

## 11 Splitting

### Goal

Splitting, ergodic averaging, gradient-descent-ascent, forward-backward, backward-backward.

### Alert 11.1: Convention

Gray boxes are not required hence can be omitted for unenthusiastic readers.

[This note is likely to be updated again soon.](#)

### Definition 11.2: The splitting/decomposition problem

Recall the familiar problem of finding a zero of a maximal monotone map  $T : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$ . We now add a small twist:

$$\text{find } \mathbf{z} \quad \text{s.t.} \quad \mathbf{0} \in T\mathbf{z}, \quad \text{where} \quad T = A + B, \quad (11.1)$$

i.e., the map  $T$  can be decomposed into the sum of two maps  $A : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$  and  $B : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$ . The catch is that we often cannot evaluate the resolvent  $J_T$  easily (so the proximal point algorithm is not directly applicable), and yet it might be possible to find a decomposition so that both  $J_A$  and  $J_B$  are readily available. Surprisingly, as we will see, many familiar algorithms are in fact instantiations of this simple but powerful idea.

We also associate the following dual with the primal problem (11.1):

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

See Example 8.9 for an explanation of the dual when both  $A$  and  $B$  are subdifferentials of convex functions.

### Theorem 11.3: Ergodic forward-backward splitting converges (Passty, 1979)

Let  $B : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$  be maximal monotone,  $A : \text{dom } B \rightrightarrows \mathbb{R}^d$  be monotone, and  $T := A + B$  be maximal monotone. Let  $\mathbf{w}_0 \in \text{dom } A$  and for all  $t \geq 0$  define

$$\begin{aligned} \mathbf{w}_{t+1} &:= J_B^{\eta_t}(\mathbf{w}_t - \eta_t \mathbf{a}_t^*), \quad \text{where} \quad \mathbf{a}_t^* \in A\mathbf{w}_t, \quad \eta_t \geq 0, \\ \mathbf{z}_t &= \sum_{k=0}^t \bar{\eta}_{t,k} \mathbf{w}_k, \quad \text{where} \quad \bar{\eta}_{t,k} := \eta_k / H_t, \quad H_t := \sum_{k=0}^t \eta_k. \end{aligned} \quad (11.2)$$

The following estimate holds for any  $(\mathbf{w}, \mathbf{w}^*) \in \text{gph } T$  and  $\mathbf{b}^* \in B\mathbf{w}$ :

$$\langle \mathbf{z}_t - \mathbf{w}, \mathbf{w}^* \rangle \leq \sum_{k=0}^t \bar{\eta}_{t,k} \langle \mathbf{w}_k - \mathbf{w}, \mathbf{a}_k^* + \mathbf{b}^* \rangle \leq \frac{\|\mathbf{w}_0 - \mathbf{w}\|_2^2 + \sum_{k=0}^t \eta_k^2 \|\mathbf{a}_k^* + \mathbf{b}^*\|_2^2}{2H_t}. \quad (11.3)$$

Moreover,

- if  $\sum_t \eta_t^2 \|\mathbf{a}_t^* + \mathbf{b}^*\|_2^2 < \infty$  and  $F := T^{-1}\mathbf{0} \neq \emptyset$ , then  $\|\mathbf{w}_t - \mathbf{w}\|_2$  converges for any  $\mathbf{w} \in F$ ;
- if  $\sum_t \eta_t^2 \|\mathbf{a}_t^* + \mathbf{b}^*\|_2^2 < \infty$  and  $H_t \rightarrow \infty$ , then either  $F = \emptyset$ , in which case  $\|\mathbf{z}_t\| \rightarrow \infty$ , or  $\mathbf{z}_t \rightarrow \mathbf{z}_\infty \in F$  (hence also follows the previous claim).

*Proof:* The assumptions guarantee that the iterates  $\{\mathbf{w}_t\}$  are well-defined. We now verify ??, starting with the last condition (III). For any  $(\mathbf{w}, \mathbf{w}^*) \in \text{gph } \mathsf{T}$  and  $\mathbf{b}^* \in \mathsf{B}\mathbf{w}$ :

$$\begin{aligned} \|\mathbf{w}_{k+1} - \mathbf{w}\|_2^2 &= \|J_{\mathsf{B}}^{\eta_k}(\mathbf{w}_k - \eta_k \mathbf{a}_k^*) - J_{\mathsf{B}}^{\eta_k}(\mathbf{w} + \eta_k \mathbf{b}^*)\|_2^2 \\ (\text{firm nonexpansiveness of } J_{\mathsf{B}}^{\eta_k}) &\leq \|\mathbf{w}_k - \mathbf{w} - \eta_k \mathbf{a}_k^* - \eta_k \mathbf{b}^*\|_2^2 - \|\mathbf{w}_k - \eta_k \mathbf{a}_k^* - \mathbf{w}_{k+1} - \eta_k \mathbf{b}^*\|_2^2 \\ &= \|\mathbf{w}_k - \mathbf{w}\|_2^2 - \|\mathbf{w}_k - \mathbf{w}_{k+1}\|_2^2 - 2\eta_k \langle \mathbf{w}_{k+1} - \mathbf{w}, \mathbf{a}_k^* + \mathbf{b}^* \rangle \\ (-\|\mathbf{x}\|_2^2 + 2\langle \mathbf{x}, \mathbf{y} \rangle \leq \|\mathbf{y}\|_2^2) &\leq \|\mathbf{w}_k - \mathbf{w}\|_2^2 + \eta_k^2 \|\mathbf{a}_k^* + \mathbf{b}^*\|_2^2 - 2\eta_k \langle \mathbf{w}_k - \mathbf{w}, \mathbf{a}_k^* + \mathbf{b}^* \rangle \end{aligned} \quad (11.4)$$

$$(\text{monotonicity of } \mathsf{A}) \leq \|\mathbf{w}_k - \mathbf{w}\|_2^2 + \eta_k^2 \|\mathbf{a}_k^* + \mathbf{b}^*\|_2^2 - 2\eta_k \langle \mathbf{w}_k - \mathbf{w}, \mathbf{w}^* \rangle. \quad (11.5)$$

Summing from  $k = 0$  to  $k = t$ , dividing by  $H_t = \sum_{k=0}^t \eta_k$ , telescoping and rearranging we obtain:

$$2\langle \mathbf{w} - \mathbf{z}_t, \mathbf{w}^* \rangle + \sum_{k=0}^t \eta_k^2 \|\mathbf{a}_k^* + \mathbf{b}^*\|_2^2 / H_t \geq (\|\mathbf{w}_{t+1} - \mathbf{w}\|_2^2 - \|\mathbf{w}_0 - \mathbf{w}\|_2^2) / H_t,$$

whence follows the estimate (11.3) (using also (11.4)). If  $\sum_t \eta_t^2 \|\mathbf{a}_t^* + \mathbf{b}^*\|_2^2 < \infty$  and  $H_t \rightarrow \infty$ , we deduce that

$$\liminf_{t \rightarrow \infty} \langle \mathbf{w} - \mathbf{z}_t, \mathbf{w}^* \rangle \geq 0,$$

whence follows from the maximality of  $\mathsf{T}$  that any limit point of  $\{\mathbf{z}_t\}$  is a zero. Note that if  $\|\mathbf{z}_t\|$  remains bounded then it admits a limit point. Therefore, from now on we assume  $\mathsf{F} \neq \emptyset$ . For any  $\mathbf{w} \in \mathsf{F}$ , from (11.5) it follows

$$\|\mathbf{w}_{t+1} - \mathbf{w}\|_2^2 \leq \|\mathbf{w}_t - \mathbf{w}\|_2^2 + \eta_t^2 \|\mathbf{a}_t^* + \mathbf{b}^*\|_2^2 - 2\eta_t \langle \mathbf{w}_t - \mathbf{w}, \mathbf{w}^* \rangle \leq \|\mathbf{w}_t - \mathbf{w}\|_2^2 + \eta_t^2 \|\mathbf{a}_t^* + \mathbf{b}^*\|_2^2.$$

If the last term is summable, then obviously  $\|\mathbf{w}_t - \mathbf{w}\|_2^2$  converges hence  $\{\mathbf{w}_t\}$  is bounded. Lastly,

$$\text{dist}(\mathbf{z}_t, W_k) \leq \left\| \mathbf{z}_t - \sum_{s=k}^t \bar{\eta}_{t,s} \mathbf{w}_s / \sum_{\kappa=k}^t \bar{\eta}_{t,\kappa} \right\|_2 \leq \sum_{\kappa=0}^{k-1} \bar{\eta}_{t,\kappa} \left[ \|\mathbf{w}_\kappa\|_2 + \left\| \sum_{s=k}^t \bar{\eta}_{t,s} \mathbf{w}_s / \sum_{\kappa=k}^t \bar{\eta}_{t,\kappa} \right\|_2 \right] \xrightarrow{t \rightarrow \infty} 0,$$

since  $\mathbf{w}_t$  is bounded and for any  $k$ ,  $\bar{\eta}_{t,k} \rightarrow 0$  as  $t \rightarrow \infty$ . ■

The special case  $\mathsf{B} = \mathcal{N}_C$  for some closed convex set  $C$  first appeared in (Bruck, 1977).

Passty, G. B. (1979). “Ergodic convergence to a zero of the sum of monotone operators in Hilbert space”. *Journal of Mathematical Analysis and Applications*, vol. 72, no. 2, pp. 383–390.

Bruck, R. E. (1977). “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, pp. 159–164.

#### Remark 11.4: Parsing the previous result

For  $\mathsf{B} = \mathcal{N}_C$  for some closed convex set  $C$ , we may take  $\mathbf{b}^* = \mathbf{0}$ , in which case, as suggested by Nemirovskii and Judin (1978), we may choose

$$\eta_t = \frac{1}{\sqrt{\|\mathbf{a}_t^*\|_2^2 + 1}} \frac{1}{(t+1)^p}, \quad p \in (\tfrac{1}{2}, 1], \quad (11.6)$$

so that obviously  $\sum_t \eta_t \|\mathbf{a}_t^*\|_2^2 < \infty$ . If there exists a zero (or  $C$  is bounded) then  $\{\mathbf{w}_t\}$  is bounded. If  $\mathsf{A}$  is also bounded on bounded sets (so that  $\sup_t \|\mathbf{a}_t^*\|_2 < \infty$ ), then letting  $H_t \rightarrow \infty$  the estimate (11.3) goes to 0 while  $\{\mathbf{z}_t\}$  converges to a zero.

It is clear that the proximal gradient Algorithm 2.17, the subgradient Algorithm 4.14 and the gradient-descent-ascent (GDA) Algorithm 8.22 are all special cases of the so-called forward-backward splitting in (11.2). In fact, Theorem 4.17 for the convergence of the subgradient Algorithm 4.14 is strictly contained in

Theorem 11.3, and now we have a similar result for GDA. Indeed, let  $A = (\partial_{\mathbf{x}}f, \partial_{\mathbf{y}}f)$  as suggested in ??, for any  $\mathbf{w} = (\mathbf{x}, \mathbf{y}) \in C$  we have from (11.3):

$$\begin{aligned} \frac{\|\mathbf{w}_0 - \mathbf{w}\|_2^2 + \sum_{k=0}^t \|\eta_k \mathbf{a}_k^*\|_2^2}{2H_t} &\geq \sum_{k=0}^t \bar{\eta}_{t,k} \langle \mathbf{w}_k - \mathbf{w}, \mathbf{a}_k^* \rangle \\ &= \sum_{k=0}^t \bar{\eta}_{t,k} [\langle \mathbf{x}_k - \mathbf{x}, \partial_{\mathbf{x}}f(\mathbf{x}_k, \mathbf{y}_k) \rangle - f(\mathbf{x}_k, \mathbf{y}_k) + f(\mathbf{x}_k, \mathbf{y}_k) + \langle \mathbf{y}_k - \mathbf{y}, \partial_{\mathbf{y}}f(\mathbf{x}_k, \mathbf{y}_k) \rangle] \\ &\geq \sum_{k=0}^t \bar{\eta}_{t,k} [-f(\mathbf{x}, \mathbf{y}_k) + f(\mathbf{x}_k, \mathbf{y})] \\ &\geq -f(\mathbf{x}, \bar{\mathbf{y}}_t) + f(\bar{\mathbf{x}}_t, \mathbf{y}), \quad \text{where } (\bar{\mathbf{x}}_t, \bar{\mathbf{y}}_t) := \sum_{k=0}^t \bar{\eta}_{t,k} \mathbf{w}_k. \end{aligned} \quad (11.7)$$

We can make the following conclusions:

- If  $C$  is bounded and  $A$  is bounded on  $C$ , then maximizing w.r.t.  $\mathbf{w} = (\mathbf{x}, \mathbf{y}) \in C$  on both sides we have

$$\frac{\text{diam}(C)^2 + L^2 S_t^2}{2H_t} \geq \underbrace{\mathfrak{d}^* - f(\bar{\mathbf{y}}_t)}_{\text{dual gap} \geq 0} + \underbrace{\bar{f}(\bar{\mathbf{x}}_t) - \mathfrak{p}_*}_{\text{primal gap} \geq 0} + \underbrace{\mathfrak{p}_* - \mathfrak{d}^*}_{\text{strong duality} = 0} \geq 0,$$

where  $\text{diam}(C)$  is the diameter of  $C$ ,  $L := \sup_t \|\mathbf{a}_t^*\|_2 < \infty$  and  $S_t^2 = \sum_t \eta_t^2$ . Thus, the primal and dual gaps go to 0 if  $H_t \rightarrow \infty$  and  $\eta_t \rightarrow 0$ , in which case any limit point of  $\{\mathbf{z}_k\}$  is a saddle point while convergence of the whole sequence requires the stronger condition  $\sum_t \eta_t^2 < \infty$ . In particular, setting  $\eta_t = O(1/\sqrt{t})$  leads to  $O((\ln t)/\sqrt{t})$  rate of convergence for the sum of gaps.

- Suppose  $C = X \times Y$  with say  $X$  bounded,  $\sum_t \|\eta_t \mathbf{a}_t^*\|_2^2 < \infty$ , there exists a saddle point, and  $A$  is bounded on bounded sets. Then, setting  $\mathbf{y} = \mathbf{y}^*$  for any  $\mathbf{y}^* \in Y^*$  we obtain:

$$\frac{\|\mathbf{x}_0 - \mathbf{x}\|_2^2 + \|\mathbf{y}_0 - \mathbf{y}^*\|_2^2 + \sum_{k=0}^t \|\eta_k \mathbf{a}_k^*\|_2^2}{2H_t} \geq -f(\mathbf{x}, \bar{\mathbf{y}}_t) + f(\bar{\mathbf{x}}_t, \mathbf{y}^*) \geq \mathfrak{p}^* - f(\mathbf{x}, \bar{\mathbf{y}}_t).$$

Maximizing w.r.t.  $\mathbf{x} \in X$  on both sides leads us to

$$\frac{\text{diam}(X)^2 + \text{dist}(\mathbf{y}_0, Y^*)^2 + \sum_{k=0}^t \|\eta_k \mathbf{a}_k^*\|_2^2}{2H_t} \geq \mathfrak{p}_* - f(\bar{\mathbf{y}}_t) \geq \mathfrak{d}^* - f(\bar{\mathbf{y}}_t) \geq 0,$$

i.e. the dual gap is bounded and converges to 0 if  $H_t \rightarrow \infty$ . And similarly for the primal gap.

- Inspecting the proof of Theorem 11.3 we realize that completely similar results still hold even with different step sizes for  $\mathbf{x}$  and  $\mathbf{y}$ , with one intriguing change: in the function estimate (11.7) we need to average  $\mathbf{x}_k$  using the step size on  $\mathbf{y}$  and vice versa. This observation is useful when only say  $X$  is bounded (such as in a Lagrangian) so that we need only use the adaptive step size (11.6) for updating  $\mathbf{y}$ .

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

### Theorem 11.5: Ergodic backward-backward splitting converges (Passty, 1979)

Let  $A, B : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$  be maximal monotone, with maximal monotone sum  $T := A + B$ . Starting with any  $\mathbf{w}_0$  and for all  $t \geq 0$  define

$$\mathbf{w}_{t+1} := J_B^{\eta_t} J_A^{\eta_t} \mathbf{w}_t, \quad \text{where } \eta_t \geq 0, \quad (11.8)$$

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

If  $\sum_t \eta_t = \infty$  and  $\eta_t \rightarrow 0$ , then either  $F := T^{-1}\mathbf{0} = \emptyset$ , in which case  $\|\mathbf{z}_t\| \rightarrow \infty$ , or  $\mathbf{z}_t \rightarrow \mathbf{z}_\infty \in F$ .

*Proof:* We simply verify ???. Let  $\mathbf{w} \in \text{dom } T$ ,  $\mathbf{a}^* \in A\mathbf{w}$  and  $\mathbf{b}^* \in B\mathbf{w}$ . Applying the firm nonexpansiveness of  $J^{\eta_k}$  (see ??):

$$\begin{aligned} \|J_A^{\eta_k} \mathbf{w}_k - \mathbf{w}\|_2^2 &= \|J_A^{\eta_k} \mathbf{w}_k - J_A^{\eta_k}(\mathbf{w} + \eta_k \mathbf{a}^*)\|_2^2 \leq \|\mathbf{w}_k - \mathbf{w} - \eta_k \mathbf{a}^*\|_2^2 - \|\mathbf{w}_k - J_A^{\eta_k} \mathbf{w}_k - \eta_k \mathbf{a}^*\|_2^2 \\ &= \|\mathbf{w}_k - \mathbf{w}\|_2^2 - \|\mathbf{w}_k - J_A^{\eta_k} \mathbf{w}_k\|_2^2 + 2\eta_k \langle \mathbf{w} - J_A^{\eta_k} \mathbf{w}_k; \mathbf{a}^* \rangle \quad (11.9) \\ \|J_B^{\eta_k} J_A^{\eta_k} \mathbf{w}_k - \mathbf{w}\|_2^2 &\leq \|J_A^{\eta_k} \mathbf{w}_k - \mathbf{w}\|_2^2 - \|J_A^{\eta_k} \mathbf{w}_k - J_B^{\eta_k} J_A^{\eta_k} \mathbf{w}_k\|_2^2 + 2\eta_k \langle \mathbf{w} - J_B^{\eta_k} J_A^{\eta_k} \mathbf{w}_k; \mathbf{b}^* \rangle. \end{aligned}$$

Summing the above two inequalities and applying the inequality  $-\|\mathbf{x}\|_2^2 + 2\langle \mathbf{x}; \mathbf{y} \rangle \leq \|\mathbf{y}\|_2^2$  repeatedly:

$$\|J_B^{\eta_k} J_A^{\eta_k} \mathbf{w}_k - \mathbf{w}\|_2^2 \leq \|\mathbf{w}_k - \mathbf{w}\|_2^2 + 2\eta_k \langle \mathbf{w} - \mathbf{w}_k; \mathbf{a}^* + \mathbf{b}^* \rangle + \eta_k^2 [\|\mathbf{a}^* + \mathbf{b}^*\|_2^2 + \|\mathbf{a}^*\|_2^2]. \quad (11.10)$$

Summing from  $k = 0$  to  $k = t$  and rearranging as in Theorem 11.3 we obtain for any  $\mathbf{w} \in \text{dom } T$ ,  $\mathbf{w}^* = \mathbf{a}^* + \mathbf{b}^* \in T\mathbf{w}$ :

$$2\langle \mathbf{w} - \mathbf{z}_t; \mathbf{w}^* \rangle + [\|\mathbf{a}^*\|_2^2 + \|\mathbf{w}^*\|_2^2] \sum_{k=0}^t \eta_k^2 / H_t \geq (\|\mathbf{w}_{t+1} - \mathbf{w}\|_2^2 - \|\mathbf{w}_0 - \mathbf{w}\|_2^2) / H_t.$$

Using the assumptions on  $\eta_t$  we thus know

$$\liminf_{t \rightarrow \infty} \langle \mathbf{w} - \mathbf{z}_t; \mathbf{w}^* \rangle \geq 0,$$

whence follows from the maximality of the sum  $T$  that any limit point of  $\{\mathbf{z}_t\}$  is a zero. If  $\{\mathbf{z}_t\}$  is bounded, then  $F \neq \emptyset$ , which we assume now. Let  $\mathbf{w} \in F$  and set  $\mathbf{w}^* = \mathbf{0}$  we know from (11.10) that  $\{\mathbf{w}_t\}$  is (uniformly) quasi-Fejér monotone w.r.t.  $F$ . Lastly, we verify condition (II) in ?? as in Theorem 11.3. ■

The special case  $B = \mathcal{N}_C$  for some closed convex set first appeared in (Lions, 1978). We may also define

$$\mathbf{z}_{t+1} := \sum_{k=0}^t \bar{\eta}_{t,k} \mathbf{w}_{k+1},$$

which will remove the constant  $\|\mathbf{a}^* + \mathbf{b}^*\|_2^2$ .

Passty, G. B. (1979). “Ergodic convergence to a zero of the sum of monotone operators in Hilbert space”. *Journal of Mathematical Analysis and Applications*, vol. 72, no. 2, pp. 383–390.

Lions, P.-L. (1978). “Une methode iterative de resolution d’une inequation variationnelle”. *Israel Journal of Mathematics*, vol. 31, no. 2, pp. 204–208.

### Alert 11.6: Comparing forward-backward and backward-backward

It is instructive to compare forward-backward with backward-backward:

$$\mathfrak{F} := J_B^\eta(\text{Id} - \eta A) \quad \text{vs.} \quad \mathfrak{B} := J_B^\eta J_A^\eta.$$

For the former, it is clear that

$$\begin{aligned} \mathbf{w} \in \mathfrak{F}\mathbf{w} &\iff [\mathbf{w} + \eta B\mathbf{w}] \cap [\mathbf{w} - \eta A\mathbf{w}] \neq \emptyset \iff \mathbf{0} \in (A + B)\mathbf{w} \\ \exists(\mathbf{w}, \mathbf{w}^*) \in \text{gph } A, \mathbf{w} \in \mathfrak{F}\mathbf{w} &\iff \exists(\mathbf{w}, \mathbf{w}^*) \in \text{gph } A, \mathbf{0} \in (A + B)\mathbf{w} \iff \mathbf{0} \in -A^{-1}(-\mathbf{w}^*) + B^{-1}\mathbf{w}^*. \end{aligned}$$

Therefore, applying the forward-backward map with any  $\eta$  at least makes sense in principle. However, the

latter, as pointed out by Bauschke et al. (2005), solves a “regularized” problem:

$$\mathbf{w} = \mathfrak{B}\mathbf{w} \iff \mathbf{0} \in (\eta\mathbf{A} + \mathbf{B})\mathbf{w}, \quad \text{where} \quad \eta\mathbf{A} := \frac{\text{Id} - J_{\mathbf{A}}^{\eta}}{\eta}.$$

In other words,

backward-backward on  $\mathbf{A} + \mathbf{B}$  is forward-backward on  $\eta\mathbf{A} + \mathbf{B}$ !

In general,  $(\eta\mathbf{A} + \mathbf{B})^{-1}\mathbf{0} \cap (\mathbf{A} + \mathbf{B})^{-1}\mathbf{0} = \emptyset$ , with one notable exception: when  $\mathbf{A}^{-1}\mathbf{0} \cap \mathbf{B}^{-1}\mathbf{0} \neq \emptyset$ , see Theorem 11.7 below. This is the reason why in all our results about backward-backward (e.g. Theorem 11.5) we require  $\eta_t \rightarrow 0$ , since then  $\eta\mathbf{A} \rightarrow \mathbf{0}\mathbf{A}$  as  $\eta \rightarrow 0$ , where recall that  $\mathbf{0}\mathbf{A}$  is the minimum-norm element in  $\mathbf{A}\mathbf{w}$ . In contrast, it is possible to use constant  $\eta$  in forward-backward (e.g. Theorem 11.13), at the expense of  $\eta$  depending on properties of  $\mathbf{A}$ . Still, it is surprising that with  $\eta_t$  decreasing to 0 slowly, (ergodic) backward-backward actually converges to a zero!

Bauschke, H. H., P. L. Combettes, and S. Reich (2005). “The asymptotic behavior of the composition of two resolvents”. *Nonlinear Analysis: Theory, Methods & Applications*, vol. 60, no. 2, pp. 283–301.

### Theorem 11.7: Backward-backward converges under a common fixed point (Tseng, 1992)

Let  $\mathbf{T}_i : \mathbb{R}^d \rightarrow \mathbb{R}^d, i = 1, \dots, m$  be  $\alpha$ -averaged with a common fixed point, i.e.  $\mathbf{F} := \cap_i \text{Fix}\mathbf{T}_i \neq \emptyset$ . Then, the (random) iterate

$$\mathbf{w}_{t+1} = (1 - \gamma_t)\mathbf{w}_t + \gamma_t \mathbf{T}_{i(t)}\mathbf{w}_t + \epsilon_t, \quad i_t \in \{1, \dots, m\}, \quad \gamma_t \in (0, \frac{1}{\alpha}), \quad \sum_t \|\epsilon_t\|_2 < \infty$$

converges to some  $\mathbf{w}_{\infty} \in \mathbf{F}$ , as long as each  $\mathbf{T}_i$  appears infinitely often and  $\liminf_t \gamma_t(\frac{1}{\alpha} - \gamma_t) > 0$ .

*Proof:* From the proof of ?? we know  $\{\mathbf{w}_t\}$  is (uniformly) quasi-Fejér monotone w.r.t.  $\mathbf{F}$  and

$$\mathbf{w}_t - \mathbf{T}_{i(t)}\mathbf{w}_t \rightarrow \mathbf{0}.$$

Let  $\mathbf{z} \in \cap_{i \in I} \text{Fix}\mathbf{T}_i$  be a limit point of  $\{\mathbf{w}_t\}$  for some  $I \neq \emptyset$  (e.g.  $I = \{i\}$  for some  $i$ ; see ??). Take a subsequence  $\mathbf{w}_{t_k} \rightarrow \mathbf{z}$  and let  $s_k = \min\{t \geq t_k : i(t) \notin I\}$ . Pass to a subsequence we may assume  $i(s_k) \equiv j$  and  $\mathbf{w}_{s_k} \rightarrow \mathbf{w}$ . Since  $\mathbf{w}_{s_k} - \mathbf{T}_j\mathbf{w}_{s_k} \rightarrow \mathbf{0}$  we have  $\mathbf{w} \in \text{Fix}\mathbf{T}_j$  (see ??). Since  $\mathbf{z} \in \cap_{i \in I} \text{Fix}\mathbf{T}_i$  and  $i(t) \in I$  for  $t \in [t_k, s_k)$  we have

$$\|\mathbf{w}_{s_k} - \mathbf{z}\|_2 \leq \|\mathbf{w}_{t_k} - \mathbf{z}\|_2 + \sum_{\kappa=t_k}^{s_k-1} \|\epsilon_{\kappa}\|_2 \rightarrow 0,$$

and hence  $\mathbf{z} = \mathbf{w} \in \text{Fix}\mathbf{T}_j$ . Since each  $\mathbf{T}_i$  appears infinitely often, we may continue the argument to conclude that any limit point  $\mathbf{z} \in \mathbf{F}$ . Applying ?? we know the whole sequence  $\mathbf{w}_t \rightarrow \mathbf{w}_{\infty} \in \mathbf{F}$ . ■

Aleyner and Reich (2009) pointed out that we only need the following weaker condition on each  $\mathbf{T}_i$ : it is continuous and there exists some  $\alpha > 0$  such that for any  $\mathbf{z} \in \text{Fix}\mathbf{T}_i$

$$\|\mathbf{T}_i\mathbf{w} - \mathbf{z}\|_2^2 + \alpha\|\mathbf{w} - \mathbf{T}_i\mathbf{w}\|_2^2 \leq \|\mathbf{w} - \mathbf{z}\|_2^2.$$

Tseng, P. (1992). “On the Convergence of the Products of Firmly Nonexpansive Mappings”. *SIAM Journal on Optimization*, vol. 2, no. 3, pp. 425–434.

Aleyner, A. and S. Reich (2009). “Random Products of Quasi-Nonexpansive Mappings in Hilbert Space”. *Journal of Convex Analysis*, vol. 16, no. 3, pp. 633–640.

**Example 11.8: Method of barycenter (Cimmino, 1938)**

Let  $H_i := \{\mathbf{w} : \langle \mathbf{w}, \mathbf{a}_i \rangle = b_i\}$  be a hyperplane and  $P_i$  the orthogonal projection onto it. Cimmino (1938) proposed the method of barycenter for finding a point in the intersection  $H = \cap_i H_i$ :

$$\mathbf{w}_{t+1} \leftarrow \frac{1}{n} \sum_i P_i \mathbf{w}_t,$$

which is exactly a backward-backward algorithm for the reformulation:

$$\min_{\mathbf{w}=(\mathbf{w}_1, \dots, \mathbf{w}_n)} \sum_i \iota_{H_i}(\mathbf{w}_i) + \iota_L(\mathbf{w}), \quad \text{where } L := \{\mathbf{w} : \mathbf{w}_1 = \dots = \mathbf{w}_n\}.$$

Applying Theorem 11.7, we actually know the more general version

$$\mathbf{w}_{t+1} \leftarrow \text{Avg}(P_{i_1}, \dots, P_{i_{k(t)}}) \mathbf{w}_t$$

also converges to a point in  $H$ , as long as each projection appears infinitely often. Setting  $k(t) \equiv 1$  we obtain Kaczmarz's (sequential) algorithm (Kaczmarz, 1937).

Reich (1983) studied the Barycenter method for both linear and nonlinear projectors in Banach spaces.

Cimmino, G. (1938). “Calcolo Approssimato Per le Soluzioni dei Sistemi di Equazioni Lineari”. *La Ricerca Scientifica*, vol. 9, no. 1, pp. 326–333.

Kaczmarz, S. (1937). “Angenäherte Auflösung von Systemen linearer Gleichungen”. *Bulletin International de l'Académie Polonaise des Sciences et des Lettres*, vol. 35, 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.

Reich, S. (1983). “A note on the mean ergodic theorem for nonlinear semigroups”. *Journal of Mathematical Analysis and Applications*, vol. 91, no. 2, pp. 547–551.

**Theorem 11.9: (Strong) non-ergodic convergence of backward-backward (Passty, 1979)**

Let  $A$  and  $B$  be maximal monotone with maximal monotone sum  $T := A + B$  and  $F := T^{-1}\mathbf{0} \neq \emptyset$ . Choose  $\sum_t \eta_t = \infty$  and  $\eta_t \rightarrow 0$ . Suppose either

- one of  $A$  and  $B$  is strongly monotone, or
- $F$  has nonempty interior.

Then, the (non-ergodic) backward-backward iterate  $\mathbf{w}_t \rightarrow \mathbf{w}_\infty \in F$  (see (11.8)).

*Proof:* The second claim readily follows from ???. For the first claim, assume w.l.o.g. that  $A$  is  $\sigma$ -strongly monotone. We strengthen (11.9) into

$$\|J_A^{\eta_k} \mathbf{w}_k - \mathbf{w}_\infty\|_2^2 \leq \|\mathbf{w}_k - \mathbf{w}_\infty\|_2^2 - \|\mathbf{w}_k - J_A^{\eta_k} \mathbf{w}_k\|_2^2 + 2\eta_k \langle \mathbf{w}_\infty - J_A^{\eta_k} \mathbf{w}_k; \mathbf{a}^* \rangle - 2\eta_k \sigma \|J_A^{\eta_k} \mathbf{w}_k - \mathbf{w}_\infty\|_2^2,$$

leading (11.10) now to

$$\begin{aligned} \|\mathbf{w}_{k+1} - \mathbf{w}_\infty\|_2^2 &\leq \|\mathbf{w}_k - \mathbf{w}_\infty\|_2^2 + \eta_k^2 \|\mathbf{a}^*\|_2^2 - 2\eta_k \sigma \|J_A^{\eta_k} \mathbf{w}_k - \mathbf{w}_\infty\|_2^2 \implies \liminf \|J_A^{\eta_k} \mathbf{w}_k - \mathbf{w}_\infty\|_2 = 0, \text{ hence} \\ \liminf \|\mathbf{w}_{k+1} - \mathbf{w}_\infty\|_2 &= \liminf \|\mathbf{w}_{k+1} - J_B^{\eta_k}(\mathbf{w}_\infty + \eta_k \mathbf{b}^*)\|_2 \leq \liminf \|J_A^{\eta_k} \mathbf{w}_k - \mathbf{w}_\infty - \eta_k \mathbf{b}^*\|_2 = 0. \end{aligned}$$

Since  $\{\mathbf{w}_t\}$  is quasi-Fejér monotone w.r.t.  $F = \{\mathbf{w}_\infty\}$ , it follows that  $\mathbf{w}_t \rightarrow \mathbf{w}_\infty$ . ■

Passty, G. B. (1979). “Ergodic convergence to a zero of the sum of monotone operators in Hilbert space”. *Journal of Mathematical Analysis and Applications*, vol. 72, no. 2, pp. 383–390.

**Definition 11.10: Inversely strong monotonicity, a.k.a., cocoercive**

We call an operator  $T : \text{dom } T \subseteq \mathbb{R}^d \rightarrow \mathbb{R}^d$  **inversely  $\sigma$ -strongly monotone** (a.k.a.  $\sigma$ -cocoercive) if

$$\forall(\mathbf{u}, \mathbf{u}^*) \in \text{gph } T, \forall(\mathbf{v}, \mathbf{v}^*) \in \text{gph } T, \quad \langle \mathbf{u} - \mathbf{v}; \mathbf{u}^* - \mathbf{v}^* \rangle \geq \sigma \|\mathbf{u}^* - \mathbf{v}^*\|_2^2,$$

i.e.  $T^{-1}$  is  $\sigma$ -strongly monotone or **equivalently  $\sigma T$  is firmly nonexpansive** and hence  $T$  is  $\frac{1}{\sigma}$ -Lipschitz continuous. When  $T = \partial f$  for a closed (proper) convex function  $f$ , we know from Alert 3.25 that  $\partial f$  is inversely  $\sigma$ -strongly monotone iff  $\partial f$  is  $\frac{1}{\sigma}$ -Lipschitz continuous.

**Exercise 11.11: Strongly monotone + Lipschitz continuity  $\implies$  inversely strongly monotone**

Let  $T$  be  $\sigma$ -strongly monotone and  $L$ -Lipschitz continuous. Prove that  $T$  is inversely  $\frac{\sigma}{L^2}$ -strongly monotone.

If  $T = \partial f$  for some (closed proper) convex function  $f$ , we may improve the factor  $\frac{\sigma}{L^2}$  to  $\frac{1}{L}$ .

**Theorem 11.12: Non-ergodic convergence of forward-backward**

Let  $B : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$  be maximal monotone and  $A : \text{dom } B \rightarrow \mathbb{R}^d$  be **inversely  $\frac{1}{L}$ -strongly monotone**. Consider the (non-ergodic) relaxed forward-backward iterate:

$$\mathbf{w}_{t+1} := (1 - \gamma_t)\mathbf{w}_t + \gamma_t J_B^{\eta_t}(\mathbf{w}_t - \eta_t \mathbf{a}_t^*) + \epsilon_t, \quad \text{where } \mathbf{a}_t^* \in A\mathbf{w}_t, \eta_t \in [0, \frac{2}{L}], \gamma_t \geq 0.$$

Assume  $F := (A+B)^{-1}\mathbf{0} \neq \emptyset$  and  $\sum_t \|\epsilon_t\|_2 < \infty$ . If  $\gamma_t \in [0, 2 - \frac{\eta_t L}{2}]$ ,  $\liminf_t \eta_t \geq \underline{\eta} > 0$  and  $\sum_t \gamma_t (2 - \frac{\eta_t L}{2} - \gamma_t) = \infty$ , then  $\mathbf{w}_t \rightarrow \mathbf{w}_\infty \in F$  and  $A\mathbf{w}_t \rightarrow A\mathbf{w}_\infty = T^* \mathbf{0}$ , where  $T^* := -A^{-1}(-\text{Id}) + B^{-1}$ .

*Proof:* We simply analyze the forward-backward map

$$\mathfrak{F}_\eta := J_B^\eta(\text{Id} - \eta A).$$

If  $A$  is inversely  $\frac{1}{L}$ -strongly monotone, i.e.  $\frac{1}{L}A$  is firmly nonexpansive, then

$$\text{Id} - \eta A = \text{Id} - \eta L \frac{\text{Id} + N}{2} = (1 - \frac{\eta L}{2})\text{Id} + \frac{\eta L}{2}(-N)$$

is  $\frac{\eta L}{2}$ -averaged for any  $\eta \in [0, \frac{2}{L}]$ . According to ??,  $\mathfrak{F}_\eta$  is  $\frac{2}{4 - \eta L}$ -averaged. As shown in the proof of ??,  $\|\mathbf{w}_t - \mathfrak{F}_{\eta_t} \mathbf{w}_t\|_2 \rightarrow 0$ . Since  $\liminf_t \eta_t \geq \underline{\eta} > 0$ , we apply ?? to obtain

$$\limsup_t \underline{\eta} \|\mathbf{w}_t - \mathfrak{F}_{\underline{\eta}}(\mathbf{w}_t)\|_2 \leq \limsup_t \|\mathbf{w}_t - \mathfrak{F}_{\eta_t}\|_2 = 0.$$

Applying ?? and ?? we know the quasi-Fejér monotone sequence  $\mathbf{w}_t \rightarrow \mathbf{w}_\infty \in F$ .

Since  $A^{-1}$  is strongly monotone,  $T^* \mathbf{0} = A\mathbf{w}$  for any  $\mathbf{w} \in F$ , see Alert 11.6. Let  $\tilde{\mathbf{w}}_t = \mathbf{w}_t - \eta_t A\mathbf{w}_t$ :

$$\begin{aligned} \langle \mathfrak{F}_{\eta_t} \mathbf{w}_t - \mathbf{w}_\infty, \mathbf{w}_t - \mathfrak{F}_{\eta_t} \mathbf{w}_t \rangle &= \langle J_B^{\eta_t} \tilde{\mathbf{w}}_t - J_B^{\eta_t} \tilde{\mathbf{w}}_\infty, \mathbf{w}_t - J_B^{\eta_t} \tilde{\mathbf{w}}_t \rangle \\ &= \langle J_B^{\eta_t} \tilde{\mathbf{w}}_t - J_B^{\eta_t} \tilde{\mathbf{w}}_\infty, (\tilde{\mathbf{w}}_t - J_B^{\eta_t} \tilde{\mathbf{w}}_t) - (\tilde{\mathbf{w}}_\infty - J_B^{\eta_t} \tilde{\mathbf{w}}_\infty) \rangle + \eta_t \langle J_B^{\eta_t} \tilde{\mathbf{w}}_t - J_B^{\eta_t} \tilde{\mathbf{w}}_\infty, A\mathbf{w}_t - A\mathbf{w}_\infty \rangle \\ &\geq \eta_t [\langle \mathfrak{F}_{\eta_t} \mathbf{w}_t - \mathbf{w}_t, A\mathbf{w}_t - A\mathbf{w}_\infty \rangle + \langle \mathbf{w}_t - \mathbf{w}_\infty, A\mathbf{w}_t - A\mathbf{w}_\infty \rangle] \\ &\geq -\eta_t \|\mathfrak{F}_{\eta_t} \mathbf{w}_t - \mathbf{w}_t\|_2 \cdot L \|\mathbf{w}_t - \mathbf{w}_\infty\|_2 + \frac{\eta_t}{L} \|A\mathbf{w}_t - A\mathbf{w}_\infty\|_2^2. \end{aligned}$$

Since  $\liminf_t \eta_t \geq \underline{\eta} > 0$  and we already know  $\mathfrak{F}_{\eta_t} \mathbf{w}_t - \mathbf{w}_t \rightarrow \mathbf{0}$ , it follows  $A\mathbf{w}_t \rightarrow A\mathbf{w}_\infty$ . ■

The primal convergence  $\mathbf{w}_t \rightarrow \mathbf{w}_\infty$ , with  $\gamma_t \equiv 1$ ,  $\eta_t \equiv \underline{\eta}$ ,  $\epsilon_t \equiv \mathbf{0}$  and  $B = \mathcal{N}_C$ , appeared in e.g. Mercier (1979, pp 157–158) and Gabay (1983, Thm 6.1). The dual convergence  $A\mathbf{w}_t \rightarrow A\mathbf{w}_\infty$  was due to Tseng (1991) who also considered relaxation and allowed varying step size. Our proof, exploiting the monotonicity in ??, confirms that the usual argument based on Opial's ?? does suffice.

Mercier, B. (1979). “Lectures on Topics in Finite Element Solution of Elliptical Problems”. Springer.

Gabay, D. (1983). “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, pp. 299–331.

Tseng, P. (1991). “Applications of a Splitting Algorithm to Decomposition in Convex Programming and Variational Inequalities”. *SIAM Journal on Control and Optimization*, vol. 29, no. 1, pp. 119–138.

**Theorem 11.13: Linear convergence of forward-backward (Chen and Rockafellar, 1997)**

Let  $\mathbf{B} : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$  be  $\sigma_b$ -strongly maximal monotone and  $\mathbf{A} : \text{dom } \mathbf{B} \rightarrow \mathbb{R}^d$  be  $\sigma_a$ -strongly monotone. Consider the (non-ergodic) forward-backward iterate:

$$\mathbf{w}_{t+1} := J_{\mathbf{B}}^{\eta_t}(\mathbf{w}_t - \eta_t \mathbf{a}_t^*) + \epsilon_t, \quad \text{where } \mathbf{a}_t^* \in \mathbf{A}\mathbf{w}_t, \eta_t \geq 0.$$

Assume  $\bar{\mathbf{A}} := \mathbf{A} - \sigma_a \cdot \text{Id}$  is  $\bar{\mathbf{L}}$ -Lipschitz continuous, then

$$\|\mathbf{w}_{t+1} - \mathbf{w}_\infty\|_2 \leq q_t \|\mathbf{w}_t - \mathbf{w}_\infty\|_2 + \|\epsilon_t\|_2, \quad \text{where } \mathbf{w}_\infty \in \mathbf{F} \text{ and } q_t := \frac{\sqrt{[(\eta_t \sigma_a - 1)_+ + \bar{\mathbf{L}}]^2 + (1 - \eta_t \sigma_a)_+^2}}{1 + \sigma_b \eta_t}. \quad (11.11)$$

Setting  $\eta_t \equiv \eta_* = \frac{1}{\sigma_a + \bar{\mathbf{L}}^2/\sigma}$ , where  $\sigma := \sigma_a + \sigma_b$ , we obtain the optimal  $q_* = 1/\sqrt{1 + \kappa^2}$ , where  $\kappa := \sigma/\bar{\mathbf{L}}$ .

*Proof:* When  $\mathbf{B}$  is  $\sigma_b$ -strongly monotone, we know from ?? that  $J_{\mathbf{B}}^\eta$  is  $\frac{1}{1+\eta\sigma_b}$ -Lipschitz continuous. When  $\mathbf{A}$  is  $\sigma_a$ -strongly monotone and  $\bar{\mathbf{A}}$  is  $\bar{\mathbf{L}}$ -Lipschitz continuous, then

$$\begin{aligned} \|(\mathbf{w} - \eta \mathbf{A}\mathbf{w}) - (\mathbf{z} - \eta \mathbf{A}\mathbf{z})\|_2^2 &= \|(1 - \eta \sigma_a)(\mathbf{w} - \mathbf{z}) - \eta(\bar{\mathbf{A}}\mathbf{w} - \bar{\mathbf{A}}\mathbf{z})\|_2^2 \\ &= (1 - \eta \sigma_a)^2 \|\mathbf{w} - \mathbf{z}\|_2^2 - 2\eta(1 - \eta \sigma_a) \langle \mathbf{w} - \mathbf{z}, \bar{\mathbf{A}}\mathbf{w} - \bar{\mathbf{A}}\mathbf{z} \rangle + \eta^2 \|\bar{\mathbf{A}}\mathbf{w} - \bar{\mathbf{A}}\mathbf{z}\|_2^2 \\ &\leq [(1 - \eta \sigma_a)^2 + \eta^2 \bar{\mathbf{L}}^2 + 2\eta \bar{\mathbf{L}}(\eta \sigma_a - 1)_+] \cdot \|\mathbf{w} - \mathbf{z}\|_2^2. \end{aligned}$$

Combing the results for the forward and backward maps we obtain the estimate (11.11).

A case analysis as in Chen and Rockafellar (1997, p. 431) justifies the optimal choice for  $\eta_t$  and  $q_t$ . ■

Following Chen and Rockafellar (1997) we have chosen to “center” the result in terms of the Lipschitz constant  $\bar{\mathbf{L}}$  of the *barely monotonic* forward map  $\bar{\mathbf{A}}$ . Doing so reveals something fundamental:

If the step size  $\eta_*$  is set accordingly, then the convergence rate  $q_*$  does **not depend on how we split strong monotonicity** between the forward map  $\mathbf{A}$  and backward map  $\mathbf{B}$ .

Of course, splitting the sum  $\mathbf{T} = \mathbf{A} + \mathbf{B}$  into non-shifted versions of  $\mathbf{A}$  and  $\mathbf{B}$  may still lead to drastically different convergence, through changing the Lipschitz constant  $\bar{\mathbf{L}}$  and possibly the easiness of evaluating  $J_{\mathbf{B}}^\eta$ .

If  $\mathbf{A}$  is  $\mathbf{L}$ -Lipschitz, then

$$\|\bar{\mathbf{A}}\mathbf{w} - \bar{\mathbf{A}}\mathbf{z}\|_2^2 = \|\mathbf{A}\mathbf{w} - \mathbf{A}\mathbf{z}\|_2^2 - 2\sigma_a \langle \mathbf{w} - \mathbf{z}, \mathbf{A}\mathbf{w} - \mathbf{A}\mathbf{z} \rangle + \sigma_a^2 \|\mathbf{w} - \mathbf{z}\|_2^2 \leq (\mathbf{L}^2 - \sigma_a^2) \|\mathbf{w} - \mathbf{z}\|_2^2,$$

leading to the simple estimate

$$\bar{\mathbf{L}} \leq \sqrt{\mathbf{L}^2 - \sigma_a^2}, \quad \eta_* = \frac{\sigma}{\mathbf{L}^2 + \sigma_a \sigma_b}, \quad q_* = 1/\sqrt{1 + \frac{\sigma^2}{\mathbf{L}^2 - \sigma_a^2}}.$$

In particular, shifting all strong monotonicity to the forward map  $\mathbf{A}$ , i.e.  $\sigma = \sigma_a$  yields

$$\eta_* = \sigma/\mathbf{L}^2, \quad q_* = \sqrt{1 - \sigma^2/\mathbf{L}^2} \geq q_*, \text{ but it is}$$

- worse than the proximal algorithm  $\mathbf{w}_{t+1} = J_{\mathbf{A}+\mathbf{B}}^{\eta_t} \mathbf{w}_t$ , which is the **most difficult to implement but enjoys the best rate**  $\frac{1}{1+\eta_t\sigma}$ , see ??;
- worse than the reflector-based ??, which is **more difficult to implement than the forward step but enjoys the better rate**  $\sqrt{1 - \sigma/\mathbf{L}}$ , see ??.

Chen, G. H.-G. and R. T. Rockafellar (1997). “Convergence Rates in Forward–Backward Splitting”. *SIAM Journal on Optimization*, vol. 7, no. 2, pp. 421–444.



**Remark 11.14: Some refinements**

We mention some further improvements on Theorem 11.13:

- **Maximality:** Chen and Rockafellar (1997) actually showed that  $T = A + B$  is maximal monotone under the assumptions in Theorem 11.13.
- **Variable metric:** Chen and Rockafellar (1997) considered changing the norm  $\|\mathbf{w}\|_2$  to  $\|\mathbf{w}\|_H := \sqrt{\langle H\mathbf{w}, \mathbf{w} \rangle}$  for some symmetric positive definite matrix  $H$ , and adapting strong monotonicity and Lipschitz continuity to the norm  $\|\cdot\|_H$  (and its dual  $\|\cdot\|_{H^{-1}}$ ). Note that the backward step is now  $(H + \eta B)^{-1}$  while the forward step is  $H - \eta A$ , where  $H = \partial \frac{1}{2} \|\cdot\|_H^2$ . Theorem 11.13 still holds after obvious adjustments. In fact, we may even allow  $H$  to change with  $t$ .
- **Convex function:** When  $A = \partial f$  we follow the same refinement in ?? to get

$$q_t = \begin{cases} \frac{1-\eta_t\sigma_a}{1+\eta_t\sigma_b}, & \text{if } \eta_t \leq \frac{2}{\sigma_a+\bar{L}} \\ \frac{\eta_t\bar{L}-1}{1+\eta_t\sigma_b}, & \text{if } \eta_t \geq \frac{2}{\sigma_a+\bar{L}} \end{cases}, \quad \text{where } \bar{L} := \sigma_a + \bar{L} \implies \eta_\star \equiv \frac{2}{\sigma_a+\bar{L}}, \quad q_\star = \frac{1}{1+2\kappa}, \quad \kappa := \sigma/\bar{L}.$$

- **Localization:** Chen and Rockafellar (1997, Thm 4.1) noted that as long as

$$\frac{1}{\eta_t} > \frac{\sigma_a - \sigma_b}{2} + \frac{\bar{L}}{2} \left( \frac{\bar{L}}{\sigma} \vee 1 \right),$$

where the parameters  $\sigma_a, \sigma_b$  and  $\bar{L}$  are localized w.r.t. a neighborhood around the unique fixed point  $\mathbf{w}_\infty$ , then the forward-backward algorithm (with  $\epsilon_t \equiv \mathbf{0}$ ) does not leave this neighborhood. A similar albeit weaker result was already known in e.g. Dem'yanov and Pevnyi (1972, Thm 4.2).

- **Asymmetry:** Chen and Rockafellar (1997) pointed out the following change-of-variable for reducing asymmetric implementations to symmetric ones:

$$(H + L + \eta B)^{-1}(H + L - \eta A) = (H + \eta(B + L/\eta))^{-1}(H - \eta(A - L/\eta)),$$

where  $H$  is symmetric but the linear map  $L$  may not.

- **Over-relaxation:** We have set  $\gamma_t \equiv 1$  in Theorem 11.13 since under-relaxation (i.e.  $\gamma_t < 1$ ) is clearly not beneficial. However, when  $B$  is Lipschitz continuous and strongly monotone, it may be beneficial to over-relax (i.e.  $\gamma_t > 1$ ), see ??.

Chen, G. H.-G. and R. T. Rockafellar (1997). “Convergence Rates in Forward–Backward Splitting”. *SIAM Journal on Optimization*, vol. 7, no. 2, pp. 421–444.

Dem'yanov, V. F. and A. B. Pevnyi (1972). “Numerical methods for finding saddle points”. *USSR Computational Mathematics and Mathematical Physics*, vol. 12, no. 5, pp. 11–52.

**Example 11.15: Application of forward-backward splitting to VI (Tseng, 1991)**

Let  $L := \{(\mathbf{u}, \mathbf{v}) : M\mathbf{u} + N\mathbf{v} = \mathbf{b}\}$  be an affine subspace and consider the following variational inequality: find  $(\mathbf{u}, \mathbf{v}) \in (\mathcal{U} \times \mathcal{V}) \cap L$  such that

$$\forall (\bar{\mathbf{u}}, \bar{\mathbf{v}}) \in (\mathcal{U} \times \mathcal{V}) \cap L, \quad \langle \bar{\mathbf{u}} - \mathbf{u}, U\mathbf{u} \rangle + \langle \bar{\mathbf{v}} - \mathbf{v}, V\mathbf{v} \rangle + f(\bar{\mathbf{u}}) - f(\mathbf{u}) + g(\bar{\mathbf{v}}) - g(\mathbf{v}) \geq 0, \quad (11.12)$$

where  $f$  and  $g$  are convex functions,  $\mathcal{U}$  and  $\mathcal{V}$  are convex sets, and  $U$  and  $V$  are monotone maps. Under mild conditions, using subdifferential calculus we may rewrite (11.12) as: find  $(\mathbf{u}, \mathbf{v})$  such that

$$\mathbf{0} \in [S + T + \mathcal{N}_L](\mathbf{u}, \mathbf{v}), \quad \text{where } S := U + \partial f + \mathcal{N}_\mathcal{U}, \quad T := V + \partial g + \mathcal{N}_\mathcal{V}.$$

Since  $\mathcal{N}_L = \text{rge}[M, N]^\top$ , equivalently we may reduce to finding some  $\mathbf{w}$  such that

$$M^\top \mathbf{w} \in S\mathbf{u}, \quad N^\top \mathbf{w} \in T\mathbf{v}, \quad M\mathbf{u} + N\mathbf{v} = \mathbf{b} \iff \mathbf{0} \in \underbrace{-\mathbf{b} + MS^{-1}M^\top}_{\mathbf{A}} \mathbf{w} + \underbrace{NT^{-1}N^\top}_{\mathbf{B}} \mathbf{w}.$$

We can now apply forward-backward splitting to obtain the following algorithm:

---

**Algorithm:** Forward-backward splitting for VI (11.12)

---

**Input:**  $\mathbf{w}_0$

```

1 for  $t = 0, 1, \dots$  do
2   find  $\mathbf{u}_t$  s.t.  $\forall \bar{\mathbf{u}} \in \mathcal{U}, f(\bar{\mathbf{u}}) - f(\mathbf{u}_t) + \langle \bar{\mathbf{u}} - \mathbf{u}_t, \mathbf{U}\mathbf{u}_t - M^\top \mathbf{w}_t \rangle \geq 0$  //  $\mathbf{u} \in S^{-1}M^\top \mathbf{w}$ 
3    $\mathbf{w}_{t+1/2} \leftarrow \mathbf{w}_t - \eta_t(M\mathbf{u}_t - \mathbf{b})$  // forward step
   // compute  $(\text{Id} + \eta_t N T^{-1} N^\top)^{-1} \mathbf{w}_{t+1/2}$  using Sherman-Morrison, see ??
4   find  $\mathbf{v}_t$  s.t.  $\forall \bar{\mathbf{v}} \in \mathcal{V}, g(\bar{\mathbf{v}}) - g(\mathbf{v}_t) + \langle \bar{\mathbf{v}} - \mathbf{v}_t, \mathbf{V}\mathbf{v}_t - N^\top(\mathbf{w}_{t+1/2} - \eta_t N \mathbf{v}_t) \rangle \geq 0$  // backward step
5    $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_{t+1/2} - \eta_t N \mathbf{v}_t = \mathbf{w}_t - \eta_t(M\mathbf{u}_t + N\mathbf{v}_t - \mathbf{b})$ 
```

---

Assuming  $M$  has full column rank and  $S$  is strongly monotone, so that  $A$  is inversely strongly monotone and hence convergence (i.e.  $\mathbf{w}_t \rightarrow \mathbf{w}_\infty$ ,  $M\mathbf{u}_t - \mathbf{b} = A\mathbf{w}_t \rightarrow A\mathbf{w}_\infty$  and  $M\mathbf{u}_t + N\mathbf{v}_t - \mathbf{b} = (\mathbf{w}_t - \mathbf{w}_{t+1})/\eta_t \rightarrow \mathbf{0}$ ) and linear rate of convergence immediately follow from Theorem 11.12 and Theorem 11.13 (and Theorem 11.3 if we average), respectively. Note that  $\mathbf{u}_t \rightarrow \mathbf{u}_\infty$  since  $M\mathbf{u}_t$  converges and  $M$  has full column rank. See also Makler-Scheinberg et al. (1996) for inexact implementations.

Tseng, P. (1991). “Applications of a Splitting Algorithm to Decomposition in Convex Programming and Variational Inequalities”. *SIAM Journal on Control and Optimization*, vol. 29, no. 1, pp. 119–138.

Makler-Scheinberg, S., V. H. Nguyen, and J. J. Strodhot (1996). “Family of perturbation methods for variational inequalities”. *Journal of Optimization Theory and Applications*, vol. 89, pp. 423–452.

### Example 11.16: Application to separable convex program (Tseng, 1991)

Next, we consider the separable convex program:

$$\min_{\mathbf{u} \in \mathcal{U}, \mathbf{v} \in \mathcal{V}} f(\mathbf{u}) + g(\mathbf{v}), \quad \text{s.t.} \quad M\mathbf{u} + N\mathbf{v} = \mathbf{b}, \quad (11.13)$$

where  $f$  is **strongly convex** and  $g$  is convex. Let  $\tilde{f} = f + \iota_{\mathcal{U}}$  and  $\tilde{g} = g + \iota_{\mathcal{V}}$  we obtain the dual problem

$$-\min_{\mathbf{w}} \tilde{f}^*(M^\top \mathbf{w}) + \tilde{g}^*(N^\top \mathbf{w}) - \langle \mathbf{b}, \mathbf{w} \rangle. \quad (11.14)$$

Specializing the algorithm in Example 11.15 we obtain:

---

**Algorithm:** Forward-backward splitting for separable convex program (11.13)

---

**Input:**  $\mathbf{w}_0$

```

1 for  $t = 0, 1, \dots$  do
2    $\mathbf{u}_t \leftarrow \operatorname{argmin}_{\mathbf{u} \in \mathcal{U}} f(\mathbf{u}) - \langle M\mathbf{u} + N\mathbf{v}_t - \mathbf{b}, \mathbf{w}_t \rangle$  // forward step  $\nabla \tilde{f}^*(M^\top \mathbf{w}_t)$ 
3    $\mathbf{v}_t \leftarrow \operatorname{argmin}_{\mathbf{v} \in \mathcal{V}} g(\mathbf{v}) - \langle M\mathbf{u}_t + N\mathbf{v} - \mathbf{b}, \mathbf{w}_t \rangle + \frac{\eta_t}{2} \|M\mathbf{u}_t + N\mathbf{v} - \mathbf{b}\|_2^2$  // backward step
4    $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t(M\mathbf{u}_t + N\mathbf{v}_t - \mathbf{b})$ 
```

---

Amazingly, the above algorithm is a perfect interpolation between Uzawa’s Algorithm 8.21 (where **quadratic augmentations** are not present in both  $\mathbf{u}$  and  $\mathbf{v}$ ) and ADMM ?? (where **quadratic augmentations** are present in both  $\mathbf{u}$  and  $\mathbf{v}$ ). Instead, it chooses to *only* augment the update in  $\mathbf{v}$  since the corresponding function  $g$  may not be strongly convex. Here,  $\mathbf{w}_t$  converges to a dual solution while  $\mathbf{u}_t$  converges to (part of) the primal solution (and any limit point of  $\mathbf{v}_t$  consists of the other part of the primal solution). See Mouallif et al. (1991) for inexact implementations.

We also recognize that the algorithm is simply the proximal gradient Algorithm 2.17 applied to the dual (11.14), with smooth component  $\tilde{f}^*(M^\top \mathbf{w})$  and nonsmooth component  $\tilde{g}^*(N^\top \mathbf{w})$ . Indeed, the forward step simply computes the gradient  $\nabla \tilde{f}^*(M^\top \mathbf{w})$  while the backward step reduces to

$$\min_{\mathbf{w}} \langle \mathbf{w}, M\mathbf{u}_t - \mathbf{b} \rangle + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|_2^2 + \tilde{g}^*(N^\top \mathbf{w}_t) \equiv \min_{\mathbf{v}} \tilde{g}(\mathbf{v}) + \frac{\eta_t}{2} \|N\mathbf{v}\|_2^2 - \langle \mathbf{w}_t - \eta_t(M\mathbf{u}_t - \mathbf{b}), N\mathbf{v} \rangle.$$

With this interpretation we may apply Amijo’s rule (see Remark 1.20) to adapt the step size  $\eta_t$  so that

$$\tilde{f}^*(M^\top \mathbf{w}_{t+1}) \leq \tilde{f}^*(M^\top \mathbf{w}_t) + \langle \mathbf{w}_{t+1} - \mathbf{w}_t, M\mathbf{u}_t \rangle + \frac{1}{2\eta_t} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2,$$

where the function value  $\tilde{f}^*$  can already be computed in the forward step.

Tseng, P. (1991). “Applications of a Splitting Algorithm to Decomposition in Convex Programming and Variational Inequalities”. *SIAM Journal on Control and Optimization*, vol. 29, no. 1, pp. 119–138.

Mouallif, K., V. H. Nguyen, and J.-J. Strodiot (1991). “A Perturbed Parallel Decomposition Method for a Class of Nonsmooth Convex Minimization Problems”. *SIAM Journal on Control and Optimization*, vol. 29, no. 4, pp. 829–847.

### Exercise 11.17: Application to finite sum

Let us consider minimizing a convex function of the finite-sum form:

$$\min_{\mathbf{u}} f_0(\mathbf{u}) + \sum_{i=1}^k f_i(\mathbf{u}),$$

where  $f_0$  is **strongly convex** and each  $f_i$  is convex. Applying the product space trick to the latter summation term (see ??), we arrive at a special case of (11.13):

$$\min_{\mathbf{u}, \mathbf{v}_1, \dots, \mathbf{v}_k} f_0(\mathbf{u}) + \sum_{i=1}^k f_i(\mathbf{v}_i), \quad \text{s.t.} \quad \forall i, \mathbf{v}_i = \mathbf{u}.$$

Derive a splitting algorithm based on Example 11.16. Do you recognize the resulting algorithm for the special case where  $f_i = \iota_{C_i}$  and  $f_0(\mathbf{u}) = \|\mathbf{u} - \mathbf{u}_0\|_2$ , i.e. projecting  $\mathbf{u}_0$  to the intersection of convex sets  $C_i$ ?

### Exercise 11.18: Application to affine VI

Consider the (nonlinear) variational inequality:

$$\text{find } \mathbf{w} \quad \text{s.t.} \quad \forall \bar{\mathbf{w}} \in C, \langle \bar{\mathbf{w}} - \mathbf{w}, \mathbf{T}\mathbf{w} \rangle \geq 0, \text{ or more succinctly } \mathbf{0} \in (\mathbf{T} + \mathcal{N}_C)\mathbf{w}, \quad (11.15)$$

where  $\mathbf{T} : C \rightarrow \mathbb{R}^d$  is continuous and monotone, and  $C \subseteq \mathbb{R}^d$  is closed convex. We linearize  $\mathbf{T}$  iteratively:

$$\text{find } \mathbf{w}_{t+1} \quad \text{s.t.} \quad \mathbf{0} \in \underbrace{L(\mathbf{w} - \mathbf{w}_t) + \mathbf{T}\mathbf{w}_t}_{\approx \mathbf{T}\mathbf{w}} + \mathcal{N}_C\mathbf{w}, \quad \text{i.e.} \quad \mathbf{w}_{t+1} \leftarrow (L + \mathcal{N}_C)^{-1}(L - \mathbf{T})\mathbf{w}_t, \quad (11.16)$$

where  $L : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is a positive definite (but not necessarily symmetric) linear map (or equivalently a  $d \times d$  matrix). Complete the following:

- We may decompose  $L = L_s + L_a$ , where  $L_s$  is symmetric positive definite and  $L_a$  is asymmetric.
- Perform change-of-variable  $\mathbf{z} := L_s^{1/2}\mathbf{w}$  and derive from (11.16) that

$$\mathbf{z}_{t+1} \leftarrow [\text{Id} + L_s^{-1/2}(L_a + \mathcal{N}_C)L_s^{-1/2}]^{-1}[\text{Id} - L_s^{-1/2}(\mathbf{T} - L_a)L_s^{-1/2}]\mathbf{z}_t.$$

- Prove that the iterates  $\{\mathbf{w}_t\}$  are well-defined and derive conditions under which they converge to a solution of the VI (11.15).
- Suppose  $\mathbf{T}$  is **linear** and choose  $L = \lambda \text{Id} - D$  for any matrix  $D$ . Note that a **triangular**  $D$  makes the **backward step extremely efficient**. Moreover, the matrix

$$L_s^{-1/2}(\mathbf{T} - L_a)L_s^{-1/2} = \text{Id} + L_s^{-1/2}(\mathbf{T} - L)L_s^{-1/2} = \text{Id} + L_s^{-1/2}(\mathbf{T} + D - \lambda \text{Id})L_s^{-1/2}$$

is symmetric if  $\mathbf{T} + D$  is so, in which case prove that  $\mathbf{w}_t$  converges if  $\lambda$  is sufficiently large.

- Let  $\mathbf{T} = \begin{bmatrix} G & A \\ -A^\top & H \end{bmatrix}$  where  $G$  and  $H$  are symmetric PSD. Set  $L = \lambda \text{Id} - D$  with  $D = \begin{bmatrix} -D_1 & \mathbf{0} \\ 2A^\top & -D_2 \end{bmatrix}$  or  $D = \begin{bmatrix} -D_1 & -2A \\ \mathbf{0} & -D_2 \end{bmatrix}$  for some symmetric PSD  $D_1$  and  $D_2$ . Explicate (11.16) under these choices.

- Derive the underlying (affine) VI for and specialize the previous result to the quadratic program:

$$\min_{\mathbf{w} \in \mathcal{W}} \langle \mathbf{w}, \frac{1}{2} Q \mathbf{w} + \mathbf{c} \rangle, \quad \text{s.t.} \quad A \mathbf{w} = \mathbf{b}.$$

- Further specialize the previous result to the projection problem:

$$\min_{\mathbf{w} \in \cap_i C_i} \|\mathbf{w} - \mathbf{w}_0\|_2 \quad \equiv \quad \min_{\mathbf{w}_i \in C_i} \frac{1}{2} \sum_i \|\mathbf{w}_i - \mathbf{w}_0\|_2^2, \quad \text{s.t.} \quad \forall i \geq 2, \mathbf{w}_1 = \mathbf{w}_i.$$

### Exercise 11.19: Application to linear complementarity

Consider the linear complementarity problem (LCP): find  $\mathbf{w}$  such that

$$Q\mathbf{w} + \mathbf{b} \geq \mathbf{0}, \quad \mathbf{w} \geq \mathbf{0}, \quad \langle \mathbf{w}, Q\mathbf{w} + \mathbf{b} \rangle = 0,$$

where  $Q \in \mathbb{R}^{d \times d}$  is **positive definite but not necessarily symmetric**. Complete the following:

- Prove that LCP is equivalent to: find  $\mathbf{w}$  such that  $\mathbf{0} \in Q\mathbf{w} + \mathbf{b} + \mathcal{N}_{\mathbb{R}_+^d} \mathbf{w}$ .
- Derive a splitting algorithm based on Example 11.15 where we split  $Q = A + B$ .
- Argue that if  $A$  is symmetric and positive semidefinite, then the splitting algorithm converges.

In practice, we aim to find structured (e.g. tri-diagonal)  $B$  so that the backward step is easily carried out.

- Apply the result in Exercise 11.18 with  $L = \lambda \text{Id} - D$  and  $D = R^\top - S$  where  $R$  and  $S$  are the strict upper and lower triangular part of  $Q$ , respectively.

### Example 11.20: Unpacking minimization

To better appreciate the preceding results, let us first consider the special cases where  $\mathbf{T} = \partial f$  for some (closed) convex function  $f$  hence

$$\mathbf{P}_f^\eta(\mathbf{w}) = \left[ \underset{\mathbf{z}}{\operatorname{argmin}} \frac{1}{2\eta} \|\mathbf{w} - \mathbf{z}\|_2^2 + f(\mathbf{z}) \right] = (\text{Id} + \eta \cdot \partial f)^{-1} \mathbf{w}.$$

A solution of  $\text{VI}(C, \mathbf{T})$  amounts to a (global) minimizer of the constrained minimization problem (19.1):

$$\min_{\mathbf{w} \in C \subseteq \mathbb{R}^d} f(\mathbf{w}), \quad \text{or equivalently} \quad \min_{\mathbf{w}} f(\mathbf{w}) + \iota_C(\mathbf{w}).$$

In this setting a weak solution is also a solution (under mild conditions on  $C$  and  $\text{dom } f$ ), which we **assume exists** in the following (otherwise the appropriately constructed iterates will blow up).

- The iterate in Theorem 11.3 amounts to the usual projected (sub)gradient algorithm:

$$\mathbf{w}_{t+1} = \mathbf{P}_C(\mathbf{w}_t - \eta_t \partial f(\mathbf{w}_t)) = \mathbf{P}_C(\text{Id} - \eta_t \partial f) \mathbf{w}_t = \underset{\mathbf{w} \in C}{\operatorname{argmin}} f(\mathbf{w}_t) + \langle \mathbf{w} - \mathbf{w}_t, \nabla f(\mathbf{w}_t) \rangle + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|_2^2.$$

Provided that  $H_t := \sum_{k=0}^t \eta_k \rightarrow \infty$ ,  $\eta_k \rightarrow 0$  and  $f$  Lipschitz continuous, convergence of the averaged sequence  $\bar{\mathbf{w}}_t = \sum_{k=0}^t \eta_k \mathbf{w}_k / H_t$  then follows from Theorem 11.3. This result is fully complementary to Theorem 4.17, which proved convergence in function value under essentially the same assumptions.

- The iterate in Theorem 11.5 amounts to an *implicit* form of projected (sub)gradient:

$$\mathbf{w}_{t+1} = \mathbf{P}_C(\text{Id} + \eta_t \partial f)^{-1} \mathbf{w}_t = \mathbf{P}_C \mathbf{P}_f^{\eta_t}(\mathbf{w}_t), \quad \text{where} \quad \mathbf{P}_f^{\eta_t}(\mathbf{w}_t) = \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|_2^2 + f(\mathbf{w}).$$

Provided that  $H_t := \sum_{k=0}^t \eta_k \rightarrow \infty$  and  $\eta_k \rightarrow 0$ , convergence of the averaged sequence  $\bar{\mathbf{w}}_t = \sum_{k=0}^t \eta_k \mathbf{w}_k / H_t$  then follows from Theorem 11.5 but dispenses the Lipschitz assumption!

- The iterate in ?? amounts to the (exact) proximal point ??:

$$\mathbf{w}_{t+1} = (\text{Id} + \eta_t \partial f + \mathcal{N}_C)^{-1} \mathbf{w}_t = \text{P}_{f + \iota_C}^{\eta_t}(\mathbf{w}_t) = \underset{\mathbf{w} \in C}{\operatorname{argmin}} \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|_2^2 + f(\mathbf{w}),$$

where  $\mathsf{T} = \partial f + \partial \iota_C = \partial f + \mathcal{N}_C$ . Provided that  $H_t := \sum_{k=0}^t \eta_k \rightarrow \infty$ , convergence of the averaged sequence  $\bar{\mathbf{w}}_t = \sum_{k=0}^t \eta_k \mathbf{w}_k / H_t$  then follows from Theorem 11.5 while convergence of  $\mathbf{w}_t$  follows from ??. Compare also the estimate (??) with Theorem 1.17.

Needless to say, among the three algorithms, projected gradient is the easiest while proximal point is the hardest to implement. In fact, we can use projected gradient to solve the subproblems of the other two variants, although often this will not yield any improvement.

### Example 11.21: Unpacking minimax

Let us now consider  $\mathsf{T} = (\partial_{\mathbf{x}} f, \partial_{\mathbf{y}} f)$  for some function  $f$  that is convex in  $\mathbf{x}$  and concave in  $\mathbf{y}$ . We show below the connection to the minimax problem (8.1), recalled here:

$$\inf_{\mathbf{x} \in X \subseteq \mathbb{R}^d} \sup_{\mathbf{y} \in Y \subseteq \mathbb{R}^d} f(\mathbf{x}, \mathbf{y}).$$

In this setting a weak solution of  $\text{VI}(X \times Y, \mathsf{T})$  is also a solution (under mild conditions on  $X \times Y$  and  $\text{dom } f$ ), which we **assume exists** in the following (otherwise the appropriately constructed iterates will blow up).

- The iterate in Theorem 11.3 amounts to the projected (sub)gradient descent ascent Algorithm 8.22:

$$\begin{aligned} \mathbf{x}_{t+1} &= \text{P}_X(\mathbf{x}_t - \eta_t \partial_{\mathbf{x}} f(\mathbf{x}_t, \mathbf{y}_t)) \\ \mathbf{y}_{t+1} &= \text{P}_Y(\mathbf{y}_t + \eta_t \partial_{\mathbf{y}} f(\mathbf{x}_t, \mathbf{y}_t)). \end{aligned}$$

Or more explicitly,

$$(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) = \underset{\mathbf{x} \in X}{\operatorname{argmin}} \underset{\mathbf{y} \in Y}{\operatorname{argmax}} \langle \mathbf{x} - \mathbf{x}_t; \partial_{\mathbf{x}} f(\mathbf{x}_t, \mathbf{y}_t) \rangle + \langle \mathbf{y} - \mathbf{y}_t; \partial_{\mathbf{y}} f(\mathbf{x}_t, \mathbf{y}_t) \rangle + \frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}_t\|_2^2 - \frac{1}{2\eta_t} \|\mathbf{y} - \mathbf{y}_t\|_2^2,$$

Provided that  $H_t := \sum_{k=0}^t \eta_k \rightarrow \infty$ ,  $\eta_k \rightarrow 0$  and  $f$  Lipschitz continuous (in  $\mathbf{x}$  and  $\mathbf{y}$ , respectively), convergence of the averaged sequence  $(\bar{\mathbf{x}}_t, \bar{\mathbf{y}}_t) = \sum_{k=0}^t \eta_k (\mathbf{x}_k, \mathbf{y}_k) / H_t$  then follows from Theorem 11.3. More refined results have already been presented in Remark 11.4.

- The iterate in Theorem 11.5 amounts to an *implicit* form of projected (sub)GDA:

$$\begin{aligned} (\tilde{\mathbf{x}}_{t+1}, \tilde{\mathbf{y}}_{t+1}) &= \underset{\mathbf{x} \in \mathbb{R}^p}{\operatorname{argmin}} \underset{\mathbf{y} \in \mathbb{R}^d}{\operatorname{argmax}} f(\mathbf{x}, \mathbf{y}) + \frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}_t\|_2^2 - \frac{1}{2\eta_t} \|\mathbf{y} - \mathbf{y}_t\|_2^2, \\ (\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) &= (\text{P}_X(\tilde{\mathbf{x}}_{t+1}), \text{P}_Y(\tilde{\mathbf{y}}_{t+1})). \end{aligned} \tag{11.17}$$

Provided that  $H_t := \sum_{k=0}^t \eta_k \rightarrow \infty$  and  $\eta_k \rightarrow 0$ , convergence of the averaged sequence  $(\bar{\mathbf{x}}_t, \bar{\mathbf{y}}_t) = \sum_{k=0}^t \eta_k (\mathbf{x}_k, \mathbf{y}_k) / H_t$  to a solution then follows from Theorem 11.5.

- The iterate in ?? amounts to the (exact) proximal point ??:

$$(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) = \underset{\mathbf{x} \in X}{\operatorname{argmin}} \underset{\mathbf{y} \in Y}{\operatorname{argmax}} f(\mathbf{x}, \mathbf{y}) + \frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}_t\|_2^2 - \frac{1}{2\eta_t} \|\mathbf{y} - \mathbf{y}_t\|_2^2, \tag{11.18}$$

Provided that  $H_t := \sum_{k=0}^t \eta_k \rightarrow \infty$ , convergence of  $(\bar{\mathbf{x}}_t, \bar{\mathbf{y}}_t) = \sum_{k=0}^t \eta_k (\mathbf{x}_k, \mathbf{y}_k) / H_t$  follows from ?? while if  $\sum_{k=0}^{\infty} \eta_k^2 = \infty$ , then  $(\mathbf{x}_t, \mathbf{y}_t)$  also converges to a solution thanks to ??.

We remark that among the three algorithms, GDA is the easiest while proximal point is the hardest to implement. In fact, we can use GDA to solve the subproblems in both (11.17) and (11.18), a seemingly simple idea that we will revisit in ??.