

Modified Fejér Sequences and Applications*

Junhong Lin[†], Lorenzo Rosasco^{†,◦}, Silvia Villa[†], and Ding-Xuan Zhou^{*}

[†]*LCSL, Istituto Italiano di Tecnologia and Massachusetts Institute of Technology, Cambridge, MA 02139, USA*

[◦]*DIBRIS, Università degli Studi di Genova, Genova 16146, Italy*

^{*}*Department of Mathematics, City University of Hong Kong, Kowloon, Hong Kong, China*

August 22, 2016

Abstract

In this note, we propose and study the notion of modified Fejér sequences. Within a Hilbert space setting, we show that it provides a unifying framework to prove convergence rates for objective function values of several optimization algorithms. In particular, our results apply to forward-backward splitting algorithm, incremental subgradient proximal algorithm, and the Douglas-Rachford splitting method including and generalizing known results.

1 Introduction

The notion of Fejér monotonicity captures essential properties of the iterates generated by a wide range of optimization methods and provides a common

*This material is based upon work supported by the Center for Brains, Minds and Machines (CBMM), funded by NSF STC award CCF-1231216. The work by D. X. Zhou described in this paper is supported by a grant from the NSFC/RGC Joint Research Scheme [RGC Project No. N_CityU120/14 and NSFC Project No. 11461161006]. L. R. acknowledges the financial support of the Italian Ministry of Education, University and Research FIRB project RBFR12M3AC. S. V. is member of the Gruppo Nazionale per l'Analisi Matematica, la Probabilità e le loro Applicazioni (GNAMPA) of the Istituto Nazionale di Alta Matematica (INdAM).

framework to analyze their convergence [8]. Quasi-Fejér monotonicity is a relaxation of the above notion that allows for an additional error term [16, 10]. In this paper, we propose and study a novel, related notion to analyze the convergence of the objective function values, in addition to that of the iterates. More precisely, we modify the notion of quasi-Fejér monotonicity, by adding a term involving the objective function and say that a sequence satisfying the new requirement is modified Fejér monotone (modified Fejér for short). In this paper, we show the usefulness of this new notion of monotonicity by deriving convergence rates for several optimization algorithms in a unified way. Based on this approach, we not only recover known results, such as the sublinear convergence rate for the proximal forward-backward splitting algorithm, but also derive new results. Interestingly, our results show that for projected subgradient, incremental subgradient proximal, and Douglas-Rachford algorithm, considering the last iterate leads to essentially the same convergence rate as considering the best iterate selection rule [27, 26], or ergodic means [5, 28], as typically done.

2 Modified Fejér Sequences

Throughout this paper, we assume that $f : \mathcal{H} \rightarrow]-\infty, \infty]$ is a proper function on a normed vector space \mathcal{H} . Assume that the set

$$\mathcal{X} = \{z \in \mathcal{H} \mid f(z) = \min_{x \in \mathcal{H}} f(x)\}$$

is nonempty. We are interested in solving the following optimization problem

$$f_* = \min_{x \in \mathcal{H}} f(x). \tag{1}$$

Given $x \in \mathcal{H}$ and a subset $S \subset \mathcal{H}$, $d(x, S)$ denotes the distance between x and S , i.e., $d(x, S) = \inf_{x' \in S} \|x - x'\|$. \mathbb{R}_+ is the set of all non-negative real numbers and $\mathbb{N}^* = \mathbb{N} \setminus \{0\}$. For any $S \subset \mathcal{H}$, we denote by $\mathbf{1}_{\{ \cdot \}}$ the characteristic function of S .

The following definition introduces the key notion we propose in this paper.

Definition 1. A sequence $\{x_t\}_{t \in \mathbb{N}} \subset \mathcal{H}$ is modified Fejér monotone with respect to the target function f and the sequence $\{(\eta_t, \xi_t)\}_{t \in \mathbb{N}}$ in \mathbb{R}_+^2 , if

$$(\forall x \in \text{dom}f) \quad \|x_{t+1} - x\|^2 \leq \|x_t - x\|^2 - \eta_t(f(x_{t+1}) - f(x)) + \xi_t. \quad (2)$$

Remark 1.

(i) Choosing $x \in \mathcal{X}$ in (2), we get

$$\eta_t f(x_{t+1}) \leq \xi_t + \eta_t f_* + \|x_t - x\|^2 < \infty.$$

This implies that $\{x_t\}_{t \in \mathbb{N}} \subset \text{dom}f$.

(ii) By setting $x = x_t$ in (2), a direct consequence is that, for all $t \in \mathbb{N}$,

$$\|x_{t+1} - x_t\|^2 \leq \xi_t. \quad (3)$$

(iii) All the subsequent results hold if condition (2) is replaced by the following weaker condition

$$(\forall x \in \mathcal{X} \cup \{x_t\}_{t \in \mathbb{N}}) \quad \|x_{t+1} - x\|^2 \leq \|x_t - x\|^2 - \eta_t(f(x_{t+1}) - f(x)) + \xi_t. \quad (4)$$

However, in the proposed applications, condition (2) is always satisfied for every $x \in \text{dom}f$.

- Here we highlight that the inequality (2) has been proved to be satisfied and implicitly used to derive convergence rate for many algorithms, considering the best iterate selection rule, e.g., [27, 26], or ergodic means [5, 28]. The main novelty of this paper is to show that considering the last iterate leads to essentially the same convergence rate as that for considering the best iterate selection rule, or ergodic means, see Theorem 2.

In the following remark we discuss the relation with classical Fejér sequences.

Remark 2 (Comparison with quasi-Fejér sequences).

If $\sum_{t \in \mathbb{N}} \xi_t < +\infty$, Definition 1 implies that the sequence $\{x_t\}_{t \in \mathbb{N}}$ is quasi-Fejér monotone with respect to \mathcal{X} [16, 10]. Indeed, since the solutions set $\mathcal{X} \in \text{dom} f$ and that $f(x_{t+1}) - f(x_*) \geq 0$ for every $x_* \in \mathcal{X}$, (2) implies

$$(\forall x \in \mathcal{X}) \quad \|x_{t+1} - x\|^2 \leq \|x_t - x\|^2 + \xi_t.$$

Note that, in the study of convergence properties of quasi-Fejér sequences corresponding to a minimization problem, the property is considered with respect to the set of solutions \mathcal{X} , while here we will consider modified Fejér monotonicity for a general constraint set or the entire space \mathcal{H} .

We next present two main results to show how modified Fejér sequences are useful to study the convergence of optimization algorithms. The first result shows that if a sequence is modified Fejér monotone, one can bound its corresponding excess function values in terms of $\{(\eta_t, \xi_t)\}_{t \in \mathbb{N}}$ explicitly.

Theorem 1. Let $\{x_t\}_{t \in \mathbb{N}} \subset \mathcal{H}$ be a modified Fejér sequence with respect to f and $\{(\eta_t, \xi_t)\}_{t \in \mathbb{N}}$ in \mathbb{R}_+^2 . Let $\{\eta_t\}_{t \in \mathbb{N}}$ be a non-increasing sequence. Let $T \in \mathbb{N}$, $T > 1$. Then

$$\eta_T(f(x_{T+1}) - f_*) \leq \frac{1}{T} d(x_1, \mathcal{X})^2 + \sum_{t=1}^{T-1} \frac{1}{T-t} \xi_t + \xi_T. \quad (5)$$

Proof. Let $\{u_j\}_{j \in \mathbb{N}}$ be a sequence in \mathbb{R} and let $k \in \{1, \dots, T-1\}$. We have

$$\begin{aligned} & \frac{1}{k} \sum_{j=T-k+1}^T u_j - \frac{1}{k+1} \sum_{j=T-k}^T u_j \\ &= \frac{1}{k(k+1)} \left\{ (k+1) \sum_{j=T-k+1}^T u_j - k \sum_{j=T-k}^T u_j \right\} \\ &= \frac{1}{k(k+1)} \sum_{j=T-k+1}^T (u_j - u_{T-k}). \end{aligned}$$

Summing over $k = 1, \dots, T-1$, and rearranging terms, we get

$$u_T = \frac{1}{T} \sum_{j=1}^T u_j + \sum_{k=1}^{T-1} \frac{1}{k(k+1)} \sum_{j=T-k+1}^T (u_j - u_{T-k}). \quad (6)$$

For any $x \in \text{dom}f$, choosing $(\forall t \in \mathbb{N}) u_t = \eta_t(f(x_{t+1}) - f(x))$ and rearranging terms, we have the following error decomposition [19]:

$$\begin{aligned} \eta_T(f(x_{T+1}) - f(x)) &= \frac{1}{T} \sum_{t=1}^T \eta_t(f(x_{t+1}) - f(x)) \\ &+ \sum_{k=1}^{T-1} \frac{1}{k(k+1)} \sum_{t=T-k+1}^T \eta_t(f(x_{t+1}) - f(x_{T-k+1})) \\ &+ \sum_{k=1}^{T-1} \frac{1}{k+1} \left[\frac{1}{k} \sum_{t=T-k+1}^T \eta_t - \eta_{T-k} \right] \{f(x_{T-k+1}) - f(x)\}. \end{aligned}$$

Let $x = x_* \in \mathcal{X}$. Since $\{\eta_t\}_{t \in \mathbb{N}}$ is a non-increasing sequence and $f(x_{T-k+1}) - f_* \geq 0$, the last term of the above inequality is non-positive. Thus, we derive that

$$\begin{aligned} \eta_T(f(x_{T+1}) - f_*) &\leq \frac{1}{T} \sum_{t=1}^T \eta_t(f(x_{t+1}) - f(x_*)) \\ &+ \sum_{k=1}^{T-1} \frac{1}{k(k+1)} \sum_{t=T-k+1}^T \eta_t(f(x_{t+1}) - f(x_{T-k+1})). \end{aligned} \quad (7)$$

For every $j \in \{1, \dots, T\}$, and for every $x \in \text{dom}f$, summing up (2) over $t = j, \dots, T$, we get

$$\sum_{t=j}^T \eta_t(f(x_{t+1}) - f(x)) \leq \|x_j - x\|^2 + \sum_{t=j}^T \xi_t. \quad (8)$$

The above inequality with $x = x_*$ and $j = 1$ implies

$$\frac{1}{T} \sum_{t=1}^T \eta_t(f(x_{t+1}) - f(x_*)) \leq \frac{1}{T} \|x_1 - x_*\|^2 + \frac{1}{T} \sum_{t=1}^T \xi_t. \quad (9)$$

Inequality (8) with $x = x_{T-k+1}$ and $j = T - k + 1$ yields

$$\begin{aligned} &\sum_{k=1}^{T-1} \frac{1}{k(k+1)} \sum_{t=T-k+1}^T \eta_t(f(x_{t+1}) - f(x_{T-k+1})) \\ &\leq \sum_{k=1}^{T-1} \frac{1}{k(k+1)} \sum_{t=T-k}^T \xi_t. \end{aligned} \quad (10)$$

Exchanging the order in the sum, we obtain

$$\begin{aligned}
\sum_{k=1}^{T-1} \frac{1}{k(k+1)} \sum_{t=T-k}^T \xi_t &= \sum_{t=1}^{T-1} \sum_{k=T-t}^{T-1} \frac{1}{k(k+1)} \xi_t + \sum_{k=1}^{T-1} \frac{1}{k(k+1)} \xi_T \\
&= \sum_{t=1}^{T-1} \left(\frac{1}{T-t} - \frac{1}{T} \right) \xi_t + \left(1 - \frac{1}{T} \right) \xi_T \\
&= \sum_{t=1}^{T-1} \frac{1}{T-t} \xi_t + \xi_T - \frac{1}{T} \sum_{t=1}^T \xi_t. \tag{11}
\end{aligned}$$

The result follows by plugging (9),(10), and (11) into (7). \square

In the special case when, for every $t \in \mathbb{N}$, $\xi_t = 0$, we derive the following result.

Corollary 1. *Let $\{x_t\}_{t \in \mathbb{N}} \subset \mathcal{H}$ be a modified Fejér sequence with respect to f and a sequence $\{(\eta_t, \xi_t)\}_{t \in \mathbb{N}}$ in \mathbb{R}_+^2 . Assume that $\xi_t = 0$ for every $t \in \mathbb{N}$, and $\{\eta_t\}_{t \in \mathbb{N}}$ is non-increasing. Then for any $T \in \mathbb{N}$ with $T > 1$,*

$$f(x_{T+1}) - f_* \leq \frac{1}{\eta_T T} d(x_1, \mathcal{X})^2.$$

The second main result shows how to derive explicit rates for the objective function values corresponding to a modified Fejér sequence with respect to polynomially decaying sequences $\{(\eta_t, \xi_t)\}_{t \in \mathbb{N}}$ in \mathbb{R}_+^2 . Interestingly, the following result (as well as the previous ones) does not require convexity of f .

Theorem 2. *Let $\{x_t\}_{t \in \mathbb{N}} \subset C$ be a modified Fejér sequence with respect to a target function f and $\{(\eta_t, \xi_t)\}_{t \in \mathbb{N}} \subset \mathbb{R}_+^2$. Let $\eta \in]0, +\infty[$, let $\theta_1 \in [0, 1[$, and set $\eta_t = \eta t^{-\theta_1}$. Let $(\theta_2, \xi) \in \mathbb{R}_+^2$ and suppose that $\xi_t \leq \xi t^{-\theta_2}$ for all $t \in \mathbb{N}$. Let $T \in \mathbb{N}$, $T \geq 3$. Then*

$$f(x_{T+1}) - f_* \leq \frac{d(x_1, \mathcal{X})^2}{\eta} T^{\theta_1-1} + \frac{\xi c_{\theta_2}}{\eta} (\log T)^{\mathbf{1}_{\{\theta_2 \leq 1\}}} T^{\theta_1 - \min\{\theta_2, 1\}}. \tag{12}$$

Here, c_{θ_2} is a positive constant depending only on θ_2 and is given by

$$c_{\theta_2} = \begin{cases} 5 + \frac{2}{1 - \theta_2} & \text{if } \theta_2 < 1, \\ 9 & \text{if } \theta_2 = 1, \\ \frac{2^{\theta_2} + 3\theta_2 - 1}{\theta_2 - 1} & \text{if } \theta_2 > 1. \end{cases} \quad (13)$$

To prove Theorem 2, we will use Theorem 1 as well as the following lemma.

Lemma 1. *Let $q \in \mathbb{R}_+$ and $T \in \mathbb{N}$, $T \geq 3$. Then*

$$\sum_{t=1}^{T-1} \frac{1}{T-t} t^{-q} \leq \begin{cases} (4 + 2/(1-q)) T^{-q} \log T, & \text{when } q < 1, \\ 8T^{-1} \log T, & \text{when } q = 1, \\ (2^q + 2q)/(q-1) T^{-1}, & \text{when } q > 1, \end{cases}$$

Proof. We split the sum into two parts

$$\begin{aligned} \sum_{t=1}^{T-1} \frac{1}{T-t} t^{-q} &= \sum_{T/2 \leq t \leq T-1} \frac{1}{T-t} t^{-q} + \sum_{1 \leq t < T/2} \frac{1}{T-t} t^{-q} \\ &\leq 2^q T^{-q} \sum_{T/2 \leq t \leq T-1} \frac{1}{T-t} + 2T^{-1} \sum_{1 \leq t < T/2} t^{-q} \\ &= 2^q T^{-q} \sum_{1 \leq t \leq T/2} t^{-1} + 2T^{-1} \sum_{1 \leq t < T/2} t^{-q}. \end{aligned}$$

Applying, for $T \geq 3$,

$$\sum_{t=1}^T t^{-\theta_2} \leq 1 + \int_1^T u^{-\theta_2} du \leq \begin{cases} T^{1-\theta_2}/(1-\theta_2), & \text{when } \theta_2 < 1, \\ 2 \log T, & \text{when } \theta_2 = 1, \\ \theta_2/(\theta_2 - 1), & \text{when } \theta_2 > 1, \end{cases}$$

we get

$$\sum_{t=1}^{T-1} \frac{1}{T-t} t^{-q} \leq 2^{q+1} T^{-q} \log T + \begin{cases} (2/(1-q)) T^{-q}, & \text{when } q < 1, \\ 4T^{-1} \log T, & \text{when } q = 1, \\ 2qT^{-1}/(q-1), & \text{when } q > 1, \end{cases}$$

which leads to the desired result by using $T^{-q+1} \log T \leq 1/(2(q-1))$ when $q > 1$. \square

Now, we are ready to prove Theorem 2.

Proof of Theorem 2. It follows from Theorem 1 that (5) holds. Substituting $\eta_t = \eta t^{-\theta_1}$, $\xi_t \leq \xi t^{-\theta_2}$,

$$\eta T^{-\theta_1}(f(x_{T+1}) - f_*) \leq \frac{1}{T}d(x_1, \mathcal{X})^2 + \xi \sum_{t=1}^{T-1} \frac{1}{T-t} t^{-\theta_2} + \xi T^{-\theta_2}.$$

Applying Lemma 1 to bound the term $\sum_{t=1}^{T-1} \frac{1}{T-t} t^{-\theta_2}$, and by a direct calculation, with c_{θ_2} given by (13),

$$\eta T^{-\theta_1}(f(x_{T+1}) - f_*) \leq \frac{1}{T}d(x_1, \mathcal{X})^2 + \xi c_{\theta_2} (\log T)^{\mathbf{1}_{\{\theta_2 \leq 1\}}} T^{-\min\{\theta_2, 1\}}.$$

The results follows dividing both sides by $\eta T^{-\theta_1}$. \square

Remark 3. *Following the proof of Theorem 2, we see that for a sequence $\{y_t\}_{t \in \mathbb{N}} \in C$ satisfying*

$$(\forall x \in \text{dom} f) \quad \|y_{t+1} - x\|^2 \leq \|y_t - x\|^2 - \eta_t(f(y_t) - f(y)) + \xi_t,$$

under the same assumptions of Theorem 2, there holds

$$f(y_T) - f_* \leq \frac{d(x_1, \mathcal{X})^2}{\eta} T^{\theta_1 - 1} + \frac{\xi c_{\theta_2}}{\eta} (\log T)^{\mathbf{1}_{\{\theta_2 \leq 1\}}} T^{\theta_1 - \min\{\theta_2, 1\}}.$$

3 Applications in Convex Optimization

In this section, we apply previous results to some convex optimization algorithms, including forward-backward splitting, projected subgradient, incremental proximal subgradient, and Douglas-Rachford splitting methods. Convergence rates for the objective function values are obtained by using Theorem 2. The key observation is that the sequences generated by these algorithms are modified Fejér monotone.

Throughout this section, we assume that \mathcal{H} is a Hilbert space, and $f : \mathcal{H} \rightarrow]-\infty, \infty]$ is a proper, lower semicontinuous convex function. Recall that the subdifferential of f at $x \in \mathcal{H}$ is

$$\partial f(x) = \{u \in \mathcal{H} : (\forall y \in \mathcal{H}) \ f(x) + \langle u, y - x \rangle \leq f(y)\}. \quad (14)$$

The elements of the subdifferential of f at x are called subgradients of f at x . More generally, for $\epsilon \in]0, +\infty[$, the ϵ -subdifferential of f at x is the set $\partial_\epsilon f(x)$ defined by

$$\partial_\epsilon f(x) = \{u \in \mathcal{H} : (\forall y \in \mathcal{H}) \ f(x) + \langle u, y - x \rangle - \epsilon \leq f(y)\}. \quad (15)$$

The proximity operator of f [21] is

$$\text{prox}_f(x) = \underset{y \in \mathcal{H}}{\text{argmin}} \left\{ f(y) + \frac{1}{2} \|y - x\|^2 \right\}. \quad (16)$$

3.1 Forward-Backward Splitting

In this subsection, we consider a forward-backward splitting algorithm for solving Problem (1), with objective function

$$f = l + r \quad (17)$$

where $r: \mathcal{H} \rightarrow]-\infty, \infty]$ and $l: \mathcal{H} \rightarrow \mathbb{R}$ are proper, lower semicontinuous, and convex. Since l is real-valued, we have $\text{dom } \partial l = \mathcal{H}$ [2, Proposition 16.14].

Algorithm 1. *Given $x_1 \in \mathcal{H}$, a sequence of stepsizes $\{\alpha_t\}_{t \in \mathbb{N}} \subset]0, +\infty[$, and a sequence $\{\epsilon_t\}_{t \in \mathbb{N}} \subset [0, +\infty[$ set, for every $t \in \mathbb{N}$,*

$$x_{t+1} = \text{prox}_{\alpha_t r}(x_t - \alpha_t g_t) \quad (18)$$

with $g_t \in \partial_{\epsilon_t} l(x_t)$.

The forward-backward splitting algorithm has been well studied [29, 7, 9, 6] and a review of this algorithm can be found in [11] under the assumption that l is differentiable with a Lipschitz continuous gradient. Convergence is proved using arguments based on Fejér monotonicity of the generated sequences [10]. Under the assumption that l is a differentiable function with Lipschitz continuous gradient, the algorithm exhibits a sublinear convergence rate $O(T^{-1})$ on the objective f [3]. If l is not smooth, the algorithm has been studied first in [25], and has a convergence rate $O(T^{-1/2})$, considering the best point selection rule [28]. Our objective here is to provide a convergence rate for the algorithm considering the last iterate, which shares the same rate

(up-to logarithmic factors) and to allow the use of ϵ -subgradients, instead of subgradients. Before stating our main results, we introduce the following novel lemma for the forward-backward splitting for a (possibly) non-smooth l . It recovers previous result (e.g. [3]) when l is smooth.

Lemma 2. *Let $\{x_t\}_{t \in \mathbb{N}^*}$ be the sequence generated by Algorithm 1. Then for all $t \in \mathbb{N}^*$, there holds*

$$2\alpha_t[f(x_{t+1}) - f(x)] \leq \|x_t - x\|^2 - \|x_{t+1} - x\|^2 - \|x_t - x_{t+1}\|^2 + 2\alpha_t[\langle x_{t+1} - x_t, g_{t+1} - g_t \rangle + \epsilon_{t+1} + \epsilon_t]. \quad (19)$$

Proof. Let $t \in \mathbb{N}^*$. By Fermat's rule (see e.g. [2, Theorem 16.2]),

$$0 \in x_{t+1} - x_t + \alpha_t g_t + \alpha_t \partial r(x_{t+1}).$$

Thus, there exists $q_{t+1} \in \partial r(x_{t+1})$, such that x_{t+1} in (18) can be written as

$$x_{t+1} = x_t - \alpha_t g_t - \alpha_t q_{t+1}. \quad (20)$$

Let $x \in \text{dom} f$. The convexity of r implies

$$r(x_{t+1}) - r(x) \leq \langle x_{t+1} - x, q_{t+1} \rangle.$$

Multiplying both sides by $2\alpha_t$, and combining with (20), we get

$$\begin{aligned} 2\alpha_t[r(x_{t+1}) - r(x)] &\leq 2\alpha_t \langle x_{t+1} - x, q_{t+1} \rangle \\ &= 2\langle x_{t+1} - x, x_t - x_{t+1} - \alpha_t g_t \rangle \\ &= 2\langle x_{t+1} - x, x_t - x_{t+1} \rangle + 2\alpha_t \langle x - x_{t+1}, g_t \rangle. \end{aligned}$$

A direct computation yields

$$\begin{aligned} 2\langle x_{t+1} - x, x_t - x_{t+1} \rangle &= 2\langle x_{t+1} - x, x_t - x \rangle - 2\|x_{t+1} - x\|^2 \\ &= \|x_t - x\|^2 - \|x_{t+1} - x\|^2 - \|x_t - x_{t+1}\|^2. \end{aligned} \quad (21)$$

Therefore,

$$\begin{aligned} &2\alpha_t[r(x_{t+1}) - r(x)] \\ &\leq \|x_t - x\|^2 - \|x_{t+1} - x\|^2 - \|x_t - x_{t+1}\|^2 + 2\alpha_t \langle x - x_{t+1}, g_t \rangle. \end{aligned}$$

Moreover, by (15), we have

$$\langle x - x_t, g_t \rangle \leq l(x) - l(x_t) + \epsilon_t,$$

and

$$\langle x_t - x_{t+1}, g_{t+1} \rangle \leq l(x_t) - l(x_{t+1}) + \epsilon_{t+1},$$

and thus

$$\begin{aligned} \langle x - x_{t+1}, g_t \rangle &= \langle x - x_t, g_t \rangle + \langle x_t - x_{t+1}, g_{t+1} \rangle + \langle x_t - x_{t+1}, g_t - g_{t+1} \rangle \\ &\leq l(x) - l(x_t) + \epsilon_t + l(x_t) - l(x_{t+1}) + \epsilon_{t+1} + \langle x_t - x_{t+1}, g_t - g_{t+1} \rangle \\ &= l(x) - l(x_{t+1}) + \langle x_{t+1} - x_t, g_{t+1} - g_t \rangle + \epsilon_{t+1} + \epsilon_t. \end{aligned}$$

Consequently, we get

$$\begin{aligned} 2\alpha_t[r(x_{t+1}) - r(x)] &\leq \|x_t - x\|^2 - \|x_{t+1} - x\|^2 - \|x_t - x_{t+1}\|^2 \\ &\quad + 2\alpha_t[l(x) - l(x_{t+1}) + \langle x_{t+1} - x_t, g_{t+1} - g_t \rangle + \epsilon_{t+1} + \epsilon_t]. \end{aligned}$$

Rearranging terms and recalling that $f = l + r$, the desired result thus follows. \square

Theorem 3. *Let $\alpha \in]0, +\infty[$, let $\theta \in [0, 1[$, and let, for every $t \in \mathbb{N}^*$, $\alpha_t = \alpha t^{-\theta}$. Let $\epsilon \in]0, +\infty[$, $\{\epsilon_t\}_{t \in \mathbb{N}^*} \subset [0, +\infty[$, and assume that $\epsilon_t \leq \epsilon \alpha_t$. Let $\{x_t\}_{t \in \mathbb{N}^*}$ be the sequence generated by Algorithm 1. Let $T \in \mathbb{N}$ with $T > 3$. Assume that there exists $B \in]0, +\infty[$ such that*

$$(\forall 1 \leq t \leq T) \quad \|g_t\| \leq B, \tag{22}$$

and let c be defined as in (13). Then

$$f(x_{T+1}) - f_* \leq \frac{d(x_1, \mathcal{X})^2}{2\alpha} T^{\theta-1} + 2\alpha(B^2 + \epsilon)c_{2\theta}(\log T)^{\mathbf{1}_{\{2\theta \leq 1\}}} T^{-\min\{\theta, 1-\theta\}}.$$

Proof. Let $t \in \mathbb{N}^*$. By (19) and Cauchy-Schwartz inequality

$$\begin{aligned} &\|x_{t+1} - x\|^2 - \|x_t - x\|^2 \\ &\leq -\|x_t - x_{t+1}\|^2 + 2\alpha_t[\langle x_{t+1} - x_t, g_{t+1} - g_t \rangle + \epsilon_{t+1} + \epsilon_t] - 2\alpha_t[f(x_{t+1}) - f(x)] \\ &\leq -\|x_t - x_{t+1}\|^2 + 2\alpha_t[\|x_{t+1} - x_t\|\|g_{t+1} - g_t\| + \epsilon_{t+1} + \epsilon_t] - 2\alpha_t[f(x_{t+1}) - f(x)] \\ &\leq \alpha_t^2\|g_{t+1} - g_t\|^2 + 2\alpha_t[\epsilon_{t+1} + \epsilon_t] - 2\alpha_t[f(x_{t+1}) - f(x)]. \end{aligned}$$

Using the assumptions $\|g_t\| \leq B$ and $\epsilon_t \leq \epsilon\alpha_t$,

$$\begin{aligned} & \|x_{t+1} - x\|^2 - \|x_t - x\|^2 \\ & \leq 4B^2\alpha_t^2 + 2\epsilon\alpha_t[\alpha_t + \alpha_{t+1}] - 2\alpha_t[f(x_{t+1}) - f(x)] \\ & \leq 4(B^2 + \epsilon)\alpha_t^2 - 2\alpha_t[f(x_{t+1}) - f(x)]. \end{aligned}$$

Thus, $\{x_t\}_{t \in \mathbb{N}^*}$ is a modified Fejér sequence with respect to the target function f and $\{(2\alpha_t, 4(B^2 + \epsilon)\alpha_t^2)\}_{t \in \mathbb{N}^*}$. The statement follows from Theorem 2, applied with $\theta_1 = \theta$, $\theta_2 = 2\theta$, $\eta = 2\alpha$ and $\xi = 4(B^2 + \epsilon)\alpha^2$. \square

The following remark collects some comments on the previous result.

Remark 4.

1. Setting $\theta = 1/2$, we get a convergence rate $O(T^{-1/2} \log T)$ for forward-backward algorithm with nonsummable diminishing stepsizes, considering the last iterate.
2. In Theorem 3, the assumption on bounded approximate subgradients, which implies Lipschitz continuity of l , is satisfied for some practical optimization problems. For example, when r is the indicator function of a closed, bounded, and convex set $D \subset \mathbb{R}^N$, it follows that $\{x_t\}_{t \in \mathbb{N}}$ is bounded, which implies $\{g_t\}_{t \in \mathbb{N}}$ is bounded as well [1]. For general cases, similar results may be obtained by imposing a growth condition on ∂f , using a similar approach to that in [19] to bound the sequence of subgradients.
3. Theorem 3 improves [12, Corollary 2.4] in two aspects. First, the assumption (22) is weaker than the assumption, i.e., $\|g_t + u_t\| \leq B$ for some $u_t \in \partial r(x_t)$, in [12]. Second, [12] shows convergence rate only for the best point, i.e, the one with smallest function value:

$$(\forall T \in \mathbb{N}^*) \quad b_T = \underset{1 \leq t \leq T}{\operatorname{argmin}} f(x_t). \quad (23)$$

where our result holds for any last iterate.

If the function l in (17) is differentiable, with a Lipschitz differentiable gradient, we recover the well-known $O(1/T)$ convergence rate for the objective function values.

Proposition 1. [3, Theorem 3.1] Let $\beta \in [0, +\infty[$ and assume that ∇l is β -Lipschitz continuous. Consider Algorithm 1 with $\epsilon = 0$ and $\alpha_t = 1/\beta$ for all $t \in \mathbb{N}^*$. Then, for every $T \in \mathbb{N}$, $T > 1$

$$f(x_{T+1}) - f_* \leq \frac{\beta d(x_1, \mathcal{X})^2}{2T} \quad (24)$$

Proof. Following from (19) and that ∇l is β -Lipschitz continuous, with $\epsilon_t = 0$,

$$2\alpha_t[f(x_{t+1}) - f(x)] \leq \|x_t - x\|^2 - \|x_{t+1} - x\|^2 - \|x_t - x_{t+1}\|^2 + 2\alpha_t\beta\|x_{t+1} - x_t\|^2,$$

which leads to [3, Equation 3.6]

$$(\forall t \in \mathbb{N}^*) \quad \frac{2}{\beta}(f(x_{t+1}) - f_*) \leq \|x_{t+1} - x_*\|^2 - \|x_t - x_*\|^2. \quad (25)$$

Thus, $\{x_t\}_{t \in \mathbb{N}^*}$ is a modified Fejér sequence with respect to the target function f and the sequence $\{(\eta_t, \xi_t)\}_{t \in \mathbb{N}^*}$ with $(\forall t \in \mathbb{N}) \eta_t = 2/\beta$ and $\xi_t = 0$. The statement follows from Corollary 1. \square

3.2 Projected Approximate Subgradient Method

Let D be a convex and closed subset of \mathcal{H} , and let ι_D be the indicator function of D . In this subsection, we consider Problem (1) with objective function given by

$$f = l + \iota_D \quad (26)$$

where $l: \mathcal{H} \rightarrow \mathbb{R}$ is proper, lower semicontinuous, and convex. It is clear that (26) is a special case of (17) corresponding to a given choice of r . The forward-backward algorithm in this case reduces to the following projected subgradient method (see e.g. [27, 26, 5] and references therein), which allows to use ϵ -subgradients, see [1, 8].

Algorithm 2. Given $x_1 \in \mathcal{H}$, a sequence of stepsizes $\{\alpha_t\}_{t \in \mathbb{N}} \subset]0, +\infty[$, and a sequence $\{\epsilon_t\}_{t \in \mathbb{N}} \subset [0, +\infty[$, set, for every $t \in \mathbb{N}$,

$$x_{t+1} = P_D(x_t - \alpha_t g_t) \quad (27)$$

with $g_t \in \partial_{\epsilon_t} l(x_t)$.

The algorithm has been studied using different rules for choosing the stepsizes. Here, as a corollary of Theorem 3, we derive the convergence rate for the objective function values, for a nonsummable diminishing stepsize.

Theorem 4. *For some $\alpha_1 > 0$, $\epsilon \geq 0$ and $\theta \in [0, 1)$, let $\alpha_t = \eta_1 t^{-\theta}$ and $\epsilon_t \leq \epsilon \alpha_t$ for all $t \in \mathbb{N}^*$. Let $\{x_t\}_{t \in \mathbb{N}}$ be a sequence generated by Algorithm 2. Assume that for all $t \in \mathbb{N}^*$, $\|g_t\| \leq B$. Then, for every $T \in \mathbb{N}$, $T > 3$*

$$f(x_{T+1}) - f^* \leq \frac{d(x_1, \mathcal{X})^2}{2\alpha_1} T^{\theta-1} + \alpha_1 (B^2 + 2\epsilon) \tilde{c}_{2\theta} (\log T)^{\mathbf{1}_{\{2\theta \leq 1\}}} T^{-\min(\theta, 1-\theta)}$$

Choosing $\theta = 1/2$, we get a convergence rate of order $O(T^{-1/2} \log T)$ for projected approximate subgradient methods with nonsummable diminishing stepsizes, which is optimal up to a log factor without any further assumption on f [13, 24]. Since the subgradient method is not a descent method, a common approach keeps track of the best point found so far, i.e., (23). Projected subgradient method with diminishing stepsizes of the form $\{\alpha t^{-\theta}\}_t$, with $\theta \in]0, 1]$, satisfies $b_T - f_* = O(T^{-1/2})$. Our result shows that considering the last iterate for projected approximate subgradient method essentially leads to the same convergence rate, up to a logarithmic factor, as the one corresponding to the best iterate, even in the cases that the function value may not decrease at each iteration. To the best of our knowledge, our result is the first of this kind, without any assumption on strong convexity of f , or on a conditioning number with respect to subgradients (as in [17] using stepsizes $\{\gamma_t/\|g_t\|\}_t$). Note that, using nonsummable diminishing stepsizes, convergence rate $O(T^{-1/2})$ was shown, but only for a subsequence of $\{x_t\}_{t \in \mathbb{N}^*}$ [1]. Finally, let us mention that using properties of quasi-Fejér sequences, convergence properties were proved in [8].

3.3 Incremental Subgradient Proximal Algorithm

In this subsection, we consider an incremental subgradient proximal algorithm [4, 22] for solving (1), with objective function f given by, for some $m \in \mathbb{N}^*$,

$$\sum_{i=1}^m (l_i + r_i),$$

where for each i , both $l_i : \mathcal{H} \rightarrow \mathbb{R}$ and $r_i : \mathcal{H} \rightarrow]-\infty, +\infty]$ are convex, proper, and lower semicontinuous. The algorithm is similar to the proximal subgradient method, the main difference being that at each iteration, x_t is updated incrementally, through a sequence of m steps.

Algorithm 3. *Let $t \in \mathbb{N}^*$. Given $x_t \in \mathcal{H}$, an iteration of the incremental proximal subgradient algorithm generates x_{t+1} according to the recursion,*

$$x_{t+1} = \psi_t^m, \quad (28)$$

where ψ_t^m is obtained at the end of a cycle, namely as the last step of the recursion

$$\psi_t^0 = x_t, \quad \psi_t^i = \text{prox}_{\alpha_t r_i}(\psi_t^{i-1} - \alpha_t g_t^i), \quad \forall g_t^i \in \partial l_i(\psi_t^{i-1}), \quad i = 1, \dots, m \quad (29)$$

for a suitable sequence of stepsizes $\{\alpha_t\}_{t \in \mathbb{N}^*} \subset]0, +\infty[$.

Several versions of incremental subgradient proximal algorithms have been studied in [4], where convergence results for various stepsizes rules and both for stochastic or cyclic selection of the components are given. Concerning the function values, the results are stated in terms of the best iteration, i.e., (23). See also [23] for the study of the special case of incremental subgradient methods under different stepsizes rules. The paper [18] provides convergence results using approximate subgradients instead of gradients.

In this section, we derive a sublinear convergence rate for the incremental subgradient proximal algorithm in a straightforward way, relying on the properties of modified Fejér sequences assuming a boundedness assumption on the subdifferentials, already used in [23].

Theorem 5. *Let $\alpha \in]0, +\infty[$, let $\theta \in [0, 1[$, and let, for every $t \in \mathbb{N}^*$, $\alpha_t = \alpha t^{-\theta}$. Let $\{x_t\}_{t \in \mathbb{N}^*}$ be the sequence generated by Algorithm 3. Let $B \in]0, +\infty[$ be such that*

$$(\forall t \in \mathbb{N}^*)(\forall g \in \partial l_i(x_t) \cup \partial r_i(x_t)) \quad \|g\| \leq B,$$

and let c be defined as in (13). Then, for every $T \in \mathbb{N}^*$,

$$f(x_T) - f_* \leq \frac{d(x_1, \mathcal{X})^2}{2\alpha} T^{\theta-1} + \frac{\alpha(4m+5)mB^2}{2} c_{2\theta} (\log T)^{\mathbf{1}_{\{2\theta \leq 1\}}} T^{-\min\{\theta, 1-\theta\}}.$$

Proof. It was shown in [4, Proposition 3 (Equation 27)] that,

$$\|x_{t+1} - x\|^2 \leq \|x_t - x\|^2 - 2\alpha_t[f(x_t) - f(x)] + \alpha_t^2 (4m + 5) mB^2.$$

Thus, $\{x_t\}_{t \in \mathbb{N}^*}$ is a modified Fejér sequence with respect to the target function f , and $\{(2\alpha_t, \alpha_t^2 (4m + 5) mB^2)\}_{t \in \mathbb{N}^*}$. The proof is concluded by applying Remark 3 with $\theta_1 = \theta, \theta_2 = 2\theta, \eta = 2\alpha$ and $\xi = \alpha^2 (4m + 5) mB^2$. \square

Remark 5.

1. *An immediate consequence of Theorem 5, is that the choice $\theta = 1/2$ yields a convergence rate of order $O(T^{-1/2} \log T)$.*
2. *In contrast to [4, Proposition 5] where convergence rate of order $T^{-1/2}$ is derived for the best iterate (23) using a fixed stepsize, our result holds for any last iterate, considering both the fixed and diminishing stepsize setting.*

Similar to Theorem 5, we can derive convergence rates for the projected incremental subgradient method. Analogously to what we have done for the forward-backward algorithm in Section 3.1, Theorem 5 can be extended to analyze convergence of the approximate and incremental subgradient method in [18].

3.4 Douglas-Rachford splitting method

In this subsection, we consider Douglas-Rachford splitting algorithm for solving (1). Given $l: \mathcal{H} \rightarrow \mathbb{R}$ and $r: \mathcal{H} \rightarrow \mathbb{R}$ proper, convex, and lower semicontinuous functions, we assume that $f = l + r$ in (1).

Algorithm 4. *Let $\{\alpha_t\}_{t \in \mathbb{N}^*} \subset]0, +\infty[$. Let $t \in \mathbb{N}^*$. Given $x_t \in \mathcal{H}$, an iteration of Douglas-Rachford algorithm generates x_{t+1} according to*

$$\begin{cases} y_{t+1} = \text{prox}_{\alpha_t l}(x_t) \\ z_{t+1} = \text{prox}_{\alpha_t r}(2y_{t+1} - x_t), \\ x_{t+1} = x_t + z_{t+1} - y_{t+1}. \end{cases} \quad (30)$$

The algorithm has been introduced in [15] to minimize the sum of two convex functions, and then has been extended to monotone inclusions involving the sum of two nonlinear operators [20]. A review of this algorithm can be found in [11]. The convergence of the iterates is established using the theory of Fejér sequences [10]. Our objective here is to establish a new result, namely a convergence rate for the objective function values.

Theorem 6. *Let $\alpha \in]0, +\infty[$, and let $\theta \in [0, 1[$. For every $t \in \mathbb{N}^*$, let $\alpha_t = \alpha t^{-\theta}$. Let $\{(y_t, x_t, z_t)\}_{t \in \mathbb{N}^*}$ be the sequences generated by Algorithm 4. Assume that there exists $B \in]0, +\infty[$ such that*

$$(\forall t \in \mathbb{N}^*)(\forall g \in \partial l(y_t)) \quad \|g\| \leq B \quad \text{and} \quad (\exists g' \in \partial r(x_t)) \quad \|g'\| \leq B. \quad (31)$$

Let c be the function defined in (13). Then, for every $T \in \mathbb{N}$, $T > 3$,

$$f(x_{T+1}) - f_* \leq \frac{d(x_1, \mathcal{X})^2}{2\alpha} T^{\theta-1} + \frac{5\alpha B^2 c_{2\theta}}{2} (\log T)^{1_{\{2\theta \leq 1\}}} T^{-\min\{\theta, 1-\theta\}}.$$

Proof. Let $t \in \mathbb{N}^*$, set $v = (x_t - y_{t+1})/\alpha_t$ and $w = (2y_{t+1} - x_t - z_{t+1})/\alpha_t$. By Fermat's rule,

$$v \in \partial l(y_{t+1}) \quad \text{and} \quad w \in \partial r(z_{t+1}). \quad (32)$$

We can rewrite (30) as

$$\begin{cases} y_{t+1} = x_t - \alpha_t v, \\ z_{t+1} = (2y_{t+1} - x_t) - \alpha_t w, \\ x_{t+1} = x_t + z_{t+1} - y_{t+1}, \end{cases} \quad (33)$$

By (14), we have for any $x \in \text{dom} f$,

$$l(y_{t+1}) - l(x) \leq \langle y_{t+1} - x, v \rangle.$$

Multiplying both sides by $2\alpha_t$, and introducing with $v = (x_t - y_{t+1})/\alpha_t$,

$$2\alpha_t[l(y_{t+1}) - l(x)] \leq 2\alpha_t \langle y_{t+1} - x, v \rangle = 2 \langle y_{t+1} - x, x_t - y_{t+1} \rangle.$$

Similarly, we have

$$2\alpha_t[r(z_{t+1}) - r(x)] \leq 2\alpha_t \langle z_{t+1} - x, w \rangle = 2 \langle z_{t+1} - x, 2y_{t+1} - x_t - z_{t+1} \rangle.$$

Combining the above two estimates, we get

$$\begin{aligned} & 2\alpha_t[l(y_{t+1}) + r(z_{t+1}) - l(x) - r(x)] \\ & \leq 2\langle y_{t+1} - x, x_t - y_{t+1} \rangle + 2\langle z_{t+1} - x, 2y_{t+1} - x_t - z_{t+1} \rangle. \end{aligned}$$

Plugging with $z_{t+1} = x_{t+1} - x_t + y_{t+1}$ (implied by the third equality of (33)),

$$\begin{aligned} & 2\alpha_t[l(y_{t+1}) + r(z_{t+1}) - l(x) - r(x)] \\ & \leq 2\langle y_{t+1} - x, x_t - y_{t+1} \rangle + 2\langle x_{t+1} - x_t + y_{t+1} - x, 2y_{t+1} - x_{t+1} - y_{t+1} \rangle \\ & = 2\langle x_t - x_{t+1}, x_{t+1} - x \rangle. \end{aligned}$$

Introducing with (21),

$$\begin{aligned} & 2\alpha_t[l(y_{t+1}) + r(z_{t+1}) - l(x) - r(x)] \\ & \leq \|x_t - x\|^2 - \|x_{t+1} - x\|^2 - \|x_t - x_{t+1}\|^2. \end{aligned}$$

Adding both sides by $2\alpha_t[l(x_{t+1}) + r(x_{t+1}) - l(y_{t+1}) - r(z_{t+1})]$, and recalling that $f = l + r$,

$$\begin{aligned} 2\alpha_t[f(x_{t+1}) - f(x)] & \leq \|x_t - x\|^2 - \|x_{t+1} - x\|^2 - \|x_t - x_{t+1}\|^2 \\ & \quad + 2\alpha_t[l(x_{t+1}) + r(x_{t+1}) - l(y_{t+1}) - r(z_{t+1})]. \end{aligned} \quad (34)$$

Let $u \in \partial l(x_{t+1})$ and $s \in \partial r(x_{t+1})$ such that $\|u\| \leq B$ and $\|s\| \leq B$. Then using the convexity of l and r , and (33),

$$\begin{aligned} l(x_{t+1}) - l(y_{t+1}) & \leq \langle x_{t+1} - y_{t+1}, u \rangle \\ & = \langle x_{t+1} - x_t, u \rangle + \langle x_t - y_{t+1}, u \rangle \\ & = \langle x_{t+1} - x_t, u \rangle + \alpha_t \langle v, u \rangle \\ & \leq \|x_{t+1} - x_t\| \|u\| + \alpha_t \|v\| \|u\| \\ & \leq \|x_{t+1} - x_t\| B + \alpha_t B^2 \\ & \leq \|x_{t+1} - x_t\|^2 / (2\alpha_t) + B^2 \alpha_t / 2 + \alpha_t B^2, \end{aligned}$$

and

$$\begin{aligned} r(x_{t+1}) - r(z_{t+1}) & \leq \langle x_{t+1} - z_{t+1}, s \rangle \\ & = \alpha_t \langle v, s \rangle \leq \alpha_t \|v\| \|s\| \leq \alpha_t B^2. \end{aligned}$$

Introducing the last two estimates into (34), and by a direct calculation,

$$2\alpha_t[f(x_{t+1}) - f(x)] \leq \|x_t - x\|^2 - \|x_{t+1} - x\|^2 + 5B^2\alpha_t^2.$$

Thus, $\{x_t\}_{t \in \mathbb{N}^*}$ is a Super Quasi-Fejér sequence with respect to the target function f and $\{(2\alpha_t, 5\alpha_t^2 B^2)\}_{t \in \mathbb{N}^*}$. The statement follows from Theorem 2 with $\theta_1 = \theta$ and $\theta_2 = 2\theta$. \square

Again, choosing $\theta = 1/2$, we get a convergence rate $O(T^{-1/2} \log T)$ for the algorithm with nonsummable diminishing stepsizes. Nonergodic convergence rates for the objective function values corresponding to the Douglas-Rachford iteration can be derived by [14, Corollary 3.5], under the additional assumption that l is the indicator function of a linear subspace of \mathcal{H} .

Remark 6. *Theorem 6 still holds when the assumption (31) is replaced by*

$$(\forall t \in \mathbb{N}^*)(\forall g \in \partial r(y_t)) \quad \|g\| \leq B \quad \text{and} \quad (\exists g' \in \partial l(x_t)) \quad \|g'\| \leq B,$$

where the proof is essentially the same.

References

- [1] Alber, Y. I., Iusem, A. N., and Solodov, M. V.: On the projected subgradient method for nonsmooth convex optimization in a Hilbert space. *Math. Program.*, 81, 23-35 (1998).
- [2] Bauschke, H. H., and Combettes, P. L.: *Convex analysis and monotone operator theory in Hilbert spaces*, xvi+468. Springer, New York (2011).
- [3] Beck, A., and Teboulle, M.: A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sciences*, 2, 183-202 (2009).
- [4] Bertsekas, D. P.: Incremental proximal methods for large scale convex optimization. *Math. Program., Ser. B* 129, 163-195 (2009).
- [5] Boyd, S., Xiao, L., and Mutapcic, A.: Subgradient methods. http://web.stanford.edu/class/ee392o/subgrad_method.pdf. Accessed 14 October 2015.

- [6] Bredies, K. and Lorenz, D. A.: Linear convergence of iterative soft-thresholding. *J. Fourier Anal. Appl.*, 14, 813-837 (2008).
- [7] Chen, G. H., and Rockafellar, R. T.: Convergence rates in forward-backward splitting. *SIAM J. Optim.*, 7, 421-444 (1997).
- [8] Combettes, P. L.: Quasi-Fejérian analysis of some optimization algorithms. *Stud. Comput. Math.* 8, 115-152 (2001).
- [9] Combettes, P.L. and Wajs, V.R.: Signal recovery by proximal forward-backward splitting. *Multiscale Model. Simul.* 4, 1168-1200 (2005).
- [10] Combettes, P. L.: Fejér monotonicity in convex optimization. In: C. A. Floudas and P. M. Pardalos (eds.) *Encyclopedia of Optimization* (pp. 1016-1024). Springer, New York (2009).
- [11] Combettes, P. L., and Pesquet, J. C.: Proximal splitting methods in signal processing. In: *Fixed-point algorithms for inverse problems in science and engineering* (pp. 185-212). Springer, New York (2011).
- [12] Cruz, J.Y.B.: On proximal subgradient splitting method for minimizing the sum of two nonsmooth convex functions, *Set-Valued Var. Anal.*, 1-19 (2016).
- [13] Darzentas, J.: Problem complexity and method efficiency in optimization. *J. Oper. Res. Soc.*, 35(5), 455-455 (1984).
- [14] Davis, D.: Convergence rate analysis of the forward-Douglas-Rachford splitting scheme, *SIAM J. Optim* (to appear).
- [15] Douglas, J., and Rachford, H. H.: On the numerical solution of heat conduction problems in two and three space variables. *Trans. Amer. Math. Soc.*, 421-439 (1956).
- [16] Ermol'ev, Yu. M. and Tuniev, A. D.: Random Fejér and quasi-Fejér sequences, *Theory of Optimal Solutions – Akademiya Nauk Ukrainskoi SSR Kiev* 2, 76–83 (1968) ; translated in: *American Mathematical Society Selected Translations in Mathematical Statistics and Probability* 13 (1973) 143–148.

- [17] Goffin, J. L.: On convergence rates of subgradient optimization methods. *Math. Program.*, 13, 329-347 (1977).
- [18] Kiwiel, K. C.: Convergence of approximate and incremental subgradient methods for convex optimization. *SIAM J. Optim.*, 14, 807-840 (2004).
- [19] Lin J., Rosasco L., and Zhou D. X.: Iterative regularization for learning with convex loss functions. *arXiv:1503.08985* (2015)
- [20] Lions, P. L., and Mercier, B.: Splitting algorithms for the sum of two nonlinear operators. *SIAM J. Numer. Anal.*, 16, 964-979, (1979).
- [21] Moreau, J.J.: Fonctions convexes duales et points proximaux dans un espace hilbertien, *C. R. Acad. Sci. Paris* 255, 2897–2899, (1962).
- [22] Nedic, A., and Bertsekas, D. P.: Incremental subgradient methods for nondifferentiable optimization. *SIAM J. Optim.*, 12, 109-138, (2001).
- [23] Nedic, A., and Bertsekas, D.: Convergence rate of incremental subgradient algorithms. In *Stochastic optimization: algorithms and applications* (pp. 223-264). Springer, New York (2001).
- [24] Nesterov, Y.: *Introductory lectures on convex optimization*. Springer, New York (2004).
- [25] Passty, G. B.: Ergodic convergence to a zero of the sum of monotone operators in Hilbert space. *J. Math. Anal. Appl.*, 72, 383–390 (1979).
- [26] Polyak, B. T.: *Introduction to optimization*. Optimization Software, New York (1987).
- [27] Shor, N. Z.: *Minimization Methods for Non-Differentiable Functions*. Springer, New York, (1979).
- [28] Singer, Y., and Duchi, J. C.: Efficient learning using forward-backward splitting. In *Advances in Neural Information Processing Systems* (pp. 495-503) (2009).

- [29] Tseng, P.: Applications of a splitting algorithm to decomposition in convex programming and variational inequalities. *SIAM J. Control Optim.*, 29, 119-138 (1991).