows extension to larger class,
- ◇ Fewer algorithm constraints allows extension to new algorithms,
- ◇ Fewer class constraints allows extension to new algorithms.

Example of Backtracking line-search [C. Park et al., 2021]:

- $f_{t+1} \leq f_t - \frac{1}{2L} \|g_t\|^2$ is verifiable,

Extensions of a certificate

- ◇ Fewer class constraints allows extension to larger class,
- ◇ Fewer algorithm constraints allows extension to new algorithms,
- ◇ Fewer class constraints allows extension to new algorithms.

Example of Backtracking line-search [C. Park et al., 2021]:

- $f_{t+1} \leq f_t - \frac{1}{2L} \|g_t\|^2$ is verifiable,
- $f_{t+1} \leq f_t + \langle g_{t+1}, x_{t+1} - x_t \rangle - \frac{1}{2L} \|g_{t+1} - g_t\|^2$ is not.

Greedy first order method (GFOM)

◇ Update:

$$x_{t+1} = \arg \min_{x \in x_0 + \text{span}(\nabla f(x_0), \nabla f(x_1), \dots, \nabla f(x_t))} f(x) \quad (\text{GFOM})$$

Greedy first order method (GFOM)

- ◇ Update:

$$x_{t+1} = \arg \min_{x \in x_0 + \text{span}(\nabla f(x_0), \nabla f(x_1), \dots, \nabla f(x_t))} f(x) \quad (\text{GFOM})$$

- ◇ Additional constraints implied by GFOM:

$$\begin{aligned} \langle g_i, g_j \rangle &= 0 && \text{for all } j < i \\ \langle g_i, x_j - x_0 \rangle &= 0 && \text{for all } j \leq i \end{aligned}$$

Proof of worst case guarantee for GFOM

- ◇ Structure of the certificate proof:

$$\begin{aligned} \text{Performance metric} - \tau \text{initialization} &\leq \sum_i \lambda_i \underbrace{\text{Constraint}_i}_{\leq 0} \\ &+ \sum_i \sum_{j < i} \beta_{i,j} \underbrace{\langle \mathbf{g}_i, \mathbf{g}_j \rangle}_{=0} \\ &+ \sum_i \sum_{j \leq i} \gamma_{i,j} \underbrace{\langle \mathbf{g}_i, \mathbf{x}_j - \mathbf{x}_0 \rangle}_{=0} \end{aligned}$$

Proof of worst case guarantee for GFOM

- ◇ Structure of the certificate proof:

$$\begin{aligned} \text{Performance metric} - \tau \text{initialization} &\leq \sum_i \lambda_i \underbrace{\text{Constraint}_i}_{\leq 0} \\ &+ \sum_i \sum_{j < i} \beta_{i,j} \underbrace{\langle \mathbf{g}_i, \mathbf{g}_j \rangle}_{=0} \\ &+ \sum_i \sum_{j \leq i} \gamma_{i,j} \underbrace{\langle \mathbf{g}_i, \mathbf{x}_j - \mathbf{x}_0 \rangle}_{=0} \end{aligned}$$

- ◇ For all i ,

$$\langle \mathbf{g}_i, \sum_{j < i} \beta_{i,j} \mathbf{g}_j + \sum_{j \leq i} \gamma_{i,j} (\mathbf{x}_j - \mathbf{x}_0) \rangle = 0$$

would be enough to obtain the same worst case guarantee.

Automating search for Lyapunov functions

- ◇ We look for proof under the form $V_{t+1} - \tau V_t \leq \sum_i \lambda_i \text{Constraint}_i \leq 0$ where V is a linear combination of a finite family of expressions of the iterates sequence. Fixing τ , this problem is an SDP.

Automating search for Lyapunov functions

- ◇ We look for proof under the form $V_{t+1} - \tau V_t \leq \sum_i \lambda_i \text{Constraint}_i \leq 0$ where V is a linear combination of a finite family of expressions of the iterates sequence. Fixing τ , this problem is an SDP.
- ◇ In [A. Taylor et al., 2018], authors show that V does not necessarily need to be a positive form of the iterates.

Automating search for Lyapunov functions

- ◇ We look for proof under the form $V_{t+1} - \tau V_t \leq \sum_i \lambda_i \text{Constraint}_i \leq 0$ where V is a linear combination of a finite family of expressions of the iterates sequence. Fixing τ , this problem is an SDP.
- ◇ In [A. Taylor et al., 2018], authors show that V does not necessarily need to be a positive form of the iterates.
- ◇ But it must be positive on the path of iterates, i.e.

$$0 - V_{t+1} \leq \sum_i \tilde{\lambda}_i \text{Constraint}_i \leq 0$$

Automating search for Lyapunov functions

◇ We look for proof under the form $V_{t+1} - \tau V_t \leq \sum_i \lambda_i \text{Constraint}_i \leq 0$ where V is a linear combination of a finite family of expressions of the iterates sequence. Fixing τ , this problem is an SDP.

◇ In [A. Taylor et al., 2018], authors show that V does not necessarily need to be a positive form of the iterates.

◇ But it must be positive on the path of iterates, i.e.

$$0 - V_{t+1} \leq \sum_i \tilde{\lambda}_i \text{Constraint}_i \leq 0$$

◇ Bonus: this problem is defined as a feasibility problem, leaving free the objective to minimize!

Automating search for Lyapunov functions

◇ We look for proof under the form $V_{t+1} - \tau V_t \leq \sum_i \lambda_i \text{Constraint}_i \leq 0$ where V is a linear combination of a finite family of expressions of the iterates sequence. Fixing τ , this problem is an SDP.

◇ In [A. Taylor et al., 2018], authors show that V does not necessarily need to be a positive form of the iterates.

◇ But it must be positive on the path of iterates, i.e.

$$0 - V_{t+1} \leq \sum_i \tilde{\lambda}_i \text{Constraint}_i \leq 0$$

◇ Bonus: this problem is defined as a feasibility problem, leaving free the objective to minimize!

◇ Bonus 2: the search for Lyapunov can be combined with the SSEP design procedure.

Conclusion

- ◇ Finding the worst case guarantee of a First-order method can be solved through convex optimization.
 - Requires 3 ingredients: Interpolation, Homogenization and SDP lifting.
 - Can be applied for various function class, algorithms and performance metrics.
- ◇ Modelization is heavy \Rightarrow PEPit package.
- ◇ Dual view enables to understand the systematic structure of convergence proofs.
 - Always linear sum of inequalities.
 - From a function class point of view, only interpolation conditions matter.
 - We can design algorithms.
 - Search for Lyapunov can be automatized as well and combined with the algorithm design.
- ◇ Given a class, running SSEP with Lyapunov search provides an algorithm with its rate, and a Lyapunov based proof of convergence. All the information are stored in the dual variables.
- ◇ Finally, the zeros tell us which inequalities must be cancel to prove tightness of the result.

In summary, many insights are hidden in the dual variables!

Application / Open problem

Using [N. Bouselmi et al., 2023]'s interpolation constraints for the classes of quadratic functions, determine a PEP-based proof of guarantee of Heavy-ball.

Application / Open problem

Using [N. Bouselmi et al., 2023]'s interpolation constraints for the classes of quadratic functions, determine a PEP-based proof of guarantee of Heavy-ball.

By removing useless class constraints, can we identify a larger class that quadratic functions on which HB accelerates?

Thanks!

Questions?