

On Decomposing the Proximal Map

Yao-Liang Yu

University of Waterloo

Splitting Algorithms, Modern Operator Theory, and Applications

Oaxaca, Mexico

September 19, 2017

Table of Contents

1 Motivation

2 Setup

3 A Naive Sufficient Condition

4 Three Case Studies

5 More Examples

Regularized loss minimization

Generic form for many ML problems:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \ell(\mathbf{w}) + f(\mathbf{w})$$

- ℓ is the loss function;
- f is the regularizer, usually a (semi)norm;

Special interest:

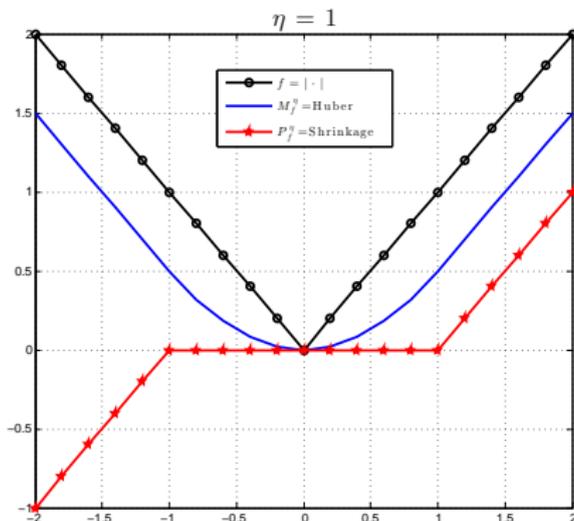
- sparsity;
- computational efficiency.

Moreau envelop and proximal map

Definition (Moreau'65)

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

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



Proximal gradient (Fukushima & Mine'81)

$$\min_{\mathbf{w} \in \mathbb{R}^d} \ell(\mathbf{w}) + f(\mathbf{w})$$

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

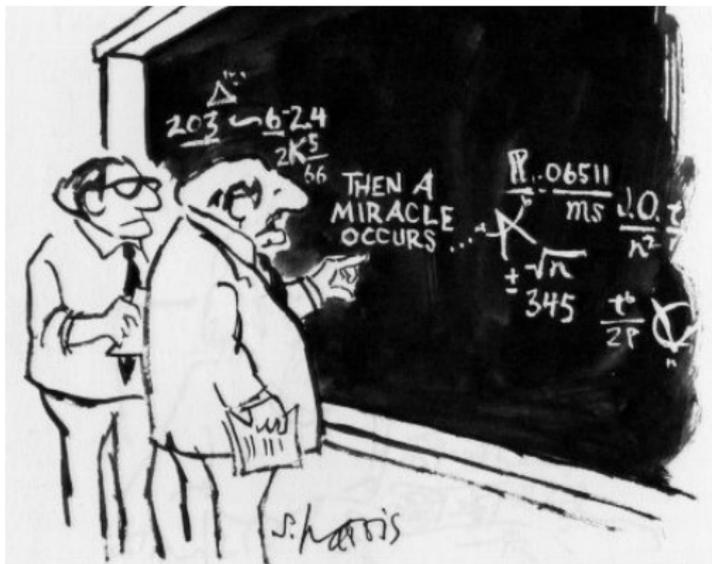
For $f = \|\cdot\|_1$, obtain the shrinkage operator

$$[P_{\|\cdot\|_1}(\mathbf{y})]_i = \text{sign}(y_i)(|y_i| - 1)_+.$$

- guaranteed convergence, 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.

Then A Miracle Occurs...



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

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

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

How to deal with sum?

- Typical structured sparse regularizers:

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

Theorem (Parallel Sum)

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

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

Table of Contents

- 1 Motivation
- 2 Setup**
- 3 A Naive Sufficient Condition
- 4 Three Case Studies
- 5 More Examples

Two previous results

Theorem (Friedman et al.'07)

$$P_{\|\cdot\|_1 + \|\cdot\|_{\text{TV}}} = P_{\|\cdot\|_1} \circ P_{\|\cdot\|_{\text{TV}}}, \quad \text{where} \quad \|\mathbf{w}\|_{\text{TV}} = \sum_{i=1}^{d-1} |w_i - w_{i+1}|.$$

Theorem (Jenatton et al.'11)

Assuming the groups $\{\mathbf{g}_i\}$ form a laminar system ($\mathbf{g}_i \cap \mathbf{g}_j \in \{\mathbf{g}_i, \mathbf{g}_j, \emptyset\}$), then, if appropriately ordered,

$$P_{\sum_{i=1}^k \|\cdot\|_{\mathbf{g}_i}} = P_{\|\cdot\|_{\mathbf{g}_1}} \circ \cdots \circ P_{\|\cdot\|_{\mathbf{g}_k}},$$

where $\|\cdot\|_{\mathbf{g}_i}$ is the restriction of $\|\cdot\|_p$, $p \in \{1, 2, \infty\}$ to the group \mathbf{g}_i .

(Wild) Generalization

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

Product of Prox's

- Long line of work: von Neumann, Halperin, Amemiya and Ando, Stiles, Dye, Reich, Bruck, Tseng, Brézis and Lions, etc., etc.
- interest was in the asymptotic behaviour
- in some sense, we want one-step convergence of such algs

Bad news

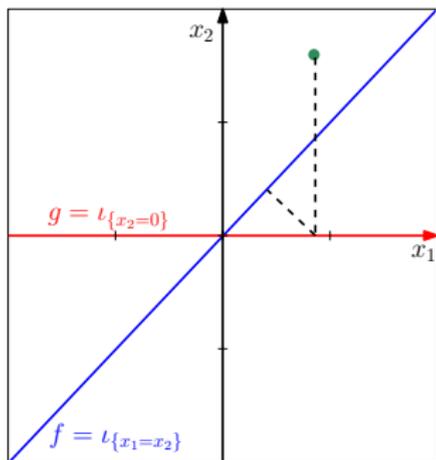
Theorem

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

- Not necessarily $h = f + g$, though

Example (A simple counterexample)

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



$$P_f = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}, \quad P_g = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}.$$

$$\text{But } P_f \circ P_g = \begin{bmatrix} 0.5 & 0 \\ 0.5 & 0 \end{bmatrix}$$

no h such that $P_h = P_f \circ P_g$

Table of Contents

- 1 Motivation
- 2 Setup
- 3 A Naive Sufficient Condition**
- 4 Three Case Studies
- 5 More Examples

Nevertheless

- Not possible to always have the decomposition — too ambitious
- More modest goal: decomposition to hold for certain functions
- Manipulating the optimality conditions:

$$P_{f+g}(z) = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{z} - \mathbf{w}\|^2 + (f + g)(\mathbf{w})$$

$$P_g(z) = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{z} - \mathbf{w}\|^2 + g(\mathbf{w})$$

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

Theorem

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

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

- “Proof” works as long as $f + g$ is convex

Nevertheless

- Not possible to always have the decomposition — too ambitious
- More modest goal: decomposition to hold for certain functions
- Manipulating the optimality conditions:

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

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

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

Theorem

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

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

- “Proof” works as long as $f+g$ is convex

Nevertheless

- Not possible to always have the decomposition — too ambitious
- More modest goal: decomposition to hold for certain functions
- Manipulating the optimality conditions:

$$\begin{aligned}P_{f+g}(\mathbf{z}) - \mathbf{z} + \partial(f+g)(P_{f+g}(\mathbf{z})) &\ni 0 \\P_f(P_g(\mathbf{z})) - \mathbf{z} + \partial g(P_g(\mathbf{z})) + \partial f(P_f(P_g(\mathbf{z}))) &\ni 0.\end{aligned}$$

Theorem

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

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

- “Proof” works as long as $f+g$ is convex

Nevertheless

- Not possible to always have the decomposition — too ambitious
- More modest goal: decomposition to hold for certain functions
- Manipulating the optimality conditions:

$$\begin{aligned}P_{f+g}(\mathbf{z}) - \mathbf{z} + \partial(f + g)(P_{f+g}(\mathbf{z})) &\ni 0 \\P_f(P_g(\mathbf{z})) - \mathbf{z} + \partial g(P_g(\mathbf{z})) + \partial f(P_f(P_g(\mathbf{z}))) &\ni 0.\end{aligned}$$

Theorem

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

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

- “Proof” works as long as $f + g$ is convex

The rest is easy



- Find f and g that clinch our sufficient condition.

Recent Results

- More sufficient conditions in (Bauschke and Combettes, 2017)
- (Adly et al., 2017) removes any condition by re-defining one prox

Table of Contents

- 1 Motivation
- 2 Setup
- 3 A Naive Sufficient Condition
- 4 Three Case Studies**
- 5 More Examples

“Trivialities”

Theorem

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

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

Asymmetry.

Theorem

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

- Reassuring the impossibility to always have $P_{f+g} = P_f \circ P_g$;
- Still hope to get interesting results!

Scaling Invariant \Leftrightarrow Positive Homogeneous

$$\partial g(P_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$$

g positive homogeneous $\Leftrightarrow \forall \lambda > 0, \partial g(\lambda \mathbf{w}) = \partial g(\mathbf{w}) \Rightarrow \forall \mathbf{z}, P_f(\mathbf{z}) \propto \mathbf{z}$

Theorem

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

- i). $\exists \lambda \in \mathcal{H} \setminus \{0\}$ for which the scaling function $\lambda \mapsto P_f(\lambda \mathbf{z})$ is linear.
- ii). For all perpendicular $\mathbf{x} \perp \mathbf{y}, P_f(\mathbf{x} + \mathbf{y}) \perp P_f(\mathbf{x})$.
- 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). $\exists \mathbf{w} \in \mathcal{H}$ and $\lambda > 0$ such that $P_f(\lambda \mathbf{w}) = \mathbf{w}$ for all positive homogeneous g .
($\mathbf{w} = \text{argmin}_{\mathbf{z}} \|\mathbf{z}\|$) \Rightarrow \mathbf{w} is the proximal map of $\mathbf{0}$ to \mathcal{H} .)

Scaling Invariant \Leftrightarrow Positive Homogeneous

$$\partial g(P_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$$

g positive homogeneous $\Leftrightarrow \forall \lambda > 0, \partial g(\lambda \mathbf{w}) = \partial g(\mathbf{w}) \Rightarrow \forall \mathbf{z}, P_f(\mathbf{z}) \propto \mathbf{z}$

Theorem

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

- i). $\exists \mathbf{w} \in \mathcal{H} \setminus \{0\}$ for which the equality $P_f(\lambda \mathbf{w}) = \lambda P_f(\mathbf{w})$ holds for all $\lambda > 0$.
- ii). For all perpendicular $\mathbf{z} \perp \mathbf{w}$, $P_f(\mathbf{z} + \mathbf{w}) \perp \mathbf{z}$.
- 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). $\exists \mathbf{w} \in \mathcal{H} \setminus \{0\}$ and $\lambda_{\mathbf{w}} > 0$ such that $P_f(\lambda \mathbf{w}) = \lambda_{\mathbf{w}} \mathbf{w}$ for all positive homogeneous g (i.e. $\partial g(\lambda \mathbf{w}) = \partial g(\mathbf{w})$) comes to hold.

Scaling Invariant \Leftrightarrow Positive Homogeneous

$$\partial g(P_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$$

g positive homogeneous $\Leftrightarrow \forall \lambda > 0, \partial g(\lambda \mathbf{w}) = \partial g(\mathbf{w}) \Rightarrow \forall \mathbf{z}, P_f(\mathbf{z}) \propto \mathbf{z}$

Theorem

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

- i). $\exists \mathbf{w} \in \mathcal{H} \setminus \{0\}$ for which $\partial g(\mathbf{w}) = \text{cone}(\mathbf{w}) = \{\lambda \mathbf{w} \mid \lambda \geq 0\}$.
- ii). For all perpendicular $\mathbf{z} \perp \mathbf{w}, \partial g(\mathbf{z}) = \{0\} = \partial g(\mathbf{w}) \perp \mathbf{z}$.
- 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). $\partial g(\mathbf{z}) = \text{cone}(\mathbf{z})$ and $\partial g(\mathbf{w}) = \text{cone}(\mathbf{w})$ for all positive homogeneous g .

Scaling Invariant \Leftrightarrow Positive Homogeneous

$$\partial g(P_f(y)) \supseteq \partial g(y)$$

g positive homogeneous $\Leftrightarrow \forall \lambda > 0, \partial g(\lambda \mathbf{w}) = \partial g(\mathbf{w}) \Rightarrow \forall \mathbf{z}, P_f(\mathbf{z}) \propto \mathbf{z}$

Theorem

Fix f . 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). For all perpendicular $x \perp y$, $f(x+y) \geq f(y)$;
 - 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). $0 \in \text{dom } f$ and $P_{f+\kappa} = P_f \circ P_{\kappa}$ for all positive homogeneous κ .
- If $\dim(\mathcal{H}) = 1$, only ii) \Rightarrow i) ceases to hold.

Scaling Invariant \Leftrightarrow Positive Homogeneous

$$\partial g(P_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$$

g positive homogeneous $\Leftrightarrow \forall \lambda > 0, \partial g(\lambda \mathbf{w}) = \partial g(\mathbf{w}) \Rightarrow \forall \mathbf{z}, P_f(\mathbf{z}) \propto \mathbf{z}$

Theorem

Fix f . 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). For all perpendicular $\mathbf{x} \perp \mathbf{y}$, $f(\mathbf{x} + \mathbf{y}) \geq f(\mathbf{y})$;
 - 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 positive homogeneous κ .
- If $\dim(\mathcal{H}) = 1$, only ii) \Rightarrow i) ceases to hold.

Scaling Invariant \Leftrightarrow Positive Homogeneous

$$\partial g(P_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$$

g positive homogeneous $\Leftrightarrow \forall \lambda > 0, \partial g(\lambda \mathbf{w}) = \partial g(\mathbf{w}) \Rightarrow \forall \mathbf{z}, P_f(\mathbf{z}) \propto \mathbf{z}$

Theorem

Fix f . 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). For all perpendicular $\mathbf{x} \perp \mathbf{y}$, $f(\mathbf{x} + \mathbf{y}) \geq f(\mathbf{y})$;
 - 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** positive homogeneous κ .
- If $\dim(\mathcal{H}) = 1$, only ii) \implies i) ceases to hold.

Some Implications

Theorem

Fix f . 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). For all perpendicular $\mathbf{x} \perp \mathbf{y}$, $f(\mathbf{x} + \mathbf{y}) \geq f(\mathbf{y})$;
 - 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 positive homogeneous κ .
- If $\dim(\mathcal{H}) = 1$, only ii) \implies i) ceases to hold.

i) \iff ii)

- Characterizing representer theorem (Dinuzzo & Schölkopf'12)

$$\operatorname{argmin} \ell(\langle \mathbf{w}, \mathbf{x}_1 \rangle, \dots, \langle \mathbf{w}, \mathbf{x}_n \rangle) + f(\mathbf{w}) \in \operatorname{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

Some Implications

Theorem

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

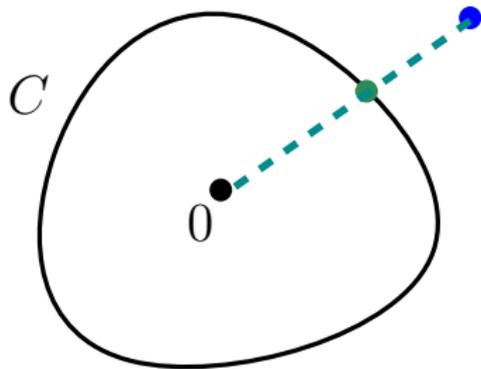
- i). $f = h(\|\cdot\|)$ for some increasing function $h : \mathbb{R}_+ \rightarrow \mathbb{R} \cup \{\infty\}$;
 - ii). For all perpendicular $\mathbf{x} \perp \mathbf{y}$, $f(\mathbf{x} + \mathbf{y}) \geq f(\mathbf{y})$;
 - iii). For all $\mathbf{z} \in \mathcal{H}$, $P_f(\mathbf{z}) = \lambda_{\mathbf{z}} \cdot \mathbf{z}$ 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** positive homogeneous κ .
- If $\dim(\mathcal{H}) = 1$, only ii) \implies i) ceases to hold.

i) \iff ii)

- Characterizing representer theorem (Dinuzzo & Schölkopf'12)

$$\operatorname{argmin} \ell(\langle \mathbf{w}, \mathbf{x}_1 \rangle, \dots, \langle \mathbf{w}, \mathbf{x}_n \rangle) + f(\mathbf{w}) \in \operatorname{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

Characterizing the Ball



Some Implications

Theorem

Fix f . 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). For all perpendicular $x \perp y$, $f(x + y) \geq f(y)$;
- iii). For all $z \in \mathcal{H}$, $P_f(z) = \lambda_z \cdot z$ for some $\lambda_z \in [0, 1]$;
- iv). $\mathbf{0} \in \text{dom } f$ and $P_{f+\kappa} = P_f \circ P_\kappa$ for **all** positive homogeneous κ .

If $\dim(\mathcal{H}) = 1$, only ii) \implies i) ceases to hold.

i) \implies iv)

$$P_{\lambda\|\cdot\|^2 + \kappa} = P_{\lambda\|\cdot\|^2} \circ P_\kappa = \frac{1}{\lambda+1} P_\kappa$$

- Double shrinkage;
- $\kappa = \|\cdot\|_1$: Elastic net (Zou & Hastie'05);
- Adding an l_2 -ish regularizer, computationally, is free.

Some Implications

Theorem

Fix f . 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). For all perpendicular $\mathbf{x} \perp \mathbf{y}$, $f(\mathbf{x} + \mathbf{y}) \geq f(\mathbf{y})$;
 - 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** positive homogeneous κ .
- If $\dim(\mathcal{H}) = 1$, only ii) \implies i) ceases to hold.

i) \implies iv)

$$P_{\lambda\|\cdot\|^2 + \kappa} = P_{\lambda\|\cdot\|^2} \circ P_{\kappa} = \frac{1}{\lambda+1} P_{\kappa}$$

- Double shrinkage;
- $\kappa = \|\cdot\|_1$: Elastic net (Zou & Hastie'05);
- Adding an l_2 -ish regularizer, computationally, is free.

Some Implications

Theorem

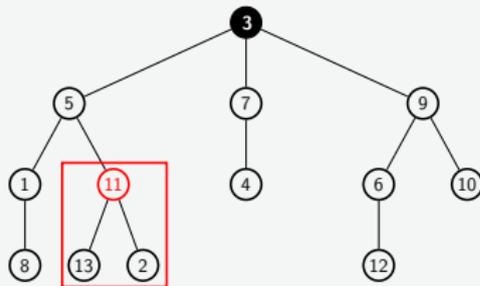
Fix f . The following are equivalent (provided $\dim(\mathcal{H}) \geq 2$):

- $f = h(\|\cdot\|)$ for some increasing function $h : \mathbb{R}_+ \rightarrow \mathbb{R} \cup \{\infty\}$;
- For all perpendicular $x \perp y$, $f(x + y) \geq f(y)$;
- For all $z \in \mathcal{H}$, $P_f(z) = \lambda_z \cdot z$ for some $\lambda_z \in [0, 1]$;
- $\mathbf{0} \in \text{dom } f$ and $P_{f+\kappa} = P_f \circ P_\kappa$ for **all** positive homogeneous κ .

i) \implies iv)

Tree-structured group norms
(Jenatton et al.'11)

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



Some Implications

Theorem

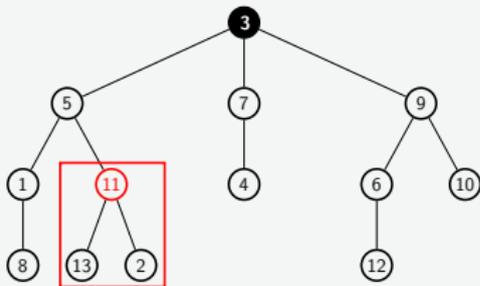
Fix f . The following are equivalent (provided $\dim(\mathcal{H}) \geq 2$):

- i). $f = h(\|\cdot\|)$ for some increasing function $h : \mathbb{R}_+ \rightarrow \mathbb{R} \cup \{\infty\}$;
- ii). For all perpendicular $\mathbf{x} \perp \mathbf{y}$, $f(\mathbf{x} + \mathbf{y}) \geq f(\mathbf{y})$;
- 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** positive homogeneous κ .

i) \implies iv)

Tree-structured group norms
(Jenatton et al.'11)

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



Choquet Integral (a.k.a. Lovász Extension)

For an increasing set function $\mu : 2^{[d]} \rightarrow \mathbb{R}$:

$$g(\mathbf{w}) := \int_0^\infty \mu(\llbracket \mathbf{w} \geq t \rrbracket) dt + \int_{-\infty}^0 [\mu(\llbracket \mathbf{w} \geq t \rrbracket) - \mu([d])] dt,$$

where we treat $\mathbf{w} : \{1, \dots, d\} \rightarrow \mathbb{R}$.

- g is positive homogeneous
- $g(\mathbf{w} + \mathbf{z}) \neq g(\mathbf{w}) + g(\mathbf{z})$ in general
- $g(\mathbf{w} + \mathbf{z}) \leq g(\mathbf{w}) + g(\mathbf{z})$ iff μ is submodular:

$$\mu(A \cap B) + \mu(A \cup B) \leq \mu(A) + \mu(B)$$

- if $\forall i, j, (w_i - w_j)(z_i - z_j) \geq 0$, then $g(\mathbf{w} + \mathbf{z}) = g(\mathbf{w}) + g(\mathbf{z})$
- $\min_{A \subseteq [d]} \mu(A) = \min_{\mathbf{w} \in [0,1]^d} g(\mathbf{w})$.

Further Properties of Choquet Integral

Theorem

Let g be the Choquet integral of some submodular function. If for all i and j ,

- $(w_i - w_j)(z_i - z_j) \geq 0$, then $\partial g(\mathbf{w}) \cap \partial g(\mathbf{z}) \neq \emptyset$
- $w_i \geq w_j \implies z_i \geq z_j$, then $\partial g(\mathbf{w}) \subseteq \partial g(\mathbf{z})$.

Theorem (Schmeidler'86)

If g is comonotone additive and increasing/continuous, then g is a Choquet integral of some set function.

TV is a Choquet integral

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

Permutation Invariant \Leftrightarrow Choquet Integral

$$\partial g(P_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$$

- For permutation-invariant f , recall

$$P_f(\mathbf{y}) = \operatorname{argmin}_x \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2 + f(\mathbf{x}).$$

By rearrangement inequality

$$y_i \geq y_j \implies [P_f(\mathbf{y})]_i \geq [P_f(\mathbf{y})]_j$$

Theorem

Let f be permutation invariant and g be the Choquet integral of some submodular set function μ . Then, $P_{f+g} = P_f \circ P_g$.

Permutation Invariant \Leftrightarrow Choquet Integral

$$\partial g(P_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$$

- For permutation-invariant f , recall

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

By rearrangement inequality

$$y_i \geq y_j \implies [P_f(\mathbf{y})]_i \geq [P_f(\mathbf{y})]_j$$

Theorem

Let f be permutation invariant and g be the Choquet integral of some submodular set function μ . Then, $P_{f+g} = P_f \circ P_g$.

Permutation Invariant \Leftrightarrow Choquet Integral

$$\partial g(P_f(y)) \supseteq \partial g(y)$$

- For permutation-invariant f , recall

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

By rearrangement inequality

$$y_i \geq y_j \implies [P_f(\mathbf{y})]_i \geq [P_f(\mathbf{y})]_j$$

Theorem

Let f be permutation invariant and g be the Choquet integral of some submodular set function μ . Then, $P_{f+g} = P_f \circ P_g$.

Some Implications

Theorem

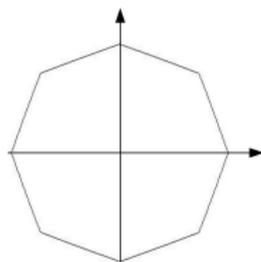
Let f be permutation invariant and g be the Choquet integral of some submodular set function. Then, $P_{f+g} = P_f \circ P_g$.

- Special case $f = \|\cdot\|_1$ in (Bach'11);
- $P_{\|\cdot\|_1 + \|\cdot\|_{TV}} = P_{\|\cdot\|_1} \circ P_{\|\cdot\|_{TV}}$ (Friedman et al.'07);
- $P_{\sum_{i=1}^k \|\cdot\|_{g_i}} = P_{\|\cdot\|_{g_1}} \circ \dots \circ P_{\|\cdot\|_{g_k}}$ (Jenatton et al.'11)

Some Implications

$$\|\mathbf{w}\|_{\text{oscar}} = \sum_{i < j} \max\{|w_i|, |w_j|\}.$$

- Feature grouping (Bondell & Reich'08)
- $P_{\|\cdot\|_{\text{oscar}}}$ in (Zhong & Kwok'11)



Let

$$\kappa_i(\mathbf{w}) := \sum_{j:j < i} \max\{|w_i|, |w_j|\}.$$

- $\|\mathbf{w}\|_{\text{oscar}} = \sum_{i=2}^d \kappa_i(\mathbf{w})$
- $P_{\|\cdot\|_{\text{oscar}}} = P_{\kappa_d} \circ \dots \circ P_{\kappa_2}$
- Given P_{κ_i} , constant time for $P_{\kappa_{i+1}}$.

Table of Contents

- 1 Motivation
- 2 Setup
- 3 A Naive Sufficient Condition
- 4 Three Case Studies
- 5 More Examples**

Projection to intersection

Theorem (Barty, Roy and Strugarek, MOR'07, Proposition 3.1)

Let $L \cap C \neq \emptyset$, where C is a closed convex set and L is a subspace. If $P_C(L) \subseteq L$, then $P_{L \cap C} = P_C \circ P_L$.

$f = \iota_C$ and $g = \iota_L$, follows from $\partial g \equiv L^\perp$.

One Solution for All, I

Theorem (Chambolle and Darbon, 2009)

Let $\varphi_i : \mathbb{R} \rightarrow \mathbb{R} \cup \{\infty\}$, $i = 1, \dots, d$, be closed convex univariate functions and $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be the Choquet integral of the set function μ . Let

$$\mathbf{u} \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + \sum_{i=1}^d \varphi_i(w_i), \quad (1)$$

whose existence is assumed. For any $t \in \bigcap_i \operatorname{dom} \partial \varphi_i$, consider the discrete problem:

$$\min_{A \subseteq [d]} F(A) + \sum_{i \in A} \varphi'_i(t). \quad (2)$$

- If for all i , $\varphi'_i(t)$ is the smallest element in the subdifferential $\partial \varphi_i(t)$ (existence assumed), then the set $[\mathbf{u} \geq t]$ solves (2).
- If for all i , $\varphi'_i(t)$ is the largest element in the subdifferential $\partial \varphi_i(t)$ (existence assumed), then the set $[\mathbf{u} > t]$ solves (2).

One Solution for All, II

Theorem (extending (Barlow and Brunk, 1972))

Let f be univariate convex and differentiable, with the induced Bregman divergence $D_f(x, y) := f(x) - f(y) - f'(y)(x - y)$. For any Choquet integral g , the following problem

$$\min_{\mathbf{x} \in \mathbb{R}^p} \sum_{i=1}^p w_i D_f(x_i, y_i) + g(\mathbf{x}) \quad (3)$$

can be solved in two steps:

$$\begin{aligned} \textcircled{1} \quad \mathbf{z} &= \operatorname{argmin}_{\mathbf{x}} \frac{1}{2} \sum_{i=1}^p w_i (x_i - f'(y_i))^2 + g(\mathbf{x}) \\ \textcircled{2} \quad \mathbf{y}^* &= (f')^{-1}(\mathbf{z}), \end{aligned}$$

Summary

- Posed the question: $P_{f+g} \stackrel{?}{=} P_f \circ P_g \stackrel{?}{=} P_g \circ P_f$;
- Presented a sufficient condition: $\partial g(P_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$;
- Identified two major cases;
- Immediately useful if plugged into splitting algs;