

robust control for analysis and design of large-scale optimization algorithms

Laurent Lessard

University of Wisconsin–Madison

Joint work with Ben Recht and Andy Packard

LCCC Workshop on Large-Scale and Distributed Optimization
Lund University, June 15, 2017

1. Many algorithms can be viewed as dynamical systems with feedback (control systems!).

algorithm convergence \iff system stability

2. By solving a small convex program, we can recover state-of-the-art convergence results for these algorithms, automatically and efficiently.
3. The ultimate goal: to move from analysis to design.

Unconstrained optimization:

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{subject to} & x \in \mathbb{R}^N \end{array}$$

- need algorithms that are *fast* and *simple*
- currently favored family: *first-order methods*

Gradient method

$$x_{k+1} = x_k - \alpha \nabla f(x_k)$$

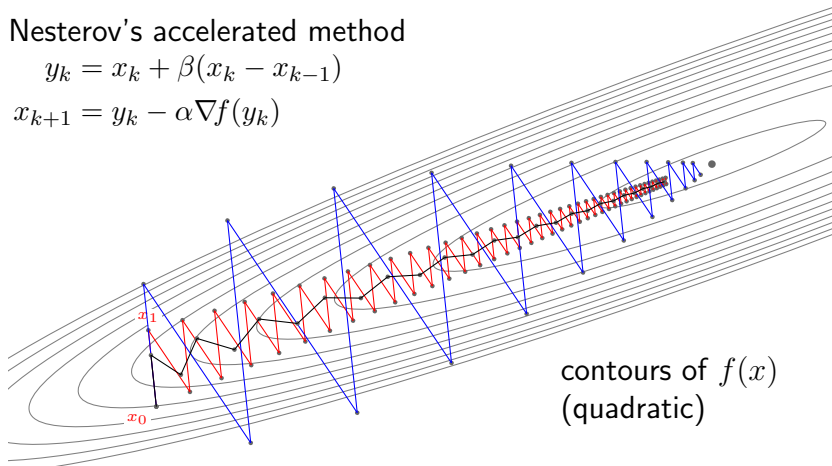
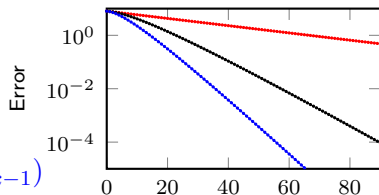
Heavy ball method

$$x_{k+1} = x_k - \alpha \nabla f(x_k) + \beta(x_k - x_{k-1})$$

Nesterov's accelerated method

$$y_k = x_k + \beta(x_k - x_{k-1})$$

$$x_{k+1} = y_k - \alpha \nabla f(y_k)$$



Robust algorithm selection

$G \in \mathcal{G}$: algorithm we're going to use

$f \in \mathcal{S}$: function we'd like to minimize

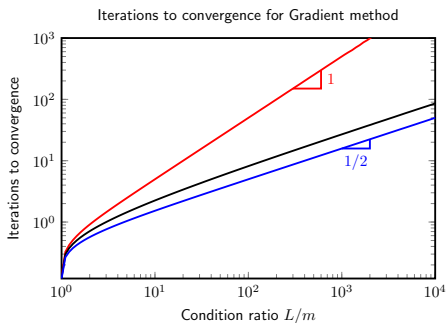
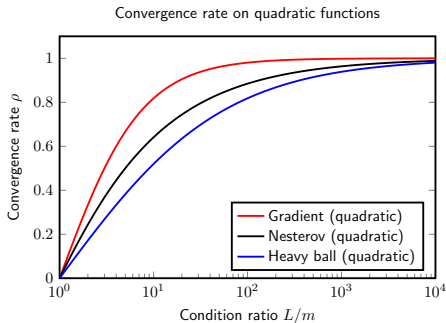
$$G_{\text{opt}} = \arg \min_{G \in \mathcal{G}} \left(\max_{f \in \mathcal{S}} \text{cost}(f, G) \right)$$

Similar problem for a finite number of iterations:

- Drori, Teboulle (2012)
- Taylor, Hendrickx, Glineur (2016)

$$G \in \mathcal{G} \left\{ \begin{array}{l} \text{Gradient method} \\ x_{k+1} = x_k - \alpha \nabla f(x_k) \\ \\ \text{Heavy ball method} \\ x_{k+1} = x_k - \alpha \nabla f(x_k) + \beta(x_k - x_{k-1}) \\ \\ \text{Nesterov's accelerated method} \\ x_{k+1} = x_k - \alpha \nabla f(x_k + \beta(x_k - x_{k-1})) + \beta(x_k - x_{k-1}) \end{array} \right.$$

$$f \in \mathcal{S} \left\{ \begin{array}{l} \text{Analytically solvable:} \\ \text{Quadratic functions: } f(x) = \frac{1}{2}x^\top Qx - p^\top x \\ \text{with the constraint: } mI \preceq Q \preceq LI \end{array} \right.$$



Convergence rate : $\|x_k - x_\star\| \leq C\rho^k \|x_0 - x_\star\|$

Iterations to convergence $\propto -\frac{1}{\log \rho}$

Robust algorithm selection

$G \in \mathcal{G}$: algorithm we're going to use

$f \in \mathcal{S}$: function we'd like to minimize

$$G_{\text{opt}} = \arg \min_{G \in \mathcal{G}} \left(\max_{f \in \mathcal{S}} \text{cost}(f, G) \right)$$

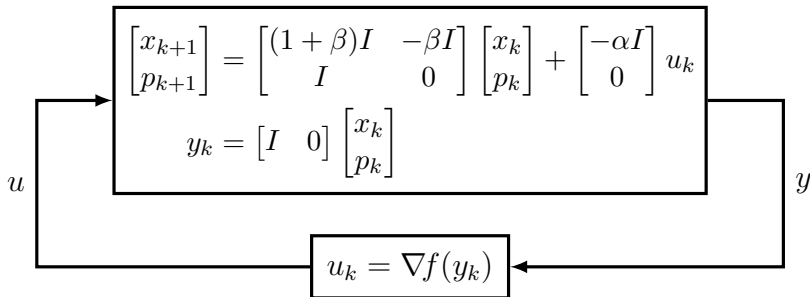
1. mathematical representation for \mathcal{G}
2. mathematical representation for \mathcal{S}
3. main robustness result

Dynamical system interpretation

Heavy ball: $x_{k+1} = x_k - \alpha \nabla f(x_k) + \beta(x_k - x_{k-1})$

Define $u_k := \nabla f(x_k)$ and $p_k := x_{k-1}$

algorithm (linear, known, decoupled)



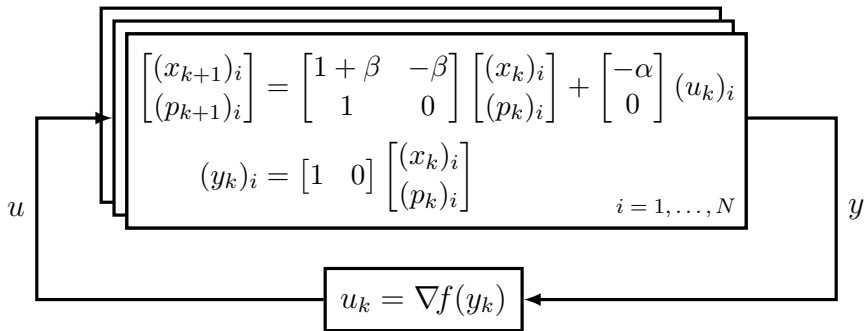
function (nonlinear, uncertain, coupled)

Dynamical system interpretation

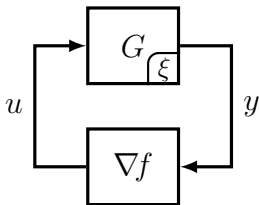
Heavy ball: $x_{k+1} = x_k - \alpha \nabla f(x_k) + \beta(x_k - x_{k-1})$

Define $u_k := \nabla f(x_k)$ and $p_k := x_{k-1}$

algorithm (linear, known, **decoupled**)



function (nonlinear, uncertain, **coupled**)

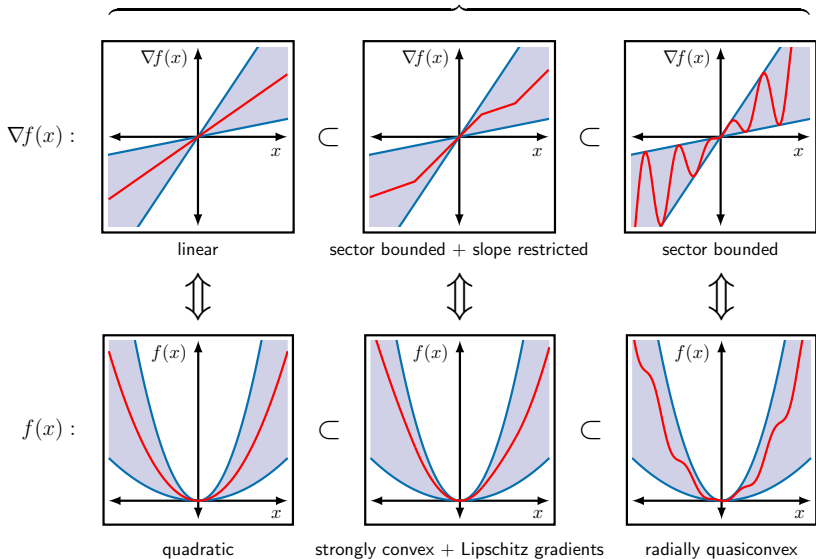
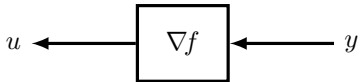


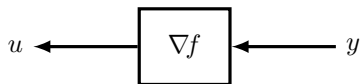
$$\xi_{k+1} = A\xi_k + Bu_k$$

$$y_k = C\xi_k$$

$$u_k = \nabla f(y_k)$$

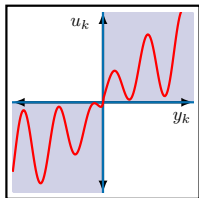
$$\left[\begin{array}{c|c} A & B \\ \hline C & 0 \end{array} \right] = \left\{ \begin{array}{l} \left[\begin{array}{cc|c} 1 & -\alpha \\ \hline 1 & 0 \end{array} \right] \quad \text{Gradient} \\ \left[\begin{array}{cc|c} 1+\beta & -\beta & -\alpha \\ \hline 1 & 0 & 0 \\ \hline 1 & 0 & 0 \end{array} \right] \quad \text{Heavy ball} \\ \left[\begin{array}{cc|c} 1+\beta & -\beta & -\alpha \\ \hline 1 & 0 & 0 \\ \hline 1+\beta & -\beta & 0 \end{array} \right] \quad \text{Nesterov} \end{array} \right.$$





Representing function classes

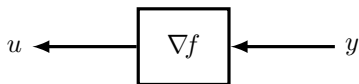
express as quadratic constraints on (y, u)



sector bounded

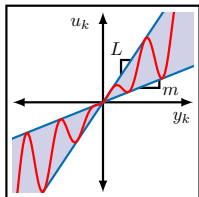
∇f is a **passive** function:

$$u_k y_k \geq 0$$



Representing function classes

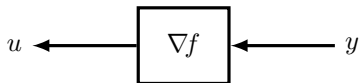
express as quadratic constraints on (y, u)



sector bounded

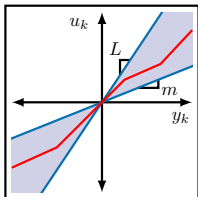
∇f is **sector-bounded**:

$$\begin{bmatrix} y_k \\ u_k \end{bmatrix}^T \begin{bmatrix} -2mL & m+L \\ m+L & -2 \end{bmatrix} \begin{bmatrix} y_k \\ u_k \end{bmatrix} \geq 0$$



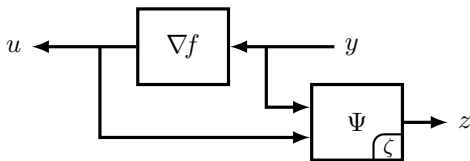
Representing function classes

express as quadratic constraints on (y, u)



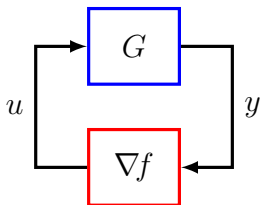
sector bounded + slope restricted

∇f is **sector-bounded** + **slope-restricted**:
 constraint on (y_k, u_k) depends on history
 $(y_0, \dots, y_{k-1}, u_0, \dots, u_{k-1})$.



Introduce extra dynamics

- Design dynamics Ψ and multiplier matrix M .
- Instead of using $q(u_k, y_k)$, use $z_k^T M z_k$.
- Systematic way of doing this for strong convexity via Zames-Falb multipliers (1968).
- General theory: Integral Quadratic Constraints (Megretski & Rantzer 1997)



$$\left[\begin{array}{c|c} 1 & -\alpha \\ \hline 1 & 0 \end{array} \right]$$

Gradient

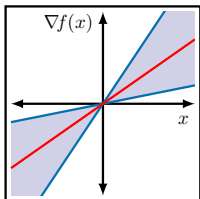
$$\left[\begin{array}{cc|c} 1+\beta & -\beta & -\alpha \\ 1 & 0 & 0 \\ \hline 1 & 0 & 0 \end{array} \right]$$

Heavy ball

$$\left[\begin{array}{cc|c} 1+\beta & -\beta & -\alpha \\ 1 & 0 & 0 \\ \hline 1+\beta & -\beta & 0 \end{array} \right]$$

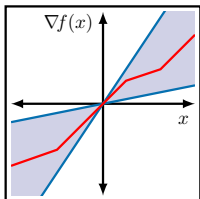
Nesterov

$$\left. \begin{array}{l} \text{Gradient} \\ \text{Heavy ball} \\ \text{Nesterov} \end{array} \right\} \left[\begin{array}{c|c} A & B \\ \hline C & 0 \end{array} \right]$$



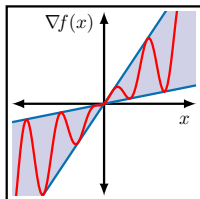
f is quadratic

\subset



f is strongly convex

\subset



f is quasiconvex

(Ψ, M)

Main result

Problem data:

- G (the algorithm)
- Ψ (what we know about f)

Auxiliary quantities:

- Compute $(\hat{A}, \hat{B}, \hat{C}, \hat{D})$ matrices from (G, Ψ)
- Choose a candidate rate $0 < \rho < 1$.

Size of LMI does **not** grow with problem dimension!
e.g. $P \in \mathbf{S}^{3 \times 3}$, LMI $\in \mathbf{S}^{4 \times 4}$

If there exists $P \succ 0$ such that

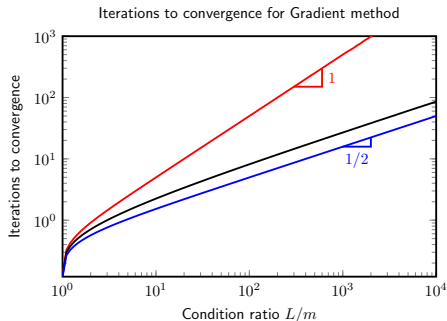
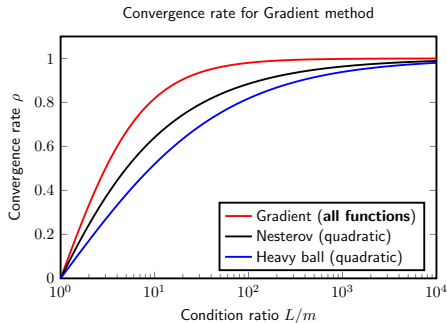
$$\begin{bmatrix} \hat{A}^\top P \hat{A} - \rho^2 P & \hat{A}^\top P \hat{B} \\ \hat{B}^\top P \hat{A} & \hat{B}^\top P \hat{B} \end{bmatrix} + \begin{bmatrix} \hat{C} & \hat{D} \end{bmatrix}^\top M \begin{bmatrix} \hat{C} & \hat{D} \end{bmatrix} \preceq 0$$

then $\|x_k - x_\star\| \leq \sqrt{\text{cond}(P)} \rho^k \|x_0 - x_\star\|$ for all k .

main results:
analytic and numerical

Gradient method

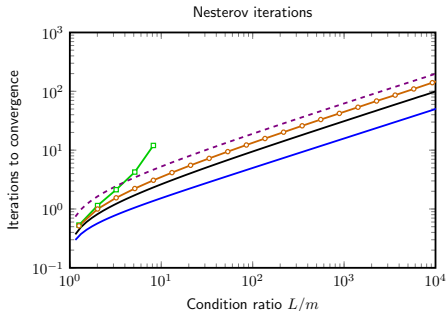
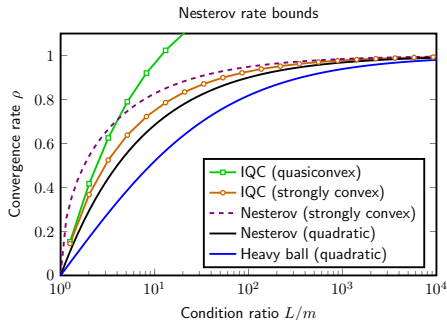
$$x_{k+1} = x_k - \alpha \nabla f(x_k)$$



analytic solution! Same rate for: quadratics, strongly convex, or quasiconvex functions.

Nesterov's method

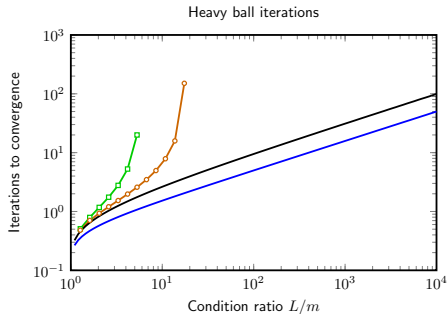
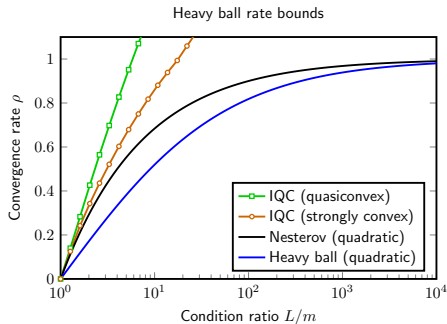
$$x_{k+1} = x_k - \alpha \nabla f(x_k + \beta(x_k - x_{k-1})) + \beta(x_k - x_{k-1})$$



- Cannot certify stability for quasiconvex functions
- IQC bound **improves** upon best known bound!

Heavy ball method

$$x_{k+1} = x_k - \alpha \nabla f(x_k) + \beta(x_k - x_{k-1})$$

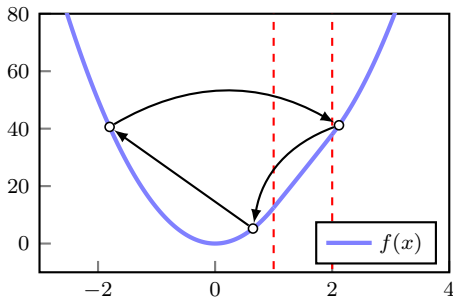


- Cannot certify stability for quasiconvex functions
- Cannot certify stability for strongly convex functions

The heavy ball method is **not** stable!

counterexample: $f(x) = \begin{cases} \frac{25}{2}x^2 & x < 1 \\ \frac{1}{2}x^2 + 24x - 12 & 1 \leq x < 2 \\ \frac{25}{2}x^2 - 24x + 36 & x \geq 2 \end{cases}$

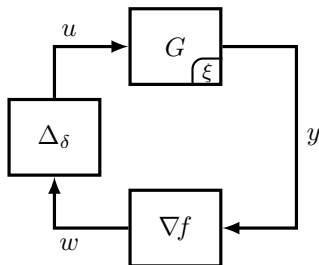
and start the heavy ball iteration at $x_0 = x_1 \in [3.07, 3.46]$.



- $L/m = 25$
- heavy ball iterations converge to a limit cycle
- simple counterexample to the Aizerman (1949) and Kalman (1957) conjectures

uncharted territory:
noise robustness and algorithm design

Noise robustness

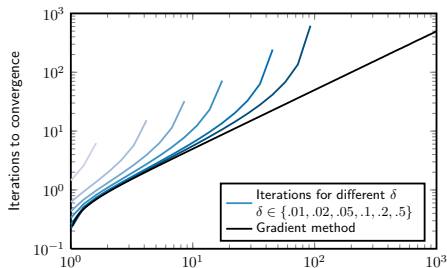
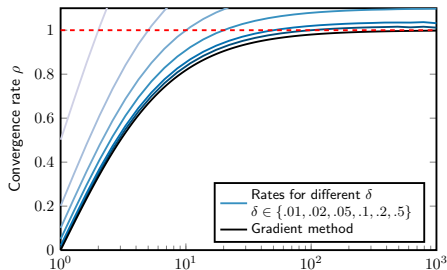


The Δ_δ block is uncertain multiplicative noise:

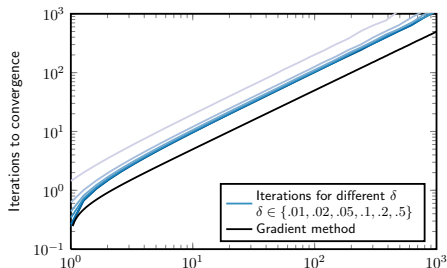
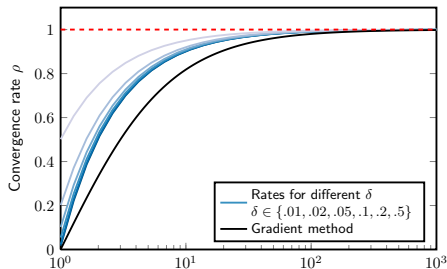
$$\|u_k - w_k\| \leq \delta \|w_k\|$$

How does an algorithm perform in the presence of noise?

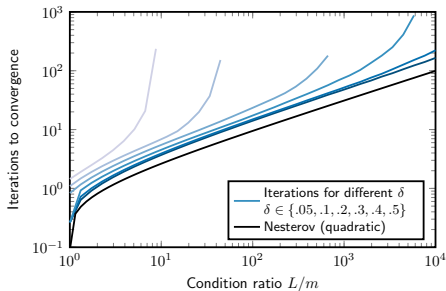
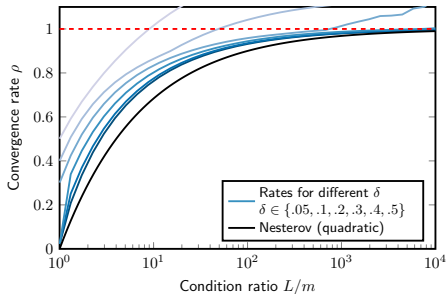
Gradient method, $\alpha = \frac{2}{L+m}$ (optimal stepsize with no noise)



Gradient method, $\alpha = \frac{1}{L}$ (more conservative stepsize)



Nesterov's method (strongly convex f , with noise)



- Nesterov's method is not robust to noise.

can we have it all? (robustness AND performance)

Brute force approach

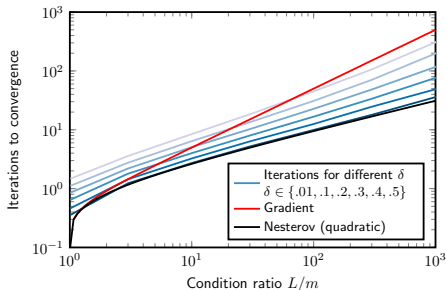
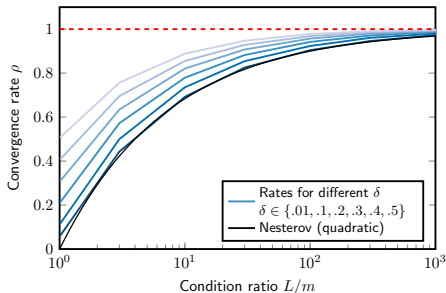
- test all strictly proper G of degree 2
- parameterization in terms of (α, β, η) :

$$\begin{aligned}x_{k+1} &= x_k - \alpha \nabla f(y_k) + \beta(x_k - x_{k-1}) \\y_k &= x_k + \eta(x_k - x_{k-1})\end{aligned}$$

Special cases:

$$(\alpha, \beta, \eta) = \begin{cases} (\alpha, 0, 0) & \text{Gradient} \\ (\alpha, \beta, 0) & \text{Heavy ball} \\ (\alpha, \beta, \beta) & \text{Nesterov} \end{cases}$$

Optimal designs over (α, β, η)



- Faster than the gradient method **and** more robust to noise than Nesterov's method
- automatic algorithm design is possible!

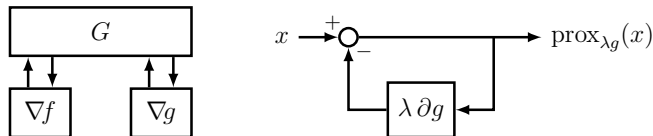
What we have (so far!)

L, Recht, Packard (SIOPT'16)

- unified framework for algorithm analysis
- read this one first!

Nishihara, L, Recht, Packard, Jordan (ICML'15)

- operator splitting methods
- application to ADMM tuning



Recent works

Boczar, L, Packard, Recht ([arXiv:1706.01337](https://arxiv.org/abs/1706.01337))

- control theory treatment
- certifying exponential convergence with IQCs

Hu, L ([arXiv:1706.04381](https://arxiv.org/abs/1706.04381))

- (energy) dissipation inequalities
- prove linear rates, $1/k$, and $1/k^2$ rates
- Lyapunov function for (time-varying) Nesterov's method

Thank you!

- Manuscripts + code available:
www.laurentlessard.com
- If you're interested, come talk to me!