

CS794/CO673: Optimization for Data Science

Lec 16: Splitting

Yaoliang Yu



UNIVERSITY OF
WATERLOO

FACULTY OF MATHEMATICS
DAVID R. CHERITON SCHOOL
OF COMPUTER SCIENCE

November 11, 2022

Problem

Find zero of a maximal monotone map $T : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$:

$$\text{find } \mathbf{z} \text{ s.t. } \mathbf{0} \in T\mathbf{z}, \quad \text{where } T = A + B$$

- $A, B : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$ can both be multi-valued
- Allow a unified treatment of many algorithms
- Dual problem:

$$\text{find } \mathbf{z}^* \text{ s.t. } \mathbf{0} \in T^*\mathbf{z}^*, \quad \text{where } T^* := [-A^{-1} \circ (-\text{Id}) + B^{-1}]$$

Some Definitions

- A multi-valued map $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is monotone iff for all $\mathbf{w}, \mathbf{z} \in \text{dom } T$, and all $\mathbf{w}^* \in T\mathbf{w}, \mathbf{z}^* \in T\mathbf{z}$,

$$\langle \mathbf{w} - \mathbf{z}, \mathbf{w}^* - \mathbf{z}^* \rangle \geq 0$$

- A monotone map is maximal if its graph is not contained in that of another monotone map:

$$\text{gph } T := \{(\mathbf{z}, \mathbf{z}^*) : \mathbf{z}^* \in T\mathbf{z}\}$$

- Inverse: $\mathbf{z}^* \in T\mathbf{z} \iff \mathbf{z} \in T^{-1}\mathbf{z}^*$
- Resolvent: $J_T^\eta := (\text{Id} + \eta T)^{-1}$, i.e. $\mathbf{z}^* \in J_T^\eta \mathbf{z} \iff \mathbf{z} \in \mathbf{z}^* + \eta T\mathbf{z}^*$
- Sum: $(A + B)(\mathbf{z}) = \{\mathbf{u}^* + \mathbf{v}^* : \mathbf{u}^* \in A\mathbf{z}, \mathbf{v}^* \in B\mathbf{z}\}$

Example: Convex minimization

Subdifferential ∂f of a convex function is maximal monotone.

$$\mathbf{0} \in \partial f(\mathbf{z}) \iff \mathbf{z} \in \operatorname{argmin} f.$$

$$J_{\partial f}^{\eta}(\mathbf{z}) = P_f^{\eta}(\mathbf{z}) = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2\eta} \|\mathbf{w} - \mathbf{z}\|_2^2 + f(\mathbf{w})$$

Example: Convex composite minimization

For convex f and g , let $A = \partial f$, $B = \partial g$ and $T = A + B = \partial f + \partial g$.

$$\mathbf{0} \in T\mathbf{z} \iff \mathbf{z} \in \operatorname{argmin} f + g$$

Example: Convex-concave minimax

For $f(\mathbf{x}, \mathbf{y})$ convex in \mathbf{x} and concave in \mathbf{y} , $T = (\partial_{\mathbf{x}}f, \partial_{\mathbf{y}}-f)$ is maximal monotone.

$\mathbf{0} \in T(\mathbf{x}, \mathbf{y}) \iff (\mathbf{x}, \mathbf{y}) \in \operatorname{argmin} \operatorname{argmax} f$, i.e. a Nash equilibrium.

- $\mathbf{0} \in \partial_{\mathbf{x}}f(\mathbf{x}, \mathbf{y}) \iff$ fixing \mathbf{y} , $\mathbf{x} \in \operatorname{argmin} f(\cdot, \mathbf{y})$
- $\mathbf{0} \in \partial_{\mathbf{y}}-f(\mathbf{x}, \mathbf{y}) \iff$ fixing \mathbf{x} , $\mathbf{y} \in \operatorname{argmax} f(\mathbf{x}, \cdot)$

Example: Constrained convex-concave minimax

For $f(\mathbf{x}, \mathbf{y})$ convex in \mathbf{x} and concave in \mathbf{y} , $g(\mathbf{x})$ convex and $h(\mathbf{y})$ concave, let $A = (\partial_{\mathbf{x}}f, \partial_{\mathbf{y}}-f)$ and $B = (\partial_{\mathbf{x}}g, \partial_{\mathbf{y}}-h)$.

$\mathbf{0} \in T(\mathbf{x}, \mathbf{y}) \iff (\mathbf{x}, \mathbf{y}) \in \operatorname{argmin}_{\mathbf{x}} \operatorname{argmax}_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) + g(\mathbf{x}) + h(\mathbf{y})$, i.e. a Nash equilibrium.

A Product Space Trick

$$\min_{\mathbf{w}} \sum_{i=1}^m f_i(\mathbf{w})$$

- m users, each having an objective f_i to optimize
- Users share the same model \mathbf{w}
- Introduce local copies $\mathbf{w}_1, \dots, \mathbf{w}_m$, and add consensus constraint:

$$\min_{\mathbf{w}_1, \dots, \mathbf{w}_m} \sum_{i=1}^m f_i(\mathbf{w}_i) + \iota_H(\mathbf{w}_1, \dots, \mathbf{w}_m)$$

$$- H := \{(\mathbf{w}_1, \dots, \mathbf{w}_m) : \mathbf{w}_1 = \mathbf{w}_2 = \dots = \mathbf{w}_m\}$$

- $J_{\partial \iota_H}(\mathbf{w}_1, \dots, \mathbf{w}_m) = \text{P}_H(\mathbf{w}_1, \dots, \mathbf{w}_m) = (\bar{\mathbf{w}}, \bar{\mathbf{w}}, \dots, \bar{\mathbf{w}})$, where $\bar{\mathbf{w}} := \frac{1}{m} \sum_{i=1}^m \mathbf{w}_i$

Algorithm 1: Forward-Backward Splitting

Input: $\mathbf{w}_0 \in \text{dom } A$

```
1 for  $t = 0, 1, 2, \dots$  do
2   choose any  $\mathbf{a}_t^* \in A\mathbf{w}_t$ 
3    $\mathbf{w}_{t+1} \leftarrow J_B^{\eta_t}(\mathbf{w}_t - \eta_t \mathbf{a}_t^*)$  //  $\eta_t \geq 0$  is the step size
4    $\mathbf{z}_t \leftarrow \sum_{k=0}^t \bar{\eta}_{t,k} \mathbf{w}_k$  // ergodic averaging,  $\bar{\eta}_{t,k} := \eta_k / H_t$ ,  $H_t := \sum_{k=0}^t \eta_k$ 
```

G. B. Passty. "Ergodic convergence to a zero of the sum of monotone operators in Hilbert space". *Journal of Mathematical Analysis and Applications*, vol. 72, no. 2 (1979), pp. 383–390, R. E. Bruck. "On the weak convergence of an ergodic iteration for the solution of variational inequalities for monotone operators in Hilbert space". *Journal of Mathematical Analysis and Applications*, vol. 61, no. 1 (1977), pp. 159–164.

Theorem:

For any $(\mathbf{w}, \mathbf{w}^*) \in \text{gph } T$ and $\mathbf{b}^* \in B\mathbf{w}$:

$$\langle \mathbf{z}_t - \mathbf{w}, \mathbf{w}^* \rangle \leq \sum_{k=0}^t \bar{\eta}_{t,k} \langle \mathbf{w}_k - \mathbf{w}, \mathbf{a}_k^* + \mathbf{b}^* \rangle \leq \frac{\|\mathbf{w}_0 - \mathbf{w}\|_2^2 + \sum_{k=0}^t \eta_k^2 \|\mathbf{a}_k^* + \mathbf{b}^*\|_2^2}{2H_t}.$$

- If $\sum_t \|\eta_t \mathbf{a}_t^* + \mathbf{b}^*\|_2^2 < \infty$ and $H_t \rightarrow \infty$, then either no solution, in which case $\|\mathbf{z}_t\| \rightarrow \infty$; or \mathbf{z}_t converges to a solution

Gradient-Descent-Ascent (GDA)

For simplicity, suppose we can choose $\mathbf{b}^* = \mathbf{0}$:

$$\begin{aligned} \frac{\|\mathbf{w}_0 - \mathbf{w}\|_2^2 + \sum_{k=0}^t \|\eta_k \mathbf{a}_k^*\|_2^2}{2H_t} &\geq \sum_{k=0}^t \bar{\eta}_{t,k} \langle \mathbf{w}_k - \mathbf{w}, \mathbf{a}_k^* \rangle \\ &= \sum_{k=0}^t \bar{\eta}_{t,k} [\langle \mathbf{x}_k - \mathbf{x}, \partial_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k) \rangle + \langle \mathbf{y}_k - \mathbf{y}, \partial_{\mathbf{y}} f(\mathbf{x}_k, \mathbf{y}_k) \rangle] \\ &\geq \sum_{k=0}^t \bar{\eta}_{t,k} [-f(\mathbf{x}, \mathbf{y}_k) + f(\mathbf{x}_k, \mathbf{y})] \\ &\geq -f(\mathbf{x}, \bar{\mathbf{y}}_t) + f(\bar{\mathbf{x}}_t, \mathbf{y}), \quad \text{where } (\bar{\mathbf{x}}_t, \bar{\mathbf{y}}_t) := \sum_{k=0}^t \bar{\eta}_{t,k} \mathbf{w}_k \end{aligned}$$

To satisfy $\sum_t \|\eta_t \mathbf{a}_t^*\|_2^2 < \infty$ we may choose $\eta_t = \frac{1}{\sqrt{\|\mathbf{a}_t^*\|_2^2 + 1}} \frac{1}{(t+1)^p}$, $p \in (\frac{1}{2}, 1]$

A. S. Nemirovskii and D. B. Judin. "Cesari convergence of the gradient method of approximating saddle points of convex-concave functions". *Soviet Mathematics Doklady*, vol. 19, no. 2 (1978), pp. 482–486.

Algorithm 2: Backward-Backward Splitting

Input: $\mathbf{w}_0 \in \text{dom} A$

1 for $t = 0, 1, 2, \dots$ do

2 $\mathbf{w}_{t+1} \leftarrow J_B^{\eta_t} J_A^{\eta_t} \mathbf{w}_t$ // $\eta_t \geq 0$ is the step size

3 $\mathbf{z}_t \leftarrow \sum_{k=0}^t \bar{\eta}_{t,k} \mathbf{w}_k$ // ergodic averaging, $\bar{\eta}_{t,k} := \eta_k / H_t$, $H_t := \sum_{k=0}^t \eta_k$

Theorem:

If $\sum_t \eta_t = \infty$ and $\eta_t \rightarrow 0$, then either there is no solution, in which case $\|\mathbf{z}_t\| \rightarrow \infty$, or \mathbf{z}_t converges to a solution.

G. B. Passty. "Ergodic convergence to a zero of the sum of monotone operators in Hilbert space". *Journal of Mathematical Analysis and Applications*, vol. 72, no. 2 (1979), pp. 383–390, P.-L. Lions. "Une methode iterative de resolution d'une inequtation variationnelle". *Israel Journal of Mathematics*, vol. 31, no. 2 (1978), pp. 204–208.

Comparing FB with BB

$$\mathfrak{F} := J_B^\eta(\text{Id} - \eta A) \quad \text{vs.} \quad \mathfrak{B} := J_B^\eta J_A^\eta, \quad \text{where} \quad J_T^\eta := (\text{Id} + \eta T)^{-1}$$

- Let us define ${}^\eta A := \frac{\text{Id} - J_A^\eta}{\eta}$ such that ${}^\eta A \rightarrow {}^0 A \subseteq A$ if $\eta \rightarrow 0$
- Then, backward-backward on $A + B$ is forward-backward on ${}^\eta A + B!$
 - $\text{Id} - \eta \cdot {}^\eta A = J_A^\eta$
- Consider $A = \partial f$ for a convex function f :
 - $J_A^\eta = P_f^\eta$
 - ${}^\eta A = \frac{\text{Id} - P_f^\eta}{\eta} = \nabla M_f^\eta$
 - backward-backward on $f + g$ is the same as forward-backward on $M_f^\eta + g!$
- For small η , $J_A^\eta = (\text{Id} + \eta A)^{-1} \approx \text{Id} - \eta A$ when A is linear

Algorithm 3: The Barycenter Method

Input: \mathbf{w}_0

1 **for** $t = 0, 1, 2, \dots$ **do**
2 $\mathbf{w}_{t+1} \leftarrow \text{Avg}(P_{i_1}, \dots, P_{i_{k(t)}})\mathbf{w}_t$

- Let $H_i := \{\mathbf{w} : \langle \mathbf{w}, \mathbf{a}_i \rangle = b_i\}$ and P_i be the orthogonal projection onto it
- Kaczmarz's method: $k(t) \equiv 1$
- Cimmino's method: $k(t) \equiv m$
- A randomized version of backward-backward!

S. Kaczmarz. "Angenäherte Auflösung von Systemen linearer Gleichungen". *Bulletin International de l'Académie Polonaise des Sciences et des Lettres*, vol. 35 (1937), pp. 355–357. "Approximate solution of systems of linear equations", English translation in *International Journal of Control*, 1993, vol. 57, no.6, pp. 1269–1271, G. Cimmino. "Calcolo Approssimato Per le Soluzioni dei Sistemi di Equazioni Lineari". *La Ricerca Scientifica*, vol. 9, no. 1 (1938), pp. 326–333.

$$\begin{cases} x + y = 0 \\ x - y = 0 \end{cases}$$

- $P_1(x, y) = \left(\frac{x-y}{2}, \frac{y-x}{2} \right)$
- $P_2(x, y) = \left(\frac{x+y}{2}, \frac{y+x}{2} \right)$
- $\text{Avg}(P_1, P_2)(x, y) = \left(\frac{x}{2}, \frac{y}{2} \right) \rightarrow (0, 0)$

Algorithm 4: Extragradient for finding a zero of $T = A + B$

Input: $\mathbf{w}_0 \in \text{dom } A \subseteq \text{dom } B$

```
1 for  $t = 0, 1, 2, \dots$  do  
2   choose step size  $\eta_t > 0$   
3    $\tilde{\mathbf{w}}_t = J_A^{\eta_t}(\mathbf{w}_t - \eta_t B \mathbf{w}_t)$  // peek  
4    $\mathbf{w}_{t+1} = J_A^{\eta_t}(\mathbf{w}_t - \eta_t B \tilde{\mathbf{w}}_t)$  // update with peeked information
```

Theorem: Convergence of extra-gradient

Let $B : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$ be L -Lipschitz and $\text{dom } A \subseteq \text{dom } B$. Set $\eta_t \in [0, 1/L]$. Then, for all $(\mathbf{w}, \mathbf{w}^*) \in \text{gph } T$, $\mathbf{a}^* \in A\mathbf{w}$,

$$\langle \tilde{\mathbf{z}}_t - \mathbf{w}, \mathbf{w}^* \rangle \leq \sum_{k=0}^t \bar{\eta}_{t,k} \langle \tilde{\mathbf{w}}_k - \mathbf{w}, B\tilde{\mathbf{w}}_k + \mathbf{a}^* \rangle \leq \frac{\|\mathbf{w}_0 - \mathbf{w}\|_2^2}{2H_t}, \quad \text{where}$$

$$\tilde{\mathbf{z}}_t := \sum_{k=0}^t \bar{\eta}_{t,k} \tilde{\mathbf{w}}_k, \quad \bar{\eta}_{t,k} := \eta_k / H_t, \quad H_t := \sum_{k=0}^t \eta_k.$$

- If $H_t \rightarrow \infty$, then either no solution and $\|\tilde{\mathbf{z}}_t\| \rightarrow \infty$, or $\tilde{\mathbf{z}}_t$ converges to a solution
- If $0 < \liminf_t \eta_t \leq \limsup_t \eta_t < 1/L$ and assume a solution exists, then $\mathbf{w}_t - \tilde{\mathbf{w}}_t \rightarrow 0$ and $\tilde{\mathbf{w}}_t$ converges to a solution
- With $\eta_t \equiv \eta$ we obtain $O(1/t)$ rate of convergence for the averaged sequence $\tilde{\mathbf{z}}_t$
- Significantly faster than FB and BB (at best $O(1/\sqrt{t})$), under the additional assumption that B is Lipschitz continuous
- The direct sequence \mathbf{w}_t converges at the slower $O(1/\sqrt{t})$ rate!

Algorithm 5: Khobotov's linear search

```
1  $\eta_t \leftarrow 2\bar{\eta}$ 
2  $\mathbf{w}_t^* \leftarrow \mathbf{B}\mathbf{w}_t$ 
3 repeat
4    $\eta_t \leftarrow \eta_t/2$ 
5    $\tilde{\mathbf{w}}_t \leftarrow J_A^{\eta_t}(\mathbf{w}_t - \eta_t\mathbf{w}_t^*)$ 
6 until  $\eta_t \leq \gamma \frac{\|\mathbf{w}_t - \tilde{\mathbf{w}}_t\|_2}{\|\mathbf{w}_t^* - \mathbf{B}\tilde{\mathbf{w}}_t\|_2}$ 
```

- Inspecting the proof to see where L is used
- A constant number of line searches \implies same convergence rate

Algorithm 6: Past extragradient for solving a smooth monotone VI

Input: $\mathbf{w}_0 = \tilde{\mathbf{w}}_{-1} \in C \subseteq \text{dom } T$

```
1 for  $t = 0, 1, 2, \dots$  do
2   choose step size  $\eta_t > 0$ 
3    $\tilde{\mathbf{w}}_t = P_C(\mathbf{w}_t - \eta_t T \tilde{\mathbf{w}}_{t-1})$ 
4    $\mathbf{w}_{t+1} = P_C(\mathbf{w}_t - \eta_t T \tilde{\mathbf{w}}_t)$ 
```

- Only requires **1** evaluation of the operator T
- But still **2** projections per step
- Compared to the extragradient, simply recycle the past evaluation $T \tilde{\mathbf{w}}_{t-1}$ to replace $T \mathbf{w}_t$, saving us 1 evaluation of T

Algorithm 7: Tseng's modified forward-backward splitting

Input: $\mathbf{w}_0 \in C \subseteq \text{dom } T$

```
1 for  $t = 0, 1, 2, \dots$  do
2   choose step size  $\eta_t > 0$ 
3    $\tilde{\mathbf{w}}_t = P_C(\mathbf{w}_t - \eta_t T \mathbf{w}_t)$ 
4    $\mathbf{w}_{t+1} = \tilde{\mathbf{w}}_t - \eta_t (T \tilde{\mathbf{w}}_t - T \mathbf{w}_t)$ 
```

- Only requires 1 projection
- But still 2 evaluations of T per step
- This variant requires say $\text{dom } T = \mathbb{R}^d$

Algorithm 8: Optimistic extragradient for solving a smooth monotone VI

Input: $\mathbf{w}_0 = \tilde{\mathbf{w}}_{-1} \in C \subseteq \text{dom } T$

```
1 for  $t = 0, 1, 2, \dots$  do
2   choose step size  $\eta_t > 0$ 
3    $\tilde{\mathbf{w}}_t = P_C(\mathbf{w}_t - \eta_t T \tilde{\mathbf{w}}_{t-1})$ 
4    $\mathbf{w}_{t+1} = \tilde{\mathbf{w}}_t - \eta_t (T \tilde{\mathbf{w}}_t - T \tilde{\mathbf{w}}_{t-1})$ 
```

- Obviously, we can combine the previous two ideas
- Obtain a variant that only requires **1** projection and **1** evaluation of T per step!

Algorithm 9: Reflected extragradient for solving a smooth monotone VI

Input: $\mathbf{w}_0 = \mathbf{w}_{-1} \in C \subseteq \text{dom } T$

```
1 for  $t = 0, 1, 2, \dots$  do
2   choose step size  $\eta_t > 0$ 
3    $\tilde{\mathbf{w}}_t = 2\mathbf{w}_t - \mathbf{w}_{t-1}$ 
4    $\mathbf{w}_{t+1} = P_C(\mathbf{w}_t - \eta_t T \tilde{\mathbf{w}}_t)$ 
```

- Uses reflection and also enjoys 1 projection and 1 evaluation of T per step
- Requires say $\text{dom } T \supseteq 2C - C$

Reflectors

- $R_{\top}^{\eta} := 2J_{\top}^{\eta} - \text{Id} \subseteq (\text{Id} - \eta\top)(\text{Id} + \eta\top)^{-1}$
 - a backward step, followed by a forward step
- $(\text{Id} - \eta\top)(\text{Id} + \eta\top)^{-1} \neq (\text{Id} + \eta\top)^{-1}(\text{Id} - \eta\top) =$
 $(\text{Id} + \eta\top)^{-1}(\text{Id} - \eta\top)(\text{Id} + \eta\top)^{-1}(\text{Id} + \eta\top)$
- $\mathbf{w} \in R_{\top}^{\eta}\mathbf{w} \iff \mathbf{w} \in J_{\top}^{\eta}(\mathbf{w}) \iff \mathbf{0} \in \top\mathbf{w}$

Algorithm 10: A general splitting algorithm based on reflectors

Input: \mathbf{w}_0

```
1 for  $t = 0, 1, \dots$  do
2   choose step size  $\eta_t \geq 0$  and relaxation size  $\gamma_t \in [0, 1]$ 
3    $\mathbf{w}_{t+1} = (1 - \gamma_t)\mathbf{w}_t + \gamma_t R_B^{\eta_t} R_A^{\eta_t}(\mathbf{w}_t) + \epsilon_t$  // allow error  $\epsilon_t$ 
4 return  $\mathbf{z} \in J_A^\eta \mathbf{w}$  // assuming the for-loop returns  $\mathbf{w}$  and  $\eta_t \equiv \eta$ 
```

- $\gamma_t \equiv 1$: Peaceman-Rachford (PR) splitting
- $\gamma_t \equiv \frac{1}{2}$: Douglas-Rachford (DR) splitting
- $\gamma_t \in [\frac{1}{2}, 1]$: over-relaxation; $\gamma_t \in [0, \frac{1}{2}]$: under-relaxation

General Convergence

Let $A, B : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$ be maximal monotone and $T = A + B$. Suppose $\sum_t \gamma_t(1 - \gamma_t) = \infty$, $\eta_t \equiv \eta > 0$, and $\sum_t \|\epsilon_t\|_2 < \infty$.

- If $T^{-1}\mathbf{0} = (A + B)^{-1}\mathbf{0} = \emptyset$, then $\{\mathbf{w}_t\}$ is unbounded.
- If $T^{-1}\mathbf{0} \neq \emptyset$, then $\mathbf{w}_t \rightarrow \mathbf{w}_\infty = R_B^\eta R_A^\eta \mathbf{w}_\infty$, $\mathbf{z}_t := J_A^\eta \mathbf{w}_t \rightarrow J_A^\eta \mathbf{w}_\infty \in T^{-1}\mathbf{0}$ and $\mathbf{z}_t^* := {}^\eta A \mathbf{w}_t \rightarrow {}^\eta A \mathbf{w}_\infty \in T^{*-1}\mathbf{0}$, where recall that $T^* := -B^{-1} \circ (-\text{Id}) + A^{-1}$ and ${}^\eta A := (\text{Id} - J_A^\eta)/\eta$.

Convergence of PR

Let $A, B : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$ be maximal monotone and $T = A + B$. Consider PR with $\gamma_t \equiv 1$, $\eta_t \equiv \eta > 0$, and $\sum_t \|\epsilon_t\|_2 < \infty$. Assume A is strictly monotone:

- There exists at most one zero \mathbf{z} of T
- If $\mathbf{z} = T^{-1}\mathbf{0} = (A + B)^{-1}\mathbf{0}$ exists, then $\mathbf{z}_t := J_A^\eta \mathbf{w}_t \rightharpoonup \mathbf{z}$
- If $T^{-1}\mathbf{0} = (A + B)^{-1}\mathbf{0} = \emptyset$, then $\{\mathbf{w}_t\}$ is unbounded

Non-convergence of PR

Consider the rotation in \mathbb{R}^2 :

$$A = B = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}, \quad \text{where } (A + B)\mathbf{z} = \mathbf{0} \iff \mathbf{z} = \mathbf{0}.$$

Since $\langle \mathbf{w}, A\mathbf{w} \rangle = 0$ for all \mathbf{w} , we know A is maximal monotone. We have

$$J_A^\eta = \frac{1}{(1 + \eta)^2} \begin{bmatrix} 1 & \eta \\ -\eta & 1 \end{bmatrix}, \quad R_A^\eta = \frac{1}{(1 + \eta)^2} \begin{bmatrix} 1 - \eta^2 & 2\eta \\ -2\eta & 1 - \eta^2 \end{bmatrix}.$$

Since both J_A^η and R_A^η are rotations (i.e. $\det = 1$), with $\gamma_t \equiv 1$ and $\epsilon_t \equiv \mathbf{0}$, $\|\mathbf{z}_t\|_2 \equiv \|\mathbf{w}_0\|_2$ hence may not converge to any point (hence also $\mathbf{0}$, the unique zero of $A + B$). Moreover, \mathbf{w}_t may not converge (to any point) either.

We verify that A is not strictly monotone and convince yourself that \mathbf{z}_t and \mathbf{w}_t do converge if we choose say $\gamma_t \equiv \gamma \in (0, 1)$.

Alternating Direction Method of Multipliers (ADMM)

Consider the generic minimization problem

$$\inf_{\mathbf{a}} g(L\mathbf{a}) + h(\mathbf{a}), \text{ or equivalently } \inf_{\mathbf{a}, \mathbf{b}} g(\mathbf{b}) + h(\mathbf{a}), \quad \text{s.t. } L\mathbf{a} = \mathbf{b},$$

and its Fenchel-Rockafellar dual

$$-\inf_{\boldsymbol{\mu}} h^*(-L^\top \boldsymbol{\mu}) + g^*(\boldsymbol{\mu}),$$

where $g : \mathbb{R}^p \rightarrow \mathbb{R} \cup \{\infty\}$ and $h : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ are closed proper convex and $L : \mathbb{R}^d \rightarrow \mathbb{R}^p$ is linear.

Introducing the Lagrangian multiplier $\boldsymbol{\mu}$ we obtain the **Lagrangian**:

$$\inf_{\mathbf{a}, \mathbf{b}} \sup_{\boldsymbol{\mu}} g(\mathbf{b}) + h(\mathbf{a}) + \langle L\mathbf{a} - \mathbf{b}; \boldsymbol{\mu} \rangle = \sup_{\boldsymbol{\mu}} \underbrace{\inf_{\mathbf{a}, \mathbf{b}} g(\mathbf{b}) + h(\mathbf{a}) + \langle L\mathbf{a} - \mathbf{b}; \boldsymbol{\mu} \rangle}_{\underline{L}(\boldsymbol{\mu})}.$$

- May apply Uzawa's algorithm to maximize the dual function $\underline{L}(\boldsymbol{\mu})$
- Nonsmooth hence requires diminishing step sizes

Instead, consider the **augmented** Lagrangian, where the penalty parameter η need **not** increase to ∞ :

$$\sup_{\boldsymbol{\mu}} \left[\inf_{\mathbf{a}, \mathbf{b}} g(\mathbf{b}) + h(\mathbf{a}) + \langle L\mathbf{a} - \mathbf{b}; \boldsymbol{\mu} \rangle + \frac{\eta}{2} \|L\mathbf{a} - \mathbf{b}\|_2^2 \right] \equiv \sup_{\boldsymbol{\mu}} \left[\sup_{\boldsymbol{\nu}} \underline{L}(\boldsymbol{\nu}) - \frac{1}{2\eta} \|\boldsymbol{\nu} - \boldsymbol{\mu}\|_2^2 \right],$$

- **smooth** inner function
- May apply Uzawa's algorithm with constant step size η :

$$\boldsymbol{\mu}_{t+1} \leftarrow \boldsymbol{\mu}_t + \eta(L\mathbf{a}_{t+1} - \mathbf{b}_{t+1})$$

- Solving \mathbf{a} and \mathbf{b} simultaneously in the augmented Lagrangian is challenging
- Fortunately, may apply **just one step of alternating minimization** to \mathbf{a} and \mathbf{b} sequentially:

$$\mathbf{a}_{t+1} \in \underset{\mathbf{a}}{\operatorname{argmin}} h(\mathbf{a}) + \langle L\mathbf{a} - \mathbf{b}_t; \boldsymbol{\mu}_t \rangle + \frac{\eta}{2} \|L\mathbf{a} - \mathbf{b}_t\|_2^2 \equiv h(\mathbf{a}) + \frac{\eta}{2} \|L\mathbf{a} - \mathbf{b}_t + \boldsymbol{\mu}_t/\eta\|_2^2$$

$$\mathbf{b}_{t+1} = \underset{\mathbf{b}}{\operatorname{argmin}} g(\mathbf{b}) + \langle L\mathbf{a}_{t+1} - \mathbf{b}; \boldsymbol{\mu}_t \rangle + \frac{\eta}{2} \|L\mathbf{a}_{t+1} - \mathbf{b}\|_2^2 \equiv g(\mathbf{b}) + \frac{\eta}{2} \|L\mathbf{a}_{t+1} - \mathbf{b} + \boldsymbol{\mu}_t/\eta\|_2^2$$

To understand the above updates, let us apply the Fenchel-Rockafellar duality again:

$$\mathbf{a}_{t+1}^* - \eta L\mathbf{a}_{t+1} = -\eta \mathbf{b}_t + \boldsymbol{\mu}_t, \quad \mathbf{a}_{t+1}^* = \underset{\mathbf{a}^*}{\operatorname{argmin}} \frac{1}{2\eta} \|\mathbf{a}^* + \eta \mathbf{b}_t - \boldsymbol{\mu}_t\|_2^2 + h^*(-L^\top \mathbf{a}^*)$$

$$\mathbf{b}_{t+1}^* + \eta \mathbf{b}_{t+1} = \eta L\mathbf{a}_{t+1} + \boldsymbol{\mu}_t, \quad \mathbf{b}_{t+1}^* = \underset{\mathbf{b}^*}{\operatorname{argmin}} \frac{1}{2\eta} \|\mathbf{b}^* - \eta L\mathbf{a}_{t+1} - \boldsymbol{\mu}_t\|_2^2 + g^*(\mathbf{b}^*)$$

It can be shown that $\boldsymbol{\mu}_t = \mathbf{b}_t^*$.

From the optimality conditions of \mathbf{b}_{t+1}^* and \mathbf{a}_{t+1} we verify that:

$$\begin{aligned} (\mathbf{b}_{t+1}^*, \mathbf{b}_{t+1}) \in \text{gph } \partial g^* &\implies \mathbf{b}_{t+1}^* = J_{\partial g^*}^\eta \mathbf{w}_{t+1}, \quad \mathbf{w}_{t+1} := \eta \mathbf{b}_{t+1} + \mathbf{b}_{t+1}^*, \\ (\mathbf{a}_{t+1}^*, \mathbf{a}_{t+1}) \in \text{gph}[\partial h^* \circ (-L^\top)] &\implies \mathbf{a}_{t+1}^* = J_{-L \circ \partial h^* \circ (-L^\top)}^\eta (-\eta L \mathbf{a}_{t+1} + \mathbf{a}_{t+1}^*) \\ (\text{since } \boldsymbol{\mu}_t = \mathbf{b}_t^*) &= J_{-L \circ \partial h^* \circ (-L^\top)}^\eta (-\eta \mathbf{b}_t + \mathbf{b}_t^*) \end{aligned}$$

Therefore, we deduce that

$$\begin{aligned} \mathbf{w}_{t+1} &:= \eta \mathbf{b}_{t+1} + \mathbf{b}_{t+1}^* = \eta L \mathbf{a}_{t+1} + \boldsymbol{\mu}_t = \mathbf{a}_{t+1}^* + \eta \mathbf{b}_t = J_{-L \circ \partial h^* \circ (-L^\top)}^\eta (-\eta \mathbf{b}_t + \mathbf{b}_t^*) + \eta \mathbf{b}_t \\ &= J_{-L \circ \partial h^* \circ (-L^\top)}^\eta (2\mathbf{b}_t^* - \mathbf{w}_t) + \mathbf{w}_t - \mathbf{b}_t^* = \frac{\text{Id} + R_{-L \circ \partial h^* \circ (-L^\top)}^\eta R_{\partial g^*}^\eta}{2} \mathbf{w}_t, \end{aligned}$$

Exactly the Douglas-Rachford algorithm applied to the dual problem, with the maximal monotone maps ∂h^* and $-L \circ \partial g^* \circ (-L^\top)$!

D. Gabay. "Applications of the Method of Multipliers to Variational Inequalities". In: *Augmented Lagrangian methods: Applications to the numerical solution of boundary-value problems*. Vol. 15. 9. 1983, pp. 299–331.

