# 19 Splitting

> **Goal**
>
> Splitting, ergodic averaging, gradient-descent-ascent, forward-backward, backward-backward.

> **Alert 19.1: Convention**
>
> Gray boxes are not required hence can be omitted for unenthusiastic readers.
> This note is likely to be updated again soon.

> **Definition 19.2: The splitting/decomposition problem**
>
> Recall the familiar problem of finding a zero of a maximal monotone map $\mathsf{T} : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$. We now add a small twist:
>
> $$\text{find } \mathbf{z} \quad \text{s.t.} \quad \mathbf{0} \in \mathsf{T}\mathbf{z}, \quad \text{where} \quad \mathsf{T} = \mathsf{A} + \mathsf{B}, \tag{19.1}$$
>
> i.e., the map $\mathsf{T}$ can be decomposed into the sum of two maps $\mathsf{A} : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$ and $\mathsf{B} : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$. The catch is that we often cannot evaluate the resolvent $J_\mathsf{T}$ easily (so the proximal point algorithm is not directly applicable), and yet it might be possible to find a decomposition so that both $J_\mathsf{A}$ and $J_\mathsf{B}$ 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 (19.1):
>
> $$\text{find } \mathbf{z}^* \quad \text{s.t.} \quad \mathbf{0} \in \mathsf{T}^*\mathbf{z}^*, \quad \text{where} \quad \mathsf{T}^* := [-\mathsf{A}^{-1} \circ (-\mathrm{Id}) + \mathsf{B}^{-1}]. \tag{19.2}$$
>
> See Example 12.9 for an explanation of the dual when both $\mathsf{A}$ and $\mathsf{B}$ are subdifferentials of convex functions.

> **Theorem 19.3: Ergodic forward-backward splitting converges (Passty 1979)**
>
> Let $\mathsf{B} : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$ be maximal monotone, $\mathsf{A} : \operatorname{dom} \mathsf{B} \rightrightarrows \mathbb{R}^d$ be monotone, and $\mathsf{T} := \mathsf{A} + \mathsf{B}$ be maximal monotone. Let $\mathbf{w}_0 \in \operatorname{dom} \mathsf{A}$ and for all $t \geq 0$ define
>
> $$\mathbf{w}_{t+1} := J_\mathsf{B}^{\eta_t}(\mathbf{w}_t - \eta_t \mathbf{a}_t^*), \quad \text{where} \quad \mathbf{a}_t^* \in \mathsf{A}\mathbf{w}_t, \ \eta_t \geq 0, \tag{19.3}$$
>
> $$\mathbf{z}_t = \sum_{k=0}^t \bar{\eta}_{t,k} \mathbf{w}_k, \quad \text{where} \quad \bar{\eta}_{t,k} := \eta_k / H_t, \qquad H_t := \sum_{k=0}^t \eta_k.$$
>
> The following estimate holds for any $(\mathbf{w}, \mathbf{w}^*) \in \operatorname{gph} \mathsf{T}$ and $\mathbf{b}^* \in \mathsf{B}\mathbf{w}$:
>
> $$\langle \mathbf{z}_t - \mathbf{w}, \mathbf{w}^* \rangle \leq \sum_{k=0}^t \bar{\eta}_{t,k} \langle \mathbf{w}_k - \mathbf{w}, \mathbf{a}_k^* + \mathbf{b}^* \rangle \leq \frac{\|\mathbf{w}_0 - \mathbf{w}\|_2^2 + \sum_{k=0}^t \eta_k^2 \|\mathbf{a}_k^* + \mathbf{b}^*\|_2^2}{2H_t}. \tag{19.4}$$
>
> Moreover,
>
> - if $\sum_t \eta_t^2 \|\mathbf{a}_t^* + \mathbf{b}^*\|_2^2 < \infty$ and $\mathsf{F} := \mathsf{T}^{-1}\mathbf{0} \neq \emptyset$, then $\|\mathbf{w}_t - \mathbf{w}\|_2$ converges for any $\mathbf{w} \in \mathsf{F}$;
>
> - if $\sum_t \eta_t^2 \|\mathbf{a}_t^* + \mathbf{b}^*\|_2^2 < \infty$ and $H_t \to \infty$, then either $\mathsf{F} = \emptyset$, in which case $\|\mathbf{z}_t\| \to \infty$, or $\mathbf{z}_t \rightharpoonup \mathbf{z}_\infty \in \mathsf{F}$ (hence also follows the previous claim).

*Proof:* The assumptions guarantee that the iterates $\{\mathbf{w}_t\}$ are well-defined. We now verify Proposition 16.2, starting with the last condition (III). For any $(\mathbf{w}, \mathbf{w}^*) \in \text{gph}\, \mathsf{T}$ and $\mathbf{b}^* \in \mathsf{B}\mathbf{w}$:

$$\|\mathbf{w}_{k+1} - \mathbf{w}\|_2^2 = \|J_{\mathsf{B}}^{\eta_k}(\mathbf{w}_k - \eta_k \mathbf{a}_k^*) - J_{\mathsf{B}}^{\eta_t}(\mathbf{w} + \eta_k \mathbf{b}^*)\|_2^2$$

$$(\text{ firm nonexpansiveness of } J_{\mathsf{B}}^{\eta_k}\,) \leq \|\mathbf{w}_k - \mathbf{w} - \eta_k \mathbf{a}_k^* - \eta_k \mathbf{b}^*\|_2^2 - \|\mathbf{w}_k - \eta_k \mathbf{a}_k^* - \mathbf{w}_{k+1} - \eta_k \mathbf{b}^*\|_2^2$$

$$= \|\mathbf{w}_k - \mathbf{w}\|_2^2 - \|\mathbf{w}_k - \mathbf{w}_{k+1}\|_2^2 - 2\eta_k \langle \mathbf{w}_{k+1} - \mathbf{w}, \mathbf{a}_k^* + \mathbf{b}^* \rangle$$

$$(\, -\|\mathbf{x}\|_2^2 + 2\langle \mathbf{x}, \mathbf{y} \rangle \leq \|\mathbf{y}\|_2^2\,) \leq \|\mathbf{w}_k - \mathbf{w}\|_2^2 + \eta_k^2 \|\mathbf{a}_k^* + \mathbf{b}^*\|_2^2 - 2\eta_k \langle \mathbf{w}_k - \mathbf{w}, \mathbf{a}_k^* + \mathbf{b}^* \rangle \tag{19.5}$$

$$(\text{ monotonicity of } \mathsf{A}\,) \leq \|\mathbf{w}_k - \mathbf{w}\|_2^2 + \eta_k^2 \|\mathbf{a}_k^* + \mathbf{b}^*\|_2^2 - 2\eta_k \langle \mathbf{w}_k - \mathbf{w}, \mathbf{w}^* \rangle. \tag{19.6}$$

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 (19.4) (using also (19.5)). If $\sum_t \eta_t^2 \|\mathbf{a}_t^* + \mathbf{b}^*\|_2^2 < \infty$ and $H_t \to \infty$, we deduce that

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

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 (19.6) 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,

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

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

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 19.4: 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

$$\eta_t = \frac{1}{\sqrt{\|\mathbf{a}_t^*\|_2^2 + 1}} \frac{1}{(t+1)^p}, \quad p \in (\tfrac{1}{2}, 1], \tag{19.7}$$

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 $H_t \to \infty$ the estimate (19.4) goes to 0 while $\{\mathbf{z}_t\}$ converges to a zero.

It is clear that the proximal gradient Algorithm 4.17, the subgradient Algorithm 5.14 and the gradient-descent-ascent (GDA) Algorithm 12.22 are all special cases of the so-called forward-backward splitting in (19.3). In fact, Theorem 5.17 for the convergence of the subgradient Algorithm 5.14 is strictly contained

in Theorem 19.3, and now we have a similar result for GDA. Indeed, let $\mathsf{A} = (\partial_{\mathbf{x}} f, \partial_{\mathbf{y}}\text{-}f)$ as suggested in Exercise 17.15, for any $\mathbf{w} = (\mathbf{x}, \mathbf{y}) \in C$ we have from (19.4):

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

$$
= \sum_{k=0}^{t} \bar{\eta}_{t,k} \left[ \langle \mathbf{x}_k - \mathbf{x}, \partial_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k) \rangle - f(\mathbf{x}_k, \mathbf{y}_k) + f(\mathbf{x}_k, \mathbf{y}_k) + \langle \mathbf{y}_k - \mathbf{y}, \partial_{\mathbf{y}}\text{-}f(\mathbf{x}_k, \mathbf{y}_k) \rangle \right]
$$

$$
\geq \sum_{k=0}^{t} \bar{\eta}_{t,k} \left[ -f(\mathbf{x}, \mathbf{y}_k) + f(\mathbf{x}_k, \mathbf{y}) \right]
$$

$$
\geq -f(\mathbf{x}, \bar{\mathbf{y}}_t) + f(\bar{\mathbf{x}}_t, \mathbf{y}), \quad \text{where} \quad (\bar{\mathbf{x}}_t, \bar{\mathbf{y}}_t) := \sum_{k=0}^{t} \bar{\eta}_{t,k} \mathbf{w}_k. \tag{19.8}
$$

We can make the following conclusions:

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

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

  where $\text{diam}(C)$ is the diameter of $C$, $\mathsf{L} := \sup_t \|\mathbf{a}_t^*\|_2 < \infty$ and $S_t^2 = \sum_t \eta_t^2$. Thus, the primal and dual gaps go to 0 if $H_t \to \infty$ and $\eta_t \to 0$, 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 $C = \mathsf{X} \times \mathsf{Y}$ with say $\mathsf{X}$ bounded, $\sum_t \|\eta_t \mathbf{a}_t^*\|_2^2 < \infty$, there exists a saddle point, and $\mathsf{A}$ is bounded on bounded sets. Then, setting $\mathbf{y} = \mathbf{y}^\star$ for any $\mathbf{y}^\star \in \mathsf{Y}^\star$ we obtain:

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

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

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

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

- Inspecting the proof of Theorem 19.3 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 (19.8) we need to average $\mathbf{x}_k$ using the step size on $\mathbf{y}$ and vice versa. This observation is useful when only say $\mathsf{X}$ is bounded (such as in a Lagrangian) so that we need only use the adaptive step size (19.7) 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.

**Alert 19.5: Never over-invest on transient steps**

In the regularization approach in Remark 18.46 and the proximal point approach in (19.19) (and also (19.18) and Uzawa's Algorithm 12.21) we need to solve some intermediate subproblems *exactly*, which is clearly wasteful, after all the exact solution will only be used for a single iteration on the next round. Intuitively, an inexact solution would probably do equally well as long as the "inexactness" is proportional to the algorithms'

progress. Indeed, Bakušinkiĭ and Poljak (1974), following Gajewski and Kluge (1970), considered the iterate

$$\mathbf{w}_{t+1} = \mathrm{P}_C(\mathbf{w}_t - \eta_t \mathsf{T}_t \mathbf{w}_t), \quad \text{where} \quad \mathsf{T}_t := \mathsf{T} + \lambda_t \mathrm{Id} - \lambda_t \mathbf{w}_0,$$

which amounts to performing 1 (projected) GDA step on (18.24) and then changing $\lambda_t$ and $\eta_t$ immediately (i.e. proceed to the next round).

Bakušinkiĭ and Poljak (1974) showed that the following condition suffices for $L$-Lipschitz continuous $\mathsf{T}$:

① $0 < \lambda_t \downarrow 0$;    ② $\sum_t \lambda_t \eta_t = \infty$;    ③ $\lambda_t/\lambda_{t+1} = 1 + o(\lambda_t \eta_t)$;    ④ $\limsup_{t\to\infty} \eta_t(L+\lambda_t)^2/\lambda_t < 2/L$;

where the last condition may be relaxed to (the familiar) $\limsup_{t\to\infty}(L+\lambda_t)\eta_t < 2$ if $\mathsf{T} = (\partial_{\mathbf{x}} f, \partial_{\mathbf{y}}\text{-}f)$. For example, if $\eta_t = c\lambda_t, \lambda_t = L/t^p, p \in (0,1/2), c \in (0, 2/L^3)$, or $\eta_t \equiv \eta \in (0, 1/L), \lambda_t = L/t^p, p \in (0,1)$ when $\mathsf{T} = (\partial_{\mathbf{x}} f, \partial_{\mathbf{y}}\text{-}f)$. In fact, convergence to the closest solution, i.e.

$$\mathbf{w}_t \to \operatorname*{argmin}_{\mathbf{w}\in C_\star} \|\mathbf{w} - \mathbf{w}_0\|_2,$$

was claimed in Bakušinkiĭ and Poljak (1974).

Bakušinkiĭ, A. B. and B. T. Poljak (1974). "On the solution of variational inequalities". *Soviet Mathematics Doklady,* vol. 15, no. 6, pp. 1705–1710.
Gajewski, H. and R. Kluge (1970). "Projektionsverfahren bei nichtlinearen Variationsungleichungen". *Mathematische Nachrichten,* vol. 46, no. 1-6, pp. 363–373.

## Theorem 19.6: Ergodic backward-backward splitting converges (Passty 1979)

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

$$\mathbf{w}_{t+1} := J_{\mathsf{B}}^{\eta_t} J_{\mathsf{A}}^{\eta_t} \mathbf{w}_t, \quad \text{where} \quad \eta_t \geq 0, \tag{19.9}$$

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

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

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

*Proof:* We simply verify Proposition 16.2. Let $\mathbf{w} \in \mathrm{dom}\,\mathsf{T}$, $\mathbf{a}^* \in \mathsf{A}\mathbf{w}$ and $\mathbf{b}^* \in \mathsf{B}\mathbf{w}$. Applying the firm nonexpansiveness of $J^{\eta_k}$ (see Exercise 16.9):

$$\|J_{\mathsf{A}}^{\eta_k}\mathbf{w}_k - \mathbf{w}\|_2^2 = \|J_{\mathsf{A}}^{\eta_k}\mathbf{w}_k - J_{\mathsf{A}}^{\eta_k}(\mathbf{w}+\eta_k\mathbf{a}^*)\|_2^2 \leq \|\mathbf{w}_k - \mathbf{w} - \eta_k\mathbf{a}^*\|_2^2 - \|\mathbf{w}_k - J_{\mathsf{A}}^{\eta_k}\mathbf{w}_k - \eta_k\mathbf{a}^*\|_2^2$$

$$= \|\mathbf{w}_k - \mathbf{w}\|_2^2 - \|\mathbf{w}_k - J_{\mathsf{A}}^{\eta_k}\mathbf{w}_k\|_2^2 + 2\eta_k \langle \mathbf{w} - J_{\mathsf{A}}^{\eta_k}\mathbf{w}_k; \mathbf{a}^* \rangle \tag{19.10}$$

$$\|J_{\mathsf{B}}^{\eta_k} J_{\mathsf{A}}^{\eta_k}\mathbf{w}_k - \mathbf{w}\|_2^2 \leq \|J_{\mathsf{A}}^{\eta_k}\mathbf{w}_k - \mathbf{w}\|_2^2 - \|J_{\mathsf{A}}^{\eta_k}\mathbf{w}_k - J_{\mathsf{B}}^{\eta_k} J_{\mathsf{A}}^{\eta_k}\mathbf{w}_k\|_2^2 + 2\eta_k \langle \mathbf{w} - J_{\mathsf{B}}^{\eta_k} J_{\mathsf{A}}^{\eta_k}\mathbf{w}_k; \mathbf{b}^* \rangle.$$

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:

$$\|J_{\mathsf{B}}^{\eta_k} J_{\mathsf{A}}^{\eta_k}\mathbf{w}_k - \mathbf{w}\|_2^2 \leq \|\mathbf{w}_k - \mathbf{w}\|_2^2 + 2\eta_k \langle \mathbf{w} - \mathbf{w}_k; \mathbf{a}^* + \mathbf{b}^* \rangle + \eta_k^2[\|\mathbf{a}^* + \mathbf{b}^*\|_2^2 + \|\mathbf{a}^*\|_2^2]. \tag{19.11}$$

Summing from $k = 0$ to $k = t$ and rearranging as in Theorem 19.3 we obtain for any $\mathbf{w} \in \mathrm{dom}\,\mathsf{T}, \mathbf{w}^* = \mathbf{a}^* + \mathbf{b}^* \in \mathsf{T}\mathbf{w}$:

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

Using the assumptions on $\eta_t$ we thus know

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

whence follows from the maximality of the sum $\mathsf{T}$ that any limit point of $\{\mathbf{z}_t\}$ is a zero. If $\{\mathbf{z}_t\}$ is bounded, then $\mathsf{F} \neq \emptyset$, which we assume now. Let $\mathbf{w} \in \mathsf{F}$ and set $\mathbf{w}^* = \mathbf{0}$ we know from (19.11) that $\{\mathbf{w}_t\}$ is (uniformly) quasi-Fejér monotone w.r.t. $\mathsf{F}$. Lastly, we verify condition (II) in Proposition 16.2 as in Theorem 19.3. ∎

The special case $\mathsf{B} = \mathcal{N}_C$ 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 19.7: Comparing forward-backward and backward-backward

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

$$\mathfrak{F} := J_{\mathsf{B}}^{\eta}(\mathrm{Id} - \eta\mathsf{A}) \qquad vs. \qquad \mathfrak{B} := J_{\mathsf{B}}^{\eta} J_{\mathsf{A}}^{\eta}.$$

For the former, it is clear that

$$\mathbf{w} \in \mathfrak{F}\mathbf{w} \iff [\mathbf{w} + \eta\mathsf{B}\mathbf{w}] \cap [\mathbf{w} - \eta\mathsf{A}\mathbf{w}] \neq \emptyset \iff \mathbf{0} \in (\mathsf{A} + \mathsf{B})\mathbf{w}$$

$$\exists(\mathbf{w}, \mathbf{w}^*) \in \mathrm{gph}\,\mathsf{A}, \mathbf{w} \in \mathfrak{F}\mathbf{w} \iff \exists(\mathbf{w}, \mathbf{w}^*) \in \mathrm{gph}\,\mathsf{A}, \mathbf{0} \in (\mathsf{A} + \mathsf{B})\mathbf{w} \iff \mathbf{0} \in -\mathsf{A}^{-1}(-\mathbf{w}^*) + \mathsf{B}^{-1}\mathbf{w}^*.$$

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:

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

In other words,

<span style="color:#d6336c">backward-backward on $\mathsf{A} + \mathsf{B}$ is forward-backward on $^{\eta}\mathsf{A} + \mathsf{B}$!</span>

In general, $(^{\eta}\mathsf{A} + \mathsf{B})^{-1}\mathbf{0} \cap (\mathsf{A} + \mathsf{B})^{-1}\mathbf{0} = \emptyset$, with one notable exception: when $\mathsf{A}^{-1}\mathbf{0} \cap \mathsf{B}^{-1}\mathbf{0} \neq \emptyset$, see Theorem 19.8 below. This is the reason why in all our results about backward-backward (e.g. Theorem 19.6) we require $\eta_t \to 0$, since then $^{\eta}\mathsf{A} \to {}^{0}\mathsf{A}$ as $\eta \to 0$, where recall that $^{0}\mathsf{A}\mathbf{w}$ is the minimum-norm element in $\mathsf{A}\mathbf{w}$. In contrast, it is possible to use constant $\eta$ in forward-backward (e.g. Theorem 19.14), at the expense of $\eta$ depending on properties of $\mathsf{A}$. 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 19.8: Backward-backward converges under a common fixed point (Tseng 1992)

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

$$\mathbf{w}_{t+1} = (1 - \gamma_t)\mathbf{w}_t + \gamma_t \mathsf{T}_{i(t)}\mathbf{w}_t + \boldsymbol{\epsilon}_t, \quad i_t \in \{1, \ldots, m\}, \quad \gamma_t \in (0, \tfrac{1}{\alpha}), \quad \sum_t \|\boldsymbol{\epsilon}_t\|_2 < \infty$$

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

*Proof:* From the proof of Theorem 16.13 we know $\{\mathbf{w}_t\}$ is (uniformly) quasi-Fejér monotone w.r.t. $\mathsf{F}$ and

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

Let $\mathbf{z} \in \cap_{i \in I}\mathsf{Fix}\mathsf{T}_i$ be a limit point of $\{\mathbf{w}_t\}$ for some $I \neq \emptyset$ (e.g. $I = \{i\}$ for some $i$; see Proposition 16.6). Take a subsequence $\mathbf{w}_{t_k} \to \mathbf{z}$ 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 $\mathbf{w}_{s_k} \to \mathbf{w}$. Since $\mathbf{w}_{s_k} - \mathsf{T}_j\mathbf{w}_{s_k} \to \mathbf{0}$ we have $\mathbf{w} \in \mathsf{Fix}\mathsf{T}_j$ (see Proposition 16.6). Since $\mathbf{z} \in \cap_{i \in I}\mathsf{Fix}\mathsf{T}_i$ and $i(t) \in I$ for $t \in [t_k, s_k)$ we have

$$\|\mathbf{w}_{s_k} - \mathbf{z}\|_2 \leq \|\mathbf{w}_{t_k} - \mathbf{z}\|_2 + \sum_{\kappa=t_k}^{s_k-1} \|\boldsymbol{\epsilon}_t\|_2 \to 0,$$

and hence $\mathbf{z} = \mathbf{w} \in \mathsf{Fix}\mathsf{T}_j$. Since each $\mathsf{T}_i$ appears infinitely often, we may continue the argument to conclude that any limit point $\mathbf{z} \in \mathsf{F}$. Applying Proposition 16.2 we know the whole sequence $\mathbf{w}_t \to \mathbf{w}_\infty \in \mathsf{F}$. ∎

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

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

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 19.9: Method of barycenter (Cimmino 1938)

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

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

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

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

Applying Theorem 19.8, we actually know the more general version

$$\mathbf{w}_{t+1} \leftarrow \mathrm{Avg}(\mathsf{P}_{i_1}, \ldots, \mathsf{P}_{i_{k(t)}})\mathbf{w}_t$$

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 Gleichunger". *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 19.10: (Strong) non-ergodic convergence of backward-backward (Passty 1979)**

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

- *one of* $\mathsf{A}$ *and* $\mathsf{B}$ *is strongly monotone, or*

- $\mathsf{F}$ *has nonempty interior.*

*Then, the (non-ergodic) backward-backward iterate* $\mathbf{w}_t \to \mathbf{w}_\infty \in \mathsf{F}$ *(see (19.9)).*

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

*Proof:* The second claim readily follows from Proposition 16.4. For the first claim, assume w.l.o.g. that $\mathsf{A}$ is $\sigma$-strongly monotone. We strengthen (19.10) into

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

leading (19.11) now to

$$\|\mathbf{w}_{k+1} - \mathbf{w}_\infty\|_2^2 \leq \|\mathbf{w}_k - \mathbf{w}_\infty\|_2^2 + \eta_k^2\|\mathbf{a}^*\|_2^2 - 2\eta_k\sigma\|J_\mathsf{A}^{\eta_k}\mathbf{w}_k - \mathbf{w}_\infty\|_2^2 \implies \liminf \|J_\mathsf{A}^{\eta_k}\mathbf{w}_k - \mathbf{w}_\infty\|_2 = 0, \text{ hence}$$
$$\liminf \|\mathbf{w}_{k+1} - \mathbf{w}_\infty\|_2 = \liminf \|\mathbf{w}_{k+1} - J_\mathsf{B}^{\eta_k}(\mathbf{w}_\infty + \eta_k\mathbf{b}^*)\|_2 \leq \liminf \|J_\mathsf{A}^{\eta_k}\mathbf{w}_k - \mathbf{w}_\infty - \eta_k\mathbf{b}^*\|_2 = 0.$$

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

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 19.11: Inversely strong monotonicity, a.k.a., cocoercive**

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

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

i.e. $\mathsf{T}^{-1}$ is $\sigma$-strongly monotone or equivalently $\sigma\mathsf{T}$ is firmly nonexpansive and hence $\mathsf{T}$ is $\frac{1}{\sigma}$-Lipschitz continuous. When $\mathsf{T} = \partial f$ for a closed (proper) convex function $f$, we know from Alert 6.25 that $\partial f$ is inversely $\sigma$-strongly monotone iff $\partial f$ is $\frac{1}{\sigma}$-Lipschitz continuous.

**Exercise 19.12: Strongly monotone + Lipschitz continuity $\implies$ inversely strongly monotone**

Let $\mathsf{T}$ be $\sigma$-strongly monotone and $\mathsf{L}$-Lipschitz continuous. Prove that $\mathsf{T}$ is inversely $\frac{\sigma}{\mathsf{L}^2}$-strongly monotone. If $\mathsf{T} = \partial f$ for some (closed proper) convex function $f$, we may improve the factor $\frac{\sigma}{\mathsf{L}^2}$ to $\frac{1}{\mathsf{L}}$.

**Theorem 19.13: Non-ergodic convergence of forward-backward**

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

$$\mathbf{w}_{t+1} := (1 - \gamma_t)\mathbf{w}_t + \gamma_t J_\mathsf{B}^{\eta_t}(\mathbf{w}_t - \eta_t\mathbf{a}_t^*) + \boldsymbol{\epsilon}_t, \quad \text{where} \quad \mathbf{a}_t^* \in \mathsf{A}\mathbf{w}_t, \ \eta_t \in [0, \tfrac{2}{\mathsf{L}}], \ \gamma_t \geq 0.$$

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

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

$$\mathfrak{F}_\eta := J_{\mathsf{B}}^\eta(\mathrm{Id} - \eta\mathsf{A}).$$

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

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

is $\frac{\eta\mathsf{L}}{2}$-averaged for any $\eta \in [0, \frac{2}{\mathsf{L}}]$. According to Exercise 16.11, $\mathfrak{F}_\eta$ is $\frac{2}{4-\eta\mathsf{L}}$-averaged. As shown in the proof of Theorem 16.13, $\|\mathbf{w}_t - \mathfrak{F}_{\eta_t}\mathbf{w}_t\|_2 \to 0$. Since $\liminf_t \eta_t \geq \underline{\eta} > 0$, we apply Theorem 17.22 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 Proposition 16.6 and Proposition 16.2 we know the quasi-Fejér monotone sequence $\mathbf{w}_t \rightharpoonup \mathbf{w}_\infty \in \mathsf{F}$.
Since $\mathsf{A}^{-1}$ is strongly monotone, $\mathsf{T}^{*-1}\mathbf{0} = \mathsf{A}\mathbf{w}$ for any $\mathbf{w} \in \mathsf{F}$, see Alert 19.7. Let $\tilde{\mathbf{w}}_t = \mathbf{w}_t - \eta_t\mathsf{A}\mathbf{w}_t$:

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

Since $\liminf_t \eta_t \geq \underline{\eta} > 0$ and we already know $\mathfrak{F}_{\eta_t}\mathbf{w}_t - \mathbf{w}_t \to \mathbf{0}$, it follows $\mathsf{A}\mathbf{w}_t \to \mathsf{A}\mathbf{w}_\infty$. ∎
The primal convergence $\mathbf{w}_t \rightharpoonup \mathbf{w}_\infty$, with $\gamma_t \equiv 1$, $\eta_t \equiv \eta$, $\boldsymbol{\epsilon}_t \equiv \mathbf{0}$ and $\mathsf{B} = \mathcal{N}_C$, appeared in e.g. Mercier (1979, pp 157-158) and Gabay (1983, Thm 6.1). The dual convergence $\mathsf{A}\mathbf{w}_t \to \mathsf{A}\mathbf{w}_\infty$ was due to Tseng (1991) who also considered relaxation and allowed varying step size. Our proof, exploiting the monotonicity in Theorem 17.22, confirms that the usual argument based on Opial's Proposition 16.6 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 19.14: Linear convergence of forward-backward (Chen and Rockafellar 1997)

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

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

*Assume $\bar{\mathsf{A}} := \mathsf{A} - \sigma_a \cdot \mathrm{Id}$ is $\bar{\mathsf{L}}$-Lipschitz continuous, then*

$$\|\mathbf{w}_{t+1} - \mathbf{w}_\infty\|_2 \leq q_t\|\mathbf{w}_t - \mathbf{w}_\infty\|_2 + \|\boldsymbol{\epsilon}_t\|_2, \quad \text{where} \quad \mathbf{w}_\infty \in \mathsf{F} \text{ and } q_t := \frac{\sqrt{[(\eta_t\sigma_a-1)_+ + \bar{\mathsf{L}}]^2 + (1-\eta_t\sigma_a)_+^2}}{1+\sigma_b\eta_t}. \quad (19.12)$$

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

*Proof:* When $\mathsf{B}$ is $\sigma_b$-strongly monotone, we know from Theorem 20.6 that $J_{\mathsf{B}}^\eta$ is $\frac{1}{1+\eta\sigma_b}$-Lipschitz continuous. When $\mathsf{A}$ is $\sigma_a$-strongly monotone and $\bar{\mathsf{A}}$ is $\bar{\mathsf{L}}$-Lipschitz continuous, then

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

Combing the results for the forward and backward maps we obtain the estimate (19.12).

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 $\bar{\mathsf{L}}$ of the *barely monotonic* forward map $\bar{\mathsf{A}}$. Doing so reveals something fundamental:

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

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

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

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

leading to the simple estimate

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

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

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

- worse than the proximal algorithm $\mathbf{w}_{t+1} = J_{\mathsf{A}+\mathsf{B}}^{\eta_t}\mathbf{w}_t$, which is the most difficult to implement but enjoys the best rate $\frac{1}{1+\eta_t\sigma}$, see Corollary 18.19;

- worse than the reflector-based Line 4, which is more difficult to implement than the forward step but enjoys the better rate $\sqrt{1 - \sigma/\mathsf{L}}$, see Theorem 20.16.

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 19.15: Some refinements**

We mention some further improvements on Theorem 19.14:

- Maximality: Chen and Rockafellar (1997) actually showed that $\mathsf{T} = \mathsf{A} + \mathsf{B}$ is maximal monotone under the assumptions in Theorem 19.14.

- 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 $(H + \eta\mathsf{B})^{-1}$ while the forward step is $H - \eta\mathsf{A}$, where $H = \partial\frac{1}{2}\|\cdot\|_H^2$. Theorem 19.14 still holds after obvious adjustments. In fact, we may even allow $H$ to change with $t$.

- Convex function: When $\mathsf{A} = \partial f$ we follow the same refinement in Alert 18.14 to get

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

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

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

where the parameters $\sigma_a, \sigma_b$ and $\bar{\mathsf{L}}$ are localized w.r.t. a neighborhood around the unique fixed point $\mathbf{w}_\infty$, then the forward-backward algorithm (with $\boldsymbol{\epsilon}_t \equiv \mathbf{0}$) 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:

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

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

- Over-relaxation: We have set $\gamma_t \equiv 1$ in Theorem 19.14 since under-relaxation (i.e. $\gamma_t < 1$) is clearly not beneficial. However, when $\mathsf{B}$ is Lipschitz continuous and strongly monotone, it may be beneficial to over-relax (i.e. $\gamma_t > 1$), see Alert 18.14.

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 19.16: 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

$$\forall(\bar{\mathbf{u}}, \bar{\mathbf{v}}) \in (\mathcal{U} \times \mathcal{V}) \cap L, \quad \langle\bar{\mathbf{u}} - \mathbf{u}, \mathsf{U}\mathbf{u}\rangle + \langle\bar{\mathbf{v}} - \mathbf{v}, \mathsf{V}\mathbf{v}\rangle + f(\bar{\mathbf{u}}) - f(\mathbf{u}) + g(\bar{\mathbf{v}}) - g(\mathbf{v}) \geq 0, \qquad (19.13)$$

where $f$ and $g$ are convex functions, $\mathcal{U}$ and $\mathcal{V}$ are convex sets, and $\mathsf{U}$ and $\mathsf{V}$ are monotone maps. Under mild conditions, using subdifferential calculus we may rewrite (19.13) as: find $(\mathbf{u}, \mathbf{v})$ such that

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

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

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

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

---
**Algorithm:** Forward-backward splitting for VI (19.13)

---
**Input:** $\mathbf{w}_0$
1 **for** $t = 0, 1, \ldots$ **do**
2     find $\mathbf{u}_t$ s.t. $\forall\bar{\mathbf{u}} \in \mathcal{U}, \ f(\bar{\mathbf{u}}) - f(\mathbf{u}_t) + \langle\bar{\mathbf{u}} - \mathbf{u}_t, \mathsf{U}\mathbf{u}_t - M^\top\mathbf{w}_t\rangle \geq 0$        `// u ∈ S⁻¹Mᵀw`
3     $\mathbf{w}_{t+1/2} \leftarrow \mathbf{w}_t - \eta_t(M\mathbf{u}_t - \mathbf{b})$                                 `// forward step`
    `// compute (Id + η_t NT⁻¹Nᵀ)⁻¹w_{t+1/2} using Sherman-Morrison, see Proposition 17.18`
4     find $\mathbf{v}_t$ s.t. $\forall\bar{\mathbf{v}} \in \mathcal{V}, \ g(\bar{\mathbf{v}}) - g(\mathbf{v}_t) + \langle\bar{\mathbf{v}} - \mathbf{v}_t, \mathsf{V}\mathbf{v}_t - N^\top(\mathbf{w}_{t+1/2} - \eta_t N\mathbf{v}_t)\rangle \geq 0$    `// 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 $\mathsf{S}$ is strongly monotone, so that $\mathsf{A}$ is inversely strongly monotone and hence convergence (i.e. $\mathbf{w}_t \rightharpoonup \mathbf{w}_\infty$, $M\mathbf{u}_t - \mathbf{b} = \mathsf{A}\mathbf{w}_t \to \mathsf{A}\mathbf{w}_\infty$ and $M\mathbf{u}_t + N\mathbf{v}_t - \mathbf{b} = (\mathbf{w}_t - \mathbf{w}_{t+1})/\eta_t \to \mathbf{0}$) and linear rate of convergence immediately follow from Theorem 19.13 and Theorem 19.14 (and Theorem 19.3 if we average), respectively. Note that $\mathbf{u}_t \to \mathbf{u}_\infty$ since $M\mathbf{u}_t$ converges and $M$ has full column rank. See also Makler-Scheimberg 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-Scheimberg, S., V. H. Nguyen, and J. J. Strodiot (1996). "Family of perturbation methods for variational inequalities". *Journal of Optimization Theory and Applications*, vol. 89, pp. 423–452.

## Example 19.17: Application to separable convex program (Tseng 1991)

Next, we consider the separable convex program:

$$\min_{\mathbf{u}\in\mathcal{U}, \mathbf{v}\in\mathcal{V}} f(\mathbf{u}) + g(\mathbf{v}), \qquad \text{s.t.} \qquad M\mathbf{u} + N\mathbf{v} = \mathbf{b}, \tag{19.14}$$

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

$$-\min_{\mathbf{w}} \ \tilde{f}^*(M^\top \mathbf{w}) + \tilde{g}^*(N^\top \mathbf{w}) - \langle \mathbf{b}, \mathbf{w} \rangle . \tag{19.15}$$

Specializing the algorithm in Line 5 we obtain:

---
**Algorithm:** Forward-backward splitting for separable convex program (19.14)

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

1   **for** $t = 0, 1, \ldots$ **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 12.21 (where quadratic augmentations are not present in both $\mathbf{u}$ and $\mathbf{v}$) and ADMM Example 20.17 (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 4.17 applied to the dual (19.15), 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

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

With this interpretation we may apply Amijo's rule (see Remark 2.20) to adapt the step size $\eta_t$ so that

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

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 19.18: 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 (19.14):

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

Derive a splitting algorithm based on Line 4. 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 19.19: Application to affine VI

Consider the (nonlinear) variational inequality:

$$\text{find } \mathbf{w} \quad \text{s.t.} \quad \forall \bar{\mathbf{w}} \in C, \; \langle \bar{\mathbf{w}} - \mathbf{w}, \mathsf{T}\mathbf{w} \rangle \geq 0, \text{ or more succinctly } \mathbf{0} \in (\mathsf{T} + \mathcal{N}_C)\mathbf{w}, \tag{19.16}$$

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

$$\text{find } \mathbf{w}_{t+1} \quad \text{s.t.} \quad \mathbf{0} \in \underbrace{L(\mathbf{w} - \mathbf{w}_t) + \mathsf{T}\mathbf{w}_t}_{\approx \mathsf{T}\mathbf{w}} + \mathcal{N}_C \mathbf{w}, \quad i.e. \quad \mathbf{w}_{t+1} \leftarrow (L + \mathcal{N}_C)^{-1}(L - \mathsf{T})\mathbf{w}_t, \tag{19.17}$$

where $L : \mathbb{R}^d \to \mathbb{R}^d$ 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 (19.17) that

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

- Prove that the iterates $\{\mathbf{w}_t\}$ are well-defined and derive conditions under which they converge to a solution of the VI (19.16).

- Suppose $\mathsf{T}$ is linear and choose $L = \lambda \mathrm{Id} - D$ for any matrix $D$. Note that a triangular $D$ makes the backward step extremely efficient. Moreover, the matrix

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

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

- Let $\mathsf{T} = \begin{bmatrix} G & A \\ -A^\top & H \end{bmatrix}$ where $G$ and $H$ are symmetric PSD. Set $L = \lambda\mathrm{Id} - D$ 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 (19.17) under these choices.

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

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

- Further specialize the previous result to the projection problem:

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

### Exercise 19.20: 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 Line 5 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 19.19 with $L = \lambda\mathrm{Id} - D$ and $D = R^\top - S$ where $R$ and $S$ are the strict upper and lower triangular part of $Q$, respectively.

---

### Example 19.21: Unpacking minimization

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

$$\mathsf{P}_f^\eta(\mathbf{w}) = \left[\operatorname*{argmin}_{\mathbf{z}} \tfrac{1}{2\eta}\|\mathbf{w} - \mathbf{z}\|_2^2 + f(\mathbf{z})\right] = (\mathrm{Id} + \eta \cdot \partial f)^{-1}\mathbf{w}.$$

A solution of $\mathrm{VI}(C, \mathsf{T})$ amounts to a (global) minimizer of the constrained minimization problem (3.1):

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

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

- The iterate in Theorem 19.3 amounts to the usual projected (sub)gradient algorithm:

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

  Provided that $H_t := \sum_{k=0}^t \eta_k \to \infty$, $\eta_k \to 0$ 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 19.3. This result is fully complementary to Theorem 5.17, which proved convergence in function value under essentially the same assumptions.

- The iterate in Theorem 19.6 amounts to an *implicit* form of projected (sub)gradient:

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

  Provided that $H_t := \sum_{k=0}^t \eta_k \to \infty$ and $\eta_k \to 0$, convergence of the averaged sequence $\bar{\mathbf{w}}_t = \sum_{k=0}^t \eta_k\mathbf{w}_k/H_t$ then follows from Theorem 19.6 but dispenses the Lipschitz assumption!

- The iterate in Theorem 18.17 amounts to the (exact) proximal point Line 3:

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

  where $\mathsf{T} = \partial f + \partial\iota_C = \partial f + \mathcal{N}_C$. Provided that $H_t := \sum_{k=0}^t \eta_k \to \infty$, convergence of the averaged sequence $\bar{\mathbf{w}}_t = \sum_{k=0}^t \eta_k\mathbf{w}_k/H_t$ then follows from Theorem 19.6 while convergence of $\mathbf{w}_t$ follows from Theorem 18.17. Compare also the estimate (18.7) with Theorem 2.17.

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 19.22: Unpacking minimax
>
> Let us now consider $\mathsf{T} = (\partial_{\mathbf{x}} f, \partial_{\mathbf{y}}\text{-} f)$ for some function $f$ that is convex in $\mathbf{x}$ and concave in $\mathbf{y}$. We show below the connection to the minimax problem (12.1), recalled here:
>
> $$\inf_{\mathbf{x} \in \mathsf{X} \subseteq \mathbb{R}^d} \sup_{\mathbf{y} \in \mathsf{Y} \subseteq \mathbb{R}^d} f(\mathbf{x}, \mathbf{y}).$$
>
> In this setting a weak solution of $\mathrm{VI}(\mathsf{X} \times \mathsf{Y}, \mathsf{T})$ is also a solution (under mild conditions on $\mathsf{X} \times \mathsf{Y}$ and $\operatorname{dom} f$), which we assume exists in the following (otherwise the appropriately constructed iterates will blow up).
>
> - The iterate in Theorem 19.3 amounts to the projected (sub)gradient descent ascent Algorithm 12.22:
>
>   $$\mathbf{x}_{t+1} = \mathrm{P}_{\mathsf{X}}(\mathbf{x}_t - \eta_t \partial_{\mathbf{x}} f(\mathbf{x}_t, \mathbf{y}_t))$$
>   $$\mathbf{y}_{t+1} = \mathrm{P}_{\mathsf{Y}}(\mathbf{y}_t + \eta_t \partial_{\mathbf{y}} f(\mathbf{x}_t, \mathbf{y}_t)).$$
>
>   Or more explicitly,
>
>   $$(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) = \operatorname*{argmin}_{\mathbf{x} \in \mathsf{X}} \operatorname*{argmax}_{\mathbf{y} \in \mathsf{Y}} \langle \mathbf{x} - \mathbf{x}_t; \partial_{\mathbf{x}} f(\mathbf{x}_t, \mathbf{y}_t) \rangle + \langle \mathbf{y} - \mathbf{y}_t; \partial_{\mathbf{y}}\text{-} f(\mathbf{x}_t, \mathbf{y}_t) \rangle + \tfrac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}_t\|_2^2 - \tfrac{1}{2\eta_t} \|\mathbf{y} - \mathbf{y}_t\|_2^2,$$
>
>   Provided that $H_t := \sum_{k=0}^t \eta_k \to \infty$, $\eta_k \to 0$ 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 19.3. More refined results have already been presented in Remark 19.4.
>
> - The iterate in Theorem 19.6 amounts to an *implicit* form of projected (sub)GDA:
>
>   $$(\tilde{\mathbf{x}}_{t+1}, \tilde{\mathbf{y}}_{t+1}) = \operatorname*{argmin}_{\mathbf{x} \in \mathbb{R}^p} \operatorname*{argmax}_{\mathbf{y} \in \mathbb{R}^d} f(\mathbf{x}, \mathbf{y}) + \tfrac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}_t\|_2^2 - \tfrac{1}{2\eta_t} \|\mathbf{y} - \mathbf{y}_t\|_2^2, \tag{19.18}$$
>   $$(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) = \left( \mathrm{P}_{\mathsf{X}}(\tilde{\mathbf{x}}_{t+1}), \mathrm{P}_{\mathsf{Y}}(\tilde{\mathbf{y}}_{t+1}) \right).$$
>
>   Provided that $H_t := \sum_{k=0}^t \eta_k \to \infty$ and $\eta_k \to 0$, 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 19.6.
>
> - The iterate in Theorem 18.9 amounts to the (exact) proximal point Line 3:
>
>   $$(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) = \operatorname*{argmin}_{\mathbf{x} \in \mathsf{X}} \operatorname*{argmax}_{\mathbf{y} \in \mathsf{Y}} f(\mathbf{x}, \mathbf{y}) + \tfrac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}_t\|_2^2 - \tfrac{1}{2\eta_t} \|\mathbf{y} - \mathbf{y}_t\|_2^2, \tag{19.19}$$
>
>   Provided that $H_t := \sum_{k=0}^t \eta_k \to \infty$, 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 Remark 18.15 while if $\sum_{k=0}^\infty \eta_k^2 = \infty$, then $(\mathbf{x}_t, \mathbf{y}_t)$ also converges to a solution thanks to Theorem 18.9.
>
> 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 (19.18) and (19.19), a seemingly simple idea that we will revisit in Remark 18.46.