

Auxiliary Material

Efficient Learning using Forward-Backward Splitting

A Batch Convergence Proofs and Corollaries

Before we begin the proofs, we provide a technical lemma.

Lemma 5 (Bounding Step Differences). *Assume that the norms of the subgradients of the functions f and r are bounded as in Eq. (5):*

$$\|\partial f(\mathbf{w})\|^2 \leq Af(\mathbf{w}) + G^2, \quad \|\partial r(\mathbf{w})\|^2 \leq Ar(\mathbf{w}) + G^2.$$

Let $\eta_{t+1} \leq \eta_{t+\frac{1}{2}} \leq \eta_t$ and suppose that $\eta_t \leq 2\eta_{t+\frac{1}{2}}$. If we use the FOBOS update of Eqs. (2) and (3), then for a constant $c \leq 4$ and any vector \mathbf{w}^* ,

$$\begin{aligned} & 2\eta_t(1 - cA\eta_t)f(\mathbf{w}_t) + 2\eta_{t+\frac{1}{2}}(1 - cA\eta_{t+\frac{1}{2}})r(\mathbf{w}_{t+1}) \\ & \leq 2\eta_t f(\mathbf{w}^*) + 2\eta_{t+\frac{1}{2}} r(\mathbf{w}^*) + \|\mathbf{w}_t - \mathbf{w}^*\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 + 7\eta_t \eta_{t+\frac{1}{2}} G^2. \end{aligned} \quad (12)$$

Proof of Lemma 5 We begin with a few simple consequences of the forward-looking subgradient steps before proceeding with the core of the proof. Note first that for some $\mathbf{g}_t^f \in \partial f(\mathbf{w}_t)$ and $\mathbf{g}_{t+1}^r \in \partial r(\mathbf{w}_{t+1})$, we have as in Eq. (4)

$$\mathbf{w}_{t+1} - \mathbf{w}_t = -\eta_t \mathbf{g}_t^f - \eta_{t+\frac{1}{2}} \mathbf{g}_{t+1}^r. \quad (13)$$

The definition of a subgradient implies that for any $\mathbf{g}_{t+1}^r \in \partial r(\mathbf{w}_{t+1})$ (and similarly for any $\mathbf{g}_t^f \in \partial f(\mathbf{w}_t)$ with $f(\mathbf{w}_t)$ and $f(\mathbf{w}^*)$) implies that

$$r(\mathbf{w}^*) \geq r(\mathbf{w}_{t+1}) + \langle \mathbf{g}_{t+1}^r, \mathbf{w}^* - \mathbf{w}_{t+1} \rangle \Rightarrow -\langle \mathbf{g}_{t+1}^r, \mathbf{w}_{t+1} - \mathbf{w}^* \rangle \leq r(\mathbf{w}^*) - r(\mathbf{w}_{t+1}). \quad (14)$$

From the Cauchy-Schwartz Inequality and Eq. (13), we obtain

$$\begin{aligned} \langle \mathbf{g}_{t+1}^r, (\mathbf{w}_{t+1} - \mathbf{w}_t) \rangle &= \langle \mathbf{g}_{t+1}^r, (-\eta_t \mathbf{g}_t^f - \eta_{t+\frac{1}{2}} \mathbf{g}_{t+1}^r) \rangle \\ &\leq \|\mathbf{g}_{t+1}^r\| \|\eta_{t+\frac{1}{2}} \mathbf{g}_{t+1}^r + \eta_t \mathbf{g}_t^f\| \leq \eta_{t+\frac{1}{2}} \|\mathbf{g}_{t+1}^r\|^2 + \eta_t \|\mathbf{g}_{t+1}^r\| \|\mathbf{g}_t^f\| \\ &\leq \eta_{t+\frac{1}{2}} (Ar(\mathbf{w}_{t+1}) + G^2) + \eta_t \max \{Af(\mathbf{w}_t) + G^2, Ar(\mathbf{w}_{t+1}) + G^2\}. \end{aligned} \quad (15)$$

We now proceed to bound the difference between \mathbf{w}^* and \mathbf{w}_{t+1} , and using a telescoping sum we will eventually bound $f(\mathbf{w}_t) + r(\mathbf{w}_t) - f(\mathbf{w}^*) - r(\mathbf{w}^*)$. First, we expand norm squared of the difference between \mathbf{w}_t and \mathbf{w}_{t+1} ,

$$\begin{aligned} \|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 &= \|\mathbf{w}_t - (\eta_t \mathbf{g}_t^f + \eta_{t+\frac{1}{2}} \mathbf{g}_{t+1}^r) - \mathbf{w}^*\|^2 \\ &= \|\mathbf{w}_t - \mathbf{w}^*\|^2 - 2 \left[\eta_t \langle \mathbf{g}_t^f, \mathbf{w}_t - \mathbf{w}^* \rangle + \eta_{t+\frac{1}{2}} \langle \mathbf{g}_{t+1}^r, \mathbf{w}_t - \mathbf{w}^* \rangle \right] + \|\eta_t \mathbf{g}_t^f + \eta_{t+\frac{1}{2}} \mathbf{g}_{t+1}^r\|^2 \\ &= \|\mathbf{w}_t - \mathbf{w}^*\|^2 - 2\eta_t \langle \mathbf{g}_t^f, \mathbf{w}_t - \mathbf{w}^* \rangle + \|\eta_t \mathbf{g}_t^f + \eta_{t+\frac{1}{2}} \mathbf{g}_{t+1}^r\|^2 \\ &\quad - 2\eta_{t+\frac{1}{2}} [\langle \mathbf{g}_{t+1}^r, \mathbf{w}_{t+1} - \mathbf{w}^* \rangle - \langle \mathbf{g}_{t+1}^r, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle]. \end{aligned} \quad (16)$$

We can bound the third term above by noting that

$$\begin{aligned} & \|\eta_t \mathbf{g}_t^f + \eta_{t+\frac{1}{2}} \mathbf{g}_{t+1}^r\|^2 \\ &= \eta_t^2 \|\mathbf{g}_t^f\|^2 + 2\eta_t \eta_{t+\frac{1}{2}} \langle \mathbf{g}_t^f, \mathbf{g}_{t+1}^r \rangle + \eta_{t+\frac{1}{2}}^2 \|\mathbf{g}_{t+1}^r\|^2 \\ &\leq \eta_t^2 Af(\mathbf{w}_t) + 2A\eta_t \eta_{t+\frac{1}{2}} \max \{f(\mathbf{w}_t), r(\mathbf{w}_{t+1})\} + \eta_{t+\frac{1}{2}}^2 Ar(\mathbf{w}_{t+1}) + 4\eta_t^2 G^2. \end{aligned}$$

We now use Eq. (15) to bound the last term of Eq. (16) and the above bound on $\eta_t \mathbf{g}_t^f + \eta_{t+\frac{1}{2}} \mathbf{g}_{t+1}^r$ to get that

$$\begin{aligned}
& \|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 \\
& \leq \|\mathbf{w}_t - \mathbf{w}^*\|^2 - 2\eta_t \langle \mathbf{g}_t^f, \mathbf{w}_t - \mathbf{w}^* \rangle - 2\eta_{t+\frac{1}{2}} \langle \mathbf{g}_{t+1}^r, \mathbf{w}_{t+1} - \mathbf{w}^* \rangle + \|\eta_t \mathbf{g}_t^f + \eta_{t+\frac{1}{2}} \mathbf{g}_{t+1}^r\|^2 \\
& \quad + 2\eta_{t+\frac{1}{2}} \left(\eta_{t+\frac{1}{2}} Ar(\mathbf{w}_{t+1}) + 2A\eta_t \max\{f(\mathbf{w}_t), r(\mathbf{w}_{t+1})\} + 2\eta_t G^2 \right) \\
& \leq \|\mathbf{w}_t - \mathbf{w}^*\|^2 + 2\eta_t (f(\mathbf{w}^*) - f(\mathbf{w}_t)) + 2\eta_{t+\frac{1}{2}} (r(\mathbf{w}^*) - r(\mathbf{w}_t)) + 7\eta_t^2 G^2 \\
& \quad + \eta_t^2 Af(\mathbf{w}_t) + 3A\eta_t \eta_{t+\frac{1}{2}} \max\{f(\mathbf{w}_t), r(\mathbf{w}_t)\} + 2\eta_{t+\frac{1}{2}}^2 Ar(\mathbf{w}_{t+1}) \tag{17}
\end{aligned}$$

$$\begin{aligned}
& \leq \|\mathbf{w}_t - \mathbf{w}^*\|^2 + 7\eta_t^2 G^2 \\
& \quad + 2\eta_t (f(\mathbf{w}^*) - (1 - cA\eta_t)f(\mathbf{w}_t)) + 2\eta_{t+\frac{1}{2}} \left(r(\mathbf{w}^*) - (1 - cA\eta_{t+\frac{1}{2}})r(\mathbf{w}_{t+1}) \right) . \tag{18}
\end{aligned}$$

To obtain Eq. (17) we used the standard convexity bounds established earlier in Eq. (14). The final bound given by Eq. (18) is due to the fact that $3A\eta_t \eta_{t+\frac{1}{2}} \leq 6A\eta_t^2$ and that for $a, b \geq 0$, $\max\{a, b\} \leq a + b$. Moving the $f(\cdot)$ and $r(\cdot)$ terms to the left hand side of the \leq gives the desired inequality. \square

Using the above lemma, the analysis for FOBOS in a batch setting is straightforward. In this setting we set $\eta_{t+\frac{1}{2}} = \eta_{t+1}$ and update \mathbf{w}_t to \mathbf{w}_{t+1} as prescribed by Eq. (2) and Eq. (3).

Proof of Theorem 1 Rearranging the $f(\mathbf{w}^*)$ and $r(\mathbf{w}^*)$ terms from the bound in lemma 5, we sum the loss terms over t from 1 through T and get a canceling telescoping sum:

$$\begin{aligned}
& \sum_{t=1}^T [\eta_t ((1 - cA\eta_t)f(\mathbf{w}_t) - f(\mathbf{w}^*)) + \eta_{t+1} ((1 - cA\eta_{t+1})r(\mathbf{w}_{t+1}) - r(\mathbf{w}^*))] \\
& \leq \|\mathbf{w}_1 - \mathbf{w}^*\|^2 - \|\mathbf{w}_{T+1} - \mathbf{w}^*\|^2 + 7G^2 \sum_{t=1}^T \eta_t^2 \leq \|\mathbf{w}_1 - \mathbf{w}^*\|^2 + 7G^2 \sum_{t=1}^T \eta_t^2 . \tag{19}
\end{aligned}$$

Now we bound the one-off $r(\mathbf{w}_{t+1})$ terms by noting that

$$\begin{aligned}
& \sum_{t=1}^T \eta_{t+1} ((1 - cA\eta_{t+1})r(\mathbf{w}_{t+1}) - r(\mathbf{w}^*)) \\
& = \sum_{t=1}^T \eta_t ((1 - cA\eta_t)r(\mathbf{w}_t) - r(\mathbf{w}^*)) + \eta_{T+1} ((1 - cA\eta_{T+1})r(\mathbf{w}_{T+1}) - r(\mathbf{w}^*)) + \eta_1 r(\mathbf{w}^*) \\
& \geq \sum_{t=1}^T \eta_t ((1 - cA\eta_t)r(\mathbf{w}_{t+1}) - r(\mathbf{w}^*)) + r(\mathbf{w}^*)(\eta_1 - \eta_{T+1}) \\
& \geq \sum_{t=1}^T \eta_t ((1 - cA\eta_t)r(\mathbf{w}_t) - r(\mathbf{w}^*)) . \tag{20}
\end{aligned}$$

Using the fact that $\|\mathbf{w}_1 - \mathbf{w}^*\| = \|\mathbf{w}^*\| \leq D$, we combine Eq. (19) with Eq. (20) to get the desired bound. \square

Corollary 6 (Convergence of decreasing step sizes). *Assume that the conditions of Thm. 1 hold and the step sizes η_t are such that $\eta_t \rightarrow 0$, and that $\sum_{t=1}^{\infty} \eta_t = \infty$. Then*

$$\liminf_{t \rightarrow \infty} f(\mathbf{w}_t) + r(\mathbf{w}_t) - (f(\mathbf{w}^*) + r(\mathbf{w}^*)) = 0 .$$

We can give tighter convergence results when we assume that f and r are Lipschitz or when we can guarantee that the $\|\partial f\|$ and $\|\partial r\|$ are bounded. In this case, we have

Corollary 7. *In addition to the conditions of Thm. 1, assume that the norm of any subgradient from ∂f and the norm of any subgradient from ∂r are bounded by G . Then*

$$\min_{t \in \{1, \dots, T\}} (f(\mathbf{w}_t) + r(\mathbf{w}_t)) - (f(\mathbf{w}^*) + r(\mathbf{w}^*)) \leq \frac{D^2 + 7G^2 \sum_{t=1}^T \eta_t^2}{2 \sum_{t=1}^T \eta_t}. \quad (21)$$

Corollary 8 (Optimal fixed step rate). *Assume that the conditions of Cor. 7 hold and that we run FOBOS for a predefined T iterations with $\eta_t = \frac{D}{\sqrt{7TG}}$. Then the following bound holds.*

$$\min_{i \in \{1, \dots, T\}} f(\mathbf{w}_t) + r(\mathbf{w}_t) - (f(\mathbf{w}^*) + r(\mathbf{w}^*)) \leq \frac{3DG}{\sqrt{T}}.$$

In corollaries 7 and 8, our assumption on bounded subgradients of the functions f and r is in practice not restrictive. If we know that an optimal \mathbf{w}^* lies in some closed and bounded set Ω and that $\Omega \subseteq \text{dom}(f + r)$, then [17, Theorem 24.7] guarantees that ∂f and ∂r are bounded for $\mathbf{w} \in \Omega$. The lingering question is thus whether we can guarantee that such a set Ω exists and that our iterates \mathbf{w}_t remain in Ω . We now describe a simple setting to show that ∂f and ∂r are indeed often bounded. If $r(\mathbf{w})$ is a norm and f is lower bounded by 0, then we know that $r(\mathbf{w}^*) \leq f(\mathbf{w}^*) + r(\mathbf{w}^*) \leq f(\mathbf{w}_1) + r(\mathbf{w}_1)$. Using standard bounds on norms, we get that for some $\gamma > 0$

$$\|\mathbf{w}^*\|_\infty \leq \gamma r(\mathbf{w}^*) \leq \gamma(f(\mathbf{w}_1) + r(\mathbf{w}_1)) = \gamma f(\mathbf{w}_1),$$

where for the last inequality we used the assumption that $r(\mathbf{w}_1) = 0$. Thus, we obtain that \mathbf{w}^* lies in a hypercube. We can easily project onto this box by truncating elements of \mathbf{w}_t lying outside it at any iteration without affecting the bounds in Eq. (21). Concretely, this follows since Euclidean projection Π_Ω to a convex set Ω with $\mathbf{w}^* \in \Omega$ satisfies $\|\Pi_\Omega(\mathbf{w}_{t+1}) - \mathbf{w}^*\| \leq \|\mathbf{w}_{t+1} - \mathbf{w}^*\|$. Further, so long as Ω is a norm ball, we know that

$$r(\Pi_\Omega(\mathbf{w}_{t+1})) \leq r(\mathbf{w}_{t+1}). \quad (22)$$

Thus, looking at Eq. (17) in our proof of Theorem 1 we notice that $r(\mathbf{w}^*) - r(\mathbf{w}_{t+1}) \leq r(\mathbf{w}^*) - r(\Pi_\Omega(\mathbf{w}_{t+1}))$ and the series of inequalities through Eq. (18) still hold (with $A = 0$). In general, so long as Eq. (22) holds and $\mathbf{w}^* \in \Omega$, we can project \mathbf{w}_{t+1} into Ω without affecting convergence guarantees.

B Online Regret Proofs and Corollaries

Proof of Theorem 3 Looking at lemma 5, we immediately see that if $\|\partial f\|$ and $\|\partial r\|$ are bounded by G ,

$$f_t(\mathbf{w}_t) - f_t(\mathbf{w}^*) + r(\mathbf{w}_{t+1}) - r(\mathbf{w}^*) \leq \frac{1}{2\eta_t} (\|\mathbf{w}_t - \mathbf{w}^*\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2) + \frac{7}{2} G^2 \eta_t. \quad (23)$$

Now we use Eq. (23) to obtain that

$$\begin{aligned} R_{f+r}(T) &= \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{w}^*) + r(\mathbf{w}_t) - r(\mathbf{w}^*)) + r(\mathbf{w}_{T+1}) - r(\mathbf{w}^*) - r(\mathbf{w}_1) + r(\mathbf{w}^*) \\ &\leq GD + \sum_{t=1}^T \frac{1}{2\eta_t} (\|\mathbf{w}_t - \mathbf{w}^*\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2) + \frac{7G^2}{2} \sum_{t=1}^T \eta_t \end{aligned}$$

since $r(\mathbf{w}) \leq r(\mathbf{0}) + G\|\mathbf{w}\| \leq GD$. We can rewrite the above bound and see

$$\begin{aligned} R_{f+r}(T) &\leq GD + \frac{1}{2\eta_1} \|\mathbf{w}_1 - \mathbf{w}^*\|^2 + \frac{1}{2} \sum_{t=2}^T \|\mathbf{w}_t - \mathbf{w}^*\|^2 \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) + \frac{7G^2}{2} \sum_{t=1}^T \eta_t \\ &\leq GD + \frac{D^2}{2\eta_1} + \frac{D^2}{2} \sum_{t=2}^T \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) + \frac{7G^2}{2} \sum_{t=1}^T \eta_t, \end{aligned}$$

where we used again the bound on the distance of each \mathbf{w}_t to \mathbf{w}^* for the last inequality. Lastly, we use the fact that the sum $\frac{1}{\eta_1} + \sum_{t=2}^T (\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}})$ telescopes and get that

$$R_{f+r}(T) \leq GD + \frac{D^2}{2\eta_T} + \frac{7G^2}{2} \sum_{t=1}^T \eta_t .$$

Setting $\eta_t = c/\sqrt{t}$ and recognizing as in [7] that $\sum_{t=1}^T \eta_t \leq 2c\sqrt{T}$ concludes the proof. \square

We assume as in Sec. 3 that we are minimizing $f(\mathbf{w}) + r(\mathbf{w})$. Suppose that on each step of FOBOS, we choose instead of some $\mathbf{g}_t^f \in \partial f(\mathbf{w}_t)$ a stochastic estimate of the gradient $\tilde{\mathbf{g}}_t^f$ where $E[\tilde{\mathbf{g}}_t^f] \in \partial f(\mathbf{w}_t)$. We assume that we still use the true r (which is generally easy, as it is simply the regularization function). It is straightforward to use Theorem 3 above as in the derivation of Theorems 2 and 3 from [14] to derive the following corollary on the convergence rate of stochastic FOBOS.

Corollary 9. *Assume that the conditions on ∂f , ∂r , and \mathbf{w}^* hold as in the previous theorems and let FOBOS be run for T iterations. Let s be an integer chosen uniformly at random from $\{1, \dots, T\}$. If $\eta_t = \frac{D}{4G\sqrt{t}}$, then*

$$E_s[f(\mathbf{w}_s) + r(\mathbf{w}_s)] \leq f(\mathbf{w}^*) + r(\mathbf{w}^*) + \frac{2GD + 4GD\sqrt{T}}{T}.$$

With probability at least $1 - \delta$, $f(\mathbf{w}_s) + r(\mathbf{w}_s) \leq f(\mathbf{w}^*) + r(\mathbf{w}^*) + \frac{2GD + 4GD\sqrt{T}}{\delta T}$.

C High-dimensional Efficiency

Proof of Proposition 4 It suffices to show that the proposition is correct for $T = 2$ and then use an inductive argument, because the proposition trivially holds for $T = 1$. We provide here a direct proof for each norm separately by examining the updates we derived in Sec. 4 and showing that $\mathbf{w}_2 = \mathbf{w}^*$.

Note that the objective functions are separable for $q = 1$. Therefore, for ℓ_1 -regularization it suffices to prove the proposition for any component of the vector \mathbf{w} . We omit the index of the component and denote by $w_0, w_1, w_2, w_3, \dots$ one coordinate of \mathbf{w} along the iterations of $\mathcal{P}.1$ and by w^* the result for the same component when solving $\mathcal{P}.2$. We need to show that $w^* = w_2$. Expanding the ℓ_1 -update of Eq. (6) over two iterations we get the following:

$$\begin{aligned} w_2 &= \text{sign}(w_1) [|w_1| - \lambda_2]_+ = \text{sign}(w_1) [|\text{sign}(w_0) [|w_0| - \lambda_1]_+| - \lambda_2]_+ \\ &= \text{sign}(w_0) [|w_0| - \lambda_1 - \lambda_2]_+ , \end{aligned}$$

where we used the positivity of $|\cdot|$. Examining $\mathcal{P}.2$ and using Eq. (6) again we get

$$w^* = \text{sign}(w_0) [|w_0| - \lambda_1 - \lambda_2]_+ .$$

Therefore, $w^* = w_2$ as claimed.

Next we prove the proposition for ℓ_2 , returning to using the entire vector for the proof. Using the explicit ℓ_2 -update from Eq. (7), we can expand the norm of the vector \mathbf{w}_1 due to the program $\mathcal{P}.1$ as follows,

$$\|\mathbf{w}_1\| = \left[1 - \frac{\lambda_1}{\|\mathbf{w}_0\|} \right]_+ \|\mathbf{w}_0\| = [\|\mathbf{w}_0\| - \lambda_1]_+ .$$

Similarly, we get that $\|\mathbf{w}_2\| = [\|\mathbf{w}_1\| - \lambda_2]_+$. Combining the norm equalities we see that the norm of \mathbf{w}_2 due to the succession of the two updates is

$$\|\mathbf{w}_2\| = [([\|\mathbf{w}_0\| - \lambda_1]_+ - \lambda_2)]_+ = [\|\mathbf{w}_0\| - \lambda_1 - \lambda_2]_+ .$$

Computing directly the norm of \mathbf{w}^* due to the update given by Eq. (7) yields

$$\|\mathbf{w}^*\| = \left[1 - \frac{\lambda_1 + \lambda_2}{\|\mathbf{w}_0\|} \right]_+ \|\mathbf{w}_0\| = [\|\mathbf{w}_0\| - \lambda_1 - \lambda_2]_+ .$$

Thus, \mathbf{w}^* and \mathbf{w}_2 have the same norm. Since the update itself retains the direction of the original vector \mathbf{w}_0 , we get that $\mathbf{w}^* = \mathbf{w}_2$ as needed.

We now turn to the most complicated update and proof of the three norms, the ℓ_∞ norm. We start by recapping the programs $\mathcal{P}.1$ and $\mathcal{P}.2$ for $T = 2$ and $q = \infty$,

$$\mathcal{P}.1 : \quad \mathbf{w}_1 = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|^2 + \lambda_1 \|\mathbf{w}\|_\infty \right\} \quad (24)$$

$$\mathbf{w}_2 = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\mathbf{w} - \mathbf{w}_1\|^2 + \lambda_2 \|\mathbf{w}\|_\infty \right\} , \quad (25)$$

$$\mathcal{P}.2 : \quad \mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\mathbf{w} - \mathbf{v}\|_2^2 + (\lambda_1 + \lambda_2) \|\mathbf{w}\|_\infty \right\} . \quad (26)$$

We prove the equivalence of the two programs in two stages. First, we examine the case $\|\mathbf{w}_0\|_1 > \lambda_1 + \lambda_2$, and then consider the complement case $\|\mathbf{w}_0\|_1 \leq \lambda_1 + \lambda_2$. For concreteness and simplicity, we assume that $\mathbf{w}_0 \succeq \mathbf{0}$, since, clearly, the objective is symmetric in \mathbf{w}_0 and $-\mathbf{w}_0$. We thus can assume that all entries of \mathbf{w}_0 are non-negative. In the proof we use the following operators: $[v]_+$ now denotes the positive component of each entry of \mathbf{v} , $\min\{\mathbf{v}, \theta\}$ denotes the component-wise minimum between the elements of \mathbf{v} and θ , and likewise $\max\{\mathbf{v}, \theta\}$ is the component-wise maximum. Starting with the case $\|\mathbf{w}_0\|_1 > \lambda_1 + \lambda_2$, we examine Eq. (24). From Lagrange duality we know that that $\mathbf{w}_1 = \mathbf{w}_0 - \boldsymbol{\alpha}_1$, where $\boldsymbol{\alpha}_1$ is the solution of

$$\min_{\boldsymbol{\alpha}} \frac{1}{2} \|\boldsymbol{\alpha} - \mathbf{w}_0\|_2^2 \quad \text{s.t.} \quad \|\boldsymbol{\alpha}\|_1 \leq \lambda_1 .$$

As described by [6] and reviewed above in Sec. 4, $\boldsymbol{\alpha}_1 = [\mathbf{w}_0 - \theta_1]_+$ for some $\theta_1 \in \mathbb{R}_+$. The form of $\boldsymbol{\alpha}_1$ readily translates to the following form for \mathbf{w}_1 : $\mathbf{w}_1 = \mathbf{w}_0 - \boldsymbol{\alpha}_1 = \min(\mathbf{w}_0, \theta_1)$. Applying similar reasoning to the second step of $\mathcal{P}.1$ yields $\mathbf{w}_2 = \mathbf{w}_1 - \boldsymbol{\alpha}_2 = \mathbf{w}_0 - \boldsymbol{\alpha}_1 - \boldsymbol{\alpha}_2$, where $\boldsymbol{\alpha}_2$ is the minimizer of

$$\frac{1}{2} \|\boldsymbol{\alpha} - \mathbf{w}_1\|_2^2 = \frac{1}{2} \|\boldsymbol{\alpha} - (\mathbf{w}_0 - \boldsymbol{\alpha}_1)\|_2^2 \quad \text{s.t.} \quad \|\boldsymbol{\alpha}\|_1 \leq \lambda_2 .$$

Again, we have $\boldsymbol{\alpha}_2 = [\mathbf{w}_1 - \theta_2]_+ = [\mathbf{w}_0 - \boldsymbol{\alpha}_1 - \theta_2]_+$ for some $\theta_2 \in \mathbb{R}_+$. The successive steps then imply that

$$\mathbf{w}_2 = \min\{\mathbf{w}_1, \theta_2\} = \min\{\min\{\mathbf{w}_0, \theta_1\}, \theta_2\} .$$

We next show that regardless of the ℓ_1 -norm of \mathbf{w}_0 , $\theta_2 \leq \theta_1$. Intuitively, if $\theta_2 > \theta_1$, the second minimization step of $\mathcal{P}.1$ would perform no shrinkage of \mathbf{w}_1 to get \mathbf{w}_2 . Formally, assume for the sake of contradiction that $\theta_2 > \theta_1$. Under this assumption, we would have that $\mathbf{w}_2 = \min\{\min\{\mathbf{w}_0, \theta_1\}, \theta_2\} = \min\{\mathbf{w}_0, \theta_1\} = \mathbf{w}_1$. In turn, we obtain that $\mathbf{0}$ belongs to the subgradient set of Eq. (25) when evaluated at $\mathbf{w} = \mathbf{w}_1$, thus,

$$\mathbf{0} \in \mathbf{w}_1 - \mathbf{w}_1 + \lambda_2 \partial \|\mathbf{w}_1\|_\infty = \lambda_2 \partial \|\mathbf{w}_1\|_\infty .$$

Clearly, the set $\partial \|\mathbf{w}_1\|_\infty$ can contain $\mathbf{0}$ only when $\mathbf{w}_1 = \mathbf{0}$. Since we assumed that $\lambda_1 < \|\mathbf{w}_0\|_1$, and hence that $\boldsymbol{\alpha}_1 \preceq \mathbf{w}_0$ and $\boldsymbol{\alpha}_1 \neq \mathbf{w}_0$, we have that $\mathbf{w}_1 = \mathbf{w}_0 - \boldsymbol{\alpha}_1 \neq \mathbf{0}$. This contradiction implies that $\theta_2 \leq \theta_1$.

We now examine the solution vectors to the dual problems of $\mathcal{P}.1$, $\boldsymbol{\alpha}_1$ and $\boldsymbol{\alpha}_2$. We know that $\|\boldsymbol{\alpha}_1\|_1 = \lambda_1$ so that $\|\mathbf{w}_0 - \boldsymbol{\alpha}_1\|_1 > \lambda_2$ and hence $\boldsymbol{\alpha}_2$ is at the boundary $\|\boldsymbol{\alpha}_2\|_1 = \lambda_2$ (see again [6]). Furthermore, the sum of these vectors is

$$\boldsymbol{\alpha}_1 + \boldsymbol{\alpha}_2 = [\mathbf{w}_0 - \theta_1]_+ + [\mathbf{w}_0 - [\mathbf{w}_0 - \theta_1]_+ - \theta_2]_+ . \quad (27)$$

Let v denote a component of \mathbf{w}_0 greater than θ_1 . For any such component the right hand side of Eq. (27) amounts to

$$[v - (v - \theta_1) - \theta_2]_+ + [v - \theta_1]_+ = [\theta_1 - \theta_2]_+ + v - \theta_1 = v - \theta_1 = [v - \theta_1]_+ ,$$

where we used the fact that $\theta_2 \leq \theta_1$ to eliminate the term $[\theta_1 - \theta_2]_+$. Next, let u denote a component of \mathbf{w}_0 smaller than θ_1 . In this case, the right hand side of Eq. (27) amounts to $[u - 0 - \theta_2]_+ + 0 =$

$[u - \theta_2]_+$. Recapping, the end result is that the vector sum $\alpha_1 + \alpha_2$ equals $[w_0 - \theta_2]_+$. Moreover, α_1 and α_2 are in \mathbb{R}_+^n as we assumed that $w_0 \succeq \mathbf{0}$, and thus

$$\| [w_0 - \theta_2]_+ \|_1 = \|\alpha_1 + \alpha_2\|_1 = \lambda_1 + \lambda_2. \quad (28)$$

We now show that $\mathcal{P}.2$ has the same dual solution as the sequential updates above. The dual of $\mathcal{P}.2$ is

$$\min_{\alpha} \frac{1}{2} \|\alpha - w_0\|_2^2 \quad \text{s.t.} \quad \|\alpha\|_1 \leq \lambda_1 + \lambda_2.$$

Denoting by α_0 the solution of the above dual problem, we have $w^* = w_0 - \alpha_0$ and $\alpha_0 = [w_0 - \theta]_+$ for some $\theta \in \mathbb{R}_+$. Examining the norm of α_0 we obtain that

$$\|\alpha_0\|_1 = \|[w_0 - \theta]_+\|_1 = \lambda_1 + \lambda_2 \quad (29)$$

because we assumed that $\|w_0\|_1 > \lambda_1 + \lambda_2$. We can view the terms $\|[w_0 - \theta_2]_+\|_1$ from Eq. (28) and $\|[w_0 - \theta]_+\|_1$ from Eq. (29) as functions of θ_2 and θ , respectively. The functions are strictly decreasing functions of θ and θ_2 over the interval $[0, \|w_0\|_\infty]$. Therefore, they are invertible for $0 < \lambda_1 + \lambda_2 < \|w_0\|_1$. Since $\|[w_0 - \theta]_+\|_1 = \|[w_0 - \theta_2]_+\|_1$, we must have $\theta_2 = \theta$. Recall that the solution of Eq. (26) is $w^* = \min\{w_0, \theta\}$, and the solution of the sequential update induced by Eq. (24) and Eq. (25) is $\min\{\min\{w_0, \theta_1\}, \theta_2\} = \min\{w_0, \theta_2\}$. The programs $\mathcal{P}.1$ and $\mathcal{P}.2$ therefore result in the same vector $\min\{w_0, \theta_2\} = \min\{w_0, \theta\}$ and their induced updates are equivalent.

We now examine the case when $\|w_0\|_1 \leq \lambda_1 + \lambda_2$. If the 1-norm of w_0 is also smaller than λ_1 , $\|w_0\|_1 \leq \lambda_1$, then the dual solution for the first step of $\mathcal{P}.1$ is $\alpha_1 = w_0$, which makes $w_1 = w_0 - \alpha_1 = \mathbf{0}$ and hence $w_2 = \mathbf{0}$. The dual solution for the combined problem is clearly $\alpha_0 = w_0$; again, $w^* = w_0 - \alpha_0 = \mathbf{0}$. We are thus left with the case $\lambda_1 < \|w_0\|_1 \leq \lambda_1 + \lambda_2$. We straightforwardly get that the solution to Eq. (26) is $w^* = \mathbf{0}$. We now prove that the iterated solution obtained by $\mathcal{P}.1$ results in the zero vector as well. First, consider the dual solution α_1 , which is the minimizer of $\|\alpha - w_0\|^2$ subject to $\|\alpha\|_1 \leq \lambda_1$. Since $\alpha_1 = [w_0 - \theta_1]_+$ for some $\theta_1 \geq 0$, we know that each component of α_1 is between zero and its corresponding component in w_0 , therefore, $\|w_0 - \alpha_1\|_1 = \|w_0\|_1 - \|\alpha_1\|_1 = \|w_0\|_1 - \lambda_1 \leq \lambda_2$. The dual of the second step of $\mathcal{P}.1$ distills to the minimization $\frac{1}{2} \|\alpha - (w_0 - \alpha_1)\|^2$ subject to $\|\alpha\|_1 \leq \lambda_2$. Since we showed that $\|w_0 - \alpha_1\|_1 \leq \lambda_2$, we get $\alpha_2 = w_0 - \alpha_1$. This means that $\theta_2 = 0$. Recall that the solution of $\mathcal{P}.1$ is $\min\{w_0, \theta_2\}$, which amounts to the zero vector when $\theta_2 = 0$. We have thus showed that both optimization problems result in the zero vector. This proves the equivalence of $\mathcal{P}.1$ and $\mathcal{P}.2$ for $q = \infty$. \square