

Chapter 6

The Primal-Dual Model of an Objective Function



In the previous chapters, we have proved that in the Black-Box framework the non-smooth optimization problems are much more difficult than the smooth ones. However, very often we know the explicit structure of the functional components. In this chapter we show how this knowledge can be used to accelerate the minimization methods and to extract a useful information about the dual counterpart of the problem. The main acceleration idea is based on the approximation of a nondifferentiable function by a differentiable one. We develop a technique for creating computable smoothed versions of non-differentiable functions and minimize them by Fast Gradient Methods. The number of iterations of the resulting methods is proportional to the square root of the number of iterations of the standard subgradient scheme. At the same time, the complexity of each iteration does not change. This technique can be used either in the primal form, or in the symmetric primal-dual form. We include in this chapter an example of application of this approach to the problem of Semidefinite Optimization. The chapter is concluded by analysis of performance of the Conditional Gradient method, which is based only on solving at each iteration an auxiliary problem of minimization of a linear function. We show that this method can also reconstruct the primal-dual solution of the problem. A similar idea is used in the second-order Trust Region Method with contraction, the first method of this type with provable global worst-case performance guarantees.

6.1 Smoothing for an Explicit Model of an Objective Function

(The minimax model of non-differentiable objective functions; The Fast Gradient Method for arbitrary norms and composite objective function; Application examples: minimax strategies for matrix games, the continuous location problem, variational inequalities with linear operator, minimization of piece-wise linear functions; Implementation issues.)

6.1.1 Smooth Approximations of Non-differentiable Functions

As we have seen in Chap. 3, subgradient methods solve the problem of Nonsmooth Convex Optimization in

$$O\left(\frac{1}{\epsilon^2}\right) \quad (6.1.1)$$

calls of the oracle, where ϵ is the desired absolute accuracy of finding the approximate solution in the function value. Moreover, we have already seen that the efficiency bound of the simplest Subgradient Method *cannot* be improved uniformly in the dimension of the space of variables (see Sect. 3.2). Of course, this statement is valid only for a Black-Box model of the objective function. However, the proof is constructive: it can be shown that the simplest problems like

$$\min_{x \in \mathbb{R}^n} \left\{ \gamma \max_{1 \leq i \leq k} x^{(i)} + \frac{\mu}{2} \|x\|^2 \right\}, \quad 1 \leq k \leq n,$$

where the norm is standard Euclidean, are difficult for all numerical schemes. The extremal simplicity of these functions possibly explains a common pessimistic belief that the actual worst-case complexity bound for finding an ϵ -approximation of the minimal value of a piece-wise linear function by gradient schemes is indeed given by (6.1.1).

In fact, this is not absolutely true. In practice, we almost never meet a pure Black-Box model. We always know something about the structure of the underlying objects (we have already discussed this in Sect. 5.1.1), and the proper use of this structure can and does help in constructing more efficient schemes.

In this section, we discuss one such possibility based on constructing a smooth approximation of a nonsmooth function. Let us look at the following situation. Consider a function f which is convex on \mathbb{E} . Assume that f satisfies the following growth condition:

$$f(x) \leq f(0) + L\|x\|, \quad \forall x \in \mathbb{R}^n, \quad (6.1.2)$$

where the Euclidean norm $\|x\| = \langle Bx, x \rangle^{1/2}$ is defined by a self-adjoint positive definite linear operator $B : \mathbb{E} \rightarrow \mathbb{E}^*$. Define the *Fenchel conjugate* of the function f as follows:

$$f_*(s) = \sup_{x \in \mathbb{E}} [\langle s, x \rangle - f(x)], \quad s \in \mathbb{E}^*. \quad (6.1.3)$$

Clearly, this function is closed and convex in view of Theorem 3.1.8. Its domain is not empty since by Theorem 3.1.20

$$\text{dom } f_* \supseteq \partial f(x), \quad \forall x \in \mathbb{E}.$$

At the same time, $\text{dom } f_*$ is bounded:

$$\|s\| \stackrel{(6.1.2)}{\leq} L \quad \forall s \in \text{dom } f_*. \quad (6.1.4)$$

Note that for all $x \in \mathbb{E}$ and $g \in \partial f(x)$, we have

$$f(x) + f_*(g) = \langle g, x \rangle. \quad (6.1.5)$$

Hence, for any $s \in \text{dom } f_*$ this implies that

$$f_*(s) \stackrel{(6.1.3)}{\geq} \langle s, x \rangle - f(x) \stackrel{(6.1.5)}{=} f_*(g) + \langle s - g, x \rangle.$$

In other words, if $g \in \partial f(x)$, then $x \in \partial f_*(g)$.

Let us prove the following relation (compare with general Theorem 3.1.16).

Lemma 6.1.1 *For all $x \in \mathbb{R}^n$, we have*

$$f(x) = \max_{s \in \text{dom } f_*} [\langle s, x \rangle - f_*(s)].$$

Proof Indeed, for any $s \in \text{dom } f_*$, we have $\langle s, x \rangle - f_*(s) \stackrel{(6.1.3)}{\leq} f(x)$, and, in view of (6.1.5), equality is achieved for $s \in \partial f(x)$. \square

Let us now look at the following smooth approximation of function f :

$$f_\mu(x) = \max_{s \in \text{dom } f_*} \left\{ \langle s, x \rangle - f_*(s) - \frac{1}{2}\mu(\|s\|^*)^2 \right\}, \quad (6.1.6)$$

where $\mu \geq 0$ is a smoothing parameter and the dual norm is defined as $\|s\|^* = \langle s, B^{-1}s \rangle^{1/2}$. In view of Lemma 6.1.1, we have

$$f(x) \geq f_\mu(x) \stackrel{(6.1.4)}{\geq} f(x) - \frac{1}{2}\mu L^2, \quad \forall x \in \mathbb{E}. \quad (6.1.7)$$

On the other hand, it appears that the function f_μ has a Lipschitz continuous gradient.

Lemma 6.1.2 *The function f_μ is differentiable on \mathbb{E} , and for any points x_1 and $x_2 \in \mathbb{E}$ we have*

$$\|\nabla f_\mu(x_1) - \nabla f_\mu(x_2)\|^* \leq \frac{1}{\mu} \|x_1 - x_2\|. \quad (6.1.8)$$

Proof Consider two points x_1 and x_2 from \mathbb{E} . Let s_i^* , $i = 1, 2$ be the optimal solutions of the corresponding optimization problems in (6.1.6). They are uniquely defined since the objective function in definition (6.1.6) is strongly concave.

Note that by Theorem 3.1.14, $s_i^* \in \partial f_\mu(x_i)$, $i = 1, 2$. On the other hand, by the first-order optimality condition of Theorem 3.1.20, there exist vectors $\tilde{x}_i \in \partial f_*(s_i^*)$ such that

$$\langle s - s_i^*, x_i - \tilde{x}_i - \mu B^{-1} s_i^* \rangle \leq 0, \quad \forall s \in \text{dom } f_*, \quad i = 1, 2.$$

Taking in this inequality $s = s_{3-i}^*$ and adding two copies of it with $i = 1, 2$, we get

$$\begin{aligned} \mu(\|s_1^* - s_2^*\|^*)^2 &\leq \langle s_1^* - s_2^*, x_1 - \tilde{x}_1 - (x_2 - \tilde{x}_2) \rangle \stackrel{(3.1.24)}{\leq} \langle s_1^* - s_2^*, x_1 - x_2 \rangle \\ &\leq \|s_1^* - s_2^*\|^* \cdot \|x_1 - x_2\|. \end{aligned}$$

Thus, $\|s_1^* - s_2^*\|^* \leq \frac{1}{\mu} \|x_1 - x_2\|$. Now, applying Lemma 3.1.10, we get $\nabla f_\mu(x_i) = s_i^*$, $i = 1, 2$. \square

Of course the smooth approximation (6.1.6) of the function f is not very practical since its internal minimization problem includes a potentially complicated function f_* . However, it already gives us some hints. Indeed, if we choose $\mu \approx \epsilon$, then the Lipschitz constant L_μ for the gradient of f_μ will be $O(\frac{1}{\epsilon})$. Therefore, Fast Gradient Methods (e.g. (2.2.20)) can find an ϵ -approximation of function f (this is f_μ) in $O\left(\sqrt{\frac{L_\mu}{\epsilon}}\right) \approx O(\frac{1}{\epsilon})$ calls of an oracle.

It remains to find a systematic and computationally inexpensive way of approximating the initial non-smooth objective function by a function with a Lipschitz continuous gradient. This can be done by exploiting a special max-representation of the objective function, which we introduce in Sect. 6.1.2.

For our goals, it is convenient to use the following notation. We often work with two finite-dimensional real vector spaces \mathbb{E}_1 and \mathbb{E}_2 . In these spaces, we use the corresponding scalar products and general norms

$$\langle s, x \rangle_{E_i}, \quad \|x\|_{\mathbb{E}_i}, \quad \|s\|_{\mathbb{E}_i}^*, \quad x \in \mathbb{E}_i, \quad s \in \mathbb{E}_i^*, \quad i = 1, 2,$$

which are not necessarily Euclidean. A *norm* of a linear operator $A : \mathbb{E}_1 \rightarrow \mathbb{E}_2^*$ is defined in the standard way:

$$\|A\|_{1,2} = \max_{x,u} \{\langle Ax, u \rangle_{\mathbb{E}_2} : \|x\|_{\mathbb{E}_1} = 1, \|u\|_{\mathbb{E}_2} = 1\}.$$

Clearly,

$$\begin{aligned} \|A\|_{1,2} &= \|A^*\|_{2,1} = \max_x \{\|Ax\|_{\mathbb{E}_2}^* : \|x\|_{\mathbb{E}_1} = 1\} \\ &= \max_u \{\|A^*u\|_{\mathbb{E}_1}^* : \|u\|_{\mathbb{E}_2} = 1\}. \end{aligned}$$

Hence, for any $x \in \mathbb{E}_1$ and $u \in \mathbb{E}_2$ we have

$$\|Ax\|_{\mathbb{E}_2}^* \leq \|A\|_{1,2} \cdot \|x\|_{\mathbb{E}_1}, \quad \|A^*u\|_{\mathbb{E}_1}^* \leq \|A\|_{1,2} \cdot \|u\|_{\mathbb{E}_2}. \quad (6.1.9)$$

6.1.2 The Minimax Model of an Objective Function

In this section, our main problem of interest is as follows:

$$\text{Find } f^* = \min_x \{f(x) : x \in Q_1\}, \quad (6.1.10)$$

where Q_1 is a bounded closed convex set in a finite-dimensional real vector space E_1 , and $f(\cdot)$ is a continuous convex function on Q_1 . We do not assume f to be differentiable.

Quite often, the *structure* of the objective function in (6.1.10) is given explicitly. Let us assume that this structure can be described by the following *model*:

$$f(x) = \hat{f}(x) + \max_u \{\langle Ax, u \rangle_{\mathbb{E}_2} - \hat{\phi}(u) : u \in Q_2\}, \quad (6.1.11)$$

where the function $\hat{f}(\cdot)$ is continuous and convex on Q_1 , Q_2 is a bounded closed convex set in a finite-dimensional real vector space E_2 , $\hat{\phi}(\cdot)$ is a continuous convex function on Q_2 , and the linear operator A maps E_1 to E_2^* . In this case, problem (6.1.10) can be written in an *adjoint* form. Indeed,

$$\begin{aligned} f^* &= \min_{x \in Q_1} \max_{u \in Q_2} \{\hat{f}(x) + \langle Ax, u \rangle_{\mathbb{E}_2} - \hat{\phi}(u)\} \\ &\stackrel{(1.3.6)}{\geq} \max_{u \in Q_2} \min_{x \in Q_1} \{\hat{f}(x) + \langle Ax, u \rangle_{\mathbb{E}_2} - \hat{\phi}(u)\}. \end{aligned}$$

Thus, the adjoint problem can be stated as follows:

$$f_* = \max_{u \in Q_2} \phi(u), \quad (6.1.12)$$

$$\phi(u) = -\hat{\phi}(u) + \min_{x \in Q_1} \{\langle Ax, u \rangle_{\mathbb{E}_2} + \hat{f}(x)\}.$$

However, the complexity of this problem is not completely identical to that of (6.1.10). Indeed, in the primal problem (6.1.10), we implicitly assume that the function $\hat{\phi}(\cdot)$ and set Q_2 are so simple that the solution of the optimization problem in (6.1.11) can be found in a closed form. This assumption may be not valid for the objects defining the function $\phi(\cdot)$.

Note that usually, for a convex function f , representation (6.1.11) is *not* uniquely defined. If we decide to use, for example, the Fenchel dual of f ,

$$\hat{\phi}(u) \equiv f_*(u) = \max_x \{ \langle u, x \rangle_{\mathbb{E}_1} - f(x) : x \in \mathbb{E}_1 \}, \quad Q_2 \equiv \mathbb{E}_2 = \mathbb{E}_1^*,$$

then we can take $\hat{f}(x) \equiv 0$, and A is equal to I_n , the identity operator. However, in this case the function $\hat{\phi}(\cdot)$ may be too complicated for our goals. Intuitively, it is clear that the bigger the dimension of the space \mathbb{E}_2 is, the simpler is the structure of the adjoint object defined by the function $\hat{\phi}(\cdot)$ and the set Q_2 . Let us demonstrate this with an example.

Example 6.1.1 Consider $f(x) = \max_{1 \leq j \leq m} |\langle a_j, x \rangle_{\mathbb{E}_1} - b^{(j)}|$. Let us choose $A = I_n$, $\mathbb{E}_2 = \mathbb{E}_1^* = \mathbb{R}^n$, and

$$\begin{aligned} \hat{\phi}(u) &= f_*(u) = \max_x \left\{ \langle u, x \rangle_{\mathbb{E}_1} - \max_{1 \leq j \leq m} |\langle a_j, x \rangle_{\mathbb{E}_1} - b^{(j)}| \right\} \\ &= \max_x \min_{s \in \mathbb{R}^m} \left\{ \langle u, x \rangle_{\mathbb{E}_1} - \sum_{j=1}^m s^{(j)} [\langle a_j, x \rangle_{\mathbb{E}_1} - b^{(j)}] : \sum_{j=1}^m |s^{(j)}| \leq 1 \right\} \\ &= \min_{s \in \mathbb{R}^m} \left\{ \langle b, s \rangle_{\mathbb{E}_2} : As = u, \sum_{j=1}^m |s^{(j)}| \leq 1 \right\}. \end{aligned}$$

It is clear that the structure of such a function can be very complicated.

Let us look at another possibility. Note that

$$\begin{aligned} f(x) &= \max_{1 \leq j \leq m} |\langle a_j, x \rangle_{\mathbb{E}_1} - b^{(j)}| \\ &= \max_{u \in \mathbb{R}^m} \left\{ \sum_{j=1}^m u^{(j)} [\langle a_j, x \rangle_{\mathbb{E}_1} - b^{(j)}] : \sum_{j=1}^m |u^{(j)}| \leq 1 \right\}. \end{aligned}$$

In this case $\mathbb{E}_2 = \mathbb{R}^m$, $\hat{\phi}(u) = \langle b, u \rangle_{\mathbb{E}_2}$ and $Q_2 = \left\{ u \in \mathbb{R}^m : \sum_{j=1}^m |u^{(j)}| \leq 1 \right\}$.

Finally, we can also represent $f(x)$ as follows:

$$f(x) = \max_{u=(u_1, u_2) \in \mathbb{R}_+^{2m}} \left\{ \sum_{j=1}^m (u_1^{(j)} - u_2^{(j)}) \cdot [\langle a_j, x \rangle_{\mathbb{E}_1} - b^{(j)}] : \sum_{j=1}^m (u_1^{(j)} + u_2^{(j)}) = 1 \right\}.$$

In this case $\mathbb{E}_2 = \mathbb{R}^{2m}$, $\hat{\phi}(u)$ is a linear function and Q_2 is a simplex. In Sect. 6.1.4.4 we will see that this representation is the easiest one. \square

Let us show that the knowledge of structure (6.1.11) can help in solving both problems (6.1.10) and (6.1.12). We are going to use this structure to construct a smooth approximation of the objective function in (6.1.10).

Consider a differentiable *prox-function* $d_2(\cdot)$ of the set Q_2 . This means that $d_2(\cdot)$ is strongly convex on Q_2 with convexity parameter one. Denote by

$$u_0 = \arg \min_u \{d_2(u) : u \in Q_2\}$$

its *prox-center*. Without loss of generality, we assume that $d_2(u_0) = 0$. Thus, for any $u \in Q_2$ we have

$$d_2(u) \stackrel{(2.2.40)}{\geq} \frac{1}{2} \|u - u_0\|_{\mathbb{E}_2}^2. \quad (6.1.13)$$

Let μ be a positive *smoothing* parameter. Consider the following function:

$$f_\mu(x) = \max_u \{\langle Ax, u \rangle_{\mathbb{E}_2} - \hat{\phi}(u) - \mu d_2(u) : u \in Q_2\}. \quad (6.1.14)$$

Denote by $u_\mu(x)$ the optimal solution of the above problem. Since the function $d_2(\cdot)$ is strongly convex, this solution is unique.

Theorem 6.1.1 *The function f_μ is well defined and continuously differentiable at any $x \in \mathbb{E}_1$. Moreover, this function is convex and its gradient*

$$\nabla f_\mu(x) = A^* u_\mu(x) \quad (6.1.15)$$

is Lipschitz continuous with constant

$$L_\mu = \frac{1}{\mu} \|A\|_{1,2}^2.$$

Proof Indeed the function $f_\mu(\cdot)$ is convex as a maximum of functions which are linear in x , and $A^* u_\mu(x) \in \partial f_\mu(x)$ (see Lemma 3.1.14). Let us prove now the existence and Lipschitz continuity of its gradient.

Consider two points x_1 and x_2 from \mathbb{E}_1 . From the first-order optimality conditions (3.1.56), we have

$$\langle Ax_i - g_i - \mu \nabla d_2(u_\mu(x_i)), u_\mu(x_{3-i}) - u_\mu(x_i) \rangle_{\mathbb{E}_2} \leq 0$$

for some $g_i \in \partial \hat{\phi}(u_\mu(x_i))$, $i = 1, 2$. Adding these inequalities, we get

$$\mu \|u_\mu(x_1) - u_\mu(x_2)\|_{\mathbb{E}_2}^2 \stackrel{(2.1.22)}{\leq} \mu \langle \nabla d_2(u_\mu(x_1)) - \nabla d_2(u_\mu(x_2)), u_\mu(x_1) - u_\mu(x_2) \rangle_{\mathbb{E}_2}$$

$$\leq \langle A(x_1 - x_2) - (g_1 - g_2), u_\mu(x_1) - u_\mu(x_2) \rangle_{\mathbb{E}_2}$$

$$\begin{aligned}
& \stackrel{(3.1.24)}{\leq} \langle A(x_1 - x_2), u_\mu(x_1) - u_\mu(x_2) \rangle_{\mathbb{E}_2} \\
& \leq \|A\|_{1,2} \cdot \|x_1 - x_2\|_{\mathbb{E}_1} \cdot \|u_\mu(x_1) - u_\mu(x_2)\|_{\mathbb{E}_2}.
\end{aligned}$$

Thus, in view of (6.1.9), we have

$$\begin{aligned}
\|A^*u_\mu(x_1) - A^*u_\mu(x_2)\|_{\mathbb{E}_1}^* & \leq \|A\|_{1,2} \cdot \|u_\mu(x_1) - u_\mu(x_2)\|_{\mathbb{E}_2}^2 \\
& \leq \frac{1}{\mu} \|A\|_{1,2}^2 \cdot \|x_1 - x_2\|_{\mathbb{E}_1}.
\end{aligned}$$

It remains to use Lemma 3.1.10. \square

Let $D_2 = \max_{u \in Q_2} d_2(u)$ and $f_0(x) = \max_{u \in Q_2} \{\langle Ax, u \rangle_{\mathbb{E}_2} - \hat{\phi}(u)\}$. Then, for any $x \in \mathbb{E}_1$ we have

$$f_0(x) \stackrel{(6.1.14)}{\geq} f_\mu(x) \stackrel{(6.1.14)}{\geq} f_0(x) - \mu D_2. \quad (6.1.16)$$

Thus, for $\mu > 0$ the function f_μ can be seen as a uniform μ -approximation of the objective function f_0 with Lipschitz constant for the gradient of the order $O(\frac{1}{\mu})$.

6.1.3 The Fast Gradient Method for Composite Minimization

Let $f(\cdot)$ be a convex differentiable function defined on a closed convex set $Q \subseteq E$. Assume that the gradient of this function is Lipschitz continuous:

$$\|\nabla f(x) - \nabla f(y)\|^* \leq L\|x - y\|, \quad \forall x, y \in Q.$$

Denote by $d(\cdot)$ a differentiable *prox-function* of the set Q . Assume that $d(\cdot)$ is strongly convex on Q with convexity parameter one. Let x_0 be the d -center of Q :

$$x_0 = \arg \min_{x \in Q} d(x).$$

Without loss of generality, assume that $d(x_0) = 0$. Thus, for any $x \in Q$ we have

$$d(x) \stackrel{(2.2.40)}{\geq} \frac{1}{2} \|x - x_0\|^2. \quad (6.1.17)$$

In this section, we present a fast gradient method for solving the following *composite* optimization problem:

$$\min_x \left\{ \tilde{f}(x) \stackrel{\text{def}}{=} f(x) + \Psi(x) : x \in Q \right\}, \quad (6.1.18)$$

where $\Psi(\cdot)$ is an arbitrary *simple* closed convex function defined on Q . Our main assumption is that the auxiliary minimization problem of the form

$$\min_{x \in Q} \{\langle s, x \rangle + \alpha d(x) + \beta \Psi(x)\}, \quad \alpha, \beta \geq 0,$$

is easily solvable. For simplicity, we assume that the constant $L > 0$ is known.

Method of Similar Triangles

0. Choose $x_0 \in Q$. Set $v_0 = x_0$ and $\phi_0(x) = Ld(x)$.

1. k th iteration ($k \geq 0$).

(a) Define $y_k = \frac{k}{k+2}x_k + \frac{2}{k+2}v_k$.

(b) Set $\phi_{k+1}(x) = \phi_k(x) + \frac{k+1}{2}[f(y_k) + \langle \nabla f(y_k), x - y_k \rangle + \Psi(x)]$.

(c) Compute $v_{k+1} = \min_{x \in Q} \phi_{k+1}(x)$.

(d) Define $x_{k+1} = \frac{k}{k+2}x_k + \frac{2}{k+2}v_{k+1}$.

(6.1.19)

In this scheme, we generate two sequences of feasible points $\{x_k\}_{k=0}^{\infty}$ and $\{y_k\}_{k=0}^{\infty}$, and a sequence of estimating functions $\{\phi_k(x)\}_{k=0}^{\infty}$. At each iteration of this method, all “events” happen in the two-dimensional plane defined by the triangle

$$\{x_k, v_k, v_{k+1}\}.$$

Note that this triangle is similar to the resulting triangle $\{x_k, y_k, x_{k+1}\}$, defining the new point of the sequence $\{x_k\}_{k=0}^{\infty}$, for which we are able to establish the rate of convergence.

Theorem 6.1.2 *Let the sequences $\{x_k\}_{k=0}^{\infty}$, $\{y_k\}_{k=0}^{\infty}$, and $\{v_k\}_{k=0}^{\infty}$ be generated by method (6.1.19). Then, for any $k \geq 0$ and $x \in Q$ we have*

$$\begin{aligned} & \frac{k(k+1)}{4} \tilde{f}(x_k) + \frac{L}{2} \|v_k - x\|^2 \\ & \leq \phi_k(x) = Ld(x) + \sum_{i=0}^{k-1} \frac{i+1}{2} [f(y_i) + \langle \nabla f(y_i), x - y_i \rangle] + \frac{k(k+1)}{4} \Psi(x). \end{aligned} \tag{6.1.20}$$

Therefore, for any $k \geq 1$, we get

$$\tilde{f}(x_k) - \tilde{f}(x^*) + \frac{2L}{k(k+1)} \|v_k - x^*\|^2 \leq \frac{4Ld(x^*)}{k(k+1)}, \quad (6.1.21)$$

where x^* is an optimal solution to problem (6.1.18).

Proof For $k \geq 0$, let

$$a_k = \frac{k}{2}, \quad A_k = \sum_{i=0}^k a_i = \frac{k(k+1)}{4}, \quad \tau_k = \frac{a_{k+1}}{A_{k+1}}.$$

Then the rules of method (6.1.19) can be written as follows:

$$y_k = (1 - \tau_k)x_k + \tau_k v_k, \quad x_{k+1} = (1 - \tau_k)x_k + \tau_k v_{k+1}. \quad (6.1.22)$$

Let us prove that

$$A_k \tilde{f}(x_k) \leq \phi_k^* \stackrel{\text{def}}{=} \min_{x \in Q} \phi_k = \phi_k(v_k), \quad k \geq 0. \quad (6.1.23)$$

Since $A_0 = 0$, this inequality is valid for $k = 0$. Assume that it is true for some $k \geq 0$. Since all functions ϕ_k are strongly convex with convexity parameter L , we have

$$\begin{aligned} \phi_{k+1}^* &= \phi_k(v_{k+1}) + a_{k+1}[f(y_k) + \langle \nabla f(y_k), v_{k+1} - y_k \rangle + \Psi(v_{k+1})] \\ &\stackrel{(2.2.40)}{\geq} \phi_k^* + \frac{L}{2} \|v_{k+1} - v_k\|^2 \\ &\quad + a_{k+1}[f(y_k) + \langle \nabla f(y_k), v_{k+1} - y_k \rangle + \Psi(v_{k+1})] \\ &\stackrel{(6.1.23)}{\geq} A_k[f(x_k) + \Psi(x_k)] + \frac{L}{2} \|v_{k+1} - v_k\|^2 \\ &\quad + a_{k+1}[f(y_k) + \langle \nabla f(y_k), v_{k+1} - y_k \rangle + \Psi(v_{k+1})] \\ &\stackrel{(2.1.2)}{\geq} A_{k+1}f(y_k) + \langle \nabla f(y_k), A_k(x_k - y_k) + a_{k+1}(v_{k+1} - y_k) \rangle \\ &\quad + \frac{L}{2} \|v_{k+1} - v_k\|^2 + A_k\Psi(x_k) + a_{k+1}\Psi(v_{k+1}). \end{aligned}$$

By the rules of the method, $A_k(x_k - y_k) + a_{k+1}(v_{k+1} - y_k) \stackrel{(6.1.22)}{=} a_{k+1}(v_{k+1} - v_k)$ and $A_k\Psi(x_k) + a_{k+1}\Psi(v_{k+1}) \geq A_{k+1}\Psi(x_{k+1})$. Therefore,

$$\begin{aligned} \phi_{k+1}^* &\geq A_{k+1}f(y_k) + a_{k+1}\langle \nabla f(y_k), v_{k+1} - v_k \rangle + \frac{L}{2}\|v_{k+1} - v_k\|^2 \\ &\quad + A_{k+1}\Psi(x_{k+1}) \\ &\stackrel{(6.1.22)}{=} A_{k+1}[f(y_k) + \langle \nabla f(y_k), x_{k+1} - y_k \rangle + \frac{LA_{k+1}}{2a_{k+1}^2}\|x_{k+1} - y_k\|^2 \\ &\quad + \Psi(x_{k+1})]. \end{aligned}$$

Since $\frac{A_{k+1}}{a_{k+1}^2} = \frac{(k+1)(k+2)}{4} \cdot \frac{4}{(k+1)^2} > 1$, we get $\phi_{k+1}^* \stackrel{(2.1.9)}{\geq} A_{k+1}f(x_{k+1})$. By strong convexity of the function ϕ_k , we have

$$\phi_k(x) \stackrel{(2.2.40)}{\geq} \phi_k^* + \frac{L}{2}\|x - v_k\|^2 \stackrel{(6.1.23)}{\geq} A_k\tilde{f}(x_k) + \frac{L}{2}\|x - v_k\|^2,$$

and this is inequality (6.1.20). Finally, inequality (6.1.21) follows from (6.1.20) in view of the convexity of the function f . \square

Remark 6.1.1 Note that method (6.1.19) generates bounded sequences of points. Indeed, by the rules of this method we have

$$x_k, y_k \in \text{Conv}\{v_0, \dots, v_k\}, \quad k \geq 0.$$

On the other hand, from inequality (6.1.21), it follows that

$$\|v_k - x^*\|^2 \leq 2d(x^*). \quad (6.1.24)$$

In the Euclidean case, $d(x) = \frac{1}{2}\|x - x_0\|^2$, and we get

$$\|v_k - x^*\| \leq \|x_0 - x^*\|, \quad k \geq 0. \quad (6.1.25)$$

6.1.4 Application Examples

Let us put the results of the previous sections together. Assume that the function $\hat{f}(\cdot)$ in (6.1.11) is differentiable and its gradient is Lipschitz-continuous with some constant $M \geq 0$. Then the smoothing technique as applied to problem (6.1.10) provides us with the following objective function:

$$\tilde{f}_\mu(x) = \hat{f}(x) + f_\mu(x) \rightarrow \min : x \in Q_1. \quad (6.1.26)$$

In view of Theorem 6.1.1, the gradient of this function is Lipschitz continuous with the constant

$$L_\mu = M + \frac{1}{\mu} \|A\|_{1,2}^2.$$

Let us choose some prox-function $d_1(\cdot)$ for the set Q_1 with convexity parameter equal to one. Recall that the set Q_1 is assumed to be bounded:

$$\max_{x \in Q_1} d_1(x) \leq D_1.$$

Theorem 6.1.3 *Let us apply method (6.1.19) to problem (6.1.26) with the following value of the smoothness parameter:*

$$\mu = \mu(N) = \frac{2\|A\|_{1,2}}{\sqrt{N(N+1)}} \cdot \sqrt{\frac{D_1}{D_2}}.$$

Then after N iterations we can generate approximate solutions to problems (6.1.10) and (6.1.12), namely,

$$\hat{x} = x_N \in Q_1, \quad \hat{u} = \sum_{i=0}^{N-1} \frac{2(i+1)}{(N+1)(N+2)} u_\mu(y_i) \in Q_2, \quad (6.1.27)$$

which satisfy the following inequality:

$$0 \leq f(\hat{x}) - \phi(\hat{u}) \leq \frac{4\|A\|_{1,2}}{\sqrt{N(N+1)}} \cdot \sqrt{D_1 D_2} + \frac{4MD_1}{N(N+1)}. \quad (6.1.28)$$

Thus, the complexity of finding an ϵ -solution to problems (6.1.10), (6.1.12) by the smoothing technique does not exceed

$$4\|A\|_{1,2} \sqrt{D_1 D_2} \cdot \frac{1}{\epsilon} + 2\sqrt{\frac{MD_1}{\epsilon}} \quad (6.1.29)$$

iterations of method (6.1.19).

Proof Let us fix an arbitrary $\mu > 0$. In view of Theorem 6.1.2, after N iterations of method (2.2.63) we can deliver a point $\hat{x} = x_N$ such that

$$\bar{f}_\mu(\hat{x}) \leq \frac{4L_\mu D_1}{N(N+1)} + \min_{x \in Q_1} \sum_{i=0}^{N-1} \frac{2(i+1)}{N(N+1)} [\bar{f}_\mu(y_i) + \langle \nabla \bar{f}_\mu(x_i), x - y_i \rangle_{\mathbb{E}_1}]. \quad (6.1.30)$$

Note that

$$\begin{aligned}
 f_\mu(y) &= \max_u \{ \langle Ay, u \rangle_{\mathbb{E}_2} - \hat{\phi}(u) - \mu d_2(u) : u \in \mathcal{Q}_2 \} \\
 &= \langle Ay, u_\mu(y) \rangle_{\mathbb{E}_2} - \hat{\phi}(u_\mu(y)) - \mu d_2(u_\mu(y)), \\
 \langle \nabla f_\mu(y), y \rangle_{\mathbb{E}_1} &= \langle A^* u_\mu(y), y \rangle_{\mathbb{E}_1}.
 \end{aligned}$$

Therefore, for $i = 0, \dots, N-1$ we have

$$f_\mu(y_i) - \langle \nabla f_\mu(y_i), y_i \rangle_{\mathbb{E}_1} = -\hat{\phi}(u_\mu(y_i)) - \mu d_2(u_\mu(y_i)). \quad (6.1.31)$$

Thus, in view of (6.1.15) and (6.1.31) we obtain

$$\begin{aligned}
 & \sum_{i=0}^{N-1} (i+1) [\bar{f}_\mu(y_i) + \langle \nabla \bar{f}_\mu(y_i), x - y_i \rangle_{\mathbb{E}_1}] \\
 & \stackrel{(2.1.2)}{\leq} \sum_{i=0}^{N-1} (i+1) [f_\mu(y_i) - \langle \nabla f_\mu(y_i), y_i \rangle_{\mathbb{E}_1}] + \frac{1}{2} N(N+1) (\hat{f}(x) + \langle A^* \hat{u}, x \rangle_{\mathbb{E}_1}) \\
 & \leq - \sum_{i=0}^{N-1} (i+1) \hat{\phi}(u_\mu(y_i)) + \frac{1}{2} N(N+1) (\hat{f}(x) + \langle A^* \hat{u}, x \rangle_{\mathbb{E}_1}) \\
 & \leq \frac{1}{2} N(N+1) [-\hat{\phi}(\hat{u}) + \hat{f}(x) + \langle Ax, \hat{u} \rangle_{\mathbb{E}_2}].
 \end{aligned}$$

Hence, using (6.1.30), (6.1.12) and (6.1.16), we get the following bound:

$$\frac{4L_\mu D_1}{N(N+1)} \geq \bar{f}_\mu(\hat{x}) - \phi(\hat{u}) \geq f(\hat{x}) - \phi(\hat{u}) - \mu D_2.$$

This is

$$0 \leq f(\hat{x}) - \phi(\hat{u}) \leq \mu D_2 + \frac{4\|A\|_{1,2}^2 D_1}{\mu N(N+1)} + \frac{4MD_1}{N(N+1)}. \quad (6.1.32)$$

Minimizing the right-hand side of this inequality in μ , we get inequality (6.1.28). \square

Note that the efficiency estimate (6.1.29) is much better than the standard bound $O\left(\frac{1}{\epsilon^2}\right)$. In accordance with the above theorem, for $M = 0$ the optimal dependence of the parameters μ , L_μ and N in ϵ is as follows:

$$\sqrt{N(N+1)} \geq 4\|A\|_{1,2} \sqrt{D_1 D_2} \cdot \frac{1}{\epsilon}, \quad \mu = \frac{\epsilon}{2D_2}, \quad L_\mu = D_2 \cdot \frac{\|A\|_{1,2}^2}{\epsilon}. \quad (6.1.33)$$

Remark 6.1.2 Inequality (6.1.28) shows that the pair of adjoint problems (6.1.10) and (6.1.12) has no *duality gap*:

$$f^* = f_*. \quad (6.1.34)$$

Let us now look at some examples.

6.1.4.1 Minimax Strategies for Matrix Games

Denote by Δ_n the standard simplex in \mathbb{R}^n :

$$\Delta_n = \left\{ x \in \mathbb{R}_+^n : \sum_{i=1}^n x^{(i)} = 1 \right\}.$$

Let $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $\mathbb{E}_1 = \mathbb{R}^n$, and $\mathbb{E}_2 = \mathbb{R}^m$. Consider the following saddle point problem:

$$\min_{x \in \Delta_n} \max_{u \in \Delta_m} \{ \langle Ax, u \rangle_{\mathbb{E}_2} + \langle c, x \rangle_{\mathbb{E}_1} + \langle b, u \rangle_{\mathbb{E}_2} \}. \quad (6.1.35)$$

From the viewpoint of players, this problem can be seen as a pair of non-smooth minimization problems:

$$\begin{aligned} \min_{x \in \Delta_n} f(x), \quad f(x) &= \langle c, x \rangle_{\mathbb{E}_1} + \max_{1 \leq j \leq m} [\langle a_j, x \rangle_{\mathbb{E}_1} + b^{(j)}], \\ \max_{u \in \Delta_m} \phi(u), \quad \phi(u) &= \langle b, u \rangle_{\mathbb{E}_2} + \min_{1 \leq i \leq n} [\langle \hat{a}_i, u \rangle_{\mathbb{E}_2} + c^{(i)}], \end{aligned} \quad (6.1.36)$$

where a_j are the rows and \hat{a}_i are the columns of matrix A . In order to solve this pair of problems using the smoothing approach, we need to find a reasonable prox-function for the simplex. Let us compare two possibilities.

1. Euclidean Distance Let us choose

$$\begin{aligned} \|x\|_{\mathbb{E}_1} &= \left[\sum_{i=1}^n (x^{(i)})^2 \right]^{1/2}, \quad d_1(x) = \frac{1}{2} \sum_{i=1}^n (x^{(i)} - \frac{1}{n})^2, \\ \|u\|_{\mathbb{E}_2} &= \left[\sum_{j=1}^m (u^{(j)})^2 \right]^{1/2}, \quad d_2(x) = \frac{1}{2} \sum_{j=1}^m (u^{(j)} - \frac{1}{m})^2. \end{aligned}$$

Then $D_1 = 1 - \frac{1}{n} < 1$, $D_2 = 1 - \frac{1}{m} < 1$ and

$$\|A\|_{1,2} = \max_u \{ \|Ax\|_2^* : \|x\|_{\mathbb{E}_1} = 1 \} = \lambda_{\max}^{1/2}(A^T A).$$

Thus, in our case the estimate (6.1.28) for the result (6.1.27) can be specified as follows:

$$0 \leq f(\hat{x}) - \phi(\hat{u}) \leq \frac{4\lambda_{\max}^{1/2}(A^T A)}{\sqrt{N(N+1)}}. \quad (6.1.37)$$

2. Entropy Distance

Let us choose

$$\|x\|_{\mathbb{E}_1} = \sum_{i=1}^n |x^{(i)}|, \quad d_1(x) = \ln n + \sum_{i=1}^n x^{(i)} \ln x^{(i)},$$

$$\|u\|_{\mathbb{E}_2} = \sum_{j=1}^m |u^{(j)}|, \quad d_2(u) = \ln m + \sum_{j=1}^m u^{(j)} \ln u^{(j)}.$$

Functions d_1 and d_2 are called the *entropy functions*.

Lemma 6.1.3 *The above prox-functions are strongly convex in an ℓ_1 -norm with convexity parameter one and $D_1 = \ln n$, $D_2 = \ln m$.*

Proof Note that the function d_1 is twice continuously differentiable in the interior of simplex Δ_n , and

$$\langle \nabla^2 d_1(x)h, h \rangle = \sum_{i=1}^n \frac{(h^{(i)})^2}{x^{(i)}}.$$

Thus, in view of Theorem 2.1.11 strong convexity of d_1 is a consequence of the following variant of Cauchy–Schwarz inequality,

$$\left(\sum_{i=1}^n |h^{(i)}| \right)^2 \leq \left(\sum_{i=1}^n x^{(i)} \right) \cdot \left(\sum_{i=1}^n \frac{(h^{(i)})^2}{x^{(i)}} \right),$$

which is valid for all positive vectors $x \in \mathbb{R}^n$. Since $d_1(\cdot)$ is a convex symmetric function of the arguments, its minimum is attained at the center of the simplex, the point $x_0 = \frac{1}{n}\bar{e}_n$. Clearly, $d_1(x_0) = 0$. On the other hand, its maximum is attained at one of the vertices of the simplex (see Corollary 3.1.2).

The reasoning for $d_2(\cdot)$ is similar. \square

Note also that now we get the following norm of the operator A :

$$\|A\|_{1,2} = \max_x \left\{ \max_{1 \leq j \leq m} |\langle a_j, x \rangle| : \|x\|_{\mathbb{E}_1} \leq 1 \right\} = \max_{i,j} |A^{(i,j)}|$$

(see Corollary 3.1.2). Thus, if we apply the entropy distance, the estimate (6.1.28) can be written as follows:

$$0 \leq f(\hat{x}) - \phi(\hat{u}) \leq \frac{4\sqrt{\ln n \ln m}}{\sqrt{N(N+1)}} \cdot \max_{i,j} |A^{(i,j)}|. \quad (6.1.38)$$

Note that typically the estimate (6.1.38) is much better than its Euclidean variant (6.1.37).

Let us write down explicitly the smooth approximation for the objective function in the first problem of (6.1.36) using the entropy distance. By definition,

$$\tilde{f}_\mu(x) = \langle c, x \rangle_{\mathbb{E}_1} + \max_{u \in \Delta_m} \left\{ \sum_{j=1}^m u^{(j)} [\langle a_j, x \rangle + b^{(j)}] - \mu \sum_{j=1}^m u^{(j)} \ln u^{(j)} - \mu \ln m \right\}.$$

Let us apply the following result.

Lemma 6.1.4 *The solution of the problem*

$$\text{Find } \phi_*(s) = \max_{u \in \Delta_m} \left\{ \sum_{j=1}^m u^{(j)} s^{(j)} - \mu \sum_{j=1}^m u^{(j)} \ln u^{(j)} \right\} \quad (6.1.39)$$

is given by the vector $u_\mu(s) \in \Delta_m$ with the following entries

$$u_\mu^{(j)}(s) = \frac{e^{s^{(j)}/\mu}}{\sum_{i=1}^m e^{s^{(i)}/\mu}}, \quad j = 1, \dots, m. \quad (6.1.40)$$

Therefore, $\phi_*(s) = \mu \ln \left(\sum_{i=1}^m e^{s^{(i)}/\mu} \right)$.

Proof Note that the gradient of the objective function in problem (6.1.39) goes to infinity as the argument approaches the boundary of the domain. Therefore, the first order necessary and sufficient optimality conditions for this problem are as follows (see (3.1.59)):

$$s^{(j)} - \mu(1 + \ln u^{(j)}) = \lambda, \quad j = 1, \dots, m,$$

$$\sum_{j=1}^m u^{(j)} = 1.$$

Clearly, they are satisfied by (6.1.40) with $\lambda = \mu \ln \left(\sum_{l=1}^m e^{s^{(l)}/\mu} \right) - \mu$. \square

Using the result of Lemma 6.1.4, we conclude that in our case the problem (6.1.26) is as follows:

$$\min_{x \in \Delta_n} \left\{ \tilde{f}_\mu(x) = \langle c, x \rangle_{\mathbb{E}_1} + \mu \ln \left(\frac{1}{m} \sum_{j=1}^m e^{[\langle a_j, x \rangle + b^{(j)}]/\mu} \right) \right\}.$$

Note that the complexity of the oracle for this problem is basically the same as that of the initial problem (6.1.36).

6.1.4.2 The Continuous Location Problem

Consider the following *location* problem. There are p cities with population m_j , which are located at points $c_j \in \mathbb{R}^n$, $j = 1, \dots, p$. We want to construct a service center at some position $x \in \mathbb{R}^n \equiv \mathbb{E}_1$, which minimizes the total social distance $f(x)$ to the center. On the other hand, this center must be constructed not too far from the origin.

Mathematically, the above problem can be posed as follows

$$\text{Find } f^* = \min_x \left\{ f(x) = \sum_{j=1}^p m_j \|x - c_j\|_{\mathbb{E}_1} : \|x\|_{\mathbb{E}_1} \leq \bar{r} \right\}. \quad (6.1.41)$$

In accordance to its interpretation, it is natural to choose

$$\|x\|_{\mathbb{E}_1} = \left[\sum_{i=1}^n (x^{(i)})^2 \right]^{1/2}, \quad d_1(x) = \frac{1}{2} \|x\|_{\mathbb{E}_1}^2.$$

Then $D_1 = \frac{1}{2} \bar{r}^2$.

Further, the structure of the adjoint space \mathbb{E}_2 is quite clear:

$$\mathbb{E}_2 = (\mathbb{E}_1^*)^p, \quad Q_2 = \left\{ u = (u_1, \dots, u_p) \in \mathbb{E}_2 : \|u_j\|_{\mathbb{E}_1}^* \leq 1, j = 1, \dots, p \right\}.$$

Let us choose

$$\|u\|_{\mathbb{E}_2} = \left[\sum_{j=1}^p m_j (\|u_j\|_{\mathbb{E}_1}^*)^2 \right]^{1/2}, \quad d_2(u) = \frac{1}{2} \|u\|_{\mathbb{E}_2}^2.$$

Then $D_2 = \frac{1}{2} P$ with $P \equiv \sum_{j=1}^p m_j$. Note that the value P may be interpreted as the total size of the population.

It remains to compute the norm of the operator A :

$$\begin{aligned} \|A\|_{1,2} &= \max_{x,u} \left\{ \sum_{j=1}^p m_j \langle u_j, x \rangle_{\mathbb{E}_1} : \sum_{j=1}^p m_j (\|u_j\|_{\mathbb{E}_1}^*)^2 = 1, \|x\|_{\mathbb{E}_1} = 1 \right\} \\ &= \max_{r_j} \left\{ \sum_{j=1}^p m_j r_j : \sum_{j=1}^p m_j r_j^2 = 1 \right\} = P^{1/2} \end{aligned}$$

(see Lemma 3.1.20).

Putting the computed values into the estimate (6.1.28), we get the following rate of convergence:

$$f(\hat{x}) - f^* \leq \frac{2P\bar{r}}{\sqrt{N(N+1)}}. \quad (6.1.42)$$

Note that the value $\tilde{f}(x) = \frac{1}{P}f(x)$ corresponds to the average individual expenses generated by the location x . Therefore,

$$\tilde{f}(\hat{x}) - \tilde{f}^* \leq \frac{2\bar{r}}{\sqrt{N(N+1)}}.$$

It is interesting that the right-hand side of this inequality is independent of any dimension. At the same time, it is clear that the reasonable accuracy for the approximate solution of our problem should not be too high. Given the low complexity of each iteration in the scheme (6.1.19), the total efficiency of the proposed technique looks quite promising.

To conclude with the location problem, let us write down explicitly a smooth approximation of the objective function.

$$\begin{aligned} f_\mu(x) &= \max_u \left\{ \sum_{j=1}^p m_j \langle u_j, x - c_j \rangle_{\mathbb{E}_1} - \mu d_2(u) : u \in Q_2 \right\} \\ &= \max_u \left\{ \sum_{j=1}^p m_j \left(\langle u_j, x - c_j \rangle_{\mathbb{E}_1} - \frac{1}{2} \mu (\|u_j\|_{\mathbb{E}_1}^*)^2 \right) : \|u_j\|_{\mathbb{E}_1}^* \leq 1, \right. \\ &\quad \left. j = 1, \dots, p \right\} \\ &= \sum_{j=1}^p m_j \psi_\mu(\|x - c_j\|_{\mathbb{E}_1}), \end{aligned}$$

where the function $\psi_\mu(\tau)$, $\tau \geq 0$, is defined as follows:

$$\psi_\mu(\tau) = \max_{\gamma \in [0,1]} \{ \gamma \tau - \frac{1}{2} \mu \gamma^2 \} = \begin{cases} \frac{\tau^2}{2\mu}, & 0 \leq \tau \leq \mu, \\ \tau - \frac{\mu}{2}, & \mu \leq \tau. \end{cases} \quad (6.1.43)$$

This is the so-called the *Huber loss function*.

6.1.4.3 Variational Inequalities with a Linear Operator

Consider a linear operator $B(w) = Bw + c: \mathbb{E} \rightarrow \mathbb{E}^*$, which is *monotone*:

$$\langle Bh, h \rangle \geq 0 \quad \forall h \in \mathbb{E}.$$

Let Q be a bounded closed convex set in \mathbb{E} . Then we can pose the following *variational inequality* problem:

$$\text{Find } w^* \in Q : \quad \langle B(w^*), w - w^* \rangle \geq 0 \quad \forall w \in Q. \quad (6.1.44)$$

Note that we can always rewrite problem (6.1.44) as an optimization problem. Indeed, define

$$\psi(w) = \max_v \{ \langle B(v), w - v \rangle : v \in Q \}.$$

In view of Theorem 3.1.8, $\psi(w)$ is a convex function. Let us show that the problem

$$\min_w \{ \psi(w) : w \in Q \} \quad (6.1.45)$$

is equivalent to (6.1.44).

Lemma 6.1.5 *A point w^* is a solution to (6.1.45) if and only if it solves variational inequality (6.1.44). Moreover, for such w^* we have $\psi(w^*) = 0$.*

Proof Indeed, at any $w \in Q$ the function ψ is non-negative. If w^* is a solution to (6.1.44), then for any $v \in Q$ we have

$$\langle B(v), v - w^* \rangle \geq \langle B(w^*), v - w^* \rangle \geq 0.$$

Hence, $\psi(w^*) = 0$ and $w^* \in \text{Arg min}_{w \in Q} \psi(w)$.

Now, consider some $w^* \in Q$ with $\psi(w^*) = 0$. Then for any $v \in Q$ we have

$$\langle B(v), v - w^* \rangle \geq 0.$$

Suppose there exists some $v_1 \in Q$ such that $\langle B(w^*), v_1 - w^* \rangle < 0$. Consider the points

$$v_\alpha = w^* + \alpha(v_1 - w^*), \quad \alpha \in [0, 1].$$

Then

$$\begin{aligned} 0 &\leq \langle B(v_\alpha), v_\alpha - w^* \rangle = \alpha \langle B(v_\alpha), v_1 - w^* \rangle \\ &= \alpha \langle B(w^*), v_1 - w^* \rangle + \alpha^2 \langle B \cdot (v_1 - w^*), v_1 - w^* \rangle. \end{aligned}$$

Hence, for α small enough we get a contradiction. \square

There are two possibilities for representing the problem (6.1.44), (6.1.45) in the form (6.1.10), (6.1.11).

1. Primal Form We take $\mathbb{E}_1 = \mathbb{E}_2 = \mathbb{E}$, $Q_1 = Q_2 = Q$, $d_1(x) = d_2(x) = d(x)$, $A = B$, and

$$\hat{f}(x) = \langle b, x \rangle_{\mathbb{E}_1}, \quad \hat{\phi}(u) = \langle b, u \rangle_{\mathbb{E}_1} + \langle Bu, u \rangle_{\mathbb{E}_1}.$$

Note that the quadratic function $\hat{\phi}(u)$ is convex. To compute the value and the gradient of the function $f_\mu(x)$, we need to solve the following problem:

$$\max_{u \in Q} \{ \langle Bx, u \rangle_{\mathbb{E}_1} - \mu d(u) - \langle b, u \rangle_{\mathbb{E}_1} - \langle Bu, u \rangle_{\mathbb{E}_1} \}. \quad (6.1.46)$$

Since in our case $M = 0$, from Theorem 6.1.3 we get the following estimate for the complexity of problem (6.1.44):

$$\frac{4D_1 \|B\|_{1,2}}{\epsilon}. \quad (6.1.47)$$

However, because of the presence of a non-trivial quadratic function in (6.1.46), the oracle for the function \hat{f} can be quite expensive. We can avoid that in the dual variant of this problem.

2. Dual Form Consider the dual variant of problem (6.1.45):

$$\min_{w \in Q} \max_{v \in Q} \langle B(v), w - v \rangle = \max_{v \in Q} \min_{w \in Q} \langle B(v), w - v \rangle = - \min_{v \in Q} \max_{w \in Q} \langle B(v), v - w \rangle.$$

Thus, we can take $\mathbb{E}_1 = \mathbb{E}_2 = \mathbb{E}$, $Q_1 = Q_2 = Q$, $d_1(x) = d_2(x) = d(x)$, $A = B$, and

$$\hat{f}(x) = \langle b, x \rangle_{\mathbb{E}_1} + \langle Bx, x \rangle_{\mathbb{E}_1}, \quad \hat{\phi}(u) = \langle b, u \rangle_{\mathbb{E}_1}.$$

Now the computation of the function value $f_\mu(x)$ becomes much simpler:

$$f_\mu(x) = \max_u \{ \langle Bx, u \rangle_{\mathbb{E}_1} - \mu d(u) - \langle b, u \rangle_{\mathbb{E}_1} : u \in Q \}.$$

Note that we pay quite a moderate cost for this. Indeed, now M becomes equal to $\|B\|_{1,2}$. Hence, the complexity estimate (6.1.47) increases up to the following level:

$$\frac{4D_1 \|B\|_{1,2}}{\epsilon} + \sqrt{\frac{D_1 \|B\|_{1,2}}{\epsilon}}.$$

In the important particular case of skew-symmetry of the operator B , that is $B + B^* = 0$, the primal and dual variant have a similar complexity.

6.1.4.4 Piece-Wise Linear Optimization

1. Maximum of Absolute Values Consider the following problem:

$$\min_{x \in Q_1} \left\{ f(x) = \max_{1 \leq j \leq m} |\langle a_j, x \rangle_{\mathbb{E}_1} - b^{(j)}| \right\}. \quad (6.1.48)$$

For simplicity, let us choose

$$\|x\|_{\mathbb{E}_1} = \left[\sum_{i=1}^n (x^{(i)})^2 \right]^{1/2}, \quad d_1(x) = \frac{1}{2} \|x\|^2.$$

Denote by A the matrix with rows a_j , $j = 1, \dots, m$. It is convenient to choose

$$\mathbb{E}_2 = \mathbb{R}^{2m}, \quad \|u\|_{\mathbb{E}_2} = \sum_{j=1}^{2m} |u^{(j)}|, \quad d_2(u) = \ln(2m) + \sum_{j=1}^{2m} u^{(j)} \ln u^{(j)}.$$

Then

$$f(x) = \max_u \{ \langle \hat{A}x, u \rangle_{\mathbb{E}_2} - \langle \hat{b}, u \rangle_{\mathbb{E}_2} : u \in \Delta_{2m} \},$$

where $\hat{A} = \begin{pmatrix} A \\ -A \end{pmatrix}$ and $\hat{b} = \begin{pmatrix} b \\ -b \end{pmatrix}$. Thus, $D_2 = \ln(2m)$, and

$$D_1 = \frac{1}{2} \bar{r}^2, \quad \bar{r} = \max_x \{ \|x\|_{\mathbb{E}_1} : x \in Q_1 \}.$$

It remains to compute the norm of the operator \hat{A} :

$$\begin{aligned} \|\hat{A}\|_{1,2} &= \max_{x,u} \{ \langle \hat{A}x, u \rangle_{\mathbb{E}_2} : \|x\|_{\mathbb{E}_1} = 1, \|u\|_{\mathbb{E}_2} = 1 \} \\ &= \max_x \{ \max_{1 \leq j \leq m} |\langle a_j, x \rangle_{\mathbb{E}_1}| : \|x\|_{\mathbb{E}_1} = 1 \} = \max_{1 \leq j \leq m} \|a_j\|_1^*. \end{aligned}$$

Putting all the computed values into the estimate (6.1.29), we see that the problem (6.1.48) can be solved in

$$2\sqrt{2} \bar{r} \max_{1 \leq j \leq m} \|a_j\|_1^* \sqrt{\ln(2m)} \cdot \frac{1}{\epsilon}$$

iterations of scheme (6.1.19). The standard subgradient schemes in this situation can count only on an

$$O \left(\left[\bar{r} \max_{1 \leq j \leq m} \|a_j\|_1^* \cdot \frac{1}{\epsilon} \right]^2 \right)$$

upper bound for the number of iterations.

Finally, the smooth version of the objective function in (6.1.48) is as follows:

$$\bar{f}_\mu(x) = \mu \ln \left(\frac{1}{m} \sum_{j=1}^m \xi \left(\frac{1}{\mu} [\langle a_j, x \rangle + b^{(j)}] \right) \right)$$

with $\xi(\tau) = \frac{1}{2}[e^\tau + e^{-\tau}]$. We leave the justification of this expression as an exercise for the reader.

2. Sum of Absolute Values Consider now the problem

$$\min_{x \in Q_1} \left\{ f(x) = \sum_{j=1}^m |\langle a_j, x \rangle_{\mathbb{E}_1} - b^{(j)}| \right\}. \quad (6.1.49)$$

The simplest representation of the function $f(\cdot)$ is as follows. Denote by A the matrix with the rows a_j . Let us choose

$$\mathbb{E}_2 = \mathbb{R}^m, \quad Q_2 = \{u \in \mathbb{R}^m : |u^{(j)}| \leq 1, j = 1, \dots, m\},$$

$$d_2(u) = \frac{1}{2} \|u\|_{\mathbb{E}_2}^2 = \frac{1}{2} \sum_{j=1}^m \|a_j\|_{\mathbb{E}_1}^* \cdot (u^{(j)})^2.$$

Then the smooth version of the objective function is as follows:

$$\begin{aligned} f_\mu(x) &= \max_u \{ \langle Ax - b, u \rangle_{\mathbb{E}_2} - \mu d_2(u) : u \in Q_2 \} \\ &= \sum_{j=1}^m \|a_j\|_{\mathbb{E}_1}^* \cdot \psi_\mu \left(\frac{|\langle a_j, x \rangle_{\mathbb{E}_1} - b^{(j)}|}{\|a_j\|_{\mathbb{E}_1}^*} \right), \end{aligned}$$

where the function $\psi_\mu(\tau)$ is defined by (6.1.43). Note that

$$\begin{aligned} \|A\|_{1,2} &= \max_{x,u} \left\{ \sum_{j=1}^m u^{(j)} \langle a_j, x \rangle_{\mathbb{E}_1} : \|x\|_{\mathbb{E}_1} \leq 1, \|u\|_{\mathbb{E}_2} \leq 1 \right\} \\ &\leq \max_u \left\{ \sum_{j=1}^m \|a_j\|_{\mathbb{E}_1}^* \cdot |u^{(j)}| : \sum_{j=1}^m \|a_j\|_{\mathbb{E}_1}^* \cdot (u^{(j)})^2 \leq 1 \right\} \\ &= D^{1/2} \equiv \left[\sum_{j=1}^m \|a_j\|_{\mathbb{E}_1}^* \right]^{1/2}. \end{aligned}$$

On the other hand, $D_2 = \frac{1}{2}D$. Therefore from Theorem 6.1.3 we get the following complexity bound:

$$\frac{2}{\epsilon} \cdot \sqrt{2D_1} \cdot \sum_{j=1}^m \|a_j\|_{\mathbb{E}_1}^*$$

iterations of method (6.1.19).

6.1.5 Implementation Issues

6.1.5.1 Computational Complexity

Let us discuss the computational complexity of the method (6.1.19) as applied to the function $\tilde{f}_\mu(\cdot)$. The main computations are performed at Steps (b) and (c) of the algorithm.

Step (b). Call of Oracle At this step we need to compute the solution of the following maximization problem:

$$\max_{u \in Q_2} \{ \langle Ay_k, u \rangle_{\mathbb{E}_2} - \hat{\phi}(u) - \mu d_2(u) : u \in Q_2 \}.$$

Note that from the origin of this problem we know that this computation for $\mu = 0$ can be done in a closed form. Thus, we can expect that with a properly chosen prox-function, computation of the smoothed version is not too difficult. In Sect. 6.1.4 we have seen three examples which confirm this belief.

Step (c). Computation of v_{k+1} This computation consists in solving the following problem:

$$\min_{x \in Q_1} \{ d_1(x) + \langle s, x \rangle_{\mathbb{E}_1} \}$$

for some fixed $s \in \mathbb{E}_1^*$. If the set Q_1 and the prox-function $d_1(\cdot)$ are simple enough, this computation can be done in a closed form (see Sect. 6.1.4). For some sets we need to solve an auxiliary equation with one variable.

6.1.5.2 Computational Stability

Our approach is based on the smoothing of non-differentiable functions. In accordance with (6.1.33), the value of the smoothness parameter μ must be of the order of ϵ . This may cause some numerical troubles in computing the function $\tilde{f}_\mu(x)$ and its gradient. Among examples of Sect. 6.1.4, only a smooth variant of the objective function in Sect. 6.1.4.2 does not involve dangerous operations; all others need a careful implementation.

In both Sects. 6.1.4.1 and 6.1.4.4 we need a stable technique for computing the values and derivatives of the function

$$\eta(u) = \mu \ln \left(\sum_{j=1}^m e^{u^{(j)}/\mu} \right) \quad (6.1.50)$$

with very small values of parameter μ . This can be done in the following way. Let

$$\bar{u} = \max_{1 \leq j \leq m} u^{(j)}, \quad v^{(j)} = u^{(j)} - \bar{u}, \quad j = 1, \dots, m.$$

Then

$$\eta(u) = \bar{u} + \eta(v).$$

Note that all components of the vector v are non-negative and one of them is zero. Therefore, the value $\eta(v)$ can be computed quite accurately. The same technique can be used to compute the gradient since $\nabla \eta(u) = \nabla \eta(v)$.

6.2 An Excessive Gap Technique for Non-smooth Convex Minimization

(Primal-dual problem structure; An excessive gap condition; Gradient mapping; Convergence analysis; Minimizing strongly convex functions.)

6.2.1 Primal-Dual Problem Structure

In this section, we give some extensions of the results presented in Sect. 6.1, where it was shown that some structured non-smooth optimization problems can be solved in $O(\frac{1}{\epsilon})$ iterations of a gradient-type scheme with ϵ being the desired accuracy of the solution. This complexity is much better than the theoretical lower complexity bound $O(\frac{1}{\epsilon^2})$ for Black-Box methods (see Sect. 3.2). This improvement, of course, is possible because of certain relaxations of the standard Black Box assumption. Instead, it was assumed that our problem has an explicit and quite simple minimax structure. However, the approach discussed in Sect. 6.1 has a certain drawback. Namely, the number of steps of the optimization scheme must be fixed in advance. It is chosen in accordance with the worst case complexity analysis and desired accuracy. Let us try to be more flexible.

Consider the same optimization problems as before:

$$\text{Find } f^* = \min_{x \in Q_1} f(x), \quad (6.2.1)$$

where Q_1 is a bounded closed convex set in a finite-dimensional real vector space \mathbb{E}_1 , and f is a continuous convex function on Q_1 . We do not assume f to be differentiable. Let the structure of the objective function be described by the following *model*:

$$f(x) = \hat{f}(x) + \max_{u \in Q_2} \{\langle Ax, u \rangle_{\mathbb{E}_2} - \hat{\phi}(u)\}, \quad (6.2.2)$$

where the function \hat{f} is continuous and convex on Q_1 , Q_2 is a closed convex bounded set in a finite-dimensional real vector space \mathbb{E}_2 , $\hat{\phi}(\cdot)$ is a continuous convex function on Q_2 , and the linear operator A maps \mathbb{E}_1 to \mathbb{E}_2^* . In this case, problem (6.2.1) can be written in an *adjoint* form:

$$\begin{aligned} f_* &= \max_{u \in Q_2} \phi(u), \\ \phi(u) &= -\hat{\phi}(u) + \min_{x \in Q_1} \{\langle Ax, u \rangle_{\mathbb{E}_2} + \hat{f}(x)\}, \end{aligned} \quad (6.2.3)$$

which has zero duality gap (see (6.1.34)).

We assume that this representation is completely similar to (6.2.1) in the following sense. All methods described in this section are implementable only if the optimization problems involved in the definitions of functions f and ϕ can be solved in a closed form. So, we assume that the structure of all objects in \hat{f} , $\hat{\phi}$, Q_1 and Q_2 is simple enough. We also assume that functions \hat{f} and $\hat{\phi}$ have Lipschitz continuous gradients with Lipschitz constants $L_1(\hat{f})$ and $L_2(\hat{\phi})$ respectively.

Let us show that the knowledge of structure (6.2.2) can help in solving problems (6.2.1) and (6.2.3). Consider a *prox-function* $d_2(\cdot)$ of the set Q_2 . This means that d_2 is continuous and strongly convex on Q_2 with a strong convexity parameter equal to one. Denote by

$$u_0 = \arg \min_{u \in Q_2} d_2(u)$$

the *prox-center* of the function d_2 . Without loss of generality we assume that $d_2(u_0) = 0$. Thus, in view of (4.2.18), for any $u \in Q_2$ we have

$$d_2(u) \geq \frac{1}{2} \|u - u_0\|_2^2. \quad (6.2.4)$$

Let μ_2 be a positive *smoothing* parameter. Consider the following function:

$$f_{\mu_2}(x) = \hat{f}(x) + \max_{u \in Q_2} \{\langle Ax, u \rangle_{\mathbb{E}_2} - \hat{\phi}(u) - \mu_2 d_2(u)\}. \quad (6.2.5)$$

Denote by $u_{\mu_2}(x)$ the optimal solution of this problem. Since the function d_2 is strongly convex, this solution is unique. In accordance with Danskin's theorem, the

gradient of f_{μ_2} is well defined as

$$\nabla f_{\mu_2}(x) = \nabla \hat{f}(x) + A^* u_{\mu_2}(x). \quad (6.2.6)$$

Moreover, this gradient is Lipschitz-continuous with constant

$$L_1(f_{\mu_2}) = L_1(\hat{f}) + \frac{1}{\mu_2} \|A\|_{1,2}^2 \quad (6.2.7)$$

(see Theorem 6.1.1).

Similarly, let us consider a prox-function $d_1(\cdot)$ of the set Q_1 , which has convexity parameter equal to one, and the prox-center x_0 with $d_1(x_0) = 0$. By (4.2.18), for any $x \in Q_1$ we have

$$d_1(x) \geq \frac{1}{2} \|x - x_0\|_1^2. \quad (6.2.8)$$

Let μ_1 be a positive smoothing parameter. Consider

$$\phi_{\mu_1}(u) = -\hat{\phi}(u) + \min_{x \in Q_1} \{\langle Ax, u \rangle_{\mathbb{E}_2} + \hat{f}(x) + \mu_1 d_1(x)\}. \quad (6.2.9)$$

Since the second term in the above definition is a minimum of linear functions, $\phi_{\mu_1}(u)$ is concave. Denote by $x_{\mu_1}(u)$ the unique optimal solution of the above problem. In accordance with Theorem 6.1.1, the gradient

$$\nabla \phi_{\mu_1}(u) = -\nabla \hat{\phi}(u) + Ax_{\mu_1}(u) \quad (6.2.10)$$

is Lipschitz-continuous with constant

$$L_2(\phi_{\mu_1}) = L_2(\hat{\phi}) + \frac{1}{\mu_1} \|A\|_{1,2}^2. \quad (6.2.11)$$

6.2.2 An Excessive Gap Condition

In view of Theorem 1.3.1, for any $x \in Q_1$ and $u \in Q_2$ we have

$$\phi(u) \leq f(x), \quad (6.2.12)$$

and our assumptions guarantee no duality gap for problems (6.2.1) and (6.2.3). However, $f_{\mu_2}(x) \leq f(x)$ and $\phi(u) \leq \phi_{\mu_1}(u)$. This opens a possibility to satisfy the following *excessive gap condition*:

$$\boxed{f_{\mu_2}(\bar{x}) \leq \phi_{\mu_1}(\bar{u})} \quad (6.2.13)$$

for certain $\bar{x} \in Q_1$ and $\bar{u} \in Q_2$. Let us show that condition (6.2.13) provides us with an upper bound on the quality of the primal-dual pair (\bar{x}, \bar{u}) .

Lemma 6.2.1 *Let $\bar{x} \in Q_1$ and $\bar{u} \in Q_2$ satisfy (6.2.13). Then*

$$\begin{aligned} 0 &\leq \max\{f(\bar{x}) - f^*, f^* - \phi(\bar{u})\} \\ &\leq f(\bar{x}) - \phi(\bar{u}) \leq \mu_1 D_1 + \mu_2 D_2, \end{aligned} \quad (6.2.14)$$

where $D_1 = \max_{x \in Q_1} d_1(x)$, and $D_2 = \max_{u \in Q_2} d_2(u)$.

Proof Indeed, for any $\bar{x} \in Q_1$, $\bar{u} \in Q_2$ we have

$$f(\bar{x}) - \mu_2 D_2 \leq f_{\mu_2}(\bar{x}) \stackrel{(6.2.13)}{\leq} \phi_{\mu_1}(\bar{u}) \leq \phi(\bar{u}) + \mu_1 D_1.$$

It remains to apply inequality (6.2.12). \square

Our goal is to justify a process for recursively updating the pair (\bar{x}, \bar{u}) , which maintains inequality (6.2.13) as μ_1 and μ_2 go to zero. Before we start our analysis, let us prove a useful inequality.

Lemma 6.2.2 *For any x and \hat{x} from Q_1 we have:*

$$f_{\mu_2}(\hat{x}) + \langle \nabla f_{\mu_2}(\hat{x}), x - \hat{x} \rangle_{\mathbb{E}_1} \leq \hat{f}(x) + \langle Ax, u_{\mu_2}(\hat{x}) \rangle_{\mathbb{E}_2} - \hat{\phi}(u_{\mu_2}(\hat{x})). \quad (6.2.15)$$

Proof Let us take arbitrary x and \hat{x} from Q_1 . Let $\hat{u} = u_{\mu_2}(\hat{x})$. Then

$$\begin{aligned} f_{\mu_2}(\hat{x}) + \langle \nabla f_{\mu_2}(\hat{x}), x - \bar{y} \rangle_{\mathbb{E}_1} &\stackrel{(6.2.5), (6.2.6)}{=} \hat{f}(\hat{x}) + \langle A\bar{y}, \hat{u} \rangle_{\mathbb{E}_2} - \hat{\phi}(\hat{u}) - \mu_2 d_2(\hat{u}) \\ &\quad + \langle \nabla \hat{f}(\hat{x}) + A^* \hat{u}, x - \hat{x} \rangle_{\mathbb{E}_1} \\ &\stackrel{(2.1.2)}{\leq} \hat{f}(x) + \langle Ax, \hat{u} \rangle_{\mathbb{E}_2} - \hat{\phi}(\hat{u}). \quad \square \end{aligned}$$

Let us justify the possibility of satisfying the excessive gap condition (6.2.13) at some starting primal-dual pair.

Lemma 6.2.3 *Let us choose an arbitrary $\mu_2 > 0$ and set*

$$\begin{aligned} \bar{x} &= \arg \min_{x \in Q_1} \{ \langle \nabla f_{\mu_2}(x_0), x - x_0 \rangle_{\mathbb{E}_1} + L_1(f_{\mu_2})d_1(x) \}, \\ \bar{u} &= u_{\mu_2}(x_0). \end{aligned} \quad (6.2.16)$$

Then the excessive gap condition is satisfied for any $\mu_1 \geq L_1(f_{\mu_2})$.

Proof Indeed, in view of (1.2.11) we have

$$\begin{aligned}
 f_{\mu_2}(\bar{x}) &\leq f_{\mu_2}(x_0) + \langle \nabla f_{\mu_2}(x_0), \bar{x} - x_0 \rangle_{\mathbb{E}_1} + \frac{1}{2} L_1(f_{\mu_2}) \|\bar{x} - x_0\|_1^2 \\
 &\stackrel{(6.2.4)}{\leq} f_{\mu_2}(x_0) + \langle \nabla f_{\mu_2}(x_0), \bar{x} - x_0 \rangle_{\mathbb{E}_1} + \frac{1}{2} L_1(f_{\mu_2}) d_1(\bar{x}) \\
 &\stackrel{(6.2.16)}{=} f_{\mu_2}(x_0) + \min_{x \in Q_1} \{ \langle \nabla f_{\mu_2}(x_0), x - x_0 \rangle_{\mathbb{E}_1} + L_1(f_{\mu_2}) d_1(x) \} \\
 &\stackrel{(6.2.15)}{\leq} \min_{x \in Q_1} \left\{ \hat{f}(x) + \langle Ax, u_{\mu_2}(x_0) \rangle_{\mathbb{E}_2} - \hat{\phi}(u_{\mu_2}(x_0)) + L_1(f_{\mu_2}) d_1(x) \right\} \\
 &\stackrel{(6.2.9)}{=} \phi_{L_1(f_{\mu_2})}(\bar{u}) \leq \phi_{\mu_1}(\bar{u}). \quad \square
 \end{aligned}$$

Thus, condition (6.2.13) can be satisfied for some primal-dual pair. Let us show how we can update the points \bar{x} and \bar{u} in order to keep it valid for smaller values of μ_1 and μ_2 . In view of the symmetry of the situation, at the first step of the process we can try to decrease only μ_1 , keeping μ_2 unchanged. After that, at the second step, we update μ_2 and keep μ_1 constant, and so on. The main advantage of such a switching strategy is that we need to find a justification only for the first step. The proof for the second one will be symmetric.

Theorem 6.2.1 *Let points $\bar{x} \in Q_1$ and $\bar{u} \in Q_2$ satisfy the excessive gap condition (6.2.13) for some positive μ_1 and μ_2 . Let us fix $\tau \in (0, 1)$ and choose $\mu_1^+ = (1 - \tau)\mu_1$,*

$$\begin{aligned}
 \hat{x} &= (1 - \tau)\bar{x} + \tau x_{\mu_1}(\bar{u}), \\
 \bar{u}_+ &= (1 - \tau)\bar{u} + \tau u_{\mu_2}(\hat{x}), \\
 \bar{x}_+ &= (1 - \tau)\bar{x} + \tau x_{\mu_1^+}(\bar{u}_+).
 \end{aligned} \tag{6.2.17}$$

Then the pair (\bar{x}_+, \bar{u}_+) satisfies condition (6.2.13) with smoothing parameters μ_1^+ and μ_2 provided that τ satisfies the following relation:

$$\boxed{\frac{\tau^2}{1-\tau} \leq \frac{\mu_1}{L_1(f_{\mu_2})}} \tag{6.2.18}$$

Proof Let $\hat{u} = u_{\mu_2}(\hat{x})$, $x_1 = x_{\mu_1}(\bar{u})$, and $\tilde{x}_+ = x_{\mu_1^+}(\bar{u}_+)$. Since $\hat{\phi}$ is convex, in view of the operation in (6.2.17), we have $\hat{\phi}(\bar{u}_+) \leq (1 - \tau)\hat{\phi}(\bar{u}) + \tau\hat{\phi}(\hat{u})$. Therefore,

$$\begin{aligned}
 \phi_{\mu_1^+}(\bar{u}_+) &= (1 - \tau)\mu_1 d_1(\tilde{x}_+) + \langle A\tilde{x}_+, (1 - \tau)\bar{u} + \tau\hat{u} \rangle_{\mathbb{E}_2} + \hat{f}(\tilde{x}_+) - \hat{\phi}(\bar{u}_+) \\
 &\geq (1 - \tau)[\mu_1 d_1(\tilde{x}_+) + \langle A\tilde{x}_+, \bar{u} \rangle_{\mathbb{E}_2} + \hat{f}(\tilde{x}_+) - \hat{\phi}(\bar{u})] \\
 &\quad + \tau[\hat{f}(\tilde{x}_+) + \langle A\tilde{x}_+, \hat{u} \rangle_{\mathbb{E}_2} - \hat{\phi}(\hat{u})] \\
 &\stackrel{(6.2.15)}{\geq} (1 - \tau)[\phi_{\mu_1}(\bar{u}) + \frac{1}{2}\mu_1 \|\tilde{x}_+ - x_1\|_1^2]_a \\
 &\quad + \tau[f_{\mu_2}(\hat{x}) + \langle \nabla f_{\mu_2}(\hat{x}), \tilde{x}_+ - \hat{x} \rangle_{\mathbb{E}_1}]_b.
 \end{aligned}$$

Note that in view of condition (6.2.13) and the first line in (6.2.17) we have

$$\begin{aligned}
 \phi_{\mu_1}(\bar{u}) &\geq f_{\mu_2}(\bar{x}) \geq f_{\mu_2}(\hat{x}) + \langle \nabla f_{\mu_2}(\hat{x}), \bar{x} - \hat{x} \rangle_{\mathbb{E}_1} \\
 &= f_{\mu_2}(\hat{x}) + \tau \langle \nabla f_{\mu_2}(\hat{x}), \bar{x} - x_1 \rangle_{\mathbb{E}_1}.
 \end{aligned}$$

Therefore, we can estimate the expression in the first brackets as follows:

$$[\cdot]_a \geq f_{\mu_2}(\hat{x}) + \tau \langle \nabla f_{\mu_2}(\hat{x}), \bar{x} - x_1 \rangle_{\mathbb{E}_1} + \frac{1}{2}\mu_1 \|\tilde{x}_+ - x_1\|_1^2.$$

In view of the first line in (6.2.15), for second brackets we have

$$[\cdot]_b = f_{\mu_2}(\hat{x}) + \langle \nabla f_{\mu_2}(\hat{x}), \tilde{x}_+ - x_1 + (1 - \tau)(x_1 - \bar{x}) \rangle_{\mathbb{E}_1}.$$

Thus, taking into account that $\bar{x}_+ - \hat{x} \stackrel{(6.2.17)}{=} \tau(\tilde{x}_+ - x_1)$, we finish the proof as follows:

$$\begin{aligned}
 \phi_{\mu_1^+}(\bar{u}_+) &\geq f_{\mu_2}(\hat{x}) + \tau \langle \nabla f_{\mu_2}(\hat{x}), \tilde{x}_+ - x_1 \rangle_{\mathbb{E}_1} + \frac{1}{2}(1 - \tau)\mu_1 \|\tilde{x}_+ - x_1\|_1^2 \\
 &= f_{\mu_2}(\hat{x}) + \langle \nabla f_{\mu_2}(\hat{x}), \bar{x}_+ - \hat{x} \rangle_{\mathbb{E}_1} + \frac{(1 - \tau)\mu_1}{2\tau^2} \|\bar{x}_+ - \hat{x}\|_1^2 \\
 &\stackrel{(6.2.18)}{\geq} f_{\mu_2}(\hat{x}) + \langle \nabla f_{\mu_2}(\hat{x}), \bar{x}_+ - \hat{x} \rangle_{\mathbb{E}_1} + \frac{1}{2}L_1(f_{\mu_2}) \|\bar{x}_+ - \hat{x}\|_1^2 \\
 &\stackrel{(1.2.11)}{\geq} f_{\mu_2}(\bar{x}_+). \quad \square
 \end{aligned}$$

6.2.3 Convergence Analysis

In Sect. 6.2.2, we have seen that the smoothness parameters μ_1 and μ_2 can be decreased by a switching strategy. Thus, in order to transform the result of Theorem 6.2.1 into an algorithmic scheme, we need to point out a strategy for updating these parameters, which is compatible with the growth condition (6.2.18). In this section, we do this for an important case $L_1(\hat{f}) = L_2(\hat{\phi}) = 0$.

It is convenient to represent the smoothness parameters as follows:

$$\mu_1 = \lambda_1 \cdot \|A\|_{1,2} \cdot \sqrt{\frac{D_2}{D_1}}, \quad \mu_2 = \lambda_2 \cdot \|A\|_{1,2} \cdot \sqrt{\frac{D_1}{D_2}}. \quad (6.2.19)$$

Then the estimate (6.2.14) for the duality gap becomes symmetric:

$$f(\bar{x}) - \phi(\bar{u}) \leq (\lambda_1 + \lambda_2) \cdot \|A\|_{1,2} \cdot \sqrt{D_1 D_2}. \quad (6.2.20)$$

Since by (6.2.7), $L_1(f_{\mu_2}) = \frac{1}{\mu_2} \|A\|_{1,2}^2$, condition (6.2.18) becomes problem independent:

$$\frac{\tau^2}{1-\tau} \leq \mu_1 \mu_2 \cdot \frac{1}{\|A\|_{1,2}^2} = \lambda_1 \lambda_2. \quad (6.2.21)$$

Let us write down the corresponding switching algorithmic scheme in an explicit form. It is convenient to have a permanent iteration counter. In this case, at even iterations we apply the primal update (6.2.17), and at odd iterations the corresponding dual update is used. Since at even iterations λ_2 does not change and at odd iterations λ_1 does not change it is convenient to put their new values in the same sequence $\{\alpha_k\}_{k=-1}^\infty$. Let us fix the following relations between the sequences:

$$k = 2l : \lambda_{1,k} = \alpha_{k-1}, \quad \lambda_{2,k} = \alpha_k, \quad (6.2.22)$$

$$k = 2l + 1 : \lambda_{1,k} = \alpha_k, \quad \lambda_{2,k} = \alpha_{k-1}.$$

Then the corresponding parameters τ_k (see the rule (6.2.1)) define the reduction rate of the sequence $\{\alpha_k\}_{k=-1}^\infty$.

Lemma 6.2.4 *For all $k \geq 0$ we have $\alpha_{k+1} = (1 - \tau_k)\alpha_{k-1}$.*

Proof Indeed, in accordance with (6.2.22), if $k = 2l$, then

$$\alpha_{k+1} = \lambda_{1,k+1} = (1 - \tau_k)\lambda_{1,k} = (1 - \tau_k)\alpha_{k-1}.$$

And if $k = 2l + 1$, then $\alpha_{k+1} = \lambda_{2,k+1} = (1 - \tau_k)\lambda_{2,k} = (1 - \tau_k)\alpha_{k-1}$. \square

Corollary 6.2.1 *In terms of the sequence $\{\alpha_k\}_{k=-1}^\infty$, condition (6.2.21) is as follows:*

$$(\alpha_{k+1} - \alpha_{k-1})^2 \leq \alpha_{k+1} \alpha_k \alpha_{k-1}^2, \quad k \geq 0. \quad (6.2.23)$$

Proof In view of (6.2.22), we always have $\lambda_{1,k} \lambda_{2,k} = \alpha_k \alpha_{k-1}$. Since $\tau_k = 1 - \frac{\alpha_{k+1}}{\alpha_{k-1}}$, we get (6.2.23). \square

Clearly, condition (6.2.23) is satisfied by

$$\alpha_k = \frac{2}{k+2}, \quad k \geq -1. \quad (6.2.24)$$

Then

$$\tau_k = 1 - \frac{\alpha_{k+1}}{\alpha_{k-1}} = \frac{2}{k+3}, \quad k \geq 0. \quad (6.2.25)$$

Now we are ready to write down an algorithmic scheme. Let us do this for the rule (6.2.17). In this scheme, we use the sequences $\{\mu_{1,k}\}_{k=-1}^\infty$ and $\{\mu_{2,k}\}_{k=-1}^\infty$, generated in accordance with rules (6.2.19), (6.2.22) and (6.2.24).

1. Initialization: Choose \bar{x}_0 and \bar{u}_0 in accordance with (6.2.16) taking $\mu_1 = \mu_{1,0}$ and $\mu_2 = \mu_{2,0}$.

2. Iterations ($k \geq 0$):

- (a) Set $\tau_k = \frac{2}{k+3}$.
- (b) If k is even, then generate $(\bar{x}_{k+1}, \bar{u}_{k+1})$ from (\bar{x}_k, \bar{u}_k) using (6.2.17).
- (c) If k is odd, then generate $(\bar{x}_{k+1}, \bar{u}_{k+1})$ from (\bar{x}_k, \bar{u}_k) using the symmetric dual variant of (6.2.17).

(6.2.26)

Theorem 6.2.2 *Let the sequences $\{\bar{x}_k\}_{k=0}^\infty$ and $\{\bar{u}_k\}_{k=0}^\infty$ be generated by method (6.2.26). Then each pair of points (\bar{x}_k, \bar{u}_k) satisfy the excessive gap condition. Therefore,*

$$f(\bar{x}_k) - \phi(\bar{u}_k) \leq \frac{4\|A\|_{1,2}}{k+1} \sqrt{D_1 D_2}. \quad (6.2.27)$$

Proof In accordance with our choice of parameters,

$$\mu_{1,0} \mu_{2,0} = \lambda_{1,0} \lambda_{2,0} \cdot \|A\|_{1,2}^2 = 2\mu_{2,0} L_1(f_{\mu_{2,0}}) > \mu_{2,0} L_1(f_{\mu_{2,0}}).$$

Hence, in view of Lemma 6.2.3 the pair (\bar{x}_0, \bar{u}_0) satisfies the excessive gap condition. We have already checked that the sequence $\{\tau_k\}_{k=0}^\infty$ defined by (6.2.25) satisfies

the conditions of Theorem 6.2.1. Therefore, excessive gap conditions will be valid for the sequences generated by (6.2.26). It remains to use inequality (6.2.20). \square

6.2.4 Minimizing Strongly Convex Functions

Consider now the model (6.2.2), which satisfies the following assumption.

Assumption 6.2.1 *In representation (6.2.2) the function \hat{f} is strongly convex with convexity parameter $\hat{\sigma} > 0$.*

Let us prove the following variant of Danskin's theorem.

Lemma 6.2.5 *Under Assumption 6.2.1 the function ϕ defined by (6.2.3) is concave and differentiable. Moreover, its gradient*

$$\nabla \phi(u) = -\nabla \hat{\phi}(u) + Ax_0(u), \quad (6.2.28)$$

where $x_0(u)$ is defined by (6.2.9), is Lipschitz-continuous with constant

$$L_2(\phi) = \frac{1}{\hat{\sigma}} \|A\|_{1,2}^2 + L_2(\hat{\phi}). \quad (6.2.29)$$

Proof Let $\tilde{\phi}(u) = \min_{x \in Q_1} \{ \langle Ax, u \rangle_{\mathbb{E}_2} + \hat{f}(x) \}$. This function is concave as a minimum of linear functions. Since \hat{f} is strongly convex, the solution of the latter minimization problem is unique. Therefore, $\tilde{\phi}(\cdot)$ is differentiable and $\nabla \tilde{\phi}(u) = Ax_0(u)$.

Consider two points u_1 and u_2 . From the first-order optimality conditions for (6.2.3) we have

$$\langle A^*u_1 + \nabla \hat{f}(x_0(u_1)), x_0(u_2) - x_0(u_1) \rangle_{\mathbb{E}_1} \geq 0,$$

$$\langle A^*u_2 + \nabla \hat{f}(x_0(u_2)), x_0(u_1) - x_0(u_2) \rangle_{\mathbb{E}_1} \geq 0.$$

Adding these inequalities and using the strong convexity of $\hat{f}(\cdot)$, we continue as follows:

$$\begin{aligned} & \langle Ax_0(u_2) - Ax_0(u_1), u_1 - u_2 \rangle_{\mathbb{E}_2} \\ & \geq \langle \nabla \hat{f}(x_0(u_1)) - \nabla \hat{f}(x_0(u_2)), x_0(u_1) - x_0(u_2) \rangle_{\mathbb{E}_1} \\ & \stackrel{(2.1.22)}{\geq} \hat{\sigma} \|x_0(u_1) - x_0(u_2)\|_{\mathbb{E}_1}^2 \stackrel{(6.1.9)}{\geq} \frac{\hat{\sigma}}{\|A\|_{1,2}^2} \left(\|\nabla \tilde{\phi}(u_1) - \nabla \tilde{\phi}(u_2)\|_{\mathbb{E}_2}^* \right)^2. \end{aligned}$$

Thus, $\|\nabla \tilde{\phi}(u_1) - \nabla \tilde{\phi}(u_2)\|_{\mathbb{E}_2}^* \leq \frac{1}{\hat{\sigma}} \|A\|_{1,2}^2 \cdot \|u_1 - u_2\|_{\mathbb{E}_2}$, and (6.2.29) follows. \square

Lemma 6.2.6 *For any u and \hat{u} from Q_2 , we have:*

$$\phi(\hat{u}) + \langle \nabla \phi(\hat{u}), u - \hat{u} \rangle_{\mathbb{E}_2} \geq -\hat{\phi}(u) + \langle Ax_0(\hat{u}), u \rangle_{\mathbb{E}_2} + \hat{f}(x_0(\hat{u})). \quad (6.2.30)$$

Proof Let us take arbitrary u and \hat{u} from Q_2 . Define $\hat{x} = x_0(\hat{u})$. Then

$$\begin{aligned} & \phi(\hat{u}) + \langle \nabla \phi(\hat{u}), u - \hat{u} \rangle_{\mathbb{E}_2} \\ &= -\hat{\phi}(\hat{u}) + \langle A\hat{x}, \hat{u} \rangle_{\mathbb{E}_2} + \hat{f}(\hat{x}) + \langle -\nabla \hat{\phi}(\hat{u}) + A\hat{x}, u - \hat{u} \rangle_{\mathbb{E}_2} \\ &\stackrel{(2.1.2)}{\geq} -\hat{\phi}(u) + \langle A\hat{x}, u \rangle_{\mathbb{E}_2} + \hat{f}(\hat{x}). \quad \square \end{aligned}$$

In this section, we derive an optimization scheme from the following variant of excessive gap condition:

$$\boxed{f_{\mu_2}(\bar{x}) \leq \phi(\bar{u})} \quad (6.2.31)$$

for some $\bar{x} \in Q_1$ and $\bar{u} \in Q_2$.

This condition can be seen as a variant of condition (6.2.13) with $\mu_1 = 0$. However, we prefer not to use the results of the previous sections since our assumptions will be slightly different. For example, we no longer need the set Q_1 to be bounded.

Lemma 6.2.7 *Let points \bar{x} from Q_1 and \bar{u} from Q_2 satisfy condition (6.2.31). Then*

$$0 \leq f(\bar{x}) - \phi(\bar{u}) \leq \mu_2 D_2. \quad (6.2.32)$$

Proof Indeed, for any $x \in Q_1$, we have $f_{\mu_2}(x) \geq f(x) - \mu_2 D_2$. \square

Define the adjoint gradient mapping as follows:

$$V(u) = \arg \max_{v \in Q_2} \left\{ \langle \nabla \phi(u), v - u \rangle_{\mathbb{E}_2} - \frac{1}{2} L_2(\phi) \|v - u\|_{\mathbb{E}_2}^2 \right\}. \quad (6.2.33)$$

Lemma 6.2.8 *The excessive gap condition (6.2.31) is valid for $\mu_2 = L_2(\phi)$ and*

$$\bar{x} = x_0(u_0), \quad \bar{u} = V(u_0). \quad (6.2.34)$$

Proof Indeed, in view of Lemma 6.2.5 and (1.2.11), we get the following relations:

$$\begin{aligned}
 \phi(V(u_0)) &\geq \phi(u_0) + \langle \nabla \phi(u_0), V(u_0) - u_0 \rangle_{\mathbb{E}_2} - \frac{1}{2} L_2(\phi) \|V(u_0) - u_0\|_2^2 \\
 &\stackrel{(6.2.33)}{=} \max_{u \in Q_2} \left\{ \phi(u_0) + \langle \nabla \phi(u_0), u - u_0 \rangle_{\mathbb{E}_2} - \frac{1}{2} L_2(\phi) \|u - u_0\|_2^2 \right\} \\
 &\stackrel{(6.2.3), (6.2.28)}{=} \max_{u \in Q_2} \left\{ -\hat{\phi}(u_0) + \langle Ax_0(u_0), u_0 \rangle_{\mathbb{E}_2} + \hat{f}(x_0(u_0)) \right. \\
 &\quad \left. + \langle Ax_0(u_0) - \nabla \hat{\phi}(u_0), u - u_0 \rangle_{\mathbb{E}_2} - \frac{1}{2} \mu_2 \|u - u_0\|_2^2 \right\} \\
 &\stackrel{(6.2.4)}{\geq} \max_{u \in Q_2} \left\{ -\hat{\phi}(u) + \hat{f}(x_0(u_0)) + \langle Ax_0(u_0), u \rangle_{\mathbb{E}_2} - \mu_2 d_2(u) \right\} \\
 &\stackrel{(6.2.5)}{=} f_{\mu_2}(x_0(u_0)). \quad \square
 \end{aligned}$$

Theorem 6.2.3 *Let points $\bar{x} \in Q_1$ and $\bar{u} \in Q_2$ satisfy the excessive gap condition (6.2.31) for some positive μ_2 . Let us fix $\tau \in (0, 1)$ and choose $\mu_2^+ = (1 - \tau)\mu_2$,*

$$\begin{aligned}
 \hat{u} &= (1 - \tau)\bar{u} + \tau u_{\mu_2}(\bar{x}), \\
 \bar{x}_+ &= (1 - \tau)\bar{x} + \tau x_0(\hat{u}), \\
 \bar{u}_+ &= V(\hat{u}).
 \end{aligned} \tag{6.2.35}$$

Then the pair (\bar{x}_+, \bar{u}_+) satisfies condition (6.2.31) with smoothness parameter μ_2^+ , provided that τ satisfies the following growth relation:

$$\frac{\tau^2}{1 - \tau} \leq \frac{\mu_2}{L_2(\phi)}. \tag{6.2.36}$$

Proof Let $\hat{x} = x_0(\hat{u})$ and $u_2 = u_{\mu_2}(\bar{x})$. In view of the second rule in (6.2.35), and (6.2.5), we have:

$$\begin{aligned}
 f_{\mu_2^+}(\bar{x}_+) &= \hat{f}(\bar{x}_+) + \max_{u \in Q_2} \left\{ \langle A((1-\tau)\bar{x} + \tau\hat{x}), u \rangle_{\mathbb{E}_2} - \hat{\phi}(u) \right. \\
 &\quad \left. - (1-\tau)\mu_2 d_2(u) \right\} \\
 &\stackrel{(3.1.2)}{\leq} \max_{u \in Q_2} \left\{ (1-\tau) \left[\hat{f}(\bar{x}) + \langle A\bar{x}, u \rangle_{\mathbb{E}_2} - \hat{\phi}(u) - \mu_2 d_2(u) \right] \right. \\
 &\quad \left. + \tau [\hat{f}(\hat{x}) + \langle A\hat{x}, u \rangle_{\mathbb{E}_2} - \hat{\phi}(u)] \right\} \\
 &\stackrel{(4.2.18)}{\leq} \max_{u \in Q_2} \left\{ (1-\tau) \left[f_{\mu_2}(\bar{x}) - \frac{1}{2}\mu_2 \|u - u_2\|_2^2 \right] \right. \\
 &\quad \left. + \tau [\phi(\hat{u}) + \langle \nabla\phi(\hat{u}), u - \hat{u} \rangle_{\mathbb{E}_2}] \right\},
 \end{aligned}$$

where we used (6.2.30) in the last line. Since ϕ is concave, by (6.2.31) we obtain

$$f_{\mu_2}(\bar{x}) \leq \phi(\bar{u}) \leq \phi(\hat{u}) + \langle \nabla\phi(\hat{u}), \bar{u} - \hat{u} \rangle_{\mathbb{E}_2}$$

$$\text{Line 1 in (6.2.35)} \stackrel{=}{=} \phi(\hat{u}) + \tau \langle \nabla\phi(\hat{u}), \bar{u} - u_2 \rangle_{\mathbb{E}_2}.$$

Hence, we can finish the proof as follows:

$$\begin{aligned}
 f_{\mu_2^+}(\bar{x}_+) &\leq \max_{u \in Q_2} \left\{ \phi(\hat{u}) + \tau \langle \nabla\phi(\hat{u}), u - u_2 \rangle_{\mathbb{E}_2} - \frac{1}{2}(1-\tau)\mu_2 \|u - u_2\|_2^2 \right\} \\
 &\stackrel{(6.2.36)}{\leq} \max_{u \in Q_2} \left\{ \phi(\hat{u}) + \tau \langle \nabla\phi(\hat{u}), u - u_2 \rangle_{\mathbb{E}_2} - \frac{1}{2}\tau^2 L_2(\phi) \|u - u_2\|_2^2 \right\}.
 \end{aligned}$$

Defining now $v = \bar{u} + \tau(u - \bar{u})$ with $u \in Q_2$, we continue:

$$\begin{aligned}
 f_{\mu_2^+}(\bar{x}_+) &\leq \max_{v \in \bar{u} + \tau(Q_2 - \bar{u})} \left\{ \phi(\hat{u}) + \langle \nabla\phi(\hat{u}), v - \hat{u} \rangle_{\mathbb{E}_2} - \frac{1}{2}L_2(\phi) \|v - \hat{u}\|_2^2 \right\} \\
 (Q_2 \text{ is convex}) &\leq \max_{v \in Q_2} \left\{ \phi(\hat{u}) + \langle \nabla\phi(\hat{u}), v - \hat{u} \rangle_{\mathbb{E}_2} - \frac{1}{2}L_2(\phi) \|v - \hat{u}\|_2^2 \right\} \\
 &\stackrel{(6.2.33)}{\leq} \phi(\hat{u}) + \langle \nabla\phi(\hat{u}), \bar{u}_+ - \hat{u} \rangle_{\mathbb{E}_2} - \frac{1}{2}L_2(\phi) \|\bar{u}_+ - \hat{u}\|_2^2 \\
 &\stackrel{(1.2.11)}{\leq} \phi(\bar{u}_+). \quad \square
 \end{aligned}$$

Now we can justify the following minimization scheme.

1. Initialization:

Set $\mu_{2,0} = 2L_2(\phi)$, $\bar{x}_0 = x_0(u_0)$ and $\bar{u}_0 = V(u_0)$.

2. For $k \geq 0$ iterate:

Set $\tau_k = \frac{2}{k+3}$ and $\hat{u}_k = (1 - \tau_k)\bar{u}_k + \tau_k u_{\mu_{2,k}}(\bar{x}_k)$. (6.2.37)

Update $\mu_{2,k+1} = (1 - \tau_k)\mu_{2,k}$,

$$\bar{x}_{k+1} = (1 - \tau_k)\bar{x}_k + \tau_k x_0(\hat{u}_k),$$

$$\bar{u}_{k+1} = V(\hat{u}_k).$$

Theorem 6.2.4 *Let problem (6.2.1) satisfy Assumption 6.2.1. Then the pairs (\bar{x}_k, \bar{u}_k) generated by scheme (6.2.37) satisfy the following inequality:*

$$f(\bar{x}_k) - \phi(\bar{u}_k) \leq \frac{4L_2(\phi)D_2}{(k+1)(k+2)}, \quad (6.2.38)$$

where $L_2(\phi)$ is given by (6.2.29).

Proof Indeed, in view of Theorem 6.2.3 and Lemma 6.2.8 we need only to justify that the sequences $\{\mu_{2,k}\}_{k=0}^{\infty}$ and $\{\tau_k\}_{k=0}^{\infty}$ satisfy relation (6.2.36). This is straightforward because of the following relation:

$$\mu_{2,k} = \frac{4L_2(\phi)}{(k+1)(k+2)},$$

which is valid for all $k \geq 0$. \square

Let us conclude this section with an example. Consider the problem

$$f(x) = \frac{1}{2}\|x\|_{\mathbb{E}_1}^2 + \max_{1 \leq j \leq m} [f_j + \langle g_j, x - x_j \rangle_{\mathbb{E}_1}] \rightarrow \min : x \in \mathbb{E}_1. \quad (6.2.39)$$

Let $\mathbb{E}_1 = \mathbb{R}^n$ and choose

$$\|x\|_1^2 = \sum_{i=1}^n (x^{(i)})^2, \quad x \in \mathbb{E}_1.$$

Then this problem can be solved by the method (6.2.37).

Indeed, we can represent the objective function in (6.2.39) in the form (6.2.2) using the following objects:

$$\mathbb{E}_2 = \mathbb{R}^m, \quad Q_2 = \Delta_m = \{u \in \mathbb{R}_+^m : \sum_{j=1}^m u^{(j)} = 1\},$$

$$\hat{f}(x) = \frac{1}{2} \|x\|_1^2, \quad \hat{\phi}(u) = \langle b, u \rangle_{\mathbb{E}_2}, \quad b^{(j)} = \langle g_j, x_j \rangle_{\mathbb{E}_1} - f_j, \quad j = 1, \dots, m,$$

$$A^T = (g_1, \dots, g_m).$$

Thus, $\hat{\sigma} = 1$ and $L_2(\hat{\phi}) = 0$. Let us choose for \mathbb{E}_2 the following norm:

$$\|u\|_{\mathbb{E}_2} = \sum_{j=1}^m |u^{(j)}|.$$

Then we can use the entropy distance function,

$$d_2(u) = \ln m + \sum_{j=1}^m u^{(j)} \ln u^{(j)}, \quad u_0 = (\frac{1}{m}, \dots, \frac{1}{m}),$$

for which the convexity parameter is one and $D_2 = \ln m$. Note that in this case

$$\|A\|_{1,2} = \max_{1 \leq j \leq m} \|g_j\|_1^*.$$

Thus, method (6.2.37) as applied to problem (6.2.39) converges with the following rate:

$$f(\bar{x}_k) - \phi(\bar{u}_k) \leq \frac{4 \ln m}{(k+1)(k+2)} \cdot \max_{1 \leq j \leq m} (\|g_j\|_1^*)^2.$$

Let us study the complexity of method (6.2.37) for our example. At each iteration, we need to compute the following objects.

1. **Computation of $u_{\mu_2}(\bar{x})$.** This is the solution of the following problem:

$$\max_u \left\{ \sum_{j=1}^m u^{(j)} s^{(j)}(\bar{x}) - \mu_2 d_2(u) : u \in Q_2 \right\}$$

with $s^{(j)}(\bar{x}) = f_j + \langle g_j, \bar{x} - x_j \rangle$, $j = 1, \dots, m$. As we have seen several times, this solution can be found in a closed form:

$$u_{\mu_2}^{(j)}(\bar{x}) = e^{s^{(j)}(\bar{x})/\mu_2} \cdot \left[\sum_{l=1}^m e^{s^{(l)}(\bar{x})/\mu_2} \right]^{-1}, \quad j = 1, \dots, m.$$

2. **Computation of $x_0(\hat{u})$.** In our case, this is a solution to the problem

$$\min_x \left\{ \langle Ax, \hat{u} \rangle_{\mathbb{E}_2} + \frac{1}{2} \|x\|_{\mathbb{E}_1}^2 : x \in \mathbb{E}_1 \right\}.$$

Hence, the answer is very simple: $x_0(\hat{u}) = -A^T \hat{u}$.

3. **Computation of $V(\hat{u})$.** In our case,

$$\begin{aligned} \phi(\bar{u}) &= \min_{x \in \mathbb{E}_1} \left\{ \sum_{j=1}^m u^{(j)} [f_j + \langle g_j, x - x_j \rangle_{\mathbb{E}_1}] + \frac{1}{2} \|x\|_{\mathbb{E}_1}^2 \right\} \\ &= -\langle b, u \rangle_{\mathbb{E}_2} - \frac{1}{2} \left(\|A^T \hat{u}\|_{\mathbb{E}_1}^* \right)^2. \end{aligned}$$

Thus, $\nabla \phi(\bar{u}) = -b - AA^T \hat{u}$. Now we can compute $V(\hat{u})$ by (6.2.33). It can be easily shown that the complexity of finding $V(\hat{u})$ is of the order $O(m \ln m)$, which comes from the necessity to sort the components of a vector in \mathbb{R}^m .

Thus, we have seen that all computations at each iteration of method (6.2.37) as applied to problem (6.2.39) are very cheap. The most expensive part of the iteration is the multiplication of matrix A by a vector. In a straightforward implementation, we need three such multiplications per iteration. However, a simple modification of the order of operations can reduce this amount to two.

6.3 The Smoothing Technique in Semidefinite Optimization

(Smooth symmetric functions of eigenvalues; Minimizing the maximal eigenvalue of a symmetric matrix.)

6.3.1 Smooth Symmetric Functions of Eigenvalues

In Sects. 6.1 and 6.2, we have shown that a proper use of the structure of nonsmooth convex optimization problems leads to very efficient gradient schemes, whose performance is significantly better than the lower complexity bounds derived from the Black Box assumptions. However, this observation leads to implementable algorithms only if we are able to form a computable smooth approximation of the objective function of our problem. In this case, applying to this approximation an optimal method (6.1.19) for minimizing smooth convex functions, we can easily obtain a good solution to our initial problem.

Our previous results are related mainly to piece-wise linear functions. In this section, we extend them to the problems of Semidefinite Optimization (SO).

For that, we introduce computable smooth approximation for one of the most important nonsmooth functions of symmetric matrices, its *maximal eigenvalue*. Our approximation is based on entropy smoothing.

In what follows, we denote by \mathbb{M}_n the space of real $n \times n$ -matrices, and by $\mathbb{S}_n \subset \mathbb{M}_n$ the space of symmetric matrices. A particular matrix is always denoted by a capital letter. In the spaces \mathbb{R}^n and \mathbb{M}_n we use the standard inner products

$$\langle x, y \rangle = \sum_{i=1}^n x^{(i)} y^{(i)}, \quad x, y \in \mathbb{R}^n,$$

$$\langle X, Y \rangle_F = \sum_{i,j=1}^n X^{(i,j)} Y^{(i,j)}, \quad X, Y \in \mathbb{M}_n.$$

For $X \in \mathbb{S}_n$, we denote by $\lambda(X) \in \mathbb{R}^n$ the vector of its eigenvalues. We assume that the eigenvalues are ordered in a decreasing order:

$$\lambda^{(1)}(X) \geq \lambda^{(2)}(X) \geq \dots \geq \lambda^{(n)}(X), \quad X \in \mathbb{S}_n.$$

Thus, $\lambda_{\max}(X) = \lambda^{(1)}(X)$. The notation $D(\lambda) \in \mathbb{S}_n$ is used for a diagonal matrix with vector $\lambda \in \mathbb{R}^n$ on the main diagonal. Note that any $X \in \mathbb{S}_n$ admits an eigenvalue decomposition

$$X = U(X)D(\lambda(X))U(X)^T$$

with $U(X)U(X)^T = I_n$, where $I_n \in \mathbb{S}_n$ is the identity matrix.

Let us mention some notations with different meanings for vectors and matrices. For a vector $\lambda \in \mathbb{R}^n$, we denote by $|\lambda| \in \mathbb{R}^n$ the vector with entries $|\lambda^{(i)}|$, $i = 1, \dots, n$. The notation $\lambda^k \in \mathbb{R}^n$ is used for the vector with components $(\lambda^{(i)})^k$, $i = 1, \dots, n$. However, for $X \in \mathbb{S}_n$ we define

$$|X| \stackrel{\text{def}}{=} U(X)D(|\lambda(X)|)U(X)^T \geq 0,$$

and the notation X^k is used for the standard matrix power. Since the power $k \geq 0$ does not change the ordering of nonnegative components, for any $X \geq 0$ we have

$$\lambda^k(X) = \lambda(X^k). \quad (6.3.1)$$

Further, in \mathbb{R}^n , we use a standard notation for ℓ_p -norms:

$$\|x\|_{(p)} = \left[\sum_{i=1}^n |x^{(i)}|^p \right]^{1/p}, \quad x \in \mathbb{R}^n,$$

where $p \geq 1$, and $\|x\|_{(\infty)} = \max_{1 \leq i \leq n} |x^{(i)}|$. The corresponding norms in \mathbb{S}_n are introduced by

$$\|X\|_{(p)} = \|\lambda(X)\|_{(p)} = \|\lambda(|X|)\|_{(p)}, \quad X \in \mathbb{S}_n. \quad (6.3.2)$$

For $k \geq 1$, consider the following function:

$$\pi_k(X) = \langle X^k, I_n \rangle_F = \sum_{i=1}^n (\lambda^{(i)}(X))^k, \quad X \in \mathbb{S}_n.$$

Let us derive an upper bound for its second derivative. Note that this bound is nontrivial only for $k \geq 2$.

The derivatives of this function along a direction $H \in \mathbb{S}_n$ are defined as follows:

$$\begin{aligned} \langle \nabla \pi_k(X), H \rangle_F &= k \langle X^{k-1}, H \rangle_F, \\ \langle \nabla^2 \pi_k(X) H, H \rangle_F &= k \sum_{p=0}^{k-2} \langle X^p H X^{k-2-p}, H \rangle_F. \end{aligned} \quad (6.3.3)$$

We need the following result.

Lemma 6.3.1 *For any $p, q \geq 0$, and X, H from \mathbb{S}_n we have*

$$\begin{aligned} \langle X^p H X^q + X^q H X^p, H \rangle_F &\leq 2 \langle |X|^{p+q}, H^2 \rangle_F \\ &\leq 2 \langle \lambda^{p+q}(|X|), \lambda^2(|H|) \rangle. \end{aligned} \quad (6.3.4)$$

Proof Indeed, let $\lambda = \lambda(X)$, $D = D(\lambda)$, $U = U(X)$ and $\hat{H} = U^T H U$. Then

$$\begin{aligned} \langle X^p H X^q + X^q H X^p, H \rangle_F &= \langle U D^p U^T H U D^q U^T + U D^q U^T H U D^p U^T, H \rangle_F \\ &= \langle D^p \hat{H} D^q + D^q \hat{H} D^p, \hat{H} \rangle_F \\ &= \sum_{i,j=1}^n (\hat{H}^{(i,j)})^2 \left((\lambda^{(i)})^p (\lambda^{(j)})^q + (\lambda^{(i)})^q (\lambda^{(j)})^p \right) \\ &\leq \sum_{i,j=1}^n (\hat{H}^{(i,j)})^2 \left(|\lambda^{(i)}|^p |\lambda^{(j)}|^q + |\lambda^{(i)}|^q |\lambda^{(j)}|^p \right). \end{aligned}$$

Note that for arbitrary non-negative values a and b we always have

$$0 \leq (a^p - b^p)(a^q - b^q) = (a^{p+q} + b^{p+q}) - (a^p b^q + a^q b^p).$$

Thus, we can continue as follows:

$$\begin{aligned}
 \langle X^p H X^q + X^q H X^p, H \rangle_F &\leq \sum_{i,j=1}^n (\hat{H}^{(i,j)})^2 (|\lambda^{(i)}|^{p+q} + |\lambda^{(j)}|^{p+q}) \\
 &= 2 \sum_{i,j=1}^n (\hat{H}^{(i,j)})^2 |\lambda^{(i)}|^{p+q} = 2 \langle D(|\lambda|)^{p+q} \hat{H}, \hat{H} \rangle_F \\
 &= 2 \langle D^{p+q}(|\lambda|), \hat{H}^2 \rangle_F = 2 \langle |X|^{p+q}, H^2 \rangle_F.
 \end{aligned}$$

Hence, we get the first inequality in (6.3.4). Further, by von Neumann's inequality

$$\langle |X|^{p+q}, H^2 \rangle_F \leq \langle \lambda(|X|^{p+q}), \lambda(H^2) \rangle \stackrel{(6.3.1)}{=} \langle \lambda^{p+q}(|X|), \lambda^2(|H|) \rangle,$$

and this proves the remaining part of (6.3.4). \square

Corollary 6.3.1 *For any $k \geq 2$, we have*

$$\langle \nabla^2 \pi_k(X) H, H \rangle_F \leq k(k-1) \langle \lambda^{k-2}(|X|), \lambda^2(|H|) \rangle. \quad (6.3.5)$$

Proof For $k = 2$, the bound is trivial. For $k \geq 3$, in representation (6.3.3) we can unify the terms in the expression $\sum_{p=0}^{k-2} \langle X^p H X^{k-2-p}, H \rangle_F$ in symmetric pairs

$$\langle X^p H X^{k-2-p} + X^{k-2-p} H X^p, H \rangle_F.$$

Applying inequality (6.3.4) to each pair, we get the estimate (6.3.5). \square

Let $f(\cdot)$ be a function of a real variable, defined by a power series

$$f(\tau) = a_0 + \sum_{k=1}^{\infty} a_k \tau^k$$

with $a_k \geq 0$ for $k \geq 2$. We assume that its domain $\text{dom } f = \{\tau : |\tau| < R\}$ is nonempty. For $X \in \mathbb{S}_n$, consider the following symmetric function of eigenvalues:

$$F(X) = \sum_{i=1}^n f(\lambda^{(i)}(X)).$$

Clearly, $\text{dom } F = \{X \in \mathbb{S}_n : \lambda^{(1)}(X) < R, \lambda^{(n)}(X) > -R\}$.

Theorem 6.3.1 *For any $X \in \text{dom } F$ and $H \in \mathbb{S}_n$ we have*

$$\langle \nabla^2 F(X) H, H \rangle \leq \sum_{i=1}^n \nabla^2 f(\lambda^{(i)}(|X|)) (\lambda^{(i)}(|H|))^2.$$

Proof Indeed,

$$\begin{aligned} F(X) &= n \cdot a_0 + \sum_{i=1}^n \sum_{k=1}^{\infty} a_k (\lambda^{(i)}(X))^k \\ &= n \cdot a_0 + \sum_{k=1}^{\infty} a_k \sum_{i=1}^n (\lambda^{(i)}(X))^k = n \cdot a_0 + \sum_{k=1}^{\infty} a_k \pi_k(X). \end{aligned}$$

Thus, in view of inequality (6.3.5),

$$\begin{aligned} \langle \nabla^2 F(X)H, H \rangle_F &= \sum_{k=2}^{\infty} a_k \langle \nabla^2 \pi_k(X)H, H \rangle_F \\ &\leq \sum_{k=2}^{\infty} k(k-1) a_k \langle \lambda^{k-2}(|X|), \lambda^2(|H|) \rangle \\ &= \sum_{i=1}^n \sum_{k=2}^{\infty} k(k-1) a_k (\lambda^{(i)}(|X|))^{k-2} (\lambda^{(i)}(|H|))^2 \\ &= \sum_{i=1}^n \nabla^2 f(\lambda^{(i)}(|X|)) (\lambda^{(i)}(|H|))^2. \quad \square \end{aligned}$$

Let us consider now two important examples of symmetric functions of eigenvalues.

1. Squared ℓ_p -Matrix Norm. For an integer $p \geq 1$, consider the following function:

$$F_p(X) = \frac{1}{2} \|\lambda(X)\|_{(2p)}^2 = \frac{1}{2} \langle X^{2p}, I_n \rangle_F^{1/p}, \quad X \in \mathbb{S}_n. \quad (6.3.6)$$

Thus, $F_p(X) = \frac{1}{2} (\pi_{2p}(X))^{1/p}$. Therefore, in view of (6.3.5), for any $X, H \in \mathbb{S}_n$ we have

$$\begin{aligned} \langle \nabla F_p(X), H \rangle_F &= \frac{1}{2p} (\pi_{2p}(X))^{\frac{1}{p}-1} \langle \nabla \pi_{2p}(X), H \rangle_F, \\ \langle \nabla^2 F_p(X)H, H \rangle_F &= \frac{1}{2p} \cdot \left(\frac{1}{p} - 1 \right) \cdot (\pi_{2p}(X))^{\frac{1}{p}-2} \langle \nabla \pi_{2p}(X), H \rangle_F^2 \\ &\quad + \frac{1}{2p} (\pi_{2p}(X))^{\frac{1}{p}-1} \langle \nabla^2 \pi_{2p}(X)H, H \rangle_F \quad (6.3.7) \\ &\leq (2p-1) (\pi_{2p}(X))^{\frac{1}{p}-1} \langle \lambda^{2p-2}(|X|), \lambda^2(|H|) \rangle. \end{aligned}$$

Let us apply Hölder's inequality $\langle x, y \rangle \leq \|x\|_{(\beta)} \|y\|_{(\gamma)}$ with $\beta = \frac{p}{p-1}$, $\gamma = \frac{\beta}{\beta-1} = p$, and

$$x^{(i)} = (\lambda^{(i)}(|X|))^{2p-2}, \quad y^{(i)} = (\lambda^{(i)}(|H|))^2, \quad i = 1, \dots, n.$$

Then,

$$\begin{aligned} \langle x, y \rangle &\leq \left[\sum_{i=1}^n (\lambda^{(i)}(|X|))^{2p} \right]^{\frac{p-1}{p}} \cdot \left[\sum_{i=1}^n (\lambda^{(i)}(|H|))^{2p} \right]^{\frac{1}{p}} \\ &\stackrel{(6.3.2)}{=} \pi_{2p}(X)^{\frac{p-1}{p}} \cdot \|\lambda(H)\|_{(2p)}^2, \end{aligned}$$

and we can continue:

$$\langle \nabla^2 F_p(X) H, H \rangle_F \leq (2p-1) \|\lambda(H)\|_{(2p)}^2 = (2p-1) \|H\|_{(2p)}^2. \quad (6.3.8)$$

2. Entropy Smoothing of Maximal Eigenvalue. Consider the function

$$E(X) = \ln \sum_{i=1}^n e^{\lambda^{(i)}(X)} \stackrel{\text{def}}{=} \ln F(X), \quad X \in \mathbb{S}_n. \quad (6.3.9)$$

Note that

$$\begin{aligned} \langle \nabla E(X), H \rangle_F &= \frac{1}{F(X)} \langle \nabla F(X), H \rangle_F, \\ \langle \nabla^2 E(X) H, H \rangle_F &= -\frac{1}{F^2(X)} \langle \nabla F(X), H \rangle_F^2 + \frac{1}{F(X)} \langle \nabla^2 F(X) H, H \rangle_F \\ &\leq \frac{1}{F(X)} \langle \nabla^2 F(X) H, H \rangle_F. \end{aligned}$$

Let us assume first that $X \succeq 0$. The function $F(X)$ is formed by the auxiliary function $f(\tau) = e^\tau$, which satisfies the assumptions of Theorem 6.3.1. Therefore,

$$\langle \nabla^2 E(X) H, H \rangle_F \leq \left[\sum_{i=1}^n e^{\lambda^{(i)}(X)} \right]^{-1} \sum_{i=1}^n e^{\lambda^{(i)}(X)} (\lambda^{(i)}(|H|))^2 \leq \|H\|_{(\infty)}^2. \quad (6.3.10)$$

It remains to note that $E(X + \tau I_n) = E(X) + \tau$. Hence, the Hessian $\nabla^2 E(X + \tau I_n)$ does not depend on τ , and we conclude that the estimate (6.3.10) is valid for arbitrary $X \in \mathbb{S}_n$.

6.3.2 Minimizing the Maximal Eigenvalue of the Symmetric Matrix

Consider the following problem:

$$\text{Find } \phi^* = \min_{y \in Q} \{\phi(y) \stackrel{\text{def}}{=} \lambda_{\max}(C + A(y))\}, \quad (6.3.11)$$

where Q is a closed convex set in \mathbb{R}^m and $A(\cdot)$ is a linear operator from \mathbb{R}^m to \mathbb{S}_n :

$$A(y) = \sum_{i=1}^m y^{(i)} A_i \in \mathbb{S}_n, \quad y \in \mathbb{R}^m.$$

Note that the objective function in (6.3.11) is nonsmooth. Therefore, this problem can be solved either by interior-point methods (see Chap. 5), or by general methods of nonsmooth convex optimization (see Chap. 3). However, due to the very special structure of the objective function, for problem (6.3.11) it is better to develop a special scheme.

We are going to solve problem (6.3.11) by a smoothing technique discussed in Sect. 6.1. This means that we replace the function $\lambda_{\max}(X)$ by its smooth approximation $f_\mu(X) = \mu E(\frac{1}{\mu}X)$, defined by (6.3.9) with tolerance parameter $\mu > 0$. Note that

$$f_\mu(X) = \mu \ln \left[\sum_{i=1}^n e^{\lambda^{(i)}(X)/\mu} \right] \geq \lambda_{\max}(X), \quad (6.3.12)$$

$$f_\mu(X) \leq \lambda_{\max}(X) + \mu \ln n.$$

At the same time,

$$\nabla f_\mu(X) = \left[\sum_{i=1}^n e^{\lambda^{(i)}(X)/\mu} \right]^{-1} \cdot \sum_{i=1}^n e^{\lambda^{(i)}(X)/\mu} u_i(X) u_i(X)^T, \quad (6.3.13)$$

where $u_i(X)$, $i = 1, \dots, n$, are corresponding unit eigenvectors of the symmetric matrix X . Thus, at each test point X , the gradient $\nabla f_\mu(X)$ takes into account all eigenvalues of the matrix X . However, since the factors $e^{\lambda^{(i)}(X)/\mu}$ decrease very rapidly, it actually depends only on few largest eigenvalues. Their selection is made automatically by expression (6.3.13). The ranking of importance of the eigenvalues is done in a logarithmic scale controlled by the tolerance parameter μ .

Let us analyze now the efficiency of the smoothing technique as applied to problem (6.3.11). Our goal is to find an ϵ -solution $\bar{x} \in Q$ to problem (6.3.11):

$$\phi(\bar{y}) - \phi^* \leq \epsilon. \quad (6.3.14)$$

For that, we will try to find a $\frac{1}{2}\epsilon$ -solution to the smooth problem

$$\text{Find } \phi_\mu^* = \min_{y \in Q} \{\phi_\mu(y) \stackrel{\text{def}}{=} f_\mu(C + A(y))\}, \quad (6.3.15)$$

with

$$\mu = \mu(\epsilon) = \frac{\epsilon}{2 \ln n}. \quad (6.3.16)$$

Clearly, if $\phi_\mu(\bar{y}) - \phi_\mu^* \leq \frac{1}{2}\epsilon$, then in view of (6.3.12) we have

$$\phi(\bar{y}) - \phi^* \leq \phi_\mu(\bar{y}) - \phi_\mu^* + \mu \ln n \leq \epsilon.$$

Let us analyze now the complexity of finding a $\frac{1}{2}\epsilon$ -solution to problem (6.3.15) by the optimal method (6.1.19).

Let us fix some norm $\|h\|$ for $h \in \mathbb{R}^m$. Consider a prox-function $d(\cdot)$ of the set Q with prox-center $x_0 \in Q$. We assume this function to be strongly convex on Q with convexity parameter one. Define

$$\|A\| = \max_{h \in \mathbb{R}^m} \{\|A(h)\|_{(\infty)} : \|h\| = 1\}.$$

Note that this norm is quite small. Indeed,

$$\|A(h)\|_{(\infty)} = \lambda^{(1)}(|A(h)|) \leq \langle A(h), A(h) \rangle_F^{1/2}, \quad h \in \mathbb{R}^m.$$

Therefore, for example, $\|A\| \leq \|A\|_G \stackrel{\text{def}}{=} \max_{\|h\|=1} \langle A(h), A(h) \rangle_F^{1/2}$.

Let us estimate the second derivative of the function $\phi_\mu(\cdot)$. For any y and h from \mathbb{R}^m , in view of inequality (6.3.10) we have

$$\langle \nabla \phi_\mu(y), h \rangle = \langle \nabla f_\mu(C + A(y)), h \rangle = \langle \nabla E(\frac{1}{\mu}(C + A(y))), A(h) \rangle_F,$$

$$\begin{aligned} \langle \nabla^2 \phi_\mu(y) h, h \rangle &= \frac{1}{\mu} \langle \nabla^2 E(C + A(y)) A(h), A(h) \rangle_F \\ &\leq \frac{1}{\mu} \|A(h)\|_{(\infty)}^2 \leq \frac{1}{\mu} \|A\|^2 \cdot \|h\|^2. \end{aligned}$$

Thus, by Theorem 6.1.1 the function ϕ_μ has Lipschitz continuous gradient with the constant

$$L = \frac{1}{\mu} \|A\|^2 = \frac{2 \ln n}{\epsilon} \|A\|^2.$$

Now taking into account the estimate (6.1.21), we conclude that the method (6.1.19), as applied to problem (6.3.15), has the following rate of convergence:

$$\phi_\mu(y_k) - \phi_\mu^* \leq \frac{8 \ln n \|A\|^2 d(y_\mu^*)}{\epsilon \cdot (k+1)(k+2)},$$

where $y_\mu^* \in Q$ is the solution to (6.3.15). Hence, it is able to generate a $\frac{1}{2}\epsilon$ -solution to this problem (which is an ϵ -solution to problem (6.3.11)) at most after

$$\frac{4\|A\|}{\epsilon} \sqrt{d(y_\mu^*) \ln n} \quad (6.3.17)$$

iterations.

6.4 Minimizing the Local Model of an Objective Function

(A linear optimization oracle; The method of conditional gradients; Conditional gradients with contraction; Computation of primal-dual solution; Strong convexity of the composite term; The second-order trust-region method with contraction.)

6.4.1 A Linear Optimization Oracle

In this section we consider numerical methods for solving the following *composite* minimization problem:

$$\min_x \left\{ \tilde{f}(x) \stackrel{\text{def}}{=} f(x) + \Psi(x) \right\}, \quad (6.4.1)$$

where Ψ is a *simple* closed convex function with bounded domain $Q \subset \mathbb{E}$, and f is a convex function, which is differentiable on Q . Denote by x^* one of the optimal solutions of (6.4.1), and $D \stackrel{\text{def}}{=} \text{diam}(Q)$. As usual, our assumption on the simplicity of the function Ψ means that some auxiliary optimization problems related to Ψ are easily solvable. The complexity of these problems will be always discussed for corresponding optimization schemes.

The most important examples of the function Ψ are as follows.

- Ψ is an indicator function of a closed convex set Q :

$$\Psi(x) = \text{Ind}_Q(x) \stackrel{\text{def}}{=} \begin{cases} 0, & x \in Q, \\ +\infty, & \text{otherwise.} \end{cases} \quad (6.4.2)$$

- Ψ is a self-concordant barrier for a closed convex set Q (see Sect. 5.3).
- Ψ is a nonsmooth convex function with simple structure. In this case, we need to include in Ψ an indicator function for a bounded domain. For example, it

could be

$$\Psi(x) = \begin{cases} \|x\|_{(1)}, & \text{if } \|x\|_{(1)} \leq R, \\ +\infty, & \text{otherwise.} \end{cases}$$

We assume that the function f is represented by a Black-Box oracle. If it is a *first-order oracle*, we assume its gradients satisfy the following *Hölder condition*:

$$\|\nabla f(x) - \nabla f(y)\|_* \leq G_\nu \|x - y\|^\nu, \quad x, y \in Q. \quad (6.4.3)$$

The constant G_ν is formally defined for any $\nu \in (0, 1]$. For some values of ν it can be $+\infty$. Note that for any x and y in Q we have

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{G_\nu}{1+\nu} \|y - x\|^{1+\nu}. \quad (6.4.4)$$

If this is a *second-order oracle*, we assume that its Hessians satisfy the Hölder condition

$$\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq H_\nu \|x - y\|^\nu, \quad x, y \in Q. \quad (6.4.5)$$

In this case, for any x and y in Q we have

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{1}{2} \langle \nabla^2 f(x)(y - x), y - x \rangle + \frac{H_\nu \|y - x\|^{2+\nu}}{(1+\nu)(2+\nu)}. \quad (6.4.6)$$

Our assumption on the simplicity of the function Ψ means exactly the following.

Assumption 6.4.1 For any $s \in \mathbb{E}^*$, the auxiliary problem

$$\min_{x \in Q} \{\langle s, x \rangle + \Psi(x)\} \quad (6.4.7)$$

is easily solvable. Denote by $v_\Psi(s) \in Q$ one of its optimal solutions.

Thus, for our methods we assume that we can use a *linear optimization oracle*, related to the set Q . Indeed, in the case (6.4.2), this assumption implies that we are able to solve the problem

$$\min_x \{\langle s, x \rangle : x \in Q\}.$$

For some sets (e.g. convex hulls of finite number of points), this oracle has lower complexity than the standard auxiliary problem consisting in minimizing a prox-function plus a linear term (see, for example, Sect. 6.1.3).

In view of Theorem 3.1.23 the point $v_\Psi(s)$ is characterized by the following variational principle:

$$\langle s, x - v_\Psi(s) \rangle + \Psi(x) \geq \Psi(v_\Psi(s)), \quad x \in Q. \quad (6.4.8)$$

By Definition 3.1.5, this means that $-s \in \partial\Psi(v_\Psi(s))$.

In the sequel, we often need to estimate the partial sums of different series. For that, it is convenient to use the following lemma, the proof of which we leave as an exercise for the reader.

Lemma 6.4.1 *Let the function $\xi(\tau)$, $\tau \in \mathbb{R}$, be decreasing and convex. Then, for any two integers a and b , such that $[a - \frac{1}{2}, b + 1] \subset \text{dom } \xi$, we have*

$$\int_a^{b+1} \xi(\tau) d\tau \leq \sum_{k=a}^b \xi(k) \leq \int_{a-1/2}^{b+1/2} \xi(\tau) d\tau. \quad (6.4.9)$$

For example, for any $t \geq 0$ and $p \geq -t$, we have

$$\begin{aligned} \sum_{k=t}^{2t+p} \frac{1}{k+p+1} &\stackrel{(5.4.38)}{\geq} \int_t^{2t+p+1} \frac{1}{\tau+p+1} d\tau = \ln(\tau+p+1) \Big|_t^{2t+p+1} \\ &= \ln \frac{2t+2p+2}{t+p+1} = \ln 2. \end{aligned} \quad (6.4.10)$$

On the other hand, if $t \geq 1$, then

$$\begin{aligned} \sum_{k=t}^{2t+1} \frac{1}{(k+2)^2} &\stackrel{(5.4.38)}{\leq} \int_{t-1/2}^{2t+3/2} \frac{1}{(\tau+2)^2} d\tau = -\frac{1}{\tau+2} \Big|_{t-1/2}^{2t+3/2} = \frac{1}{t+3/2} - \frac{1}{2t+7/2} \\ &= \frac{4t+8}{(2t+3)(4t+7)} \leq \frac{12}{11(2t+3)}. \end{aligned} \quad (6.4.11)$$

6.4.2 The Method of Conditional Gradients with Composite Objective

In order to solve problem (6.4.1), we apply the following method.

Conditional Gradients with Composite Objective	
1. Choose an arbitrary point $x_0 \in Q$.	(6.4.12)
2. For $t \geq 0$ iterate: (a) Compute $v_t = v_\psi(\nabla f(x_t))$.	
(b) Choose $\tau_t \in (0, 1]$ and set $x_{t+1} = (1 - \tau_t)x_t + \tau_t v_t$.	

It is clear that this method can solve only problems where the function f has continuous gradient.

Example 6.4.1 Let $\Psi(x) = \text{Ind}_Q(x)$ with $Q = \{x \in \mathbb{R}^2 : (x^{(1)})^2 + (x^{(2)})^2 \leq 1\}$. Define

$$f(x) = \max\{x^{(1)}, x^{(2)}\}.$$

Then clearly $x_* = \left(\frac{-1}{\sqrt{2}}, \frac{-1}{\sqrt{2}}\right)^T$. Let us choose in (6.4.12) $x_0 \neq x_*$.

For the function f , we can apply an oracle which returns at any $x \in Q$ a subgradient $\nabla f(x) \in \{(1, 0)^T, (0, 1)^T\}$. Then, for any feasible x , the point $v_\Psi(\nabla f(x))$ is equal either to $y_1 = (-1, 0)^T$, or to $y_2 = (0, -1)^T$. Therefore, all points of the sequence $\{x_t\}_{t \geq 0}$, generated by method (6.4.12), belong to the triangle $\text{Conv}\{x_0, y_1, y_2\}$, which does not contain the optimal point x_* . \square

In order to justify the rate of convergence of method (6.4.12) for functions with Hölder continuous gradients, we apply a variant of the estimating sequences technique (see Sects. 2.2.1 and 6.1.3). For that, it is convenient to introduce in (6.4.12) new control variables. Consider a sequence of nonnegative weights $\{a_t\}_{t \geq 0}$. Define

$$A_t = \sum_{k=0}^t a_k, \quad \tau_t = \frac{a_{t+1}}{A_{t+1}}, \quad t \geq 0. \quad (6.4.13)$$

From now on, we assume that the parameter τ_t in method (6.4.12) is chosen in accordance with the rule (6.4.13). Define

$$\begin{aligned} V_0 &= \max_x \{\langle \nabla f(x_0), x_0 - x \rangle + \Psi(x_0) - \Psi(x)\}, \\ B_{v,t} &= a_0 V_0 + \left(\sum_{k=1}^t \frac{a_k^{1+v}}{A_k^v} \right) G_v D^{1+v}, \quad t \geq 0. \end{aligned} \quad (6.4.14)$$

It is clear that

$$\begin{aligned} V_0 &\stackrel{(6.4.6)}{\leq} \max_x \left\{ f(x_0) - f(x) + \frac{G_v}{1+v} \|x - x_0\|^{1+v} + \Psi(x_0) - \Psi(x) \right\} \\ &\leq \bar{f}(x_0) - \bar{f}(x_*) + \frac{G_v D^{1+v}}{1+v} \stackrel{\text{def}}{=} \Delta(x_0) + \frac{G_v D^{1+v}}{1+v}. \end{aligned} \quad (6.4.15)$$

Theorem 6.4.1 *Let the sequence $\{x_t\}_{t \geq 0}$ be generated by method (6.4.12). Then, for any $v \in (0, 1]$ with $G_v < +\infty$, any step $t \geq 0$, and any $x \in Q$ we have*

$$A_t(f(x_t) + \Psi(x_t)) \leq \sum_{k=0}^t a_k [f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \Psi(x)] + B_{v,t}. \quad (6.4.16)$$

Proof Indeed, in view of definition (6.4.14), for $t = 0$ inequality (6.4.16) is satisfied. Assume that it is valid for some $t \geq 0$. Then

$$\begin{aligned}
 & \sum_{k=0}^{t+1} a_k [f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \Psi(x)] + B_{v,t} \\
 & \stackrel{(6.4.16)}{\geq} A_t(f(x_t) + \Psi(x_t)) + a_{t+1}[f(x_{t+1}) + \langle \nabla f(x_{t+1}), x - x_{t+1} \rangle + \Psi(x)] \\
 & \geq A_{t+1}f(x_{t+1}) + A_t\Psi(x_t) + \langle \nabla f(x_{t+1}), a_{t+1}(x - x_{t+1}) + A_t(x_t - x_{t+1}) \rangle \\
 & \quad + a_{t+1}\Psi(x) \\
 & \stackrel{(6.4.12)_b}{=} A_{t+1}f(x_{t+1}) + A_t\Psi(x_t) + a_{t+1}[\Psi(x) + \langle \nabla f(x_{t+1}), x - v_t \rangle] \\
 & \stackrel{(6.4.12)_b}{\geq} A_{t+1}(f(x_{t+1}) + \Psi(x_{t+1})) + a_{t+1}[\Psi(x) - \Psi(v_t) + \langle \nabla f(x_{t+1}), x - v_t \rangle].
 \end{aligned}$$

It remains to note that

$$\begin{aligned}
 \Psi(x) - \Psi(v_t) + \langle \nabla f(x_{t+1}), x - v_t \rangle & \stackrel{(6.4.8)}{\geq} \langle \nabla f(x_{t+1}) - \nabla f(x_t), x - v_t \rangle \\
 & \stackrel{(6.4.3)}{\geq} -\tau_t^v G_v D^{1+v}.
 \end{aligned}$$

Thus, to ensure that (6.4.16) is valid for the next iteration, it is enough to choose

$$B_{v,t+1} = B_{v,t} + \frac{a_{t+1}^{1+v}}{A_{t+1}^v} G_v D^{1+v}. \quad \square$$

Corollary 6.4.1 *For any $t \geq 0$ with $A_t > 0$, and any $v \in (0, 1]$, we have*

$$\bar{f}(x_t) - \bar{f}(x_*) \leq \frac{1}{A_t} B_{v,t}. \quad (6.4.17)$$

Let us discuss now the possible variants for choosing the weights $\{a_t\}_{t \geq 0}$.

1. *Constant weights.* Let us choose $a_t \equiv 1$, $t \geq 0$. Then $A_t = t + 1$, and for $v \in (0, 1)$ we have

$$\begin{aligned}
 B_{v,t} &= V_0 + \left(\sum_{k=1}^t \frac{1}{(1+k)^v} \right) G_v D^{1+v} \\
 & \stackrel{(6.4.9)}{\leq} V_0 + G_v D^{1+v} \frac{1}{1-v} (1 + \tau)^{1-v} \Big|_{1/2}^{t+1/2} \\
 & \stackrel{(6.4.15)}{\leq} \Delta(x_0) + G_v D^{1+v} \left[\frac{1}{1+v} + \left(\frac{3}{2} \right)^{1-v} \frac{1}{1-v} \left(\left(1 + \frac{2}{3}t \right)^{1-v} - 1 \right) \right].
 \end{aligned}$$

Thus, for $\nu \in (0, 1)$, we have $\frac{1}{A_t} B_{\nu,t} \leq O(t^{-\nu})$. For the most important case $\nu = 1$, we have $\lim_{\nu \rightarrow 1} \frac{1}{1-\nu} \left(\left(1 + \frac{2}{3}t\right)^{1-\nu} - 1 \right) = \ln(1 + \frac{2}{3}t)$. Therefore,

$$\bar{f}(x_t) - \bar{f}(x_*) \leq \frac{1}{t+1} \left(\Delta(x_0) + G_1 D^2 \left[\frac{1}{2} + \ln(1 + \frac{2}{3}t) \right] \right). \quad (6.4.18)$$

In this situation, in method (6.4.12) we take $\tau_t \stackrel{(6.4.13)}{=} \frac{1}{t+1}$.

2. *Linear weights.* Let us choose $a_t \equiv t, t \geq 0$. Then $A_t = \frac{t(t+1)}{2}$, and for $\nu \in (0, 1)$ with $t \geq 1$ we have

$$\begin{aligned} B_{\nu,t} &= \left(\sum_{k=1}^t \frac{2^{\nu} k^{1+\nu}}{k^{\nu}(1+k)^{\nu}} \right) G_{\nu} D^{1+\nu} \leq \left(\sum_{k=1}^t 2^{\nu} k^{1-\nu} \right) G_{\nu} D^{1+\nu} \\ &\stackrel{(6.4.9)}{\leq} G_{\nu} D^{1+\nu} \frac{2^{\nu}}{2^{-\nu}} \tau^{2-\nu} \Big|_{1/2}^{t+1/2} = \frac{2^{\nu}}{2^{-\nu}} \left[\left(t + \frac{1}{2}\right)^{2-\nu} - \left(\frac{1}{2}\right)^{2-\nu} \right] G_{\nu} D^{1+\nu}. \end{aligned}$$

Thus, for $\nu \in (0, 1)$, we again have $\frac{1}{A_t} B_{\nu,t} \leq O(t^{-\nu})$. For the case $\nu = 1$, we get the following bound:

$$\bar{f}(x_t) - \bar{f}(x_*) \leq \frac{4}{t+1} G_1 D^2, \quad t \geq 1. \quad (6.4.19)$$

As we can see, this rate of convergence is better than (6.4.18). In this case, in method (6.4.12) we take $\tau_t \stackrel{(6.4.13)}{=} \frac{2}{t+2}$, which is a standard recommendation for this scheme.

3. *Aggressive weights.* Let us choose, for example, $a_t \equiv t^2, t \geq 0$. Then $A_t = \frac{t(t+1)(2t+1)}{6}$. Note that for $k \geq 0$ we have $\frac{k^{2+\nu}}{(k+1)^{\nu}(2k+1)^{\nu}} \leq \frac{k^{2-\nu}}{2^{\nu}}$. Therefore, for $\nu \in (0, 1)$ with $t \geq 1$ we obtain

$$\begin{aligned} B_{\nu,t} &= \left(\sum_{k=1}^t \frac{6^{\nu} k^{2(1+\nu)}}{k^{\nu}(1+k)^{\nu}(2k+1)^{\nu}} \right) G_{\nu} D^{1+\nu} \leq \left(\sum_{k=1}^t 3^{\nu} k^{2-\nu} \right) G_{\nu} D^{1+\nu} \\ &\stackrel{(6.4.9)}{\leq} G_{\nu} D^{1+\nu} \frac{3^{\nu}}{3^{-\nu}} \tau^{3-\nu} \Big|_{1/2}^{t+1/2} = \frac{3^{\nu}}{3^{-\nu}} \left[\left(t + \frac{1}{2}\right)^{3-\nu} - \left(\frac{1}{2}\right)^{3-\nu} \right] G_{\nu} D^{1+\nu}. \end{aligned}$$

For $\nu \in (0, 1)$, we get again $\frac{1}{A_t} B_{\nu,t} \leq O(t^{-\nu})$. For $\nu = 1$, we obtain

$$\bar{f}(x_t) - \bar{f}(x_*) \leq \frac{9}{2t+1} G_1 D^2, \quad t \geq 1, \quad (6.4.20)$$

which is slightly worse than (6.4.19). The rule for choosing the coefficients τ_t in this situation is $\tau_t \stackrel{(6.4.13)}{=} \frac{6(t+1)}{(t+2)(2t+3)}$. It can be easily checked that a further increase of the rate of growth of coefficients a_t makes the rate of convergence of method (6.4.12) even worse.

Note that the above rules for choosing the coefficients $\{\tau_t\}_{t \geq 0}$ in method (6.4.12) do not depend on the smoothness parameter $\nu \in (0, 1]$. In this sense, method (6.4.12) is a *universal method* for solving the problem (6.4.1). Moreover, this method is *affine invariant*. Its behavior does not depend on the choice of norm in \mathbb{E} . Hence, its rate of convergence can be established with respect to the best norm describing the geometry of the feasible set.

6.4.3 Conditional Gradients with Contraction

In this section, we will use some special dual functions. Let $Q \subset E$ be a bounded closed convex set. For a closed convex function $F(\cdot)$ with $\text{dom } F \supseteq \text{int } Q$, we define its *restricted dual function*, (with respect to a central point $\bar{x} \in Q$), as follows:

$$F_{\bar{x}, Q}^*(s) = \max_{x \in Q} \{ \langle s, \bar{x} - x \rangle + F(\bar{x}) - F(x) \}, \quad s \in \mathbb{E}^*. \quad (6.4.21)$$

Clearly, this function is well defined for all $s \in \mathbb{E}^*$. Moreover, it is convex and nonnegative on \mathbb{E}^* .

We need to introduce in construction (6.4.21) an additional scaling parameter $\tau \in [0, 1]$, which controls the size of the feasible set. For $s \in \mathbb{E}^*$, we call the function

$$F_{\tau, \bar{x}, Q}^*(s) = \max_{x \in Q} \{ \langle s, \bar{x} - y \rangle + F(\bar{x}) - F(y) : y = (1 - \tau)\bar{x} + \tau x \} \quad (6.4.22)$$

the *scaled restricted dual* of the function F .

Lemma 6.4.2 For any $s \in \mathbb{E}^*$ and $\tau \in [0, 1]$, we have

$$F_{\bar{x}, Q}^*(s) \geq F_{\tau, \bar{x}, Q}^*(s) \geq \tau F_{\bar{x}, Q}^*(s). \quad (6.4.23)$$

Proof Since for any $x \in Q$, the point $y = (1 - \tau)\bar{x} + \tau x$ belongs to Q , the first inequality is trivial. On the other hand,

$$\begin{aligned} F_{\tau, \bar{x}, Q}^*(s) &= \max_{x \in Q} \{ \langle s, \tau(\bar{x} - x) \rangle + F(\bar{x}) - F(y) : y = (1 - \tau)\bar{x} + \tau x \} \\ &\geq \max_{x \in Q} \{ \langle s, \tau(\bar{x} - x) \rangle + F(\bar{x}) - (1 - \tau)F(\bar{x}) - \tau F(x) \} \\ &= \tau F_{\bar{x}, Q}^*(s). \end{aligned} \quad \square$$

Let us consider a variant of method (6.4.12), which takes into account the composite form of the objective function in problem (6.4.1). For $\Psi(x) \equiv \text{Ind}_Q(x)$, these

two methods coincide. Otherwise, they generate different minimization sequences.

Conditional Gradient Method with Contraction
<p>1. Choose an arbitrary point $x_0 \in Q$.</p> <p>2. For $t \geq 0$ iterate: Choose a coefficient $\tau_t \in (0, 1]$ and compute</p> $x_{t+1} = \arg \min_{x \in Q} \{ \langle \nabla f(x_t), y \rangle + \Psi(y) : y = (1 - \tau_t)x_t + \tau_t x \}.$

(6.4.24)

This method can be seen as a *Trust-Region Scheme* with a linear model of the objective function. The trust region in method (6.4.24) is formed by a contraction of the initial feasible set. In Sect. 6.4.6, we will consider a more traditional trust-region method with quadratic model of the objective.

In view of Theorem 3.1.23 the point x_{t+1} in method (6.4.24) is characterized by the following variational principle:

$$\begin{aligned}
 x_{t+1} &= (1 - \tau_t)x_t + \tau_t v_t, \quad v_t \in Q, \\
 \Psi((1 - \tau_t)x_t + \tau_t x) + \tau_t \langle \nabla f(x_t), x - x_t \rangle & \\
 &\geq \Psi(x_{t+1}) + \langle \nabla f(x_t), x_{t+1} - x_t \rangle, \quad x \in Q.
 \end{aligned}
 \tag{6.4.25}$$

Let us choose somehow the sequence of nonnegative weights $\{a_t\}_{t \geq 0}$, and define in (6.4.24) the coefficients τ_t in accordance to (6.4.13). Define now the estimating functional sequence $\{\phi_t(x)\}_{t \geq 0}$ as follows:

$$\begin{aligned}
 \phi_0(x) &= a_0 \bar{f}(x), \\
 \phi_{t+1}(x) &= \phi_t(x) + a_{t+1} [f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \Psi(x)], \quad t \geq 0.
 \end{aligned}
 \tag{6.4.26}$$

Clearly, for all $t \geq 0$ we have

$$\phi_t(x) \leq A_t \bar{f}(x), \quad x \in Q. \tag{6.4.27}$$

Define

$$C_{v,t} = a_0 \Delta(x_0) + \frac{1}{1+v} \left(\sum_{k=1}^t \frac{a_k^{1+v}}{A_k^v} \right) G_v D^{1+v}, \quad t \geq 0. \quad (6.4.28)$$

Let us introduce

$$\delta(x) \stackrel{\text{def}}{=} \max_{y \in Q} \{ \langle \nabla f(x), x - y \rangle + \Psi(x) - \Psi(y) \} \stackrel{(6.4.21)}{=} \Psi_{x,Q}^*(\nabla f(x)). \quad (6.4.29)$$

For problem (6.4.1), this value measures the level of satisfaction of the first-order optimality conditions at a point $x \in Q$. For any $x \in Q$, we have

$$\delta(x) \geq \bar{f}(x) - \bar{f}(x_*) \geq 0. \quad (6.4.30)$$

We call $\delta(x)$ the *total variation* of the linear model of the composite objective function in problem (6.4.1) over the feasible set. It justifies the first-order optimality conditions in our problem. Note that this value can be computed by a procedure for solving the auxiliary problem (6.4.7).

Theorem 6.4.2 *Let the sequence $\{x_t\}_{t \geq 0}$ be generated by method (6.4.24). Then, for any $v \in (0, 1]$ and any step $t \geq 0$, we have*

$$A_t \bar{f}(x_t) \leq \phi_t(x) + C_{v,t}, \quad x \in Q. \quad (6.4.31)$$

Moreover, for any $t \geq 0$ we have

$$\bar{f}(x_t) - \bar{f}(x_{t+1}) \geq \tau_t \delta(x_t) - \frac{G_v D^{1+v}}{1+v} \tau_t^{1+v}. \quad (6.4.32)$$

Proof Let us prove inequality (6.4.31). For $t = 0$, we have $C_{v,0} = a_0[\bar{f}(x_0) - \bar{f}(x_*)]$. Thus, in this case (6.4.31) follows from (6.4.27).

Assume now that (6.4.31) is valid for some $t \geq 0$. In view of definition (6.4.13), optimality condition (6.4.25) can be written in the following form:

$$\begin{aligned} a_{t+1} \langle \nabla f(x_t), x - x_t \rangle &\geq A_{t+1} [\Psi(x_{t+1}) - \Psi((1 - \tau_t)x_t + \tau_t x) \\ &\quad + \langle \nabla f(x_t), x_{t+1} - x_t \rangle] \end{aligned}$$

for all $x \in Q$. Therefore,

$$\begin{aligned}
 \phi_{t+1}(x) + C_{v,t} &= \phi_t(x) + C_{v,t} \\
 &\quad + a_{t+1}[f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \Psi(x)] \\
 &\stackrel{(6.4.25), (6.4.31)}{\geq} A_t[f(x_t) + \Psi(x_t)] + a_{t+1}[f(x_t) + \Psi(x)] \\
 &\quad + A_{t+1}[\Psi(x_{t+1}) - \Psi((1 - \tau_t)x_t + \tau_t x)] \\
 &\quad + \langle \nabla f(x_t), x_{t+1} - x_t \rangle \\
 &\geq A_{t+1}[f(x_t) + \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \Psi(x_{t+1})] \\
 &\stackrel{(6.4.4)}{\geq} A_{t+1} \left[\bar{f}(x_{t+1}) - \frac{1}{1+\nu} G_v \|x_{t+1} - x_t\|^{1+\nu} \right].
 \end{aligned}$$

It remains to note that $\|x_{t+1} - x_t\| = \tau_t \|x_t - v_t\| \stackrel{(6.4.13)}{\leq} \frac{a_{t+1}}{A_{t+1}} D$. Thus, we can take

$$C_{v,t+1} = C_{v,t} + \frac{1}{1+\nu} \frac{a_{t+1}^{1+\nu}}{A_{t+1}^\nu} G_v D^{1+\nu}.$$

In order to prove inequality (6.4.32), let us introduce the values

$$\begin{aligned}
 \delta_{\tau_t}(x) &\stackrel{\text{def}}{=} \max_{u \in Q} \{ \langle \nabla f(x), x - y \rangle + \Psi(x) - \Psi(y) : y = (1 - \tau)x + \tau u \} \\
 &\stackrel{(6.4.22)}{=} \Psi_{\tau, x, Q}^*(\nabla f(x)), \quad \tau \in [0, 1].
 \end{aligned}$$

Clearly,

$$\begin{aligned}
 -\delta_{\tau_t}(x_t) &= \min_{x \in Q} \{ \langle \nabla f(x_t), y - x_t \rangle + \Psi(y) - \Psi(x_t) : y = (1 - \tau_t)x_t + \tau_t x \} \\
 &= \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \Psi(x_{t+1}) - \Psi(x_t) \\
 &\stackrel{(6.4.4)}{\geq} \bar{f}(x_{t+1}) - \bar{f}(x_t) - \frac{G_v}{1+\nu} \|x_{t+1} - x_t\|^{1+\nu}.
 \end{aligned}$$

Since $\|x_{t+1} - x_t\| \leq \tau_t D$, we conclude that

$$\bar{f}(x_t) - \bar{f}(x_{t+1}) \geq \delta_{\tau_t}(x_t) - \frac{G_v D^{1+\nu}}{1+\nu} \tau_t^{1+\nu} \stackrel{(6.4.23)}{\geq} \tau_t \delta(x_t) - \frac{G_v D^{1+\nu}}{1+\nu} \tau_t^{1+\nu}. \quad \square$$

In view of (6.4.27), inequality (6.4.31) results in the following rate of convergence:

$$\bar{f}(x_t) - \bar{f}(x_*) \leq \frac{1}{A_t} C_{\nu,t}, \quad t \geq 0. \quad (6.4.33)$$

For the linearly growing weights $a_t = t$, $A_t = \frac{t(t+1)}{2}$, $t \geq 0$, we have already seen that

$$C_{\nu,t} = \frac{1}{1+\nu} B_{\nu,t} \leq \frac{2^\nu}{(1+\nu)(2-\nu)} \left[\left(t + \frac{1}{2}\right)^{2-\nu} - \left(\frac{1}{2}\right)^{2-\nu} \right] G_\nu D^{1+\nu}.$$

In the case $\nu = 1$, this results in the following rate of convergence:

$$\bar{f}(x_t) - \bar{f}(x_*) \leq \frac{2}{t+1} G_1 D^2, \quad t \geq 1. \quad (6.4.34)$$

Let us justify for this case the rate of convergence of the sequence $\{\delta(x_t)\}_{t \geq 1}$. We have $\tau_t \stackrel{(6.4.13)}{=} \frac{a_{t+1}}{A_{t+1}} = \frac{2}{t+2}$. On the other hand, for any $T \geq t$,

$$\begin{aligned} \frac{2G_1 D^2}{t+1} &\stackrel{(6.4.34)}{\geq} \bar{f}(x_t) - \bar{f}(x_*) \\ &\stackrel{(6.4.32)}{\geq} \sum_{k=t}^T \left[\tau_k \delta(x_k) - \frac{1}{2} G_1 D^2 \tau_k^2 \right] + \bar{f}(x_{T+1}) - \bar{f}(x_*). \end{aligned} \quad (6.4.35)$$

Let $\delta_T^* = \min_{0 \leq t \leq T} \delta(x_t)$. Then, choosing $T = 2t + 1$, we get

$$\begin{aligned} 2 \ln 2 \cdot \delta_T^* &\stackrel{(6.4.10)}{\leq} \left(\sum_{k=t}^T \frac{2}{k+2} \right) \delta_T^* \stackrel{(6.4.35)}{\leq} 2G_1 D^2 \left[\frac{1}{t+1} + \sum_{k=t}^T \frac{1}{(k+2)^2} \right] \\ &\stackrel{(6.4.11)}{\leq} 2G_1 D^2 \left[\frac{1}{t+1} + \frac{12}{11(2t+3)} \right] = 2G_1 D^2 \left[\frac{2}{T+1} + \frac{12}{11(T+2)} \right] \\ &\leq \frac{68}{11} \cdot \frac{G_1 D^2}{T+1}. \end{aligned}$$

Thus, in the case $\nu = 1$, for odd T , we get the following bound:

$$\delta_T^* \leq \frac{34}{11 \ln 2} \cdot \frac{G_1 D^2}{T+1}. \quad (6.4.36)$$

6.4.4 Computing the Primal-Dual Solution

Note that both methods (6.4.12) and (6.4.24) admit computable accuracy certificates. For the first method, define

$$\ell_t = \frac{1}{A_t} \min_x \left\{ \sum_{k=0}^t a_k [f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \Psi(x)] : x \in Q \right\}.$$

This value can be computed by the standard operation (6.4.7). Clearly,

$$\bar{f}(x_t) - \bar{f}(x_*) \leq \bar{f}(x_t) - \ell_t \stackrel{(6.4.16)}{\leq} \frac{1}{A_t} B_{v,t}. \quad (6.4.37)$$

For the second method, let us choose $a_0 = 0$. Then the estimating functions are linear:

$$\phi_t(x) = \sum_{k=1}^t a_k [f(x_{k-1}) + \langle \nabla f(x_{k-1}), x - x_{k-1} \rangle + \Psi(x)].$$

Therefore, defining $\hat{\ell}_t = \frac{1}{A_t} \min_x \{\phi_t(x) : x \in Q\}$, we also have

$$\bar{f}(x_t) - \bar{f}(x_*) \leq \bar{f}(x_t) - \hat{\ell}_t \stackrel{(6.4.16)}{\leq} \frac{1}{A_t} C_{v,t}, \quad t \geq 1. \quad (6.4.38)$$

Accuracy certificates (6.4.37) and (6.4.38) justify that both methods (6.4.12) and (6.4.24) are able to recover some information on the optimal dual solution. However, in order to implement this ability, we need to open the Black Box and introduce an *explicit model* of the function $f(\cdot)$.

Let us assume that the function f is representable in the following form:

$$f(x) = \max_u \{ \langle Ax, u \rangle - g(u) : u \in Q_d \}, \quad (6.4.39)$$

where $A : \mathbb{E} \rightarrow \mathbb{E}_1^*$, Q_d is a closed convex set in a finite-dimensional linear space \mathbb{E}_2 , and the function $g(\cdot)$ is p -uniformly convex on Q_d :

$$\langle \nabla g(u_1) - \nabla g(u_2), u_1 - u_2 \rangle \geq \sigma_g \|u_1 - u_2\|^p, \quad u_1, u_2 \in Q_d, \quad (6.4.40)$$

where the *convexity degree* $p \geq 2$. Denote by $u(x) \in Q_d$ the unique optimal solution to optimization problem in (6.4.39).

Lemma 6.4.3 *The function f has Hölder continuous gradient $\nabla f(x) = A^*u(x)$ with parameter $v = \frac{1}{p-1}$ and constant $G_v = \left(\frac{1}{\sigma_g}\right)^v \|A\|^{1+v}$.*

Proof Let $u_1 = u(x_1)$, $u_2 = u(x_2)$, $g'_1 = \nabla g(u_1)$, and $g'_2 = \nabla g(u_2)$. Then, in view of the optimality condition (2.2.39), we have

$$\langle Ax_1 - g'_1, u_2 - u_1 \rangle \leq 0, \quad \langle Ax_2 - g'_2, u_1 - u_2 \rangle \leq 0.$$

Adding these two inequalities, we get

$$\langle A(x_1 - x_2), u_1 - u_2 \rangle \geq \langle g'_1 - g'_2, u_1 - u_2 \rangle \stackrel{(6.4.40)}{\geq} \sigma_g \|u_1 - u_2\|^p.$$

Thus,

$$\begin{aligned} \|\nabla f(x_1) - \nabla f(x_2)\|^* &= \|A^*(u_1 - u_2)\|^* \leq \|A\| \cdot \|u_1 - u_2\| \\ &\leq \|A\| \cdot \left(\frac{1}{\sigma_g} \|A(x_1 - x_2)\| \right)^{\frac{1}{p-1}} \\ &\leq \|A\|^{\frac{p}{p-1}} \left(\frac{1}{\sigma_g} \|x_1 - x_2\| \right)^{\frac{1}{p-1}}. \quad \square \end{aligned}$$

Let us write down an *adjoint problem* to (6.4.1).

$$\begin{aligned} \min_x \{ \bar{f}(x) : x \in Q \} &\stackrel{(6.4.39)}{=} \min_x \left\{ \Psi(x) + \max_u \{ \langle Ax, u \rangle - g(u) : u \in Q_d \} \right\} \\ &\geq \max_{u \in Q_d} \left\{ -g(u) + \min_x \{ \langle A^*u, x \rangle + \Psi(x) \} \right\}. \end{aligned}$$

Thus, defining $\Phi(u) = \min_x \{ \langle A^*u, x \rangle + \Psi(x) \}$, we get the following adjoint problem:

$$\max_{u \in Q_d} \left\{ \bar{g}(u) \stackrel{\text{def}}{=} -g(u) + \Phi(u) \right\}. \quad (6.4.41)$$

In this problem, the objective function is nonsmooth and uniformly strongly concave of degree p . Clearly, we have

$$\bar{f}(x) - \bar{g}(u) \geq 0, \quad x \in Q, \quad u \in Q_d. \quad (6.4.42)$$

Let us show that both methods (6.4.12) and (6.4.24) are able to approximate the optimal solution to the problem (6.4.41).

Note that for any $\bar{x} \in Q$ we have

$$\begin{aligned} f(\bar{x}) + \langle \nabla f(\bar{x}), x - \bar{x} \rangle &\stackrel{(6.4.39)}{=} \langle A\bar{x}, u(\bar{x}) \rangle - g(u(\bar{x})) + \langle A^*u(\bar{x}), x - \bar{x} \rangle \\ &= \langle Ax, u(\bar{x}) \rangle - g(u(\bar{x})). \end{aligned}$$

Therefore, defining for the first method (6.4.12) $u_t = \frac{1}{A_t} \sum_{k=0}^t a_k u(x_k)$, we obtain

$$\begin{aligned} \ell_t &= \min_{x \in Q} \left\{ \Psi(x) + \frac{1}{A_t} \sum_{k=0}^t a_k [\langle Ax, u(x_k) \rangle - g(u(x_k))] \right\} \\ &= \Phi(u_t) - \frac{1}{A_t} \sum_{k=0}^t a_k g(u(x_k)) \leq \bar{g}(u_t). \end{aligned}$$

Thus, we get

$$0 \stackrel{(6.4.42)}{\leq} \bar{f}(x_t) - \bar{g}(u_t) \leq \bar{f}(x_t) - \ell_t \stackrel{(6.4.37)}{\leq} \frac{1}{A_t} B_{v,t}, \quad t \geq 0. \quad (6.4.43)$$

For the second method (6.4.24), we choose $a_0 = 0$ and take $u_t = \frac{1}{A_t} \sum_{k=1}^t a_k u(x_{k-1})$.

In this case, by a similar reasoning, we get

$$0 \stackrel{(6.4.42)}{\leq} \bar{f}(x_t) - \bar{g}(u_t) \leq \bar{f}(x_t) - \hat{\ell}_t \stackrel{(6.4.38)}{\leq} \frac{1}{A_t} C_{v,t}, \quad t \geq 1. \quad (6.4.44)$$

6.4.5 Strong Convexity of the Composite Term

In this section, we assume that the function Ψ in problem (6.4.1) is *strongly convex* (see Sect. 3.2.6). In view of (3.2.37), this means that there exists a positive constant σ_Ψ such that

$$\Psi(\tau x + (1 - \tau)y) \leq \tau \Psi(x) + (1 - \tau) \Psi(y) - \frac{1}{2} \sigma_\Psi \tau(1 - \tau) \|x - y\|^2 \quad (6.4.45)$$

for all $x, y \in Q$ and $\tau \in [0, 1]$. Let us show that in this case CG-methods converge much faster. We demonstrate this for method (6.4.12).

In view of the strong convexity of Ψ , the variational principle (6.4.8) characterizing the point v_t in method (6.4.12) can be strengthened:

$$\Psi(x) + \langle \nabla f(x_t), x - v_t \rangle \geq \Psi(v_t) + \frac{1}{2} \sigma_\Psi \|x - v_t\|^2, \quad x \in Q. \quad (6.4.46)$$

Let V_0 be defined as in (6.4.14). Define

$$\hat{B}_{v,t} = a_0 V_0 + \left(\sum_{k=1}^t \frac{a_k^{1+2v}}{A_k^{2v}} \right) \frac{G_v^2 D^{2v}}{2\sigma_\Psi}, \quad t \geq 0. \quad (6.4.47)$$

Theorem 6.4.3 *Let the sequence $\{x_t\}_{t \geq 0}$ be generated by method (6.4.12), and assume the function Ψ is strongly convex. Then, for any $v \in (0, 1]$, any step $t \geq 0$, and any $x \in Q$ we have*

$$A_t(f(x_t) + \Psi(x_t)) \leq \sum_{k=0}^t a_k[f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \Psi(x)] + \hat{B}_{v,t}. \quad (6.4.48)$$

Proof The beginning of the proof of this statement is very similar to that of Theorem 6.4.1. Assuming that (6.4.48) is valid for some $t \geq 0$, we get the following inequality:

$$\begin{aligned} & \sum_{k=0}^{t+1} a_k[f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \Psi(x)] + B_{v,t} \\ & \geq A_{t+1}(f(x_{t+1}) + \Psi(x_{t+1})) + a_{t+1}[\Psi(x) - \Psi(v_t) + \langle \nabla f(x_{t+1}), x - v_t \rangle]. \end{aligned}$$

Further,

$$\begin{aligned} & \Psi(x) - \Psi(v_t) + \langle \nabla f(x_{t+1}), x - v_t \rangle \\ & \stackrel{(6.4.46)}{\geq} \langle \nabla f(x_{t+1}) - \nabla f(x_t), x - v_t \rangle + \frac{1}{2}\sigma_\Psi \|x - v_t\|^2 \\ & \stackrel{(4.2.3)}{\geq} -\frac{1}{2\sigma_\Psi} \|\nabla f(x_{t+1}) - \nabla f(x_t)\|_*^2 \\ & \stackrel{(6.4.3)}{\geq} -\frac{1}{2\sigma_\Psi} \left(\frac{a_{t+1}^v}{A_{t+1}^v} G_v D^v \right)^2. \end{aligned}$$

Thus, to ensure that (6.4.48) is valid for the next iteration, it is enough to choose

$$\hat{B}_{v,t+1} = \hat{B}_{v,t} + \frac{1}{2\sigma_\Psi} \frac{a_{t+1}^{1+2v}}{A_{t+1}^{2v}} G_v^2 D^{2v}. \quad \square$$

It can be easily checked that in our situation, the linear weights strategy $a_t \equiv t$ is not the best one. Let us choose $a_t = t^2$, $t \geq 0$. Then $A_t = \frac{t(t+1)(2t+1)}{6}$, and we get

$$\begin{aligned} \hat{B}_{v,t} &= \left(\sum_{k=1}^t \frac{6^{2v} k^{2(1+2v)}}{k^{2v} (k+1)^{2v} (2k+1)^{2v}} \right) \frac{G_v^2 D^{2v}}{2\sigma_\Psi} \leq \left(3^{2v} \sum_{k=1}^t k^{2(1-v)} \right) \frac{G_v^2 D^{2v}}{2\sigma_\Psi} \\ &\stackrel{(6.4.9)}{\leq} \frac{G_v^2 D^{2v}}{2\sigma_\Psi} \cdot \frac{3^{2v}}{3-2v} \tau^{3-2v} \Big|_{1/2}^{t+1/2} = \frac{3^{2v}}{3-2v} \left[\left(t + \frac{1}{2}\right)^{3-2v} - \left(\frac{1}{2}\right)^{3-2v} \right] \frac{G_v^2 D^{2v}}{2\sigma_\Psi}. \end{aligned}$$

Thus, for $\nu \in (0, 1)$, we get $\frac{1}{A_t} \hat{B}_{\nu, t} \leq O(t^{-2\nu})$. For $\nu = 1$, we obtain

$$\tilde{f}(x_t) - \tilde{f}(x_*) \leq \frac{54}{(t+1)(2t+1)} \cdot \frac{G_1^2 D^2}{2\sigma_\Psi}, \quad (6.4.49)$$

which is much better than (6.4.19). This gives us an example of acceleration of the Conditional Gradient Method by a strong convexity assumption.

6.4.6 Minimizing the Second-Order Model

Let us assume now that in problem (6.4.1) the function f is twice continuously differentiable. Then we can apply to this problem the following method.

Composite Trust-Region Method with Contraction

1. Choose an arbitrary point $x_0 \in Q$.

2. For $t \geq 0$ iterate: Define the coefficient $\tau_t \in (0, 1]$ and choose

$$x_{t+1} \in \operatorname{Arg\,min}_y \left\{ \langle \nabla f(x_t), y - x_t \rangle + \frac{1}{2} \langle \nabla^2 f(x_t)(y - x_t), y - x_t \rangle + \Psi(y) : y \in (1 - \tau_t)x_t + \tau_t x, x \in Q \right\}.$$

(6.4.50)

Note that this scheme is well defined even if the Hessian of the function f is positive semidefinite. Of course, in general, the computational cost of each iteration of this scheme can be big. However, in one important case, when $\Psi(\cdot)$ is an indicator function of a Euclidean ball, the complexity of each iteration of this scheme is dominated by the complexity of matrix inversion. Thus, method (6.4.50) can be easily applied to problems of the form

$$\min_x \{f(x) : \|x - x_0\| \leq r\}, \quad (6.4.51)$$

where the norm $\|\cdot\|$ is Euclidean.

Let $H_\nu < +\infty$ for some $\nu \in (0, 1]$. In this section we assume that

$$\langle \nabla^2 f(x)h, h \rangle \leq L\|h\|^2, \quad x \in Q, \quad h \in \mathbb{E}. \quad (6.4.52)$$

Let us choose a sequence of nonnegative weights $\{a_t\}_{t \geq 0}$, and define in (6.4.50) the coefficients $\{\tau_t\}_{t \geq 0}$ in accordance with (6.4.13). Define the estimating functional sequence $\{\phi_t(x)\}_{t \geq 0}$ by recurrent relations (6.4.26), where the sequence $\{x_t\}_{t \geq 0}$ is generated by method (6.4.50). Finally, define

$$\hat{C}_{v,t} = a_0 \Delta(x_0) + \left(\sum_{k=1}^t \frac{a_k^{2+v}}{A_k^{1+v}} \right) \frac{H_v D^{2+v}}{(1+v)(2+v)} + \left(\sum_{k=1}^t \frac{a_k^2}{2A_k} \right) L D^2. \quad (6.4.53)$$

In our convergence results, we also estimate the *second-order optimality measure* for problem (6.4.1) at the current test points. Let us introduce

$$\theta(x) \stackrel{\text{def}}{=} \max_{y \in Q} \{ \langle \nabla f(x), x - y \rangle - \frac{1}{2} \langle \nabla^2 f(x)(y - x), y - x \rangle + \Psi(x) - \Psi(y) \}. \quad (6.4.54)$$

For any $x \in Q$ we have $\theta(x) \geq 0$. We call $\theta(x)$ the *total variation* of the quadratic model of the composite objective function in problem (6.4.1) over the feasible set. Defining

$$F_x(y) = \frac{1}{2} \langle \nabla^2 f(x)(y - x), y - x \rangle + \Psi(y),$$

we get $\theta(x) = (F_x)^*_{x,Q}(\nabla f(x))$ (see definition (6.4.21)).

Theorem 6.4.4 *Let the sequence $\{x_t\}_{t \geq 0}$ be generated by method (6.4.50). Then, for any $v \in [0, 1]$ and any step $t \geq 0$ we have*

$$A_t \bar{f}(x_t) \leq \phi_t(x) + \hat{C}_{v,t}, \quad x \in Q. \quad (6.4.55)$$

Moreover, for any $t \geq 0$ we have

$$\bar{f}(x_t) - \bar{f}(x_{t+1}) \geq \tau_t \theta(x_t) - \frac{H_v D^{2+v}}{(1+v)(2+v)} \tau_t^{2+v}. \quad (6.4.56)$$

Proof Let us prove inequality (6.4.55). For $t = 0$, $\hat{C}_{v,0} = a_0[\bar{f}(x_0) - \bar{f}(x_*)]$. Therefore, this inequality is valid.

In view of Theorem 3.1.23 the point x_{t+1} is characterized by the following variational principle:

$$x_{t+1} = (1 - \tau_t)x_t + \tau_t v_t, \quad v_t \in Q,$$

$$\Psi(y) + \langle \nabla f(x_t) + \nabla^2 f(x_t)(x_{t+1} - x_t), y - x_{t+1} \rangle \geq \Psi(x_{t+1}),$$

$$\forall y = (1 - \tau_t)x_t + \tau_t x, \quad x \in Q.$$

Therefore, in view of definition (6.4.13), for any $x \in Q$ we have

$$\begin{aligned}
 a_{t+1} \langle \nabla f(x_t), x - x_t \rangle &\geq A_{t+1} \langle \nabla f(x_t) + \nabla^2 f(x_t)(x_{t+1} - x_t), x_{t+1} - x_t \rangle \\
 &\quad + a_{t+1} \langle \nabla^2 f(x_t)(x_{t+1} - x_t), x_t - x \rangle \\
 &\quad + A_{t+1} [\Psi(x_{t+1}) - \Psi((1 - \tau_t)x_t + \tau_t x)] \\
 &\stackrel{(6.4.52)}{\geq} A_{t+1} \langle \nabla f(x_t) + \frac{1}{2} \nabla^2 f(x_t)(x_{t+1} - x_t), x_{t+1} - x_t \rangle \\
 &\quad + A_{t+1} [\Psi(x_{t+1}) - \Psi((1 - \tau_t)x_t + \tau_t x)] - \frac{a_{t+1}^2}{2A_{t+1}} LD^2.
 \end{aligned}$$

Hence,

$$\begin{aligned}
 &A_t \bar{f}(x_t) + a_{t+1} [f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \Psi(x)] \\
 &\geq A_t \Psi(x_t) + A_{t+1} [f(x_t) + \langle \nabla f(x_t) + \frac{1}{2} \nabla^2 f(x_t)(x_{t+1} - x_t), x_{t+1} - x_t \rangle] \\
 &\quad + a_{t+1} \Psi(x) + A_{t+1} [\Psi(x_{t+1}) - \Psi((1 - \tau_t)x_t + \tau_t x)] - \frac{a_{t+1}^2}{2A_{t+1}} LD^2 \\
 &\stackrel{(6.4.6)}{\geq} A_{t+1} [f(x_{t+1}) + \Psi(x_{t+1})] - A_{t+1} \frac{H_v \|x_{t+1} - x_t\|^{2+v}}{(1+v)(2+v)} - \frac{a_{t+1}^2}{2A_{t+1}} LD^2 \\
 &\geq A_{t+1} \bar{f}(x_{t+1}) - \frac{a_{t+1}^{2+v}}{A_{t+1}^{1+v}} \cdot \frac{H_v D^{2+v}}{(1+v)(2+v)} - \frac{a_{t+1}^2}{2A_{t+1}} LD^2.
 \end{aligned}$$

Thus, if (6.4.55) is valid for some $t \geq 0$, then

$$\begin{aligned}
 \phi_{t+1}(x) + \hat{C}_{v,t} &\geq A_t \bar{f}(x_t) + a_{t+1} [f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \Psi(x)] \\
 &\geq A_{t+1} \bar{f}(x_{t+1}) - \frac{a_{t+1}^{2+v}}{A_{t+1}^{1+v}} \cdot \frac{H_v D^{2+v}}{(1+v)(2+v)} - \frac{a_{t+1}^2}{2A_{t+1}} LD^2.
 \end{aligned}$$

Therefore, we can take $\hat{C}_{v,t+1} = \hat{C}_{v,t} + \frac{a_{t+1}^{2+v}}{A_{t+1}^{1+v}} \cdot \frac{H_v D^{2+v}}{(1+v)(2+v)} + \frac{a_{t+1}^2}{2A_{t+1}} LD^2$.

In order to justify inequality (6.4.56), let us introduce the values

$$\begin{aligned}
 \theta_t(\tau) &\stackrel{\text{def}}{=} \max_{x \in Q} \{ \langle \nabla f(x_t), x_t - y \rangle - \frac{1}{2} \langle \nabla^2 f(x_t)(y - x_t), y - x_t \rangle \\
 &\quad + \Psi(x_t) - \Psi(y) : y = (1 - \tau)x_t + \tau x \} \\
 &\stackrel{(6.4.22)}{=} \left(F_{x_t} \right)_{\tau, x_t, Q}^* (\nabla f(x_t)), \quad \tau \in [0, 1].
 \end{aligned}$$

Clearly,

$$\begin{aligned}
 -\theta_t(\tau_t) &= \min_{x \in Q} \{ \langle \nabla f(x_t), y - x_t \rangle - \frac{1}{2} \langle \nabla^2 f(x_t)(y - x_t), y - x_t \rangle \\
 &\quad + \Psi(y) - \Psi(x_t) : y = (1 - \tau_t)x_t + \tau_t x \} \\
 &= \langle \nabla f(x_t), x_{t+1} - x_t \rangle - \frac{1}{2} \langle \nabla^2 f(x_t)(x_{t+1} - x_t), x_{t+1} - x_t \rangle \\
 &\quad + \Psi(x_{t+1}) - \Psi(x_t) \\
 &\stackrel{(6.4.6)}{\geq} \bar{f}(x_{t+1}) - \bar{f}(x_t) - \frac{H_v}{(1+v)(2+v)} \|x_{t+1} - x_t\|^{2+v}.
 \end{aligned}$$

Since $\|x_{t+1} - x_t\| \leq \tau_t D$, we conclude that

$$\bar{f}(x_t) - \bar{f}(x_{t+1}) \geq \theta_t(\tau_t) - \frac{H_v D^{2+v}}{(1+v)(2+v)} \tau_t^{2+v} \stackrel{(6.4.23)}{\geq} \tau_t \theta(x_t) - \frac{H_v D^{2+v}}{(1+v)(2+v)} \tau_t^{2+v}. \square$$

Thus, inequality (6.4.55) ensures the following rate of convergence of method (6.4.50)

$$\bar{f}(x_t) - \bar{f}(x_*) \leq \frac{1}{A_t} \hat{C}_{v,t}. \quad (6.4.57)$$

A particular expression of the right-hand side of this inequality for different values of $v \in [0, 1]$ can be obtained in exactly the same way as it was done in Sect. 6.4.2. Here, we restrict ourselves only to the case when $v = 1$ and $a_t = t^2$, $t \geq 0$. Then $A_t = \frac{t(t+1)(2t+1)}{6}$, and

$$\begin{aligned}
 \sum_{k=1}^t \frac{a_k^3}{A_k^2} &= \sum_{k=1}^t \frac{36k^6}{k^2(k+1)^2(2k+1)^2} \leq 18t, \\
 \sum_{k=1}^t \frac{a_k^2}{2A_k} &= \sum_{k=1}^t \frac{3k^4}{k(k+1)(2k+1)} \leq \frac{3}{2} \sum_{k=1}^t k = \frac{3}{4} t(t+1).
 \end{aligned}$$

Thus, we get

$$\bar{f}(x_t) - \bar{f}(x_*) \leq \frac{18H_1 D^3}{(t+1)(2t+1)} + \frac{9LD^2}{2(2t+1)}. \quad (6.4.58)$$

Note that the rate of convergence (6.4.58) is worse than the convergence rate of cubic regularization of the Newton method (see Sect. 4.2.3). However, to the best of our knowledge, inequality (6.4.58) gives us the first global rate of convergence of an optimization scheme belonging to the family of trust-region methods. In view of inequality (6.4.55), the optimal solution of the dual problem (6.4.41) can be

approximated by method (6.4.50) with $a_0 = 0$ in the same way as it was suggested in Sect. 6.4.4 for Conditional Gradient Methods.

Let us now estimate the rate of decrease of the values $\theta(x_t)$, $t \geq 0$, in the case when $\nu = 1$. Note that $\tau_t \stackrel{(6.4.13)}{=} \frac{a_{t+1}}{A_{t+1}} = \frac{6(t+1)}{(t+2)(2t+3)}$. It is easy to see that these coefficients satisfy the following inequalities:

$$\frac{3}{t+3} \leq \tau_t \leq \frac{6}{2t+5}, \quad t \geq 0. \quad (6.4.59)$$

Therefore, choosing the total number of steps $T = 2t + 2$, we have

$$\begin{aligned} \sum_{k=t}^T \tau_k &\stackrel{(6.4.59)}{\geq} 3 \sum_{k=t}^{2t+2} \frac{1}{k+3} \stackrel{(6.4.10)}{\geq} 3 \ln 2, \\ \sum_{k=t}^T \tau_k^3 &\stackrel{(6.4.59)}{\leq} \sum_{k=t}^{2t+2} \frac{27}{(k+5/2)^3} \stackrel{(6.4.11)}{\leq} -\frac{27}{2(k+5/2)^2} \Big|_{t-1/2}^{2t+5/2} \\ &= \frac{27}{2} \left[\frac{1}{(t+2)^2} - \frac{1}{(2t+5)^2} \right] = \frac{27}{2} \left[\frac{4}{(T+2)^2} - \frac{1}{(T+3)^2} \right] \\ &= \frac{27(3T+8)(T+4)}{2(T+2)^2(T+3)^2} \leq \frac{81}{2(T+1)(T+2)}. \end{aligned} \quad (6.4.60)$$

Now we can use the same trick as at the end of Sect. 6.4.2. Define

$$\theta_T^* = \min_{0 \leq t \leq T} \theta(x_t).$$

Then

$$\begin{aligned} \frac{36H_1D^3}{T(T-1)} + \frac{9LD^2}{2(T-1)} &\stackrel{(6.4.58)}{\geq} \bar{f}(x_t) - \bar{f}(x_*) \geq \sum_{k=t}^T (\bar{f}(x_k) - \bar{f}(x_{k+1})) \\ &\stackrel{(6.4.56)}{\geq} \theta_T^* \sum_{k=t}^T \tau_k - \frac{H_1D^3}{6} \sum_{k=t}^T \tau_k^3 \\ &\stackrel{(6.4.60)}{\geq} 3\theta_T^* \ln 2 - \frac{27H_1D^3}{4(T+1)(T+2)}. \end{aligned}$$

Thus, for even T , we get the following bound:

$$\begin{aligned} \theta_T^* &\leq \frac{3}{\ln 2} \left[\frac{4H_1D^3}{T(T-1)} + \frac{3H_1D^3}{4(T+1)(T+2)} + \frac{LD^2}{2(T-1)} \right] \\ &\leq \frac{3}{\ln 2} \left[\frac{5H_1D^3}{T(T-1)} + \frac{LD^2}{2(T-1)} \right]. \end{aligned} \quad (6.4.61)$$