

# Fast gradient algorithms for structured sparsity

Yao-Liang Yu

PhD: University of Alberta  
Now: Carnegie Mellon University

June 5, CAIAC 2015 @ Dalhousie



# Statistical inference 101

To estimate unknown parameter  $\theta \in \mathbb{R}^p$ :

$$y_i = \mathbf{x}_i^\top \theta + \epsilon_i, \quad i = 1, \dots, n$$

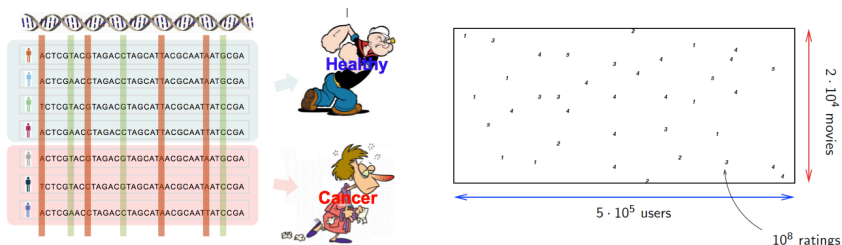
- classical setting:  $p$  fixed small,  $n \rightarrow \infty$ , lots of results.

# Statistical inference 101

To estimate unknown parameter  $\theta \in \mathbb{R}^p$ :

$$y_i = \mathbf{x}_i^\top \theta + \epsilon_i, \quad i = 1, \dots, n$$

- classical setting:  $p$  fixed small,  $n \rightarrow \infty$ , lots of results.
- modern setting:



$$p \sim 10^7, n \sim 10^3$$

$$p \sim 10^{10}, n \sim 10^8$$

# High-dimensional challenge

- More unknown parameters than observations, ill-defined.
  - ▶ **structure**: effective number of unknown parameters is moderate.
    - ★  $\theta$  is sparse:  $\text{nnz}(\theta)$  small, but do not know which is which.
    - ★  $\theta$  as a matrix is low-rank, but do not know the column/row spaces.
- Extremely large scale, takes forever to run.
  - ▶ **first order grad alg**: scales (sub)linearly with problem size.
- Ideally, want algorithm to exploit structure for faster convergence.
  - ▶ **open the blackbox**.



- ▶ contributions of this thesis lie in.

- 1 Introduction
- 2 Decomposing the Proximal Map
- 3 Approximation by the Proximal Average
- 4 Generalized Conditional Gradient
- 5 Post-PhD Extensions

# Table of Contents

- 1 Introduction
- 2 Decomposing the Proximal Map
- 3 Approximation by the Proximal Average
- 4 Generalized Conditional Gradient
- 5 Post-PhD Extensions

# Regularized loss minimization

Generic form for many ML problems:

$$\min_{\mathbf{w} \in \mathbb{R}^p} \ell(\mathbf{w}) + f(\mathbf{w}), \quad \text{where}$$

- $\ell$  is the loss/-likelihood function, usually smooth;
- $f$  is the regularizer, usually nondifferentiable;
  - ▶ **structure** inducing

Special interest:

- sparsity (structure);
- computational efficiency.

# The LASSO (Tibshirani'96)

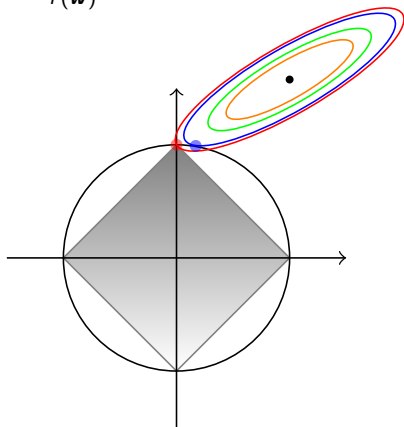
$$\min_{\mathbf{w} \in \mathbb{R}^p} \underbrace{\|A\mathbf{w} - \mathbf{b}\|^2}_{\ell(\mathbf{w})} + \lambda \underbrace{\|\mathbf{w}\|_1}_{f(\mathbf{w})}.$$

## Multiple benefits

- interpretability;
- complexity control;
- storage saving;
- perfect recovery;
- etc.

## Computationally?

- **convex** quadratic program
- but  $P \neq E$  !
- especially when  $p$  is large.



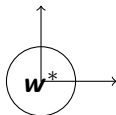


# Nonsmooth optimization

Generic subgradient descent:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta[\nabla\ell(\mathbf{w}_t) + \partial f(\mathbf{w}_t)]$$

- guaranteed convergence,  $O(1/\epsilon^2)$ ;



- dense iterates;
- weak regularizing effect;
- and slow, very slow...



Naum Zuselevich Shor  
(1937–2006)

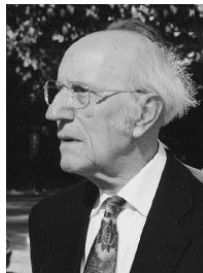
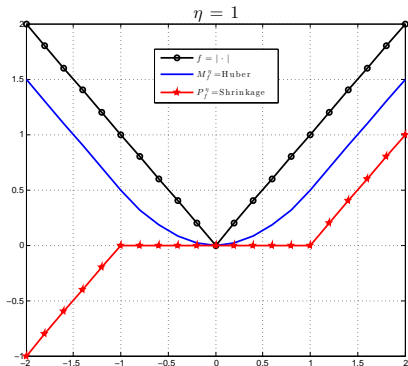
Second order methods (e.g. IPM) do not scale.

# Moreau envelope and proximal map

## Definition (Moreau'65)

$$M_f^\eta(\mathbf{y}) = \min_{\mathbf{w}} \frac{1}{2\eta} \|\mathbf{w} - \mathbf{y}\|^2 + f(\mathbf{w})$$

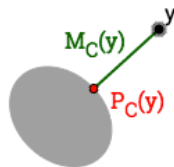
$$P_f^\eta(\mathbf{y}) = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2\eta} \|\mathbf{w} - \mathbf{y}\|^2 + f(\mathbf{w})$$



Jean Jacques Moreau, 1923–2014

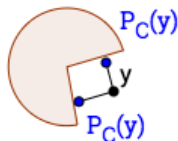
# Some properties of the proximal map

- For  $f(\mathbf{w}) = \iota_C(\mathbf{w}) := \begin{cases} 0, & \mathbf{w} \in C \\ \infty, & \text{otherwise} \end{cases}$ ,
  - ▶  $P_f^\eta(\cdot)$  is the usual Euclidean projection onto  $C$ ;
  - ▶  $M_f^\eta(\cdot)$  is the (squared) distance function;
  - ▶ Both well-defined as long as  $C$  is closed.



- For  $f$  convex (and closed),
  - ▶  $P_f^\eta(\cdot)$  is a **nonexpansion**:  $\|P_f^\eta(\mathbf{x}) - P_f^\eta(\mathbf{y})\| \leq \|\mathbf{x} - \mathbf{y}\|$ ;
  - ▶  $M_f^\eta(\cdot)$  is **continuously differentiable**;
  - ▶  $\eta \downarrow 0 \implies M_f^\eta \uparrow f$ .

- For general  $f$  (that decreases not too fast),
  - ▶  $P_f^\eta(\cdot)$  is a nonempty compact set;
  - ▶  $M_f^\eta(\cdot)$  is continuous;
  - ▶ Still  $\eta \downarrow 0 \implies M_f^\eta \uparrow f$ .



# Proximal gradient (Fukushima & Mine'81)

$$\min_{\mathbf{w} \in \mathbb{R}^m} \ell(\mathbf{w}) + f(\mathbf{w}), \quad \text{where } \ell \in \mathcal{C}^1.$$

1	$\mathbf{y}_t = \mathbf{w}_t - \eta \nabla \ell(\mathbf{w}_t);$	(forward)
2	$\mathbf{w}_{t+1} = P_f^\eta(\mathbf{y}_t).$	(backward)

For  $f = \|\cdot\|_1$ , obtain the shrinkage operator

$$[P_{\|\cdot\|_1}^\eta(\mathbf{y})]_i = \text{sign}(y_i)(|y_i| - \eta)_+.$$

- much faster,  $O(1/\epsilon)$ , can be accelerated;
- generalization of projected gradient:  $f = \iota_C$ ;
- reveals the sparsity-inducing property.

Refs: Combettes & Wajs'05; Beck & Teboulle'09; Duchi & Singer'09; Nesterov'13; etc.

## CONVEX PROGRAMMING IN HILBERT SPACE

BY A. A. GOLDSTEIN<sup>1</sup>

Communicated by V. Klee, May 1, 1964

This note gives a construction for minimizing certain twice-differentiable functions on a closed convex subset  $C$ , of a Hilbert Space,  $H$ . The algorithm assumes one can constructively "project" points onto convex sets. A related algorithm may be found in Cheney-Goldstein [1], where a constructive fixed-point theorem is employed to construct points inducing a minimum distance between two convex sets. In certain instances when such projections are not too difficult to construct, say on spheres, linear varieties, and orthants, the method can be effective. For applications to control theory, for example, see Balakrishnan [2], and Goldstein [3].

In what follows  $P$  will denote the "projection" operator for the convex set  $C$ . This operator, which is well defined and Lipschitzian, assigns to a given point in  $H$  its closest point in  $C$  (see, e.g., [1]). Take  $x \in H$  and  $y \in C$ . Then  $\|x - y, P(x) - y\| \geq \|P(x) - y\|^2$ . In the nontrivial case this inequality is a consequence of the fact that  $C$  is supported by a hyperplane through  $P(x)$  with normal  $x - P(x)$ . Let  $f$  be a real-valued function on  $H$  and  $x_0$  an arbitrary point of  $C$ . Let  $S$  denote the level set  $\{x \in C: f(x) \leq f(x_0)\}$ , and let  $\bar{S}$  be any open set containing the convex hull of  $S$ . Let  $f'(x, \cdot) = \nabla f(x, \cdot)$  signify the Fréchet derivative of  $f$  at  $x$ . A point  $z$  in  $C$  will be called stationary if  $P(z - \rho \nabla f(z)) = z$  for all  $\rho > 0$ ; equivalently, when  $f$  is convex the linear functional  $f'(z, \cdot)$  achieves a minimum on  $C$  at  $z$ .

**THEOREM.** Assume  $f$  is bounded below. For each  $x \in \bar{S}$ ,  $h$  in  $H$  and for some  $\rho_2 > 0$ , assume that  $f'(x, h)$  exists in the sense of Fréchet,  $f''(x, h, h)$  exists in the sense of Gâteaux, and  $\|f''(x, h, h)\| \leq \|h\|^2/\rho_2$ . Choose  $\sigma$  and  $\rho_1$  satisfying  $0 < \sigma \leq \rho_1$  and  $\sigma \leq \rho_2 \leq 2\rho_1 - \sigma$ . Set  $x_{k+1} = P(x_k - \rho_k \nabla f(x_k))$ . Then:

- (i) The sequence  $x_k$  belongs to  $S$ ,  $(x_{k+1} - x_k)$  converges to 0, and  $f(x_k)$  converges downward to a limit  $L$ .
- (ii) If  $S$  is compact,  $z$  is a cluster point of  $\{x_k\}$ , and  $\nabla f$  is continuous in some neighborhood of  $z$ , then  $z$  is a stationary point. If  $z$  is unique,  $x_k$  converges to  $z$ , and  $z$  minimizes  $f$  on  $C$ .
- (iii) If  $S$  is convex and  $f''(x, h, h) \geq \mu \|h\|^2$  for each  $x \in S$ ,  $h \in H$  and some  $\mu \geq 0$ , then  $L = \inf\{f(x); x \in C\}$ .
- (iv) Assume (iii) with  $S$  bounded. Weak cluster points of  $\{x_k\}$  minimize  $f$  on  $C$ .

<sup>1</sup> Present address, University of Washington, Seattle. This research was supported by grant AF-AFOSR-62-348.

(v) Assume (iii) with  $\mu$  positive and  $\nabla f$  bounded on  $S$ . Then  $f(z) = L$  for some  $z$  in  $S$ ,  $x_k$  converges to  $z$ , and  $z$  is unique.

**PROOF.** Assume  $x_0$  belongs to  $S$  and that  $x_0$  is not stationary. Let  $\nabla f(x_0) = \nabla f_0$ ,  $x(\rho) = P(x_0 - \rho \nabla f_0)$ ,  $\delta(\rho) = x(\rho) - x_0$  and  $\Delta(\rho) = f(x_0) - f(x(\rho))$ . If we notice that  $-\rho \langle \nabla f_0, \delta(\rho) \rangle \geq \|\delta(\rho)\|^2$  and invoke Taylor's theorem, we obtain  $\Delta(\rho) \geq \|\delta(\rho)\|^2 \{\rho^{-1} - f''(\xi(\rho), \delta(\rho))/2\} \|\delta(\rho)\|^2$ . Here  $\xi(\rho) = x_0 + t\delta(\rho)$  with  $t \in (0, 1)$ . For some  $\rho$  sufficiently small and positive,  $\Delta(\rho)$  is positive and continuous. Let  $\beta$  denote the least positive  $\rho$  satisfying  $\Delta(\rho) = 0$ , if such exists. If  $\beta$  exists,  $\Delta(\beta) = 0$  implies that  $\beta \geq 2\rho_0$ . Thus if  $\sigma \leq \rho \leq 2\rho_0 - \sigma$ ,  $\Delta(\rho) > 0$  and  $x(\rho) \in S$ , whence  $\Delta(\rho_0) \geq \|x_{k+1} - x_k\|^2/4\rho_0^2$ , proving (i).

The proof of (ii) being straightforward, we proceed with the proof of (iii). Suppose that  $L \neq \inf\{f(x); x \in C\}$  and choose  $z \in C$  such that  $f(z) < L$ . Then  $0 > f(z) - f(x_k) \geq \langle \nabla f, z - x_k \rangle$ . If  $\liminf \langle \nabla f, z - x_k \rangle = \beta$  were non-negative, a contradiction would be manifest. But the inequality  $\langle \rho_k \nabla f, z - x_{k+1} \rangle \geq \langle x_k - x_{k+1}, z \rangle + \langle x_{k+1}, x_{k+1} - x_k \rangle$  holds because either  $x_k - \rho_k \nabla f, x_{k+1}$  is the normal to  $C$  at  $x_{k+1}$ , or it is 0. If the sequence  $x_k$  is bounded, clearly  $\beta > 0$ ; otherwise choose a subsequence satisfying  $\|x_{k+1} - x_k\| > \|x_k\|$ . Then  $\beta \geq 0$ .

To prove (iv) we observe that  $f$  is lower semi-continuous on  $S$  if and only if the set  $S_m = \{x \in S; f(x) \leq m\}$  is closed in  $S$  for each  $m$ . Since  $f$  is convex and continuous,  $S_m$  is closed and convex, and is thus weakly closed. Hence  $f$  is weakly l.s.c. If  $x_k$  converges weakly to  $z$ , then  $\liminf f(x_k) = L \geq f(z)$ .

Assume the hypotheses of (v). If  $s > h$ , we may write that  $0 > f(x_0) - f(x_k) \geq \langle \nabla f, x_k - x_0 \rangle + (1/2)\mu \|x_k - x_0\|^2$ , whence  $\{x_k\}$  is bounded. Invoking again the supporting hyperplane at  $x_{k+1}$ ,  $\langle \rho_k \nabla f, x_k - x_k \rangle \geq \langle \rho_k \nabla f, x_{k+1} - x_k \rangle + \langle x_{k+1} - x_k, x_{k+1} - x_k \rangle$ . Thus when  $k$  is sufficiently large  $\|x_k - x_0\| < \epsilon$ . There exists therefore  $z \in S$  minimizing  $f$  on  $C$ , and  $f(x) \geq f(z) + \langle \nabla f(z), x - z \rangle + (1/2)\mu \|x - z\|^2$ . Since  $\langle \nabla f(z), x - z \rangle \geq 0$ ,  $f(x) - f(z) \geq (1/2)\mu \|x - z\|^2$ ; and therefore  $z$  is unique.

### REFERENCES

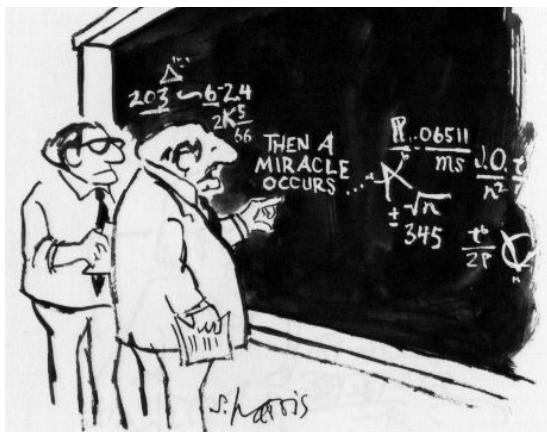
1. E. W. Cheney and A. A. Goldstein, *Proximity maps for convex sets*, Proc. Amer. Math. Soc. 10 (1959), 448-450.
2. A. V. Balakrishnan, *An operator theoretic formulation of a class of control problems and a steepest descent method of solution*, J. SIAM Control Ser. A 1 (1963), 109-127.
3. A. A. Goldstein, *Minimizing functionals on Hilbert space*, Computer methods in optimization problems, Academic Press, New York, 1964, pp. 159-165.

UNIVERSITY OF TEXAS

# Modern significance & rediscovery

- Donoho & Johnstone (90s), wavelet shrinkage;
- Starck, Donoho, and Candès (2003), astronomical image representation;
- Figueiredo & Nowak (2003), image restoration;
- Daubechies, Defrise, and De Mol (2004), inverse problem.
- Many many more...

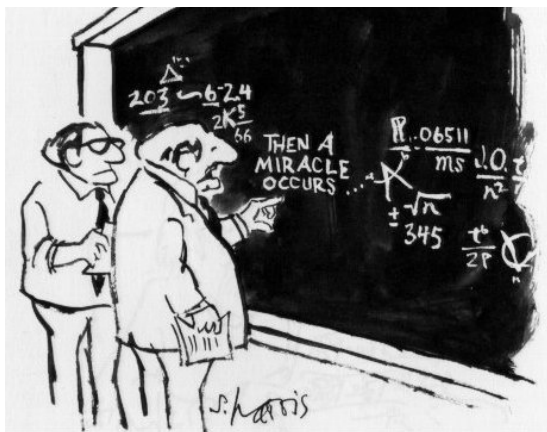
However...



"I think you should be more explicit here in step two."

from *What's so Funny about Science?* by Sidney Harris (1977)

However...



“I think you should be more explicit here in step two.”

from *What's so Funny about Science?* by Sidney Harris (1977)

Step 2: 
$$P_f^\eta(\mathbf{y}) = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2\eta} \|\mathbf{y} - \mathbf{w}\|^2 + f(\mathbf{w})$$



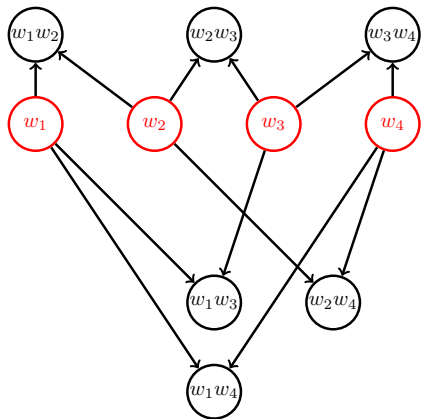
# Structured sparsity: group

Group level sparse regularizer

$$f(\mathbf{w}) = \sum_i \|\mathbf{w}\|_{g_i}.$$

For  $P_f$ , when groups are

- non-overlapping: decouple;
- tree structured: decompose;
- arbitrary?



Refs: Bakin'99; Yuan & Lin'06; Zhao et al.'09; etc.

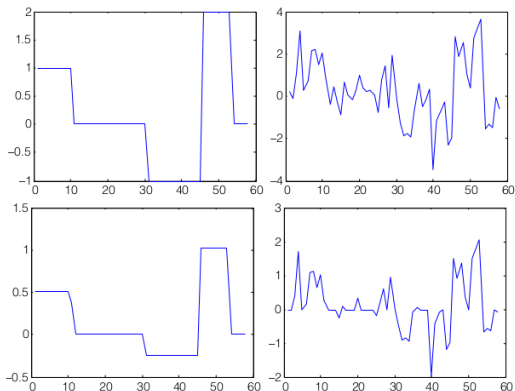
# Structured sparsity: graph

Neighborhood sparse regularizer

$$f(\mathbf{w}) = \sum_{\{i,j\} \in E} |w_i - w_j|.$$

For  $P_f$ , when graph is

- a chain: DP;
- arbitrary?
- vector valued?



Refs: Tibshirani et al.'05; Kim et al.'09; Kim & Xing'09; Hoefling'10; etc.

# Structured sparsity: matrix

- Matrix completion:

$$\min_{X \in \mathbb{R}^{m \times n}} \underbrace{\sum_{(i,j) \in \mathcal{O}} (X_{ij} - Z_{ij})^2}_{\ell(X)} + \underbrace{\lambda \|X\|_{\text{tr}}}_{f(X)}$$

- Can apply PG:

$$P_{\lambda \|\cdot\|_{\text{tr}}}^{\eta}(Y) = \sum_k (\sigma_k - \lambda \eta)_+ \mathbf{u}_k \mathbf{v}_k^{\top}$$

- Require **full** SVD in each step.



Refs: Candès & Recht'09; Cai et al.'10; Pong et al.'10; Toh & Yun'10; Ma et al.'11; etc.

# Learned so far

- Proximal gradient is simple, efficient, and structure-friendly.
  - ▶ easily parallelizable, can randomize, can block-wise.
- But backward step (proximal map) not always easy/cheap.
  - ▶ decompose;
  - ▶ approximate;
  - ▶ bypass proximal gradient;
- Constant theme: **exploit the structure of your problem!**
  - ▶ statistically;
  - ▶ and computationally.

# Table of Contents

- 1 Introduction
- 2 Decomposing the Proximal Map**
- 3 Approximation by the Proximal Average
- 4 Generalized Conditional Gradient
- 5 Post-PhD Extensions

## How to decompose?

- Typical structured sparse regularizers:

$$f(\mathbf{w}) = \sum_i f_i(\mathbf{w});$$

- ▶ Also applies to ERM, each  $i$  is a sample.
- **Key** observation: each  $P_{f_i}$  is easy to compute.
- Can we compute  $P_f = P_{\sum_i f_i}$  efficiently?

### Theorem (Folklore)

$$P_{f+g} = (P_{2f}^{-1} + P_{2g}^{-1})^{-1} \circ (2\text{Id}).$$

- Not directly useful;
- Can numerically reduce to  $P_f$  and  $P_g$  (Combettes et al.'11);
- But a two-loop routine can be as slow as subgradient (Villa et al.'13).

## Two previous results

$$\|\mathbf{w}\|_{\text{TV}} = \sum_{i=1}^p |w_i - w_{i+1}|.$$

### Theorem (Friedman et al.'07)

$$P_{\|\cdot\|_1 + \|\cdot\|_{\text{TV}}} = P_{\|\cdot\|_1} \circ P_{\|\cdot\|_{\text{TV}}}.$$

### Theorem (Jenatton et al.'11)

$$P_{\sum_{i=1}^k \|\cdot\|_{g_i}} = P_{\|\cdot\|_{g_1}} \circ \dots \circ P_{\|\cdot\|_{g_k}}.$$

### Generalization

$$P_{f+g} \stackrel{?}{=} P_f \circ P_g \stackrel{?}{=} P_g \circ P_f.$$

But, is it even sensible?

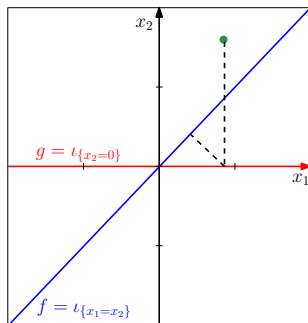
# Good news and bad news

## Theorem

On the real line,  $\exists h$  such that  $P_h = P_f \circ P_g$ .

## Example (But not so in general...)

Consider  $\mathbb{R}^2$ , and let  $f = \iota_{\{x_1=x_2\}}$ ,  $g = \iota_{\{x_2=0\}}$ .





## Nevertheless

- Can ask the decomposition to hold for many but not all cases.
- Setting the subdifferential to 0:

$$P_{f+g}(\mathbf{z}) - \mathbf{z} + \partial(f + g)(P_{f+g}(\mathbf{z})) \ni 0$$

$$P_g(\mathbf{z}) - \mathbf{z} + \partial g(P_g(\mathbf{z})) \ni 0$$

$$P_f(P_g(\mathbf{z})) - P_g(\mathbf{z}) + \partial f(P_f(P_g(\mathbf{z}))) \ni 0.$$

- Adding the last two equations we obtain

$$P_f(P_g(\mathbf{z})) - \mathbf{z} + \partial g(P_g(\mathbf{z})) + \partial f(P_f(P_g(\mathbf{z}))) \ni 0.$$

### Theorem (Y'13a)

A sufficient condition for  $P_{f+g}(\mathbf{z}) = P_f(P_g(\mathbf{z}))$  is

$$\forall \mathbf{y} \in \text{dom } g, \partial g(P_f(\mathbf{y})) \supseteq \partial g(\mathbf{y}).$$

# The rest is easy



- Find  $f$  and  $g$  that clinch our sufficient condition.

## Result I: Start with “trivialities”

### Theorem (Y'13a)

Fix  $f \in \Gamma_0$ .  $P_{f+g} = P_f \circ P_g$  for *all*  $g \in \Gamma_0$  if and only if

- $\dim(\mathcal{H}) \geq 2$ ;  $f \equiv c$  or  $f = \iota_{\{\mathbf{w}\}} + c$  for some  $c \in \mathbb{R}$  and  $\mathbf{w} \in \mathcal{H}$ ;
- $\dim(\mathcal{H}) = 1$  and  $f = \iota_C + c$  for some closed and convex set  $C$  and  $c \in \mathbb{R}$ .

Asymmetry.

### Theorem (Y'13a)

Fix  $g \in \Gamma_0$ .  $P_{f+g} = P_f \circ P_g$  for *all*  $f \in \Gamma_0$  if and only if  $g$  is a continuous affine function.

## Result II: Positive homogeneity and “roundness”

### Theorem (Y'13a)

Let  $f \in \Gamma_0$ . The following are equivalent (provided  $\dim(\mathcal{H}) \geq 2$ ):

- i).  $f = h(\|\cdot\|)$  for some increasing function  $h : \mathbb{R}_+ \rightarrow \mathbb{R} \cup \{\infty\}$ ;
- ii).  $\mathbf{x} \perp \mathbf{y} \implies f(\mathbf{x} + \mathbf{y}) \geq f(\mathbf{y})$ ;
- iii). For all  $\mathbf{z} \in \mathcal{H}$ ,  $P_f(\mathbf{z}) = \lambda_{\mathbf{z}} \cdot \mathbf{z}$  for some  $\lambda_{\mathbf{z}} \in [0, 1]$ ;
- iv).  $\mathbf{0} \in \text{dom } f$  and  $P_{f+\kappa} = P_f \circ P_{\kappa}$  for *all* p.h. functions  $\kappa \in \Gamma_0$ .

- Include and generalize many results;
- Connects to the representer theorem in kernel methods (YCSS'13).

# More implications

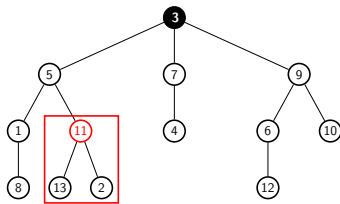
## Example (Elastic net, Zou & Hastie'05)

$$P_{\lambda \|\cdot\|_2^2 + \kappa} = P_{\lambda \|\cdot\|_2} \circ P_{\kappa} = \frac{1}{\lambda+1} P_{\kappa} \quad \text{double shrinkage}$$

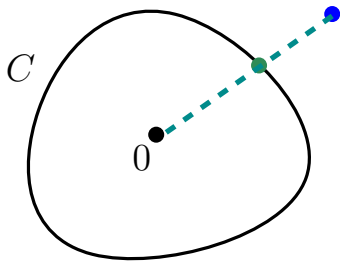
## Example (Jenatton et al.'11)

Tree-structured (laminar system)

$$P_{\sum_i \|\cdot\|_{g_i}} = P_{\|\cdot\|_{g_1}} \circ \dots \circ P_{\|\cdot\|_{g_k}}$$



# Characterizing the ball



## Result III: Comonotonicity and Choquet integral

Initially case by case for many polyhedral regularizers.

### Theorem (Y'13a)

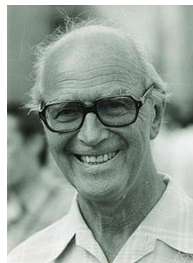
Let  $f$  be permutation invariant and  $g$  be the Choquet integral of some submodular set function.

$$P_{f+g} = P_f \circ P_g.$$

### Example (Friedman et al.'07)

$$P_{\|\cdot\|_1 + \|\cdot\|_{TV}} = P_{\|\cdot\|_1} \circ P_{\|\cdot\|_{TV}}.$$

- $\|\cdot\|_1$ : permutation invariant;
- $\|\cdot\|_{TV}$ : Choquet integral of something.



Gustave Choquet  
(1915–2006)

“Always consider a problem under the minimum structure in which it makes sense.”

# Summary

- Posed the question:  $P_{f+g} \stackrel{?}{=} P_f \circ P_g \stackrel{?}{=} P_g \circ P_f$ ;
- Presented a sufficient condition:  $\partial g(P_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$ ;
- “Trivial” case;
- Positive homogeneity and “roundness”;
- Comonotonicity and Choquet integral;
- Immediately useful if plugged into PG;

What if the sufficient condition fails?



# Table of Contents

- 1 Introduction
- 2 Decomposing the Proximal Map
- 3 Approximation by the Proximal Average**
- 4 Generalized Conditional Gradient
- 5 Post-PhD Extensions

## More generally

Recall: typical structured sparse regularizers:  $\bar{f} = \sum_i \alpha_i f_i$

- $P_{f_i}^\eta$  easy to compute;
- $f_i$  Lipschitz continuous.

### Example (Overlapping group lasso, Zhao et al.'09)

$f_i(\mathbf{w}) = \|\mathbf{w}\|_{g_i}$  where  $g_i$  is a group (subset) of variables.

- When the groups overlap arbitrarily,  $P_{\bar{f}}^\eta$  cannot be easily computed;
- Each  $f_i$  is 1-Lipschitz continuous *w.r.t.*  $\|\cdot\|$ ;
- The proximal map  $P_{f_i}^\eta$  is simply a re-scaling:

$$[P_{f_i}^\eta(\mathbf{w})]_j = \begin{cases} w_j, & j \notin g_i \\ (1 - \eta/\|\mathbf{w}\|_{g_i})_+ w_j, & j \in g_i \end{cases}.$$

## Example cont'

### Example (Graph-guided fused lasso, Kim & Xing'09)

Given some graph, we let  $f_{ij}(\mathbf{w}) = |w_i - w_j|$  for every edge  $\{i, j\}$ .

- For a general graph, the proximal map of the regularizer  $\bar{f} = \sum_{\{i,j\} \in E} \alpha_{ij} f_{ij}$  can not be easily computed;
- Each  $f_{ij}$  is 1-Lipschitz continuous *w.r.t.* the Euclidean norm;
- The proximal map  $P_{f_{ij}}^\eta$  is easy to compute:

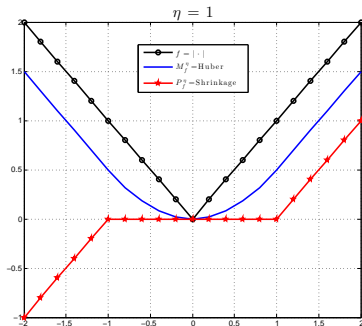
$$[P_{f_{ij}}^\eta(\mathbf{w})]_s = \begin{cases} w_s, & s \notin \{i, j\} \\ w_s - \text{sign}(w_i - w_j) \min\{\eta, |w_i - w_j|/2\}, & s \in \{i, j\} \end{cases}.$$

Other examples abound.

# Smoothing (Nesterov'05)

$$M_f^\eta(\mathbf{y}) = \min_{\mathbf{w}} \frac{1}{2\eta} \|\mathbf{w} - \mathbf{y}\|^2 + f(\mathbf{w})$$

$$P_f^\eta(\mathbf{y}) = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2\eta} \|\mathbf{w} - \mathbf{y}\|^2 + f(\mathbf{w})$$



## Proposition (Nesterov'05)

If  $f$  is  $L$ -Lipschitz continuous, then  $0 \leq f - M_f^\eta \leq \eta L^2/2$ .

## In retrospect

Suppose want:

$$\min_{\mathbf{w} \in C} \ell(\mathbf{w}).$$

Same for large  $\lambda > 0$ :

$$\min_{\mathbf{w}} \ell(\mathbf{w}) + \lambda \cdot \text{dist}(\mathbf{w}, C)$$

- $\text{dist}(\mathbf{w}, C) := \min_{\mathbf{z} \in C} \|\mathbf{w} - \mathbf{z}\|$ , nonsmooth but Lipschitz continuous.
- Can smooth dist and apply gradient descent.
- But nobody does that, overkill.
- Can just use projected gradient.

## A “naive” idea (Y’13b)

$$\bar{f} = \sum_i \alpha_i f_i$$

↓ as if have linearity?

$$P_{\bar{f}}^\eta \approx \sum_i \alpha_i P_{f_i}^\eta$$

Definition (Proximal Average, Moreau’65; Bauschke et al.’08)

There exists a unique function  $A^\eta$  such that  $P_{A^\eta}^\eta = \sum_i \alpha_i P_{f_i}^\eta$ .

# What mathematicians call a “picture”

$$\begin{array}{ccc} \Gamma_0 \ni f_i & \xrightarrow{\text{onto}} & M_{f_i}^\eta \in \text{SS}_{1/\eta} \\ & & \downarrow \text{convex} \\ \Gamma_0 \ni A^\eta & \xleftarrow{1-1} & \sum_i \alpha_i M_{f_i}^\eta \in \text{SS}_{1/\eta} \\ \downarrow \nabla & & \downarrow \nabla \\ P_{A^\eta}^\eta & \longleftrightarrow & \sum_i \alpha_i P_{f_i}^\eta \end{array}$$

- Not so easy to compute  $A^\eta$ , but existence is enough.

# The algorithm

Dream

$$\textcircled{1} \mathbf{z}_t = \mathbf{w}_t - \mu \nabla \ell(\mathbf{w}_t)$$

$$\textcircled{2} \mathbf{w}_{t+1} = P_{\bar{f}}^{\eta}(\mathbf{z}_t)$$



$$\min_{\mathbf{w}} \ell(\mathbf{w}) + \bar{f}(\mathbf{w})$$

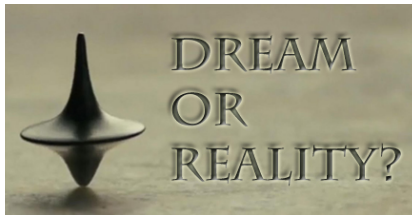
Reality

$$\textcircled{1} \mathbf{z}_t = \mathbf{w}_t - \mu \nabla \ell(\mathbf{w}_t)$$

$$\textcircled{2} \mathbf{w}_{t+1} = P_{A^{\eta}}^{\eta}(\mathbf{z}_t) = \sum_i \alpha_i P_{f_i}^{\eta}(\mathbf{z}_t)$$



$$\min_{\mathbf{w}} \ell(\mathbf{w}) + A^{\eta}(\mathbf{w})$$



When are they close?



# Nonsmooth approximation

- How good the proximal average  $A^\eta$  approximates  $\bar{f}$ ?

## Proposition (Uniform lower approximation)

Assuming  $f_i$  is  $M_i$ -Lipschitz continuous, and  $M := \sum_i \alpha_i M_i^2$ , then

$$0 \leq \bar{f} - A^\eta \leq \eta M^2 / 2.$$

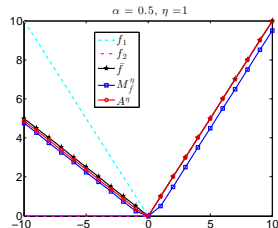
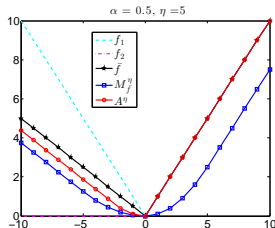
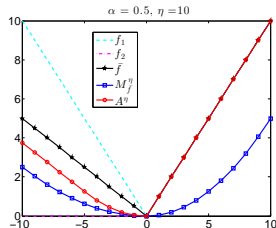
- Proximal average is a tighter approximation than smoothing:

$$\sum_i \alpha_i M_{f_i}^\eta \leq A^\eta \leq \bar{f}.$$

# An example

## Example

Consider  $f_1(x) = |x|$ , and  $f_2(x) = \max\{x, 0\}$ .



- The proximal average is smooth iff some  $f_i$  is;
- Essentially we de-smooth Nesterov's approximation.

# Convergence guarantee

## Theorem (Y'13b)

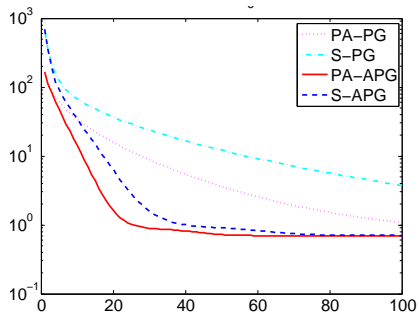
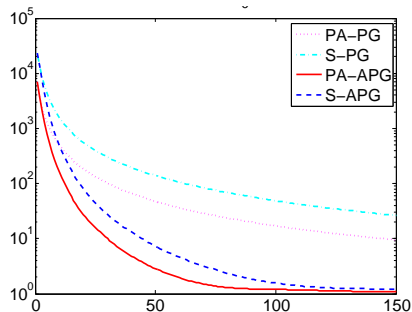
Using a suitable step size, we get an  $\epsilon$ -accurate solution in at most  $O(\sqrt{\max\{L_0, L^2/(2\epsilon)\}}\sqrt{1/\epsilon})$  steps.

- Improves Nesterov's complexity  $O(\sqrt{L_0 + L^2/(2\epsilon)}\sqrt{1/\epsilon})$  by removing secondary term;
- No overhead, same assumption, strict improvement;
- Simple update rule.

$$\text{S-PG: } \mathbf{w}_{t+1} = \frac{\eta L_0}{1 + \eta L_0} \left[ \mathbf{w}_t - \frac{1}{L_0} \nabla \ell(\mathbf{w}_t) \right] + \frac{1}{1 + \eta L_0} \sum_i \alpha_i P_{f_i}^\eta(\mathbf{w}_t),$$

$$\text{PA-PG: } \mathbf{w}_{t+1} = \sum_i \alpha_i P_{f_i}^\eta(\mathbf{w}_t - \eta \nabla \ell(\mathbf{w}_t)).$$

# Experiment



# Summary

- Linear approximation of the proximal map;
- Improved convergence guarantee;
- Retain nonsmoothness (to some extent);
- How to combine regularizers?

# Table of Contents

- 1 Introduction
- 2 Decomposing the Proximal Map
- 3 Approximation by the Proximal Average
- 4 Generalized Conditional Gradient**
- 5 Post-PhD Extensions

# Conditional gradient (Frank-Wolfe'56)

$$\min_{\mathbf{w} \in C} \ell(\mathbf{w})$$

- $C$ : compact convex;
- $\ell$ : smooth convex.

$$\begin{aligned} \textcircled{1} \quad & \mathbf{y}_t \in \operatorname{argmin}_{\mathbf{w} \in C} \langle \mathbf{w}, \nabla \ell(\mathbf{w}_t) \rangle; \\ \textcircled{2} \quad & \mathbf{w}_{t+1} = (1 - \eta) \mathbf{w}_t + \eta \mathbf{y}_t. \end{aligned}$$

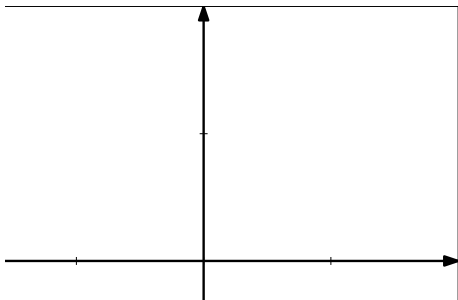
Gained much recent attention due to

- its simplicity;
- the greedy nature in step 1.

Refs: Zhang'03; Clarkson'10; Hazan'08; Jaggi-Sulovsky'10; etc.

## An example

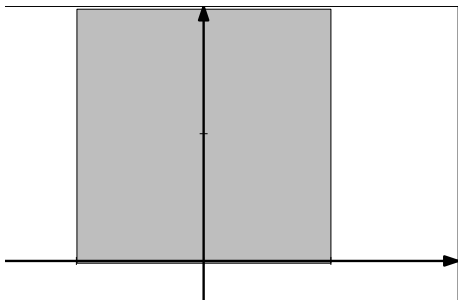
$$\min_{a,b} a^2 + (b + 1)^2, \text{ s.t. } |a| \leq 1, 2 \geq b \geq 0$$





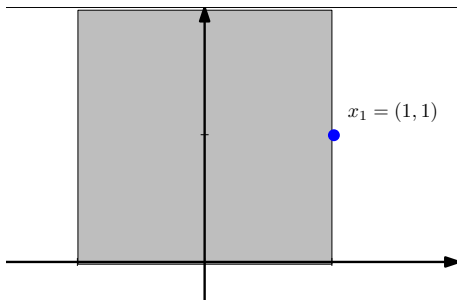
## An example

$$\min_{a,b} a^2 + (b + 1)^2, \text{ s.t. } |a| \leq 1, 2 \geq b \geq 0$$



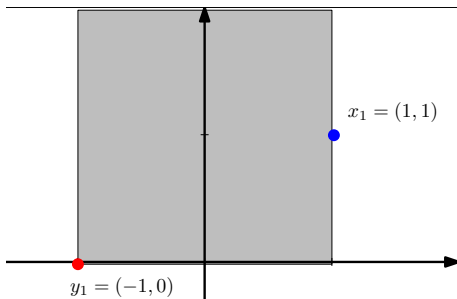
## An example

$$\min_{a,b} a^2 + (b+1)^2, \text{ s.t. } |a| \leq 1, 2 \geq b \geq 0$$



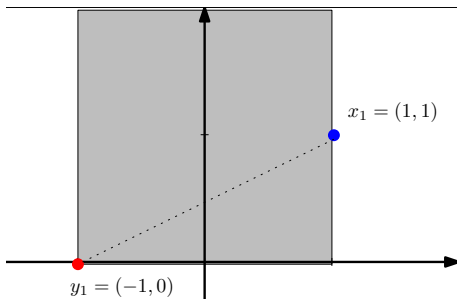
## An example

$$\min_{a,b} a^2 + (b+1)^2, \text{ s.t. } |a| \leq 1, 2 \geq b \geq 0$$



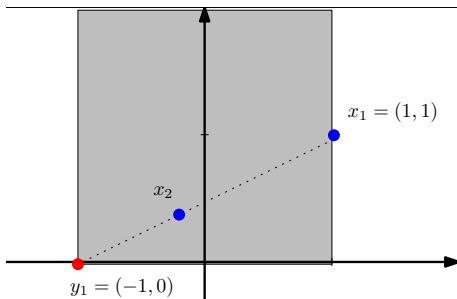
## An example

$$\min_{a,b} a^2 + (b+1)^2, \text{ s.t. } |a| \leq 1, 2 \geq b \geq 0$$



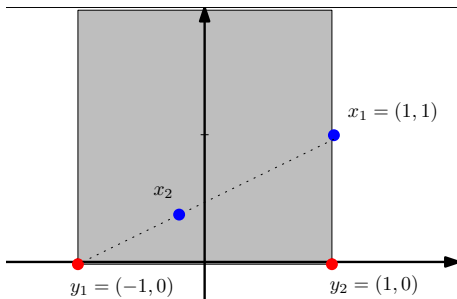
## An example

$$\min_{a,b} a^2 + (b+1)^2, \text{ s.t. } |a| \leq 1, 2 \geq b \geq 0$$



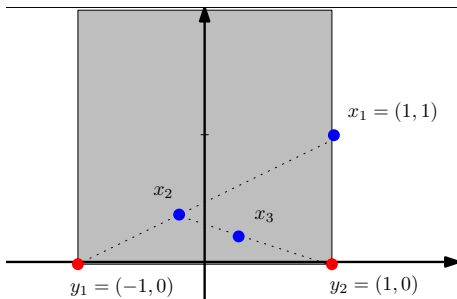
## An example

$$\min_{a,b} a^2 + (b+1)^2, \text{ s.t. } |a| \leq 1, 2 \geq b \geq 0$$



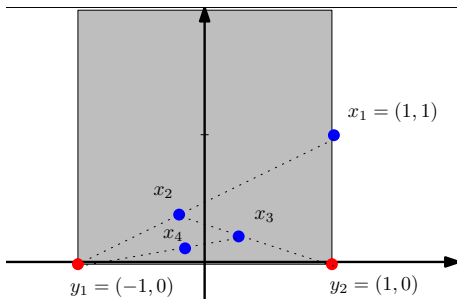
## An example

$$\min_{a,b} a^2 + (b+1)^2, \text{ s.t. } |a| \leq 1, 2 \geq b \geq 0$$



## An example

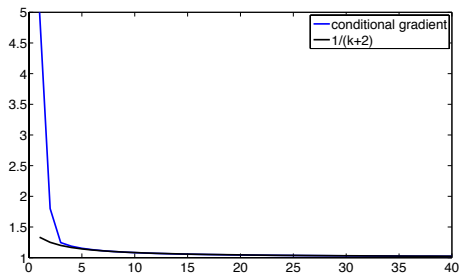
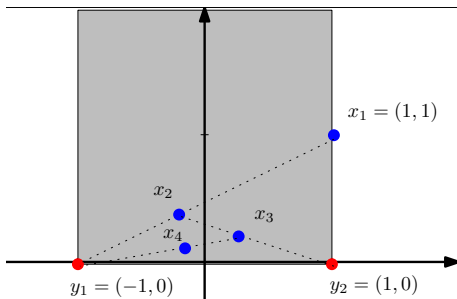
$$\min_{a,b} a^2 + (b+1)^2, \text{ s.t. } |a| \leq 1, 2 \geq b \geq 0$$





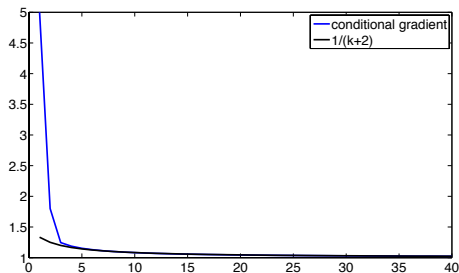
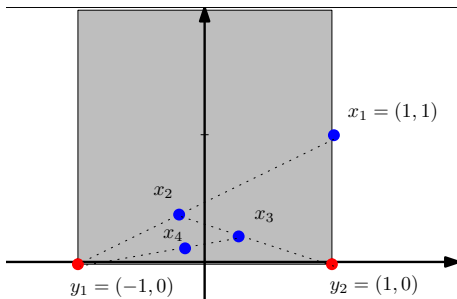
# An example

$$\min_{a,b} a^2 + (b+1)^2, \text{ s.t. } |a| \leq 1, 2 \geq b \geq 0$$



# An example

$$\min_{a,b} a^2 + (b+1)^2, \text{ s.t. } |a| \leq 1, 2 \geq b \geq 0$$

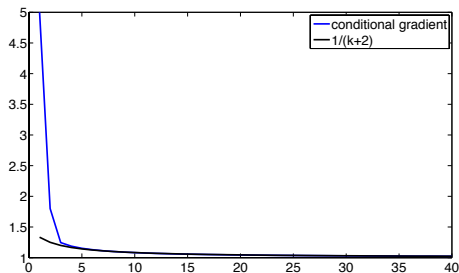
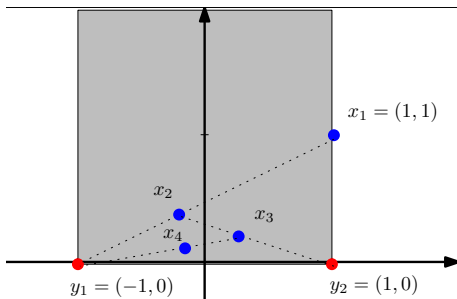


Can show  $\ell(\mathbf{w}_t) - \ell(\mathbf{w}^*) = 4/t + o(1/t)$ .

PG converges in two iterations.

## An example

$$\min_{a,b} a^2 + (b+1)^2, \text{ s.t. } |a| \leq 1, 2 \geq b \geq 0$$



Can show  $\ell(\mathbf{w}_t) - \ell(\mathbf{w}^*) = 4/t + o(1/t)$ .

PG converges in two iterations.

Refs: (Levtin-Polyak'66; Polyak'87; Beck-Teboulle'04) for faster rates.

# The revival of CG: sparsity!

The revived popularity of conditional gradient is due to (Clarkson'10; Shalev-Shwartz-Srebro-Zhang'10), both focusing on

$$\min_{\mathbf{w}: \|\mathbf{w}\|_1 \leq 1} \ell(\mathbf{w}).$$

$$\textcircled{1} \quad \mathbf{y}_t \leftarrow \underset{\|\mathbf{y}\|_1 \leq 1}{\operatorname{argmin}} \langle \mathbf{y}, \nabla \ell(\mathbf{w}_t) \rangle, \quad \text{card}(\mathbf{y}_t) = 1;$$

$$\textcircled{2} \quad \mathbf{w}_{t+1} \leftarrow (1 - \eta)\mathbf{w}_t + \eta\mathbf{y}_t, \quad \text{card}(\mathbf{w}_{t+1}) \leq \text{card}(\mathbf{w}_t) + 1.$$

Explicit control of the sparsity.

Later on, (Hazan'08; Jaggi-Sulovsky'10) generalized the idea to SDPs.

# Generalized conditional gradient

$$\min_{\mathbf{w}} \ell(\mathbf{w}) + \lambda \cdot f(\mathbf{w})$$

- composite, with a nonsmooth term;
- unconstrained, hence unbounded domain;
- first studied by Mine & Fukushima'81 and then Bredies et al.'09;
- generalizes CG.

$$\begin{aligned} \textcircled{1} \quad & \mathbf{y}_t \in \operatorname{argmin}_{\mathbf{w}} \langle \mathbf{w}, \nabla \ell(\mathbf{w}_t) \rangle + f(\mathbf{w}); \\ \textcircled{2} \quad & \mathbf{w}_{t+1} = (1 - \eta) \mathbf{w}_t + \eta \mathbf{y}_t. \end{aligned}$$

Our interest:

- $f$  p.h. (e.g., a norm);
- Step 1 undefined.

# Positive homogeneous regularizer

$$\min_{\mathbf{w}} \ell(\mathbf{w}) + \lambda \cdot \kappa(\mathbf{w})$$

- $\ell$ : smooth convex;
- $\kappa$ : positive homogeneous convex—gauge (not necessarily smooth).

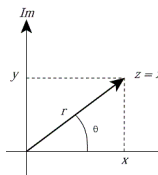
Challenges:

- composite, with a nonsmooth term;
- unconstrained, hence unbounded domain;
- $\kappa$  expensive to evaluate.

① **Polar operator:**  $\mathbf{y}_t \in \operatorname{argmin}_{\mathbf{w}: \kappa(\mathbf{w}) \leq 1} \langle \mathbf{w}, \nabla \ell(\mathbf{w}_t) \rangle$ ;

② line search:  $s_t \in \operatorname{argmin}_{s \geq 0} \ell((1 - \eta)\mathbf{w}_t + \eta s \mathbf{y}_t) + \lambda \eta s$ ;

③  $\mathbf{w}_{t+1} = (1 - \eta)\mathbf{w}_t + \eta s_t \mathbf{y}_t$ .



# Convergence guarantee

$$\min_{\mathbf{w}} \ell(\mathbf{w}) + \lambda \cdot \kappa(\mathbf{w})$$

## Theorem (ZYS'12)

*If  $\ell$  and  $\kappa$  have bounded level sets and  $\ell \in C^1$ , then GCG converges at rate  $O(1/t)$ , where the constant is independent of  $\lambda$ .*

- Proof is simple: Line search is as good as knowing  $\kappa(\mathbf{w}^*)$ ;
- Upper bound

$$\kappa((1 - \eta)\mathbf{w}_t + \eta\mathbf{s}\mathbf{y}_t) \leq (1 - \eta)\kappa(\mathbf{w}_t) + \eta\kappa(\mathbf{s}\mathbf{y}_t) \leq (1 - \eta)\kappa(\mathbf{w}_t) + \eta s;$$

- Still too slow!

# Local improvement

Assume some procedure (say LOCAL) that can *locally* solve

$$\min_{\mathbf{w}} \ell(\mathbf{w}) + \lambda \cdot \kappa(\mathbf{w}),$$

or some variation of it.

Combine LOCAL with some GLOBAL?

Three conditions:

- LOCAL cannot incur big overhead;
- cannot ruin GLOBAL;
- easy to switch between LOCAL and GLOBAL.



## Case study: matrix completion with trace norm

GLOBAL: 
$$\min_X \sum_{(i,j) \in \mathcal{O}} (X_{ij} - Z_{ij})^2 + \lambda \cdot \|X\|_{\text{tr}}$$

The only nontrivial step in GCG:

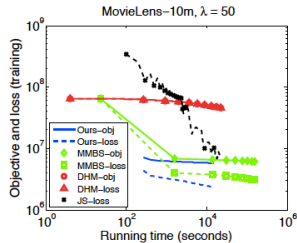
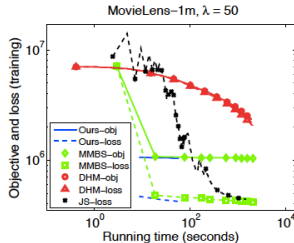
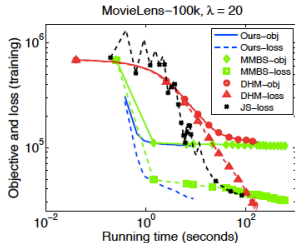
- Polar operator:  $Y_t \in \underset{\|Y\|_{\text{tr}} \leq 1}{\text{argmin}} \langle Y, G_t \rangle$ , *dominating* singular vectors.

In contrast, PG requires the *full* SVD of  $-G_t$ .

LOCAL (Srebro'05): 
$$\min_{B,W} \sum_{(i,j) \in \mathcal{O}} ((BW)_{ij} - Z_{ij})^2 + \lambda/2 \cdot (\|B\|_F^2 + \|W\|_F^2).$$

- Not jointly convex in  $B$  and  $W$ ;
- But smooth in  $B$  and  $W$ ;
- $Y_t$  in GCG is rank-1 hence  $X_t = BW$  is of rank at most  $t$ .

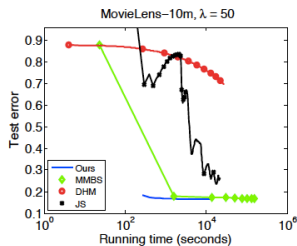
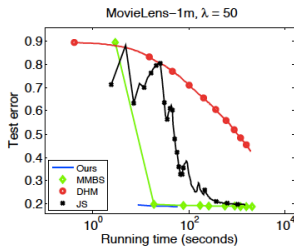
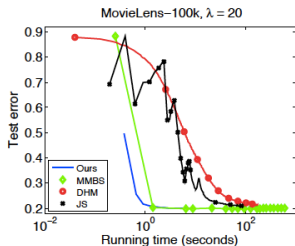
# Case study: experiment



(a) Objective & loss vs time (loglog)

(a) Objective & loss vs time (loglog)

(a) Objective & loss vs time (loglog)



(b) Test NMAE vs time (semilogx)

(b) Test NMAE vs time (semilogx)

(b) Test NMAE vs time (semilogx)

# Summary

- Generalized conditional gradient for p.h. regularizer;
- $O(1/t)$  convergence rate;
- Combined LOCAL with GCG ;
- Applied to matrix completion.

# Table of Contents

- 1 Introduction
- 2 Decomposing the Proximal Map
- 3 Approximation by the Proximal Average
- 4 Generalized Conditional Gradient
- 5 Post-PhD Extensions**

# Prox-decomposition and isotonicity

Hölder's inequality:  $\langle \mathbf{x}, \mathbf{y} \rangle \leq \|\mathbf{x}\|_r \|\mathbf{y}\|_s$ ,  $r \geq 1$ ,  $1/r + 1/s = 1$

Ky Fan's norm  $\|\mathbf{x}\|_{k,r} := \sqrt[r]{\sum_{i=1}^k |x|_{(i)}^r}$ .

$\langle \mathbf{x}, \mathbf{y} \rangle \leq \|\mathbf{x}\|_{k,r}$ ???, i.e., dual norm  $\|\mathbf{y}\|_{k,r}^\circ := \max_{\|\mathbf{x}\|_{k,r} \leq 1} \langle \mathbf{x}, \mathbf{y} \rangle = ?$

First shown in (Mudholkar et al, 1984).

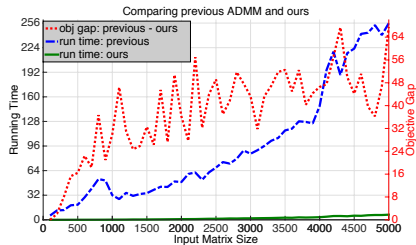
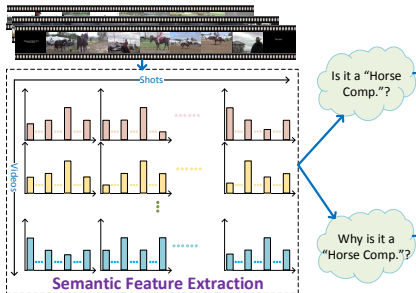
## Theorem (YYX'15)

For any  $r \geq 1$  and  $1/r + 1/s = 1$ , the dual Ky Fan norm  $\|\mathbf{y}\|_{k,r}^\circ = \|\mathbf{z}\|_s$ ,

where  $\mathbf{z} := \mathbf{P}_{\mathcal{K}}(\mathbf{m}) = \underset{w_1 \geq w_2 \geq \dots \geq w_k}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{m} - \mathbf{w}\|_2^2$  and

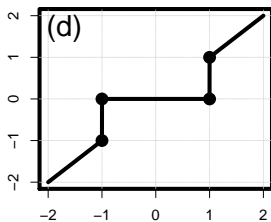
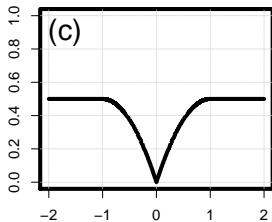
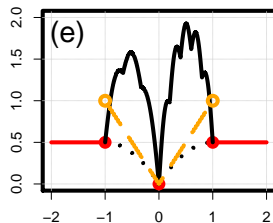
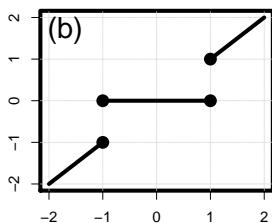
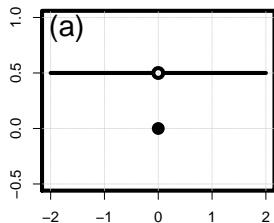
$$m_i = \begin{cases} |y|_{(i)}, & i = 1, \dots, k-1 \\ \sum_{j=k}^p |y|_{(j)}, & i = k \end{cases}.$$

# Video event detection and recounting (CYYH'15)



# Nonconvex proximal average (YZMX'14)

$$P_{\sum_i \alpha_i f_i}^\eta \stackrel{?}{\approx} \sum_i \alpha_i P_{f_i}^\eta$$



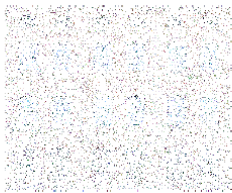
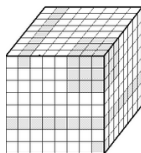
# Approximate generalized conditional gradient

Pick  $\kappa(\mathbf{y}_t) \leq 1$  such that for some  $\alpha \in (0, 1]$

$$\langle \mathbf{y}_t, \nabla \ell(\mathbf{w}_t) \rangle \leq \alpha \cdot \min_{\mathbf{y}: \kappa(\mathbf{y}) \leq 1} \langle \mathbf{y}, \nabla \ell(\mathbf{w}_t) \rangle.$$

## Theorem (YCZ'14)

Assume  $\ell \geq 0$ . Equipped with an  $\alpha$ -approximate PO, GCG “converges” to an  $\alpha$ -approximate solution at the rate  $O(1/t)$ .





# Thank you!