

The ADMM: Past, Present, and Future

Jonathan Eckstein

Rutgers University, Piscataway, NJ, USA

RUTGERS

Rutgers Business School
Newark and New Brunswick



RUTGERS

Thank You

- Thank you for inviting me to this workshop and as the Tutte speaker
- I recently read about Tutte, and it is an honor to present in a colloquium named after him

Warning/Disclaimer

- Henry requested a general-audience talk also suitable for people not attending the MOM conference
- If you are already an ADMM expert, you may not hear very much new here
- Similar to a presentation I gave at a workshop in Edinburgh in May

The ADMM: Past

- The ADMM is now considered a standard optimization algorithm
- But it has an unusual history:
 1. It was discovered empirically before it was analyzed mathematically
 2. The initial discoverers and analyzers were French applied mathematics researchers specializing in large-scale discretized PDEs
 3. Over 20 years elapsed between its initial analysis and its becoming popular

The ADMM - Background from the Standard ALM

In the mid-1970's, this group of researchers (Fortin, Glowinski, Marrocco, Gabay, Mercier) had reformulated their discretized PDEs (roughly) as follows:

- $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is a convex function
- $g : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$ is a convex function
- M is an $m \times n$ matrix

$$\boxed{\min f(x) + g(Mx)}$$

- Equivalent formulation:

$$\boxed{\begin{array}{ll} \min & f(x) + g(z) \\ \text{ST} & Mx = z \end{array}}$$

Applying the Augmented Lagrangian Method

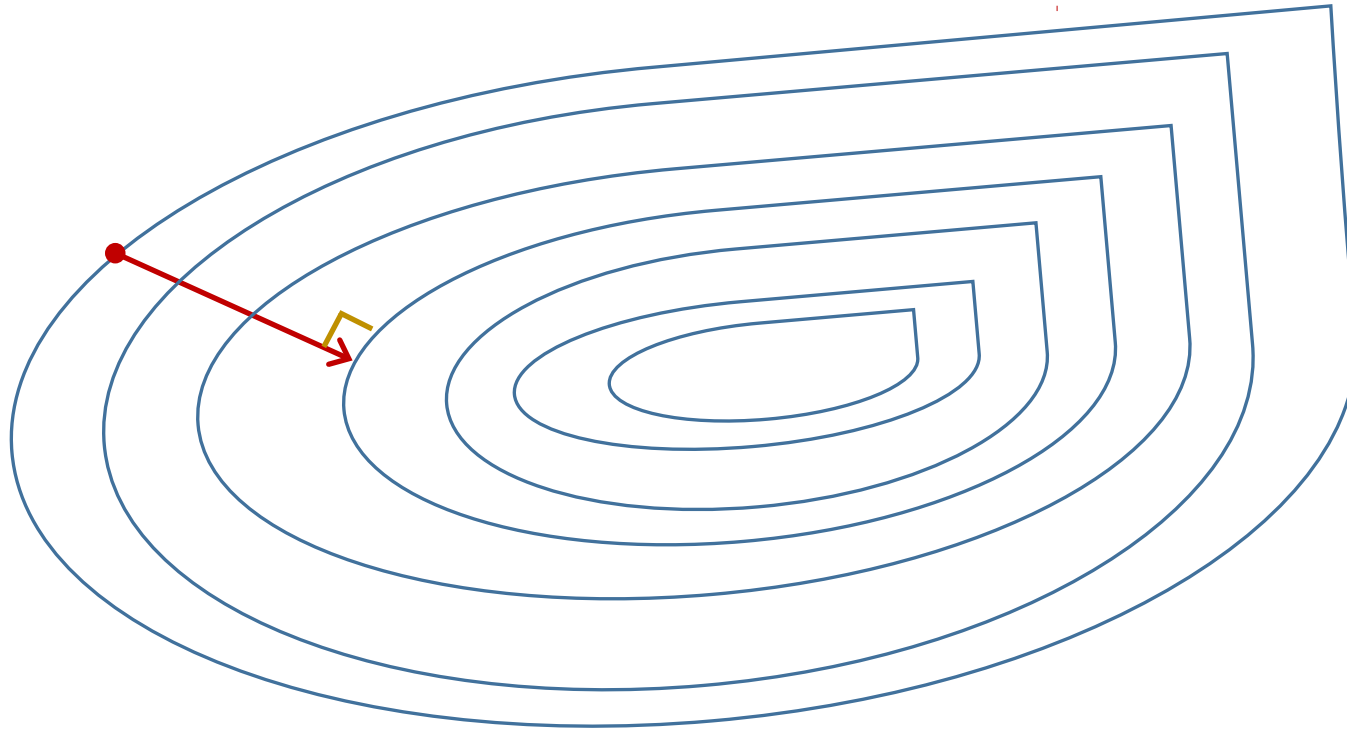
- To this formulation, they applied the augmented Lagrangian method (ALM)
 - A “hot” new method at the time
- Standard augmented Lagrangian method (ALM) for this formulation:

$$\left(x^{k+1}, z^{k+1} \right) \in \underset{x \in \mathbb{R}^n, z \in \mathbb{R}^m}{\text{Arg min}} \left\{ f(x) + g(z) + \langle p^k, Mx - z \rangle + \frac{c_k}{2} \|Mx - z\|^2 \right\}$$
$$p^{k+1} = p^k + c_k (Mx^{k+1} - z^{k+1})$$

- Although the equality constraints are gone, the cross terms in $\|Mx - z\|^2$ make the augmented Lagrangian harder to optimize than the ordinary Lagrangian (without $\|Mx - z\|^2$)
 - Cannot handle f and g independently
- But the ALM is much more stable than minimizing the ordinary Lagrangian & multiplier update (subgradient in the dual)

Interpretation of the ALM

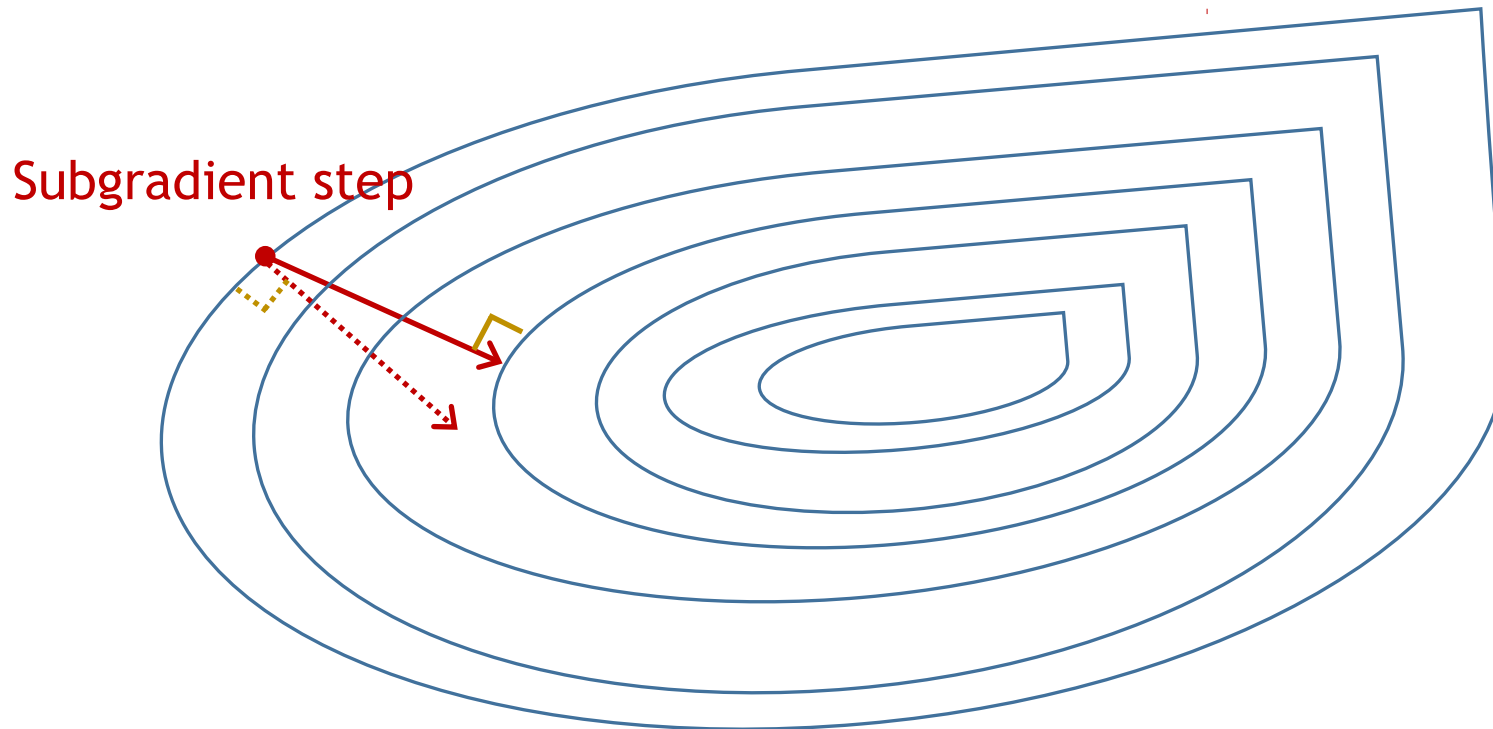
- The ALM method does an **implicit** subgradient step on the dual problem (as shown by Rockafellar; a form of “dual ascent”)



- The step direction is a subgradient of the function at the **end** of the step, not the beginning

Interpretation of the ALM

- The ALM method does an **implicit** subgradient step on the dual problem (as shown by Rockafellar; a form of “dual ascent”)



- The step direction is a subgradient of the function at the **end** of the step, not the beginning
- Much more stable, but at the cost of those cross terms

Alternating Directions

- To make the cross terms less painful, Glowinski and Marrocco (1976) suggested an alternating direction method for the subproblem:
 - Minimize over x with z held fixed
 - Then minimize over z with x held fixed
 - (There are other ways to try handle the cross terms)
- They suggested executing this loop a fixed number of times, then updating the multipliers
- This inner iteration was later shown to converge by Tseng (2001) under fairly loose assumptions
 - But not in a fixed number of steps
 - Unless you can show that you have (at least approximately) minimized the inner problem, the multiplier update is no longer dual ascent

The ADMM is Born

- Interestingly, Glowinski and Marrocco observed the best performance when making only one pass through x and z at every iteration - the ADMM:

$$\begin{aligned}x^{k+1} &\in \operatorname{Arg min}_{x \in \mathbb{R}^n} \left\{ f(x) + g(z^k) + \langle p^k, Mx - z^k \rangle + \frac{c}{2} \|Mx - z^k\|^2 \right\} \\z^{k+1} &\in \operatorname{Arg min}_{z \in \mathbb{R}^m} \left\{ f(x^{k+1}) + g(z) + \langle p^k, Mx^{k+1} - z \rangle + \frac{c}{2} \|Mx^{k+1} - z\|^2 \right\} \\p^{k+1} &= p^k + c(Mx^{k+1} - z^{k+1})\end{aligned}$$

Omitting constants from the minimands,

$$\begin{aligned}x^{k+1} &\in \operatorname{Arg min}_{x \in \mathbb{R}^n} \left\{ f(x) + \langle p^k, Mx \rangle + \frac{c}{2} \|Mx - z^k\|^2 \right\} \\z^{k+1} &\in \operatorname{Arg min}_{z \in \mathbb{R}^m} \left\{ g(z) - \langle p^k, z \rangle + \frac{c}{2} \|Mx^{k+1} - z\|^2 \right\} \\p^{k+1} &= p^k + c(Mx^{k+1} - z^{k+1})\end{aligned}$$

No Theory Yet

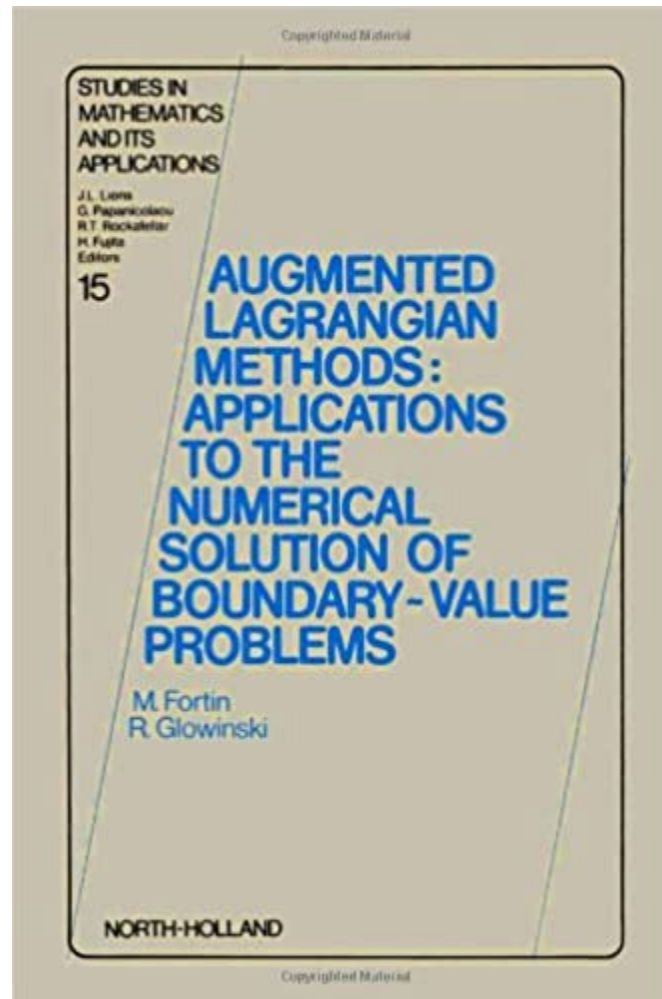
- But they did not have any theory to support this result
- The multiplier update is definitely not dual ascent now, because we are nowhere close to minimizing the augmented Lagrangian

The ADMM Does Not Approximate the ALM

- I have done experiments in which I use alternating minimization for the inner problem until some (rigorous) approximation criterion for the augmented Lagrangian is met, then update the multipliers
- This generally produces far fewer multiplier updates (usually about an order of magnitude)
- But many orders of magnitude more total inner iterations
- Alternating minimization is in general a poor algorithm for the inner problems
- So how to understand the convergence of the ADMM?

Four Years After Glowinski & Marrocco, Some Theory

- The following edited volume of papers appeared in 1983



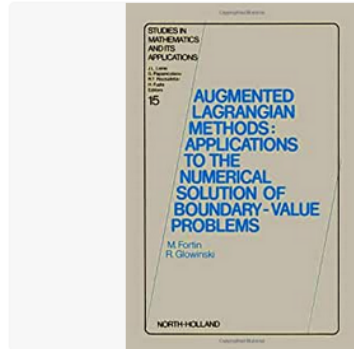
Aside: Amazon's Search Engine

- When I recently searched for this book on Amazon, I got

Aside: Amazon's Search Engine

- When I recently searched for this book on Amazon, I got

augmented lagrangian applications [See all 57 results](#)



Augmented Lagrangian Methods: Applications to the Numerical Solution of Boundary-Value Problems (Studies in Mathematics and Its Applications, V. 15) (English and French Edition)

French Edition | by [Michel Fortin](#), R. Glowinski, et al. | Jul 1, 1983

Hardcover

Best Seller



Paulas Choice--SKIN PERFECTING 2% BHA Liquid Salicylic Acid Exfoliant--Facial Exfoliant for Blackheads, Enlarged Pores, Wrinkles & Fine Lines, 4 oz Bottle

4 Fl Oz (Pack of 1)

★★★★☆ [60,366](#)

\$32⁰⁰ (\$8.00/Fl Oz)

\$30.40 with Subscribe & Save discount

[prime](#) FREE One-Day

Get it Tomorrow, May 17

[Climate Pledge Friendly](#)

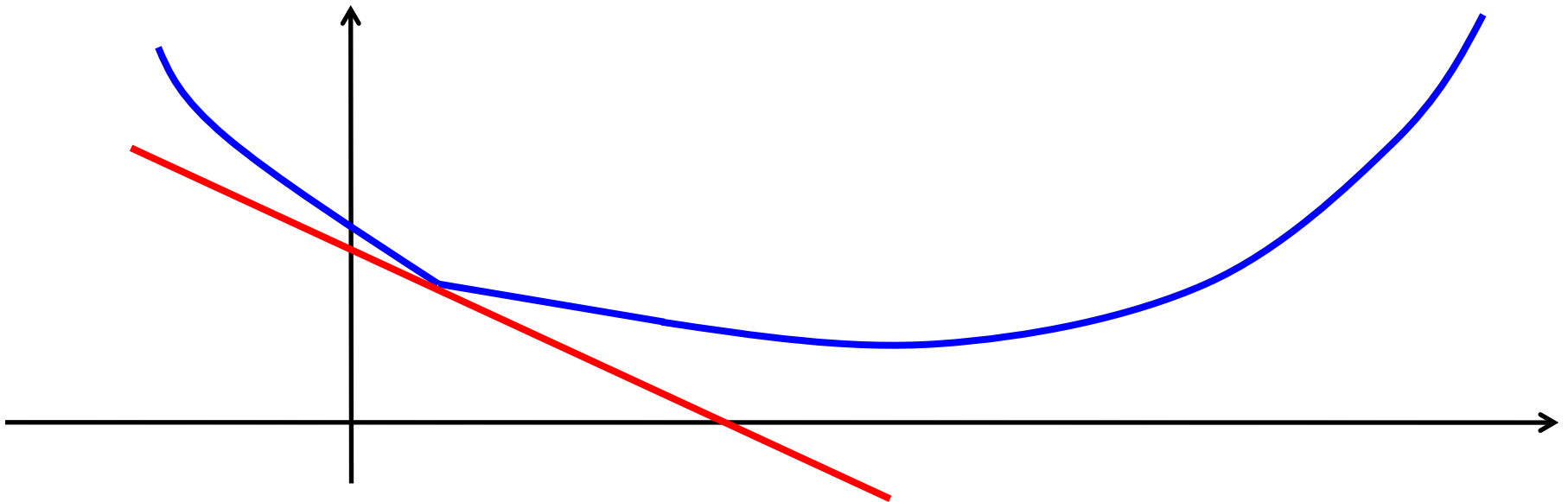
Two Proofs of Convergence in this Book

- In Fortin & Glowinski (1983), a convergence proof using a variational inequality analysis
- In Gabay (1983), a proof showing that the ADMM is an **operator splitting** method
 - The “Douglas-Rachford” splitting method for monotone (set-valued) operators analyzed by Lions and Mercier in 1979
 - Applied to the dual of $\min f(x) + g(Mx)$
 - Operator splitting methods also have their roots in the PDE world - so relatively natural for these researchers to have this insight
- I will follow the Gabay path since it is more intuitive
- The relationship between the two proofs could still use clarification
 - They lead to different forms of over-relaxation

Subgradients of a Convex Function

- Suppose that $d : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is a convex function
- d may not be smooth, but it has **subgradients**
- $\partial d(x)$ denotes the set of subgradients of d at x :

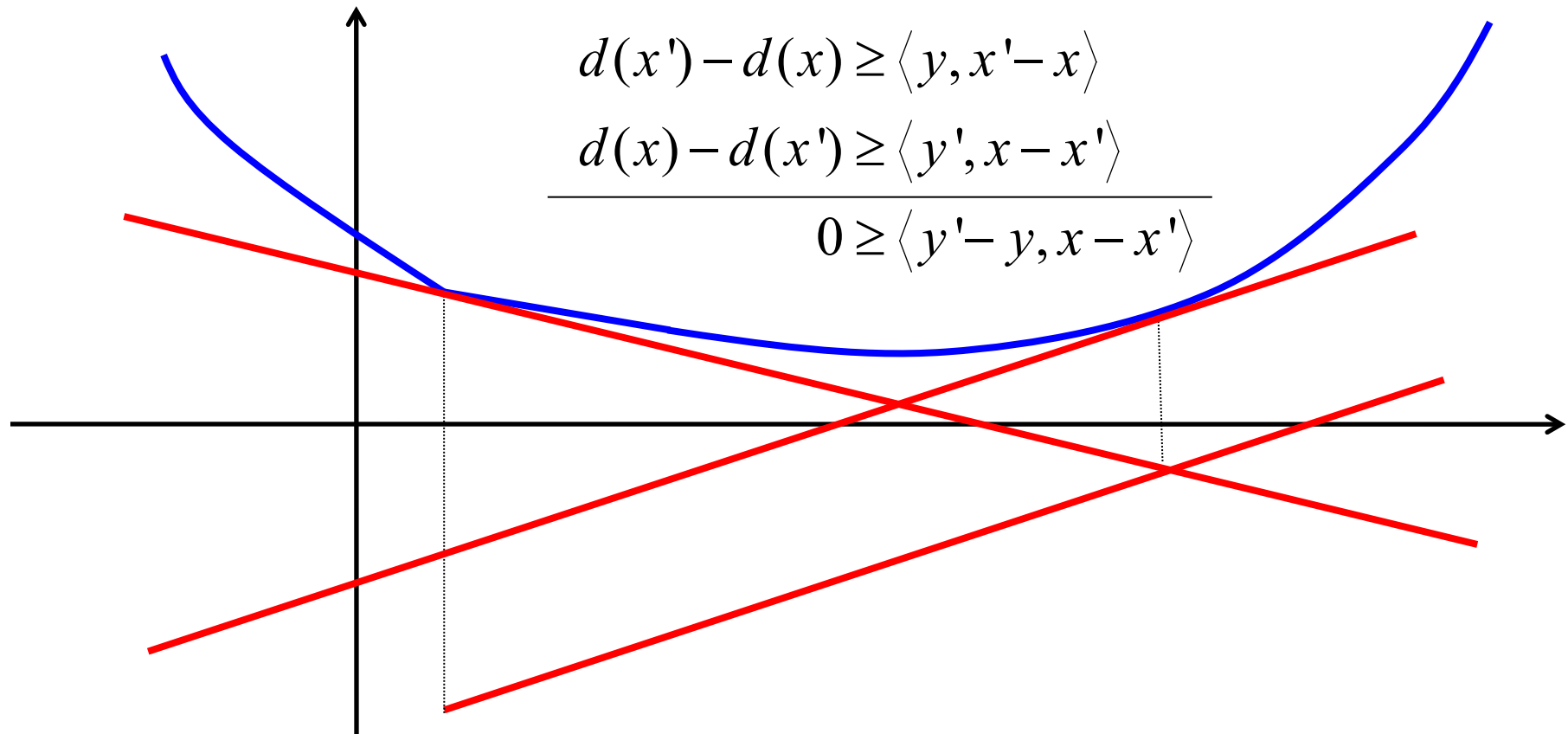
$$\partial d(x) = \left\{ y \mid d(x') \geq d(x) + \langle y, x' - x \rangle \quad \forall x' \in \mathbb{R}^m \right\}$$



Monotonicity

- Subgradient maps of convex functions are **monotone**

$$\boxed{y \in \partial d(x), y' \in \partial d(x') \Rightarrow \langle x - x', y - y' \rangle \geq 0}$$



- This condition is a natural generalization to higher dimension of a function being monotone nondecreasing

The Dual of $\min f(x) + g(Mx)$

- The dual of $\min f(x) + g(Mx)$ can be written in the form $\min_p d_1(p) + d_2(p)$, for two convex functions d_1 and d_2
- Namely, $d_1(p) = f^*(-M^\top p)$ and $d_2(p) = g^*(p)$

$$\begin{aligned}d(p) &= \min_{x \in \mathbb{R}^n, z \in \mathbb{R}^m} \{L(x, z, p)\} \\&= \min_{x \in \mathbb{R}^n, z \in \mathbb{R}^m} \{f(x) + g(z) + \langle p, Mx - z \rangle\} \\&= \min_{x \in \mathbb{R}^n} \{f(x) + \langle M^\top p, x \rangle\} + \min_{z \in \mathbb{R}^m} \{g(z) - \langle p, z \rangle\} \\&= \underbrace{f^*(-M^\top p)}_{d_1(p)} + \underbrace{g^*(p)}_{d_2(p)} \\&= d_1(p) + d_2(p)\end{aligned}$$

Splitting the Dual

- This is the same as solving

$$0 \in \partial(d_1 + d_2)(p)$$

- Unless things are really ugly (details omitted), the same as solving

$$0 \in \partial d_1(p) + \partial d_2(p)$$

where $+$ denotes the Minkowski sum of sets

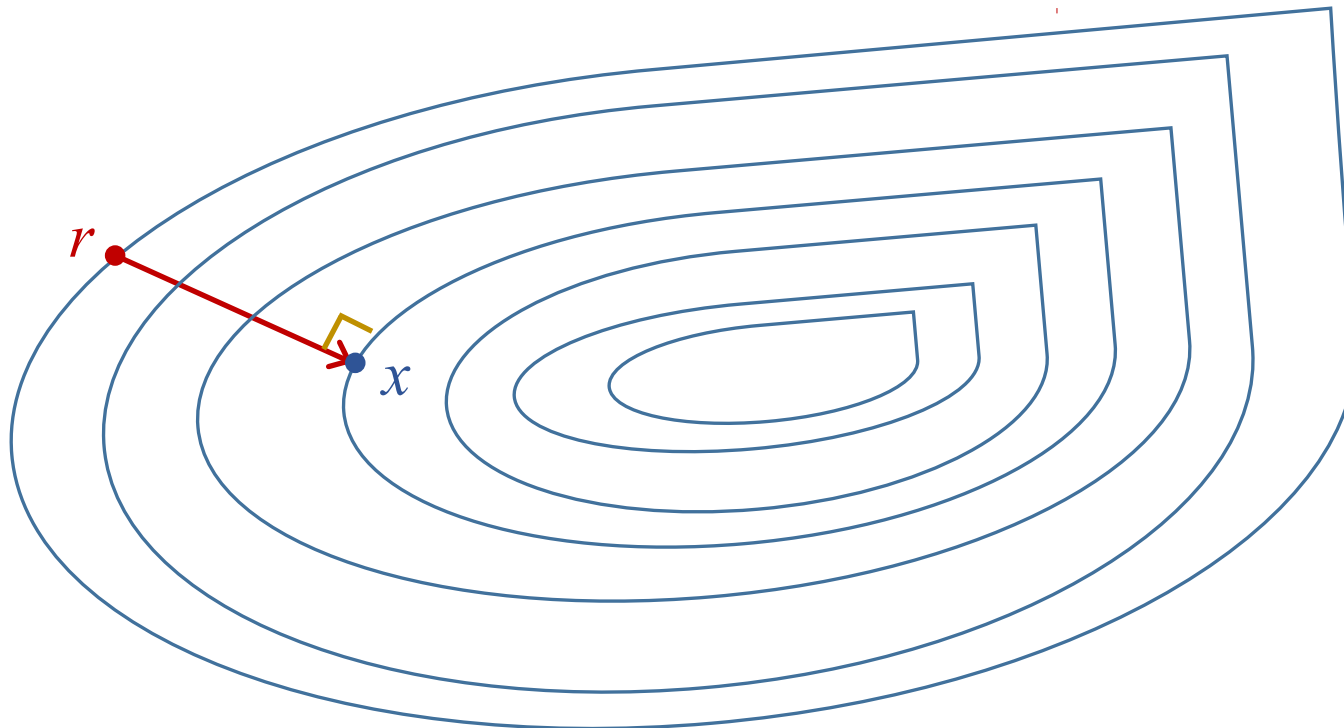
$$A + B = \{a + b \mid a \in A, b \in B\}$$

Resolvents

- Suppose that T is any point-to-set map on \mathbb{R}^n that is monotone: $y \in T(x), y' \in T(x') \Rightarrow \langle x - x', y - y' \rangle \geq 0$
- Consider any fixed scalar $c > 0$
- Then the **resolvent of T with stepsize c** is $J_{cT} = (I + cT)^{-1}$
- The same operation as an implicit step in ODE/PDE integration
- Conceptually, to evaluate $J_{cT}(r)$:
 - Find x, y such that $x + cy = r$ and $y \in T(x)$
(can only be done one way if T is monotone)
 - Return x

Interpretation of the Resolvent for Convex Functions

- When T is the subgradient map ∂f of a convex function f , ...
- The ALM method does an **implicit** (sub)gradient step on the function (as shown before for the dual function)



- A step of size c in the “ $-\nabla f(x)$ ” direction, ending at x
- Iterating such steps is called the **proximal point algorithm** (Rockafellar 1976)

Resolvents and “Reflectants”

- If T is monotone, then the resolvent $J_{cT} = (I + cT)^{-1}$ is defined everywhere, single valued, and **firmly nonexpansive**

$$(\forall x, x') \quad \|J(x) - J(x')\|^2 \leq \|x - x'\|^2 - \|(x - J(x)) - (x' - J(x'))\|^2$$

- And the “**reflectant**” $R_{cT} = 2J_{cT} - I$ is defined everywhere, single valued, and **nonexpansive**

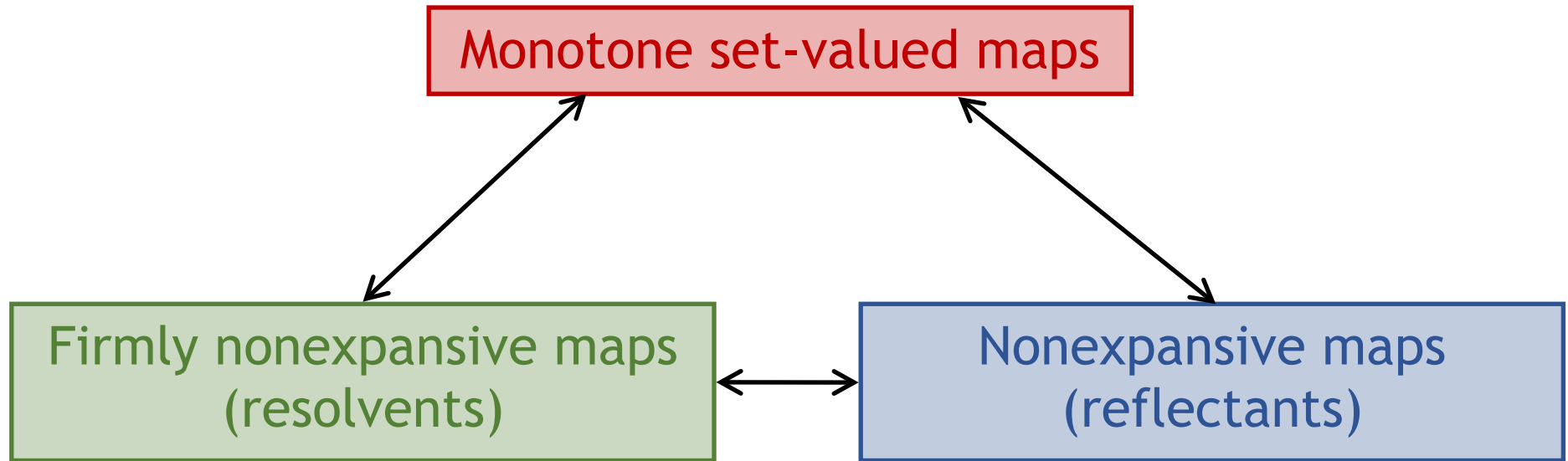
$$(\forall x, x') \quad \|R(x) - R(x')\|^2 \leq \|x - x'\|^2$$

- Conceptually, evaluating the reflectant amounts to:
 - Find x, y such that $x + cy = r$ and $y \in T(x)$
(can only be done one way if T is monotone)
 - Return $2x - r = 2x - (x + cy) = x - cy$

$$R_{cT} = 2J_{cT} - I \quad \Leftrightarrow \quad J_{cT} = \frac{1}{2}R_{cT} + \frac{1}{2}I$$

Isomorphic Relationships

- The relationships between monotone operators, resolvents, and reflectants are all isomorphisms



- A map is firmly nonexpansive if and only if it is the resolvent of some monotone operator
- A map is nonexpansive if and only if it is the reflectant of some monotone operator
- A map J is firmly nonexpansive if and only if it is of the form $J = \frac{1}{2}R + \frac{1}{2}I$, where R is nonexpansive

Convergence of the ADMM I

Fundamentally, the convergence theory of the ADMM relies on a very simple observation:

The composition of two nonexpansive mappings is nonexpansive

- Nonexpansive map $R_{c\partial d_1} : \mathbb{R}^m \rightarrow \mathbb{R}^m$ corresponding to d_1
- Nonexpansive map $R_{c\partial d_2} : \mathbb{R}^m \rightarrow \mathbb{R}^m$ corresponding to d_2
- Their composition $R_{c\partial d_1} \circ R_{c\partial d_2}$ is nonexpansive

Furthermore, the fixed points of $R_{c\partial d_1} \circ R_{c\partial d_2}$ are of the form

$$\{p + cz \mid z \in \partial d_2(p), -z \in \partial d_1(p)\}$$

Sketch of proof. $p + cz \xrightarrow{N_{c\partial d_2}} p - cz \xrightarrow{N_{c\partial d_1}} p + cz$. Easy to show this is the only possibility (two equations in two unknowns). \square

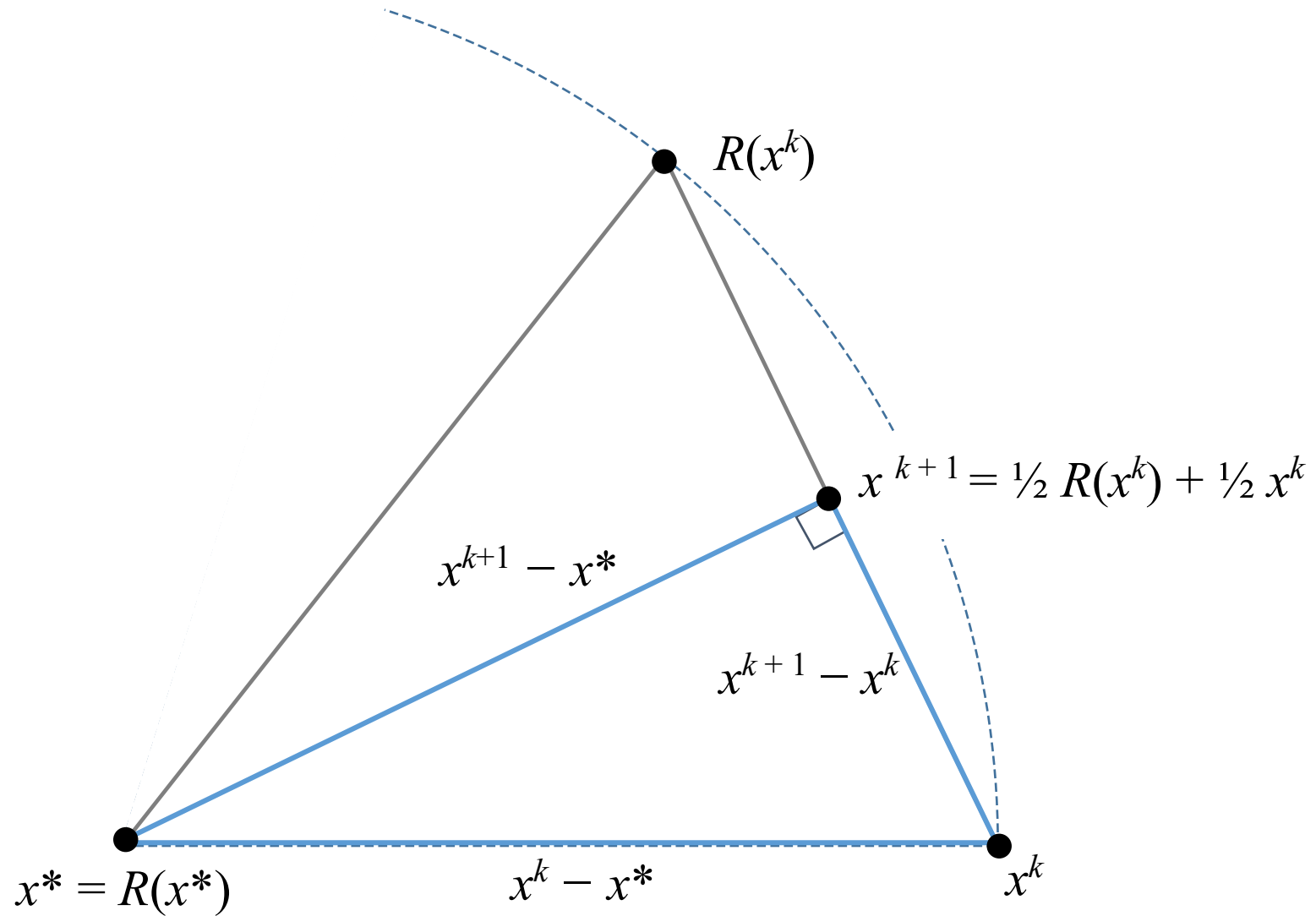
Convergence of the ADMM II

- From a fixed point t of $R_{c\hat{d}_1} \circ R_{c\hat{d}_2}$, we can find the optimal dual solution by computing $J_{c\hat{d}_2}(t)$, and we can also easily find the primal solution
- It would be nice to just iterate the map $R_{c\hat{d}_1} \circ R_{c\hat{d}_2}$ to converge to a fixed point, but since its Lipschitz constant is 1, this process might just “orbit” around the set of fixed points
- **But** (Krasnosel'skii 1955), if we blend it with the identity, it **will** converge

$$s^{k+1} = \frac{1}{2}s^k + \frac{1}{2}\left(R_{c\hat{d}_1} \circ R_{c\hat{d}_2}\right)(s^k) = \frac{1}{2}s^k + \frac{1}{2}R_{c\hat{d}_1}\left(R_{c\hat{d}_2}(s^k)\right)$$

- This is the essence of “Douglas-Rachford splitting”
- Converts the nonexpansive map $R_{c\hat{d}_1} \circ R_{c\hat{d}_2}$ to a firmly nonexpansive one (with the same fixed points)

Picture of Krasnosel'skii



What is Meant by “Operator Splitting”?

- Douglas-Rachford splitting is a kind of **operator splitting**
- We are solving the problem $\min_p d_1(p) + d_2(p)$
- Or (usually) equivalently $0 \in \partial d_1(p) + \partial d_2(p)$
- But we only deal with the individual reflectant maps $R_{c\partial d_1}$ and $R_{c\partial d_2}$ respectively associated with d_1 and d_2
- That’s the essence of “operator splitting”
 - “Forward” (gradient-like) steps are also allowed on individual functions/operators, in addition to resolvents/reflectants

Getting the History Right

- Douglas and Rachford had a different representation of the operations in their method, but equivalent
- However, the original Douglas-Rachford publication was only for linear operators
 - Applied to very specific linear operators related to the discretized 2-D heat equation
- Lions and Mercier (1979) generalized the idea from linear maps to general monotone set-valued maps (but kept the name)
- Gabay (1983) showed that the ADMM is just this idea applied to the dual of $\min f(x) + g(Mx)$
- The composition-of-nonexpansive maps interpretation may first be found (as an aside) in Lawrence & Spingarn (1987)
- E & Bertsekas (1992) contains some equivalent analysis and exploits the relationship with the proximal point algorithm to derive approximate and over-relaxed versions

Some Insights from the Convergence Analysis

- ADMM convergence based on evaluation of $R_{c\hat{d}_1} \circ R_{c\hat{d}_2}$, which is **not** an approximation of $R_{c(\hat{d}_1+\hat{d}_2)}$ (the mapping for ALM)
- Unlike the ALM, changing c between iterations in the ADMM is problematic because it shifts the set of fixed points $\{p + cv \mid v \in \hat{d}_2(p), -v \in \hat{d}_1(p)\}$ of $R_{c\hat{d}_1} \circ R_{c\hat{d}_2}$
 - There are results for variable c , but they need extra assumptions and get “messy”
- Also, the $R_{c\hat{d}_1} \circ R_{c\hat{d}_2}$ convergence theory of the ADMM does not have a “clean” extension to more than two blocks:

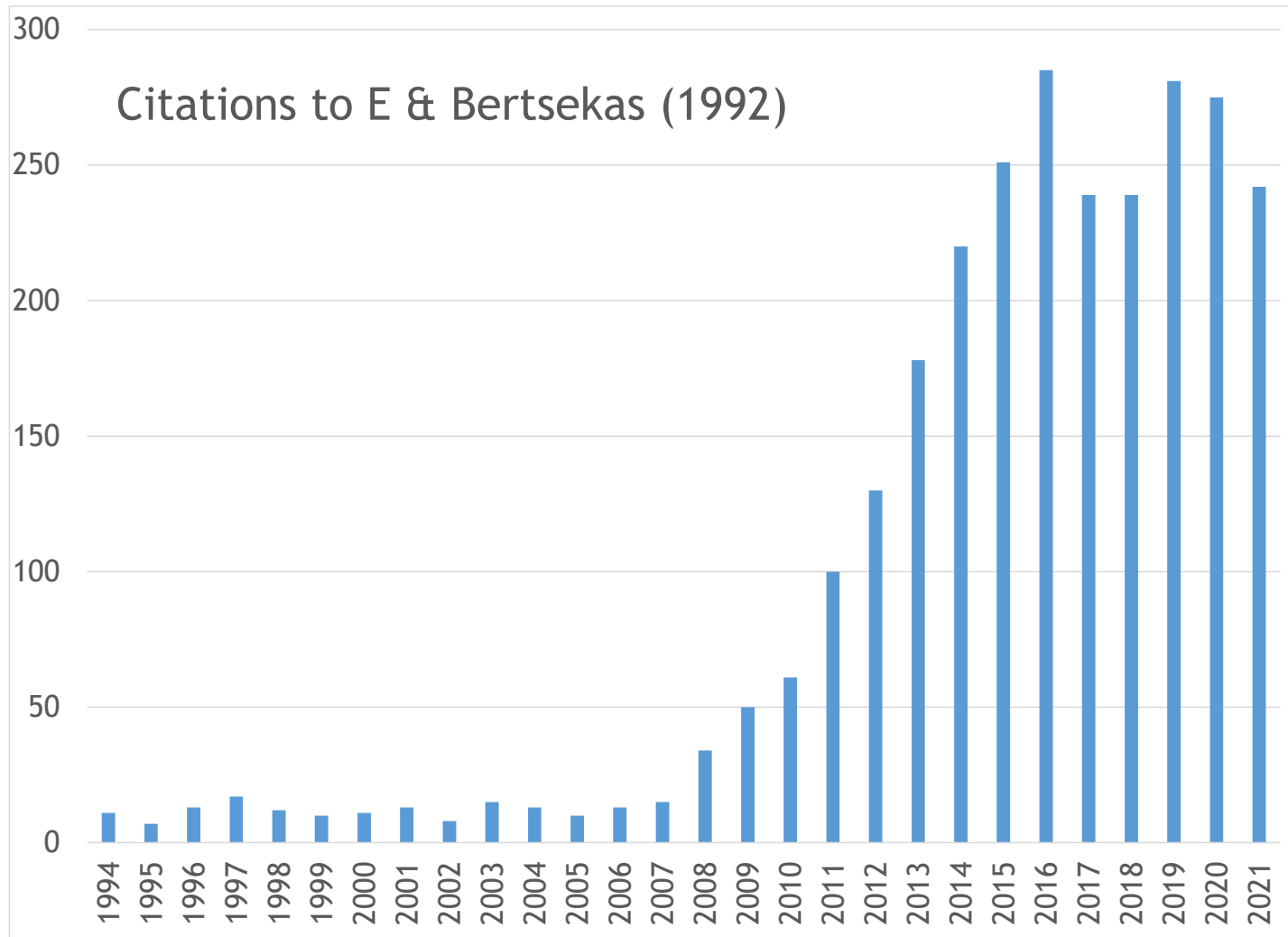
$$\min_{x \in \mathbb{R}^n} \left\{ f_1(M_1 x_1) + f_2(M_2 x_2) + \dots + f_p(M_p x_p) \right\}$$

$R_{c\hat{d}_1} \circ R_{c\hat{d}_2} \circ \dots \circ R_{c\hat{d}_p}$ does not have “nice” fixed points

- Must use a product-space reformulation, or make things “messier”

The ADMM: More Past

- The basic theory of the ADMM was settled by the early 90's
- But it remained an obscure algorithm for 15+ years
- During the period 2008-2014, things changed



The ADMM: Present

- Now the ADMM is considered part of the standard optimization “toolbox”

Some typical current applications:

- Image denoising
- Data fitting / machine learning
 - Along with other operator-splitting methods, like forward-backward
- Stochastic programming (progressive hedging)
- ...
- Even some general conic QP solvers

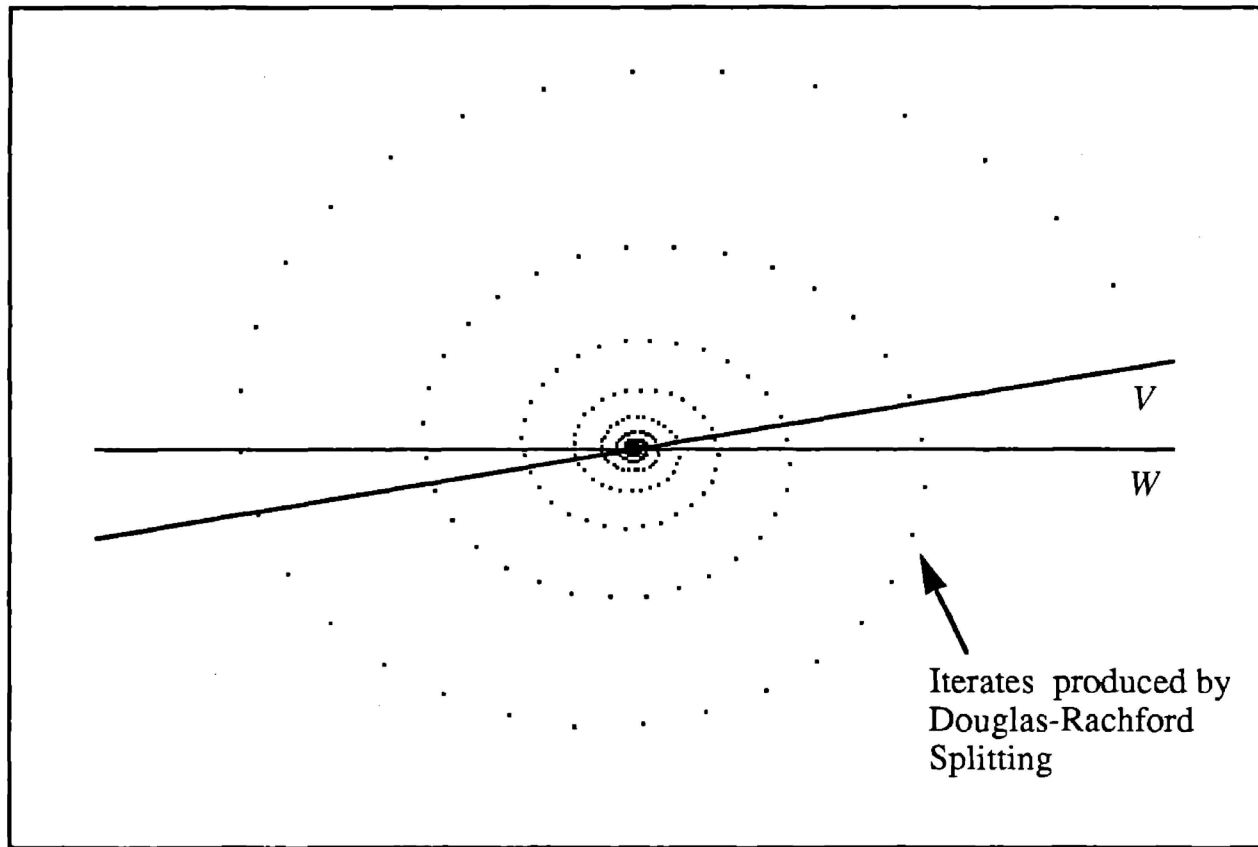
- Also, a dizzying profusion of new variants (not covered much here)

Features of Successful Applications

1. Low accuracy solutions sufficient
(fairly common knowledge - but there is work on trying to change this!)
2. Nonlinear convex objectives can often work better than linear ones
3. Do not try to “atomize” problems
(although that can be tempting for parallelism)

Low Accuracy Requirement

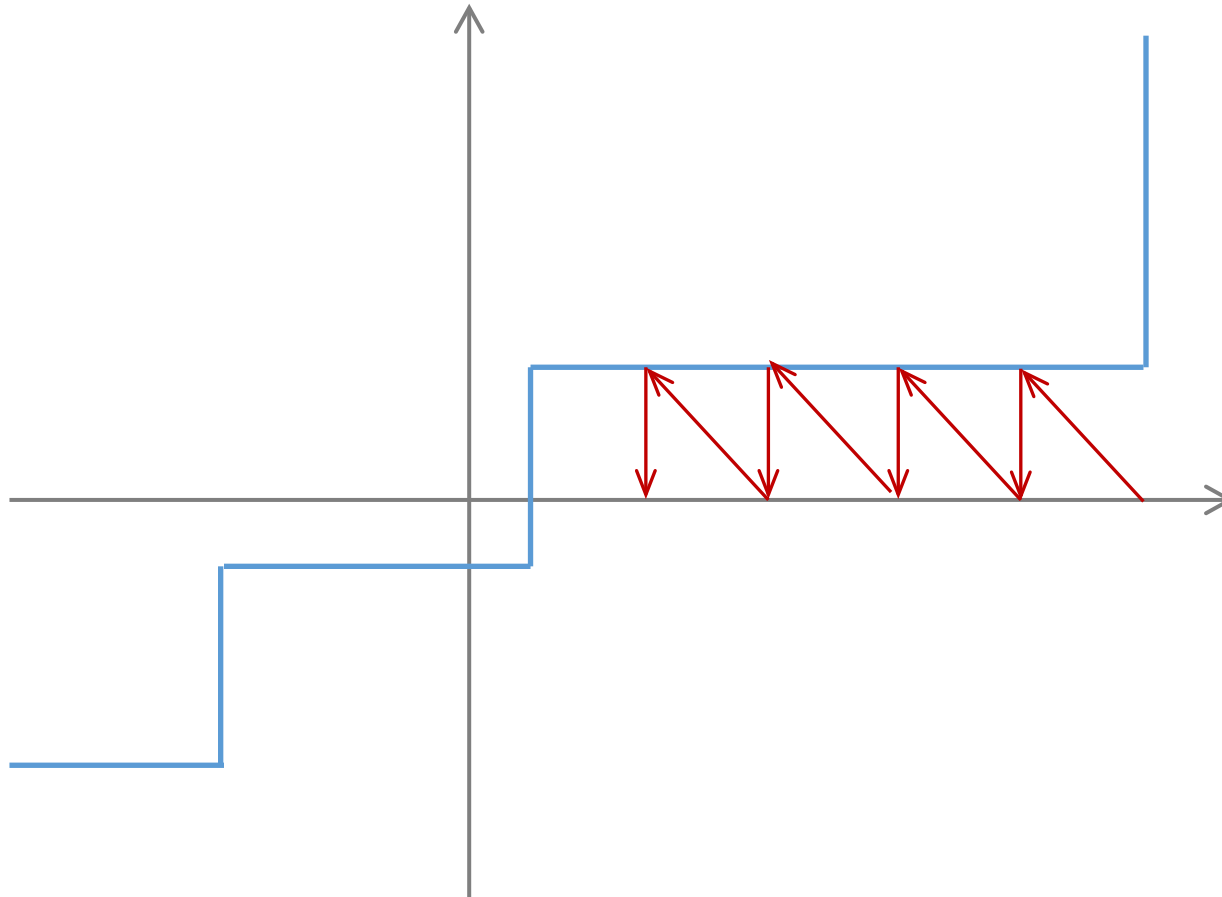
- The ADMM does not have fast asymptotic / tail convergence
- It is typically linear/geometric, but the constant can be poor



- However, applications like machine learning and image denoising typically don't require high-accuracy solutions

Nonlinear Objectives (Intuition)

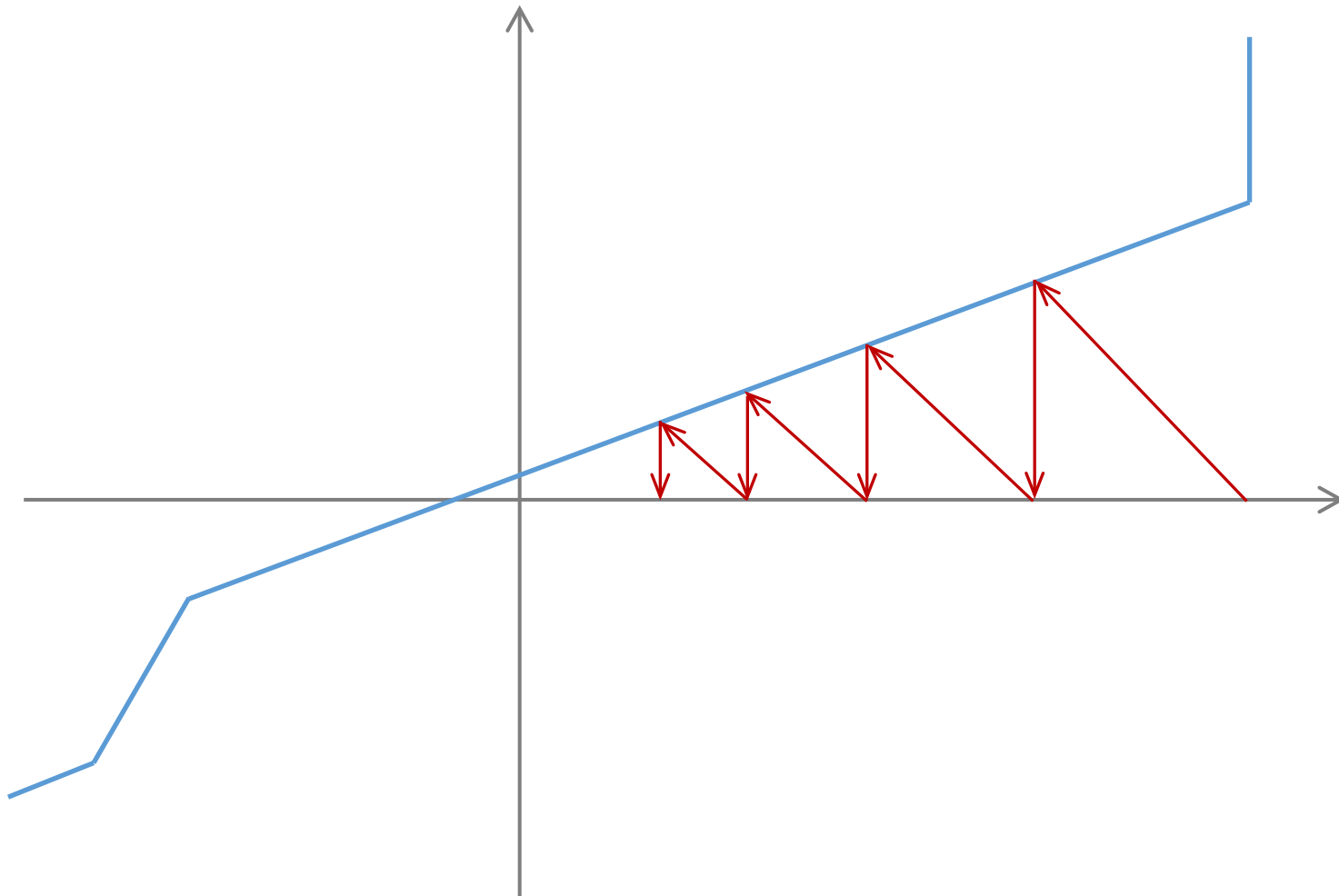
- The ADMM is a form of the proximal point algorithm (PPA) (shown for example in E & Bertsekas 1992)
- Solving an LP-like problem with the PPA



- Convergence can be sublinear

Nonlinear Objectives (Intuition)

- Solving a QP-like problem with the PPA



- Convergence is linear
(but might be slow, depending on the slopes)

Don't Atomize Problems

In successful ADMM applications...

- At least one side of the splitting (f or g) should model a substantial portion of the global interconnections between problem elements
- My former postdoc Patrick Johnstone calls this property “being meaty”
- What does that mean?

Don't Atomize Problems II

- Example: in data fitting / ML problems, one often has a structure like

$$\min_x \ell(x) + r(x)$$

- ℓ is a smooth loss function
- r is a regularizer (for example, an L_1 penalty)
- All the connections between the model parameters x and fitting the observations are contained in ℓ
- So, ℓ is “meaty” (it contains essentially all the connections between model elements) and one can set $f = \ell, g = r, M = I$
- r might not be “meaty” (usually separable by variables or blocks), but having just one kind of meat on the plate is enough

Don't Atomize Problems III

In OSQP (Stellato *et al.* 2020)

- f models all the linear relationships within the model
- g contains only conic constraints on individual subvectors
- So, f is “meaty”

In E & Ferris 1998 for optimal control problems

- f enforces a block-tridiagonal linear system capturing all the time dynamics in the model
- g enforces all the inequalities and nonsmooth elements (confined within each time step)
- So, f is “meaty” (but g is not)

Don't Atomize Models IV

In progressive hedging for stochastic programming problems (originally Rockafellar and Wets 1991)

- f contains the entire time dynamics within each scenario
- g enforces “nonanticipativity” (not seeing the future) relationships between scenarios
- So, they are both fairly “meaty”

A Tempting Example of “Non-Meatiness”: The Classic “Transportation” Problem

Given a bipartite graph (S, D, E) ,

$$\begin{array}{ll} \min & r^\top x \\ \text{ST} & \sum_{j:(i,j) \in E} x_{ij} = s_i \quad \forall i \in S \\ & \sum_{i:(i,j) \in E} x_{ij} = d_j \quad \forall i \in D \\ & x_{ij} \geq 0 \quad \forall (i, j) \in E \end{array}$$

- Illustration of how the ADMM can lead to highly parallel algorithms
- But ones that are typically not competitively efficient

Modeling Transportation in the ADMM Form

With $x, z \in \mathbb{R}^{|E|}$,

$$f(x) = \begin{cases} \frac{1}{2} r^\top x, & \text{if } x \geq 0 \text{ and } \sum_{j:(i,j) \in E} x_{ij} = s_i \quad \forall i \in S \\ +\infty & \text{otherwise} \end{cases}$$

$$g(z) = \begin{cases} \frac{1}{2} r^\top z, & \text{if } z \geq 0 \text{ and } \sum_{i:(i,j) \in E} z_{ij} = d_j \quad \forall i \in D \\ +\infty & \text{otherwise} \end{cases}$$

Then the problem reduces to

$$\begin{array}{ll} \min & f(x) + g(z) \\ \text{ST} & x - z = 0 \end{array}$$

- The x -minimization step separates by source node $i \in S$
- The z -minimization step separates by destination node $j \in D$

ADMM for Transportation

For example, the x minimization reduces to, for each $i \in S$

$$\begin{aligned} \min \quad & \sum_{j:(i,j) \in E} r_{ij} x_{ij} + \sum_{j:(i,j) \in E} p_{ij} x_{ij} + \frac{c}{2} \sum_{j:(i,j) \in E} (x_{ij} - z_{ij})^2 \\ \text{ST} \quad & \sum_{j:(i,j) \in E} x_{ij} = s_i \\ & x_{ij} \geq 0 \quad \forall j : (i, j) \in E \end{aligned}$$

This is just projection on a simplex, so it's an easy problem

- A simple implementation is $O(\delta \log \delta)$, where δ is the node degree
- Can be done in $O(\delta)$ time if one is careful (related to linear-time median finding; only matters for large δ)
- The z minimization step is similar

ADMM for Transportation – (Parallel) Implementation

High potential for parallelism:

- x minimization consists of $|S|$ independent, easy tasks
- z minimization consists of $|D|$ independent, easy tasks
- Multiplier update consists of $|E|$ independent, easy tasks
- Straightforward communication pattern between these tasks
- An example of how ADMM can lead to highly parallel algorithms
- Studied these kinds of applications in my dissertation
- Unfortunately, this ADMM application is very slow compared to network simplex etc.
- And parallelism is not enough to save it

Slowness Intuition and the Moral of the Story

- Both sides of the splitting decompose into optimizations that only “see” individual nodes
- The whole “big picture” is left to the ADMM to coordinate
- But ADMM / DR is not an outstanding linear equation solver
- So it takes a long time for all the pieces of the problem to come into alignment

- ADMM is a useful algorithm
- But don't ask too much of it

- Don't leave the entire coordination of small problem elements to the ADMM
- Keep at least some of the global connections within f or g

Approximate Iterations

- With “meatiness” of subproblems comes the potential need to solve them inexactly
- E & Bertsekas 1992 was the first publication to rigorously cover inexact solution of subproblems in the ADMM
 - Based on earlier results by Rockafellar for the proximal point algorithm
- But requires bounding the distance between the approximate iterate and the exact one
- In general, finding such a bound can be difficult
- I have repeatedly seen people using only the distance-based 1992 approximation result, not realizing that there is more recent work on the subject

Better Approximation Criteria

Two papers:

- E & Yao 2017 (COAP)
- E & Yao 2018 (*Math Programming A*) (simpler)
- Neither require estimating the distance to the exact subproblem solution - just a subgradient
- So, can be easier to implement than the 1992 criterion
- **Absolute** error criteria involve a summable sequence of allowable errors that is (formally) an external parameter
- **Relative** error criteria use a single parameter to compare two quantities generated by the algorithm
- The above two papers contain both kinds
- If you are using inexact ADMM, please look at these papers!

The ADMM: Future

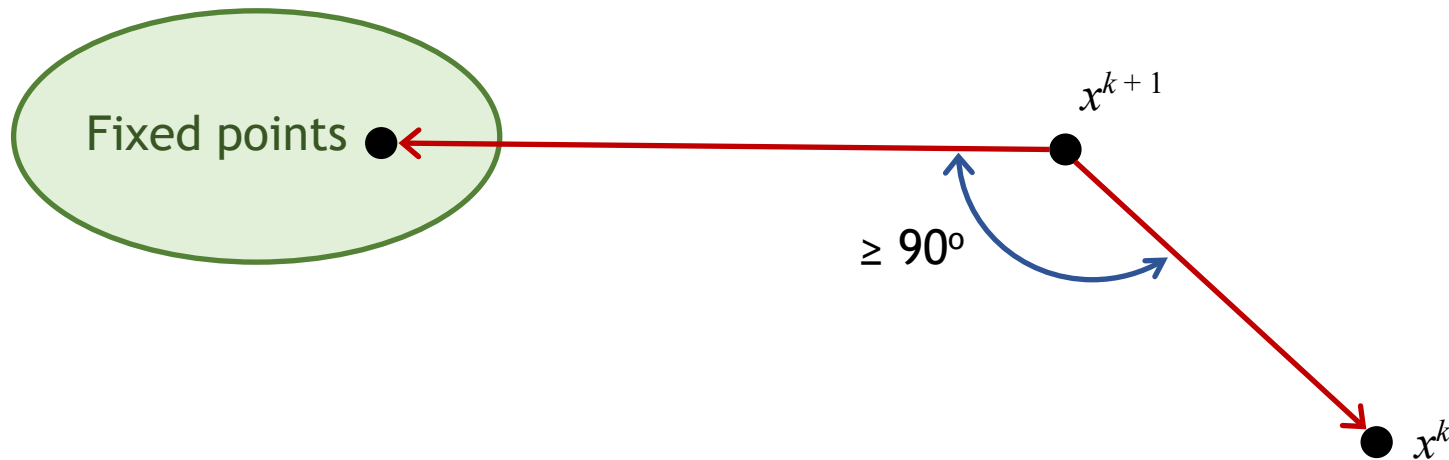
- What can be the directions for the future?
 1. Other operator splitting methods might make an impact
 2. Upper and lower bounds
 3. A wider range of applications, **if** we can solve the tail convergence issues

Other Operator Splitting Methods

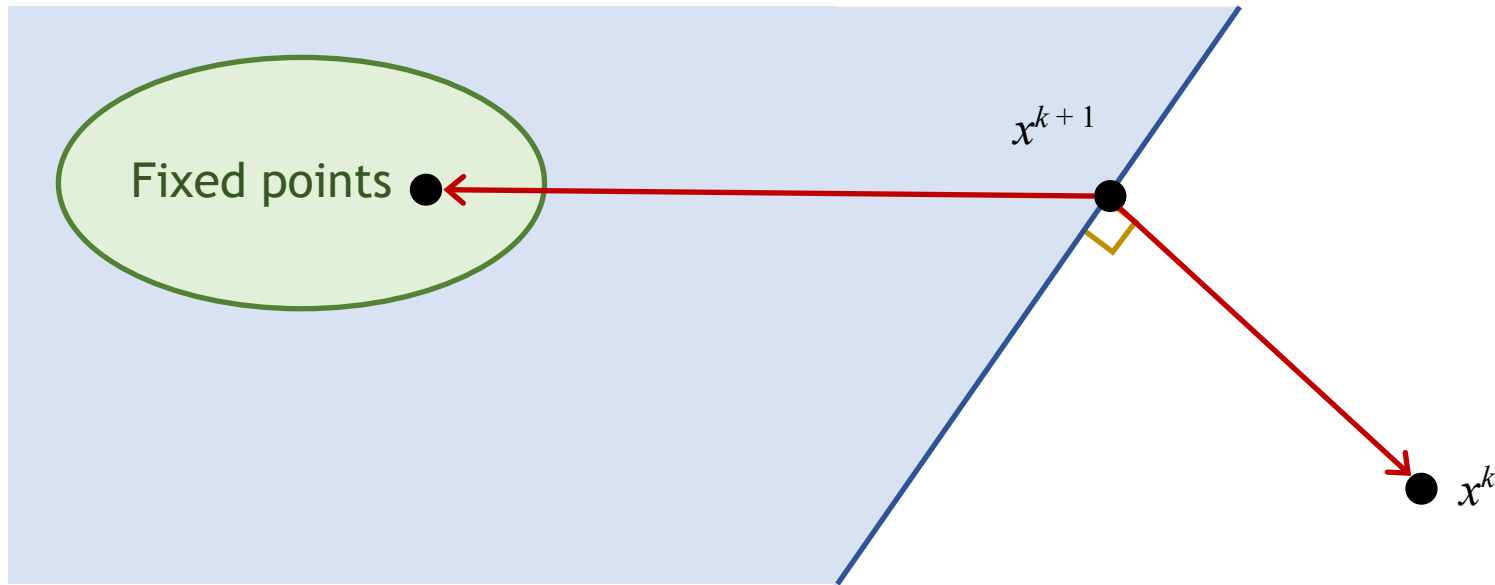
- For a long while, there were basically three classes of operator splitting method
 - Forward-backward (generalizes gradient projection)
 - Douglas-Rachford (as in the ADMM), which forms a continuous family with Peaceman-Rachford
 - Double-backward: simple to write, but solves $\min_p d_1(p) + \frac{c}{2} \|p - q\|^2 + d_2(q)$ instead of $\min_p d_1(p) + d_2(p)$
- Once the math programming and machine learning communities got interested in operator splitting, new varieties started appearing
 - Forward-backward-forward (Tseng 200)
 - Projective splitting (starting with E & Svaiter 2008)
 - Forward-reflected-backward (Malitsky & Tam 2020) ...

Projection and Nonexpansive Maps

- Iterating any firmly nonexpansive map is equivalent to:



- In turn equivalent to projecting onto a separating hyperplane:



Projection and Operator Splitting Methods

- This observation enables new operator splitting methods and variations on existing ones
- First example was Solodov and Svaiter's relative-error proximal point algorithm (1999)
- “Projective splitting” idea in E & Svaiter (2008): use separate operations on f and g (or d_1 and d_2) to construct a separator, then project onto it
 - Many variations since then
 - Also showed that Douglas-Rachford (and thus the ADMM) are a “boundary point” of a family of methods like this

Other Splitting Methods

- The picture is not clear yet
- There are so many methods and variations now
- And so many problems to apply them too
- But there are specific cases in which new operator splitting methods can outperform the ADMM

An Example

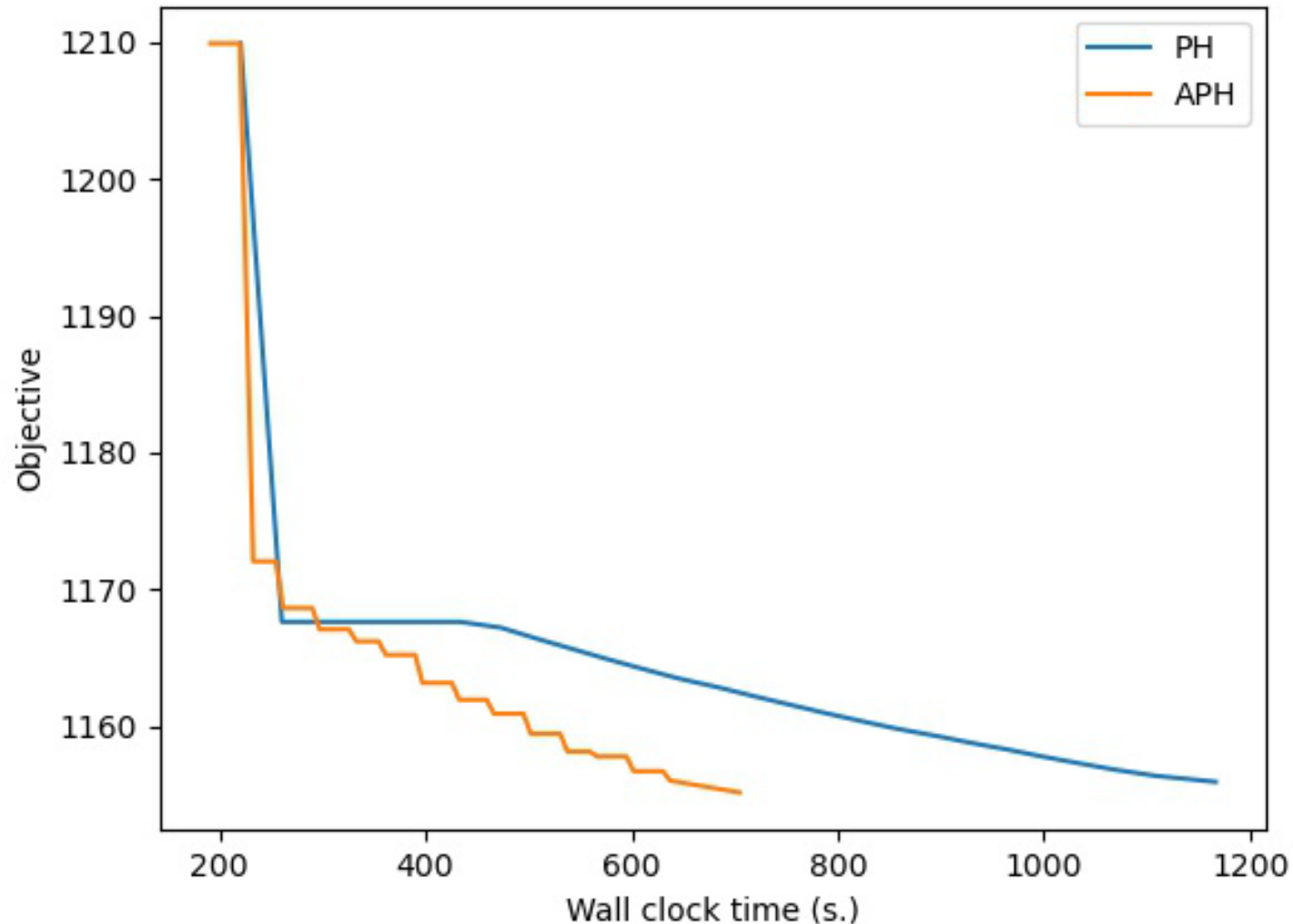
Progressive hedging (PH) for stochastic programming:

- Scenarios form a (possibly large) finite tree
- The x minimization separately optimizes a quadratic perturbation of each scenario
- The z minimization and multiplier update try to make your overall strategy non-clairvoyant

“Projective Hedging”

- Apply an “asynchronous” projective splitting variant (Combettes and E 2018, E 2017) to the same problem
- Obtain a similar algorithm to PH, except
 - You don’t have to optimize every scenario at every iteration
 - You can just optimize a subset
 - Called **asynchronous projective hedging** (APH)
- Generally, if you don’t re-optimize all the scenarios
 - The convergence slows down somewhat
 - But each iteration takes much less time
 - So overall time may be reduced, as in...

An Example Stochastic Programming Problem



- 1,000,000 scenarios! (And 5 stages)
- 600 processor cores
- Here, APH only solves 10% of the scenarios at each iteration

Upper and Lower Bounds

- ADMM, like many other operator-splitting methods...
- Converges asymptotically to solution (both primal and dual)
- But is typically never feasible: $x = z$ only in the limit
 - So in general there are no upper bounds
 - Although in ML problems everything is usually feasible, so upper bounds are easy
- And doesn't provide lower bounds
 - Since it never truly minimizes the (augmented) Lagrangian
- But if you're solving problems to low accuracy, you would often like to have upper and lower bounds
- Workarounds are generally application-specific
 - E 2020 gives a possible lower bound when all else fails...
- It would be nice to have a systematic approach

Tail Convergence

- Slow tail convergence is probably the biggest issue with ADMM-class methods
- There is a lot of recent work on “accelerated” versions
 - Nesterov/momentum (but just improves sublinear rates)
 - Anderson acceleration
 - Periodically averaging “blocks” of iterates (considerable work on this - currently “engineering-style”)
- Sometimes one can periodically test and try to “jump” to a basic solution (in simplex terms, or some generalization); see for example E & Ferris 1998
- But that’s not very satisfactory in general
- If we could speed up the tail, we could see a lot more applications

Thank you again!