

# Nearly linear-time packing and covering LP solvers

## Achieving width-independence and $O(1/\varepsilon)$ -convergence

Zeyuan Allen-Zhu<sup>1</sup> · Lorenzo Orecchia<sup>2</sup>

Received: 8 August 2016 / Accepted: 30 January 2018 / Published online: 14 February 2018  
© Springer-Verlag GmbH Germany, part of Springer Nature and Mathematical Optimization Society 2018

**Abstract** Packing and covering linear programs (PC-LPs) form an important class of linear programs (LPs) across computer science, operations research, and optimization. Luby and Nisan (in: STOC, ACM Press, New York, 1993) constructed an iterative algorithm for approximately solving PC-LPs in *nearly linear time*, where the time complexity scales nearly linearly in  $N$ , the number of nonzero entries of the matrix, and polynomially in  $\varepsilon$ , the (multiplicative) approximation error. Unfortunately, existing nearly linear-time algorithms (Plotkin et al. in Math Oper Res 20(2):257–301, 1995; Bartal et al., in: Proceedings 38th annual symposium on foundations of computer science, IEEE Computer Society, 1997; Young, in: 42nd annual IEEE symposium on foundations of computer science (FOCS'01), IEEE Computer Society, 2001; Koufogiannakis and Young in Algorithmica 70:494–506, 2013; Young in Nearly linear-time approximation schemes for mixed packing/covering and facility-location linear programs, 2014, arXiv:1407.3015; Allen-Zhu and Orecchia, in: SODA, 2015) for solving PC-LPs require time at least proportional to  $\varepsilon^{-2}$ . In this paper, we break this long-standing barrier by designing a packing solver that runs in time  $\tilde{O}(N\varepsilon^{-1})$  and covering LP solver that runs in time  $\tilde{O}(N\varepsilon^{-1.5})$ . Our packing solver can be extended to run in

---

An earlier version of this paper appeared on arXiv <http://arxiv.org/abs/1411.1124> in November 2014. A 6-paged abstract of this paper “Nearly-Linear Time Positive LP Solver with Faster Convergence Rate,” excluding Sects. 4, 5 and 6, and excluding all the proofs in Sects. 2 and 3, has been presented at the STOC 2015 conference in Portland, OR.

---

✉ Zeyuan Allen-Zhu  
zeyuan@csail.mit.edu

Lorenzo Orecchia  
orecchia@bu.edu

<sup>1</sup> Microsoft Research AI, Redmond, WA, USA

<sup>2</sup> Boston University, Boston, MA, USA

time  $\tilde{O}(N\varepsilon^{-1})$  for a class of well-behaved covering programs. In a follow-up work, Wang et al. (in: ICALP, 2016) showed that all covering LPs can be converted into well-behaved ones by a reduction that blows up the problem size only logarithmically.

**Mathematics Subject Classification** 90C05, Linear programming · 90C25, Convex programming · 65K05, Mathematical programming methods · 49M20, Methods of relaxation type

## 1 Introduction

A packing linear program (LP) takes the form  $\max\{c^T x : Ax \leq b\}$  where  $c \in \mathbb{R}_{\geq 0}^n$ ,  $b \in \mathbb{R}_{\geq 0}^m$ , and  $A \in \mathbb{R}_{\geq 0}^{m \times n}$ . A covering LP can be written as  $\min\{b^T y : A^T y \geq c\}$ , with the same requirements on  $A$ ,  $b$ , and  $c$ . We denote by  $N$  the number of non-zero elements in matrix  $A$ . We assume without loss of generality that the two LP programs are in their *standard forms*:

$$\text{Packing LP: } \max_{x \in \mathbb{R}_{\geq 0}^n} \{\mathbb{1}^T x : Ax \leq \mathbb{1}\}, \quad (1.1)$$

$$\text{Covering LP: } \min_{y \in \mathbb{R}_{\geq 0}^m} \{\mathbb{1}^T y : A^T y \geq \mathbb{1}\}. \quad (1.2)$$

The two programs are dual to each other, so we denote by  $\text{OPT} \geq 0$  their shared optimum. We say  $x$  is a  $(1 - \varepsilon)$ -approximation for the packing LP if  $Ax \leq \mathbb{1}$  and  $\mathbb{1}^T x \geq (1 - \varepsilon)\text{OPT}$ , and  $y$  a  $(1 + \varepsilon)$ -approximation for the covering LP if  $A^T y \geq \mathbb{1}$  and  $\mathbb{1}^T y \leq (1 + \varepsilon)\text{OPT}$ .

In this paper, we study first-order iterative methods for solving packing and covering linear programs (PC-LPs) efficiently.<sup>1</sup> Of course, it is possible to adopt the Interior Point or Ellipsoid methods to obtain approximate solvers with a  $\log(1/\varepsilon)$  dependence on the number of iterations. However, the computational cost of such algorithms is typically high, as each iteration requires solving a linear system, and thus is not suitable for large-scale applications.

To address this issue, researchers have developed *iterative approximate* PC-LP solvers that achieve a better dependence on the problem size (e.g., nearly linear in  $N$ ) at the cost of having a  $\text{poly}(1/\varepsilon)$  dependence on the approximation parameter  $\varepsilon$ . Such iterative solvers have been widely applied in approximation algorithms (e.g., MINSETCOVER [24], MAXSET, MAXDICT, MAX- $k$ -CSP [32], bipartite matching), probabilistic checkable proofs [32], zero-sum matrix games [29], scheduling [31], graph embedding [31], flow controls [10, 11], auction mechanisms [37], wireless sensor networks [14], and many other areas. In addition, techniques developed in this line of research have inspired important results on other fundamental algorithmic problems, such as the design of fast algorithms for multi-commodity flow problems [9, 18, 19, 25, 31] and the equivalence between QIP and PSPACE [21].

<sup>1</sup> Luby and Nisan, who originally studied iterative solvers for this class of problems [24], dubbed them *positive LPs*. However, the class of LPs with non-negative constraint matrices is slightly larger, including mixed-packing-and-covering LPs. For this reason, we prefer to stick to the PC-LP terminology.

**Table 1** Comparisons among iterative approximate solvers for packing and covering LPs

Paper	Running time	Width independent?	Nearly linear-time?
[31]	$O(N \times \frac{\rho^2 \text{OPT}^2 \log m}{\epsilon^2})$	No	No
[7]	$O(N \times \frac{\rho \text{OPT} \log m}{\epsilon^2})$	No	No
[26,29]	$O(N \times \frac{\rho \text{OPT} \log m}{\epsilon})$	No	No
[13]	$O(N \times \frac{\sqrt{Kn \log m}}{\epsilon})$	Yes	No
[15,27]: packing LP	$\tilde{O}(N \times (n + \frac{\sqrt{n}}{\epsilon}))$	Yes	No
Parallel solvers [4,8,10,11,24,34,35]	$O(N \times \frac{\log^2 N \log(1/\epsilon)}{\epsilon^2})$ at best	Yes	Yes
[35]	$O((md + N) \times \frac{\log N}{\epsilon^2})$	Yes	Almost yes
[10,11]	$O(nm \times \frac{\log N}{\epsilon^2})$	Yes	Almost yes
[36]	$O(N \times \frac{\log N}{\epsilon^2})$	Yes	Yes
[23]	$O(N + (n + m) \times \frac{\log N}{\epsilon^2})$	Yes	Yes
Theorem 3.4 packing LP	$O(N \times \frac{\log N \log \epsilon^{-1}}{\epsilon})$	Yes	Yes
Theorem 5.3 well-behaved covering LP	$O(N \times \frac{\log N \log \epsilon^{-1}}{\epsilon})$	Yes	Yes
Theorem 6.6 covering LP	$O(N \times \frac{\log N \log \epsilon^{-1}}{\epsilon^{1.5}})$	Yes	Yes

The width  $\rho \in [1/\text{OPT}, \infty)$  is defined as the largest entry of the constraint matrix  $A$ . The parameter  $d$  is the maximum number of constraints each variable is in;  $md$  may be larger than  $N$

Previous iterative approximate solvers can be divided into two classes, *width-dependent* and *width-independent* solvers (see also Table 1).

**Width-dependent solvers**<sup>2</sup> Based on multiplicative weight update ideas (a.k.a. exponentiated gradient updates), researchers have obtained solvers for PC-LPs with a running time at least  $N$  multiplied with  $\rho \text{OPT} \in [1, \infty)$ , where  $\rho$  is the *width* of the program, i.e., the largest entry of matrix  $A$ . For instance, PC-LPs can be solved in  $O(\frac{N\rho^2 \text{OPT}^2 \log m}{\epsilon^2})$ -time [31], or  $O(\frac{N\rho \text{OPT} \log m}{\epsilon^2})$ -time using some more refined analysis [7]. These algorithms only require “oracle-access” to the matrix  $A$ . When  $A$  is given explicitly like in this paper, the running time can be reduced to  $O(\frac{N\rho \text{OPT} \log m}{\epsilon})$  by deploying Nesterov’s accelerated gradient method [29], or Nemirovski’s mirror prox method [26].

<sup>2</sup> Most width-dependent solvers study the minmax problem  $\min_{x \geq 0, \mathbb{1}^T x = 1} \max_{y \geq 0, \mathbb{1}^T y = 1} y^T Ax$ , whose optimal value equals  $1/\text{OPT}$ . Their approximation guarantees are often written in terms of *additive* error. We have translated their performances to multiplicative error for a clear comparison.

Width-dependent algorithms are not polynomial time but only *pseudo-polynomial time*.

**Width-independent, but super linear-time solvers** Researchers also tried to appropriately scale the matrix so as to avoid the width penalty in the above methods. For instance, Bienstock and Iyengar [13] built on Nesterov's method [29] and obtained a running time  $O(\varepsilon^{-1}N\sqrt{Kn\log m})$  where  $K$  is the maximum number of non-zeros per row of  $A$ . This is  $O(\varepsilon^{-1}Nn\sqrt{\log m})$  in the worst case. The results of [15,27] improved this complexity (for packing LP only) to  $\tilde{O}(\varepsilon^{-1}N\sqrt{n})$ , at a cost of enduring an  $\tilde{O}(Nn)$ -time preprocessing stage.

**Width-independent, nearly linear-time solvers** Perhaps the most desirable complexity is a running time that is both independent of the width parameter  $\rho$ , and also nearly linearly scales with  $N$ .<sup>3</sup> This line of research was initiated by a seminal paper of Luby and Nisan [24], who gave an algorithm running in  $O(\frac{N\log^2 N}{\varepsilon^4})$  time with no dependence on the width  $\rho$ . This is also the first *nearly linear-time* approximate solver for PC-LPs, and also the first to run in parallel in nearly linear-work and polylogarithmic depth.

The parallel algorithm of Luby and Nisan was extended by [4,8,10,33,35]. Most notably, the algorithm of Wang *et al.* [33] runs in  $O(\frac{\log^2 N \log(1/\varepsilon)}{\varepsilon^2})$  iterations, each costing a matrix-vector multiplication that can be implemented in  $O(N)$  total work.

The ideas of Luby and Nisan also led to sequential width-independent, nearly linear-time PC-LP solvers [10,11,23,35,36]. Most notably, the algorithm of Koufogiannakis and Young [23] runs in time  $O(N + \frac{\log N}{\varepsilon^2} \times (n + m))$ .

Despite the amount of work in this area, the  $O(1/\varepsilon^2)$  convergence rate was established in 1997 [10,11] and has not been improved since then. On a separate note, Klein and Young [22] showed that all Dantzig-Wolfe type algorithms have to suffer from a  $O(1/\varepsilon^2)$  convergence rate. This lack of progress constitutes a significant limitation, as the  $\varepsilon^{-2}$ -dependence (also known as the  $1/\sqrt{T}$  convergence) on the approximation parameter  $\varepsilon$  is particularly poor.

## 1.1 Our results

**Packing LP** We present an algorithm *PacLP Solver* that runs in  $O(\frac{\log(nm/\varepsilon)\log(1/\varepsilon)}{\varepsilon}N)$  total time. This gives the first width-independent, and the first nearly linear-time solver for packing LP with an  $\varepsilon^{-1}$  convergence (i.e., an  $1/T$  convergence). In contrast, no nearly linear-time algorithm has achieved any convergence rate faster than  $\varepsilon^{-2}$  before our work.

Interestingly, the maximum (weighted) bipartite matching is just one instance of a packing LP. As a consequence, our *PacLP Solver* algorithm finds an approximate maximum bipartite matching in time  $\tilde{O}(m\varepsilon^{-1})$ . This new matching algorithm, which arises purely from convex-optimization arguments, matches the running time of the best known combinatorial algorithm for maximum weighted bipartite matching [16].

<sup>3</sup> Some of these solvers still have a  $\text{polylog}(\rho)$  dependence. Since each occurrence of  $\log(\rho)$  can be replaced with  $\log(nm)$  after slightly modifying the matrix  $A$ , we have done so in Table 1 for a fair comparisons.

**Covering LP** A symmetric design of *PacLP Solver* gives rise to an algorithm *CovLP Solver<sup>wb</sup>* with the same running time  $O(\frac{\log(nm/\varepsilon)\log(1/\varepsilon)}{\varepsilon} N)$ , but only solving *well-behaved* covering LP instances. At a high level, we say an instance is well-behaved if the constraint  $A^T y \geq \mathbb{1}$  is “never redundant”: for instance, if the optimal solution  $y^*$  satisfies  $C \cdot \mathbb{1} \geq A^T y^* \geq \mathbb{1}$  for some constant  $C > 1$  then the covering LP is well-behaved. For the *general* covering LP *without* well-behavior assumptions, we propose a different algorithm *CovLP Solver* that runs in time  $O(\frac{\log(nm/\varepsilon)\log(1/\varepsilon)}{\varepsilon^{1.5}} N)$ . Again, we emphasize that no nearly linear-time covering LP solver can achieve a convergence rate faster than  $\varepsilon^{-2}$  (or equivalently  $O(1/\sqrt{T})$ ) before our work.

REMARK After the first version of this paper appeared on arXiv in 2014, Wang, Rao and Mahoney [34] showed all covering LPs can be converted into well-behaved ones, by blowing up the problem size logarithmically. In other words, they obtained a nearly linear-time covering LP solver with  $\varepsilon^{-1}$  convergence by a reduction to *CovLP Solver<sup>wb</sup>*. Nevertheless, our *CovLP Solver*, being a direct method, may still be of practical and theoretical interests.

## 1.2 Main challenge and our approach

**Width-independence versus acceleration** Previous solvers for PC-LPs are based on standard techniques in non-smooth optimization. They first implicitly or explicitly smoothen the objective, often by the entropy regularizer. Then, they minimize the resulting convex objective either via variations of full-gradient methods, yielding parallel algorithms, or via variations of coordinate-gradient methods, yielding sequential algorithms. The *main challenge* in previous work is to show that the width dependence can sometimes be completely removed for PC-LPs, if the underlying minimization method is designed cleverly.

Of course, the slower the convergence rate is, the easier it is to design nearly linear-time solvers. The  $\varepsilon^{-4}$ -convergence solver of Awerbuch and Khandekar [8] and the  $\varepsilon^{-3}$ -convergence solver of [4] are arguably the simplest nearly linear-time solvers at this point.

In this paper, we achieve the  $\varepsilon^{-1}$  convergence that is typical for accelerated gradient descent over smoothened objectives [29], but without paying the width or any additional super-logarithmic factors. The challenge in this approach is to preserve the width-independence and the accelerated rate *at the same time*. We stress here that our algorithm is *not an instance* of any known variant of accelerated gradient descent.<sup>4</sup> Moreover, the incorporation of width-independence and Nesterov’s acceleration requires significant effort, as witnessed by the lack of progress on this problem for the last 15 years.

**Our high-level approach** Our approach is based on an improved convex formalization  $f(x)$  of the PC-LP objective, together with our linear-coupling framework for designing efficient first-order methods [5] for minimizing  $f(x)$ .

---

<sup>4</sup> This can be verified by observing that our objective  $f_\mu(x)$ , to be introduced later, is not globally Lipschitz smooth, so that one cannot apply accelerated gradient descent directly.

The improved formalization shows that our smoothed objective  $f(x)$  satisfies either the classical condition for Lipschitz smoothness or a different condition based on multiplicative change. This formalization also clarifies why width-independent algorithms exist in the first place. See Lemma 2.7 and the related discussion for more details.

The linear-coupling framework in our previous work [5] provides a different interpretation of Nesterov's acceleration for smooth optimization [28]. In a nutshell, this linear-coupling framework allows us to construct accelerated algorithms by coupling the executions of a gradient descent algorithm, yielding iterates  $\{y_k\}$  and a mirror descent step algorithm, with iterates  $\{z_k\}$ . The name "linear coupling" stems from the fact that, at iteration  $k + 1$ , the gradient of the objective is queried at a point  $x_{k+1}$ , which is a linear combination of gradient and mirror steps, i.e.,  $x_{k+1} = (1 - \tau) \cdot z_k + (1 - \tau) \cdot y_k$ .

In this paper, we apply linear coupling in a very non-trivial manner. We design a gradient and a mirror descent step, each very specific to the underlying PC-LP problem. We also perform a coupling step  $x_{k+1} = (1 - \tau) \cdot z_k + (1 - \tau) \cdot y_k$ , but need to design a different analysis to preserve width independence. None of these components has appeared in [5].

**Arithmetic precision** Throughout this paper, we assume exact arithmetic operations for presenting the cleanest proofs. If the updates are calculated within precision  $\frac{1}{\text{poly}(\varepsilon^{-1}, n, m)}$ , or equivalently when word size  $O(\log(\varepsilon^{-1} + n + m))$  is used, our results still hold.<sup>5</sup>

**Roadmap** We relax the packing LP in Sect. 2, and provide our packing LP solver in Sect. 3. We relax the covering LP in Sect. 4, and provide our covering LP solver in the well-behaved case in Sect. 5. In Sect. 6, we provide our full covering LP solver.

## 2 Relaxation of the packing linear program

To solve packing LP, we minimize a relaxed version of the original LP, where the hard constraint  $Ax \leq 1$  is regularized by entropy and replaced by an exponential penalty function.

**Notations** Recall that the packing LP in its standard form is  $\max_{x \geq 0} \{\mathbb{1}^T x : Ax \leq \mathbb{1}\}$ . Let us denote by OPT the optimal value of this linear program, and  $x^*$  any optimal solution. We say that  $x$  is a  $(1 - \varepsilon)$ -approximation for the packing LP if  $Ax \leq \mathbb{1}$  and  $\mathbb{1}^T x \geq (1 - \varepsilon)\text{OPT}$ .

Throughout this paper, we use the indices  $i \in [n]$  to denote the columns of  $A$ , and the indices  $j \in [m]$  to denote the rows of  $A$ . We let  $A_{:i}$  be the  $i$ -th column

<sup>5</sup> Due to space limitation, we quickly sketch why logarithmic word size suffices for our algorithms. On one hand, one can prove in an iteration, if  $x$  is calculated with a small additive error  $1/\text{poly}(1/\varepsilon, n, m)$ , then the objective  $f(x)$  may increase only by  $1/\text{poly}(1/\varepsilon, n, m)$  in that iteration. The proof of this relies on the fact that (1) one can assume without loss of generality all entries of  $A$  are no more than  $\text{poly}(1/\varepsilon, n, m)$  and (2) our algorithms ensure  $f(x) < \text{poly}(1/\varepsilon, n, m)$  for all iterations with high probability, so even though we are using the exponential functions,  $f(x)$  will not change additively by much. On the other hand, one can similarly prove that each  $\nabla_i f(x)$  can be calculated within an additive error  $1/\text{poly}(1/\varepsilon, n, m)$  in each iteration. They together imply that the total error incurred by arithmetic operations can be made negligible.

vector of  $A$ , and  $A_j$ : the  $j$ -th row vector of  $A$ . Given any vector  $x$ , we denote by  $\|x\|_A = \sqrt{\sum_{i \in [n]} x_i^2 \cdot \|A_{:i}\|_\infty}$  the  $A$ -norm of  $x$ . By simple scaling, we can assume without loss of generality that <sup>6</sup>

$$\min_{i \in [n]} \{\|A_{:i}\|_\infty\} = 1. \tag{2.1}$$

We restrict the domain of  $x$  and the range of  $\text{OPT}$  as follows.

**Fact 2.1** Define the bounding box  $\Delta_{\text{box}} \stackrel{\text{def}}{=} \{x \in \mathbb{R}^n : x_i \in [0, \frac{1}{\|A_{:i}\|_\infty}]\}$ . Under assumption (2.1), we have  $\text{OPT} \in [1, n]$  and  $\{x : x \geq 0 \wedge Ax \leq \mathbb{1}\} \subseteq \Delta_{\text{box}}$ .

*Proof* Suppose that  $i^*$  is the column that achieves the smallest infinite norm  $\|A_{:i}\|_\infty$  over all columns. Letting  $x$  be such that  $x_{i^*} = 1$  and  $x_i = 0$  at  $i \neq i^*$ , we claim that  $x$  is a feasible solution for the packing LP (1.1), simply because  $\|A_{:i^*}\|_\infty = 1$  according to (2.1). This feasible solution  $x$  yields an objective value  $\mathbb{1}^T x = 1$ , proving that  $\text{OPT} \geq 1$ . On the other hand, for any solution  $x \geq 0$  satisfying  $Ax \leq \mathbb{1}$ , we must have  $x_i \leq \frac{1}{\|A_{:i}\|_\infty}$  for each  $i$ . Therefore,  $\mathbb{1}^T x \leq \sum_i \frac{1}{\|A_{:i}\|_\infty} \leq n$ , proving that  $\text{OPT} \leq n$ .

The inclusion  $\{x : x \geq 0 \wedge Ax \leq \mathbb{1}\} \subseteq \Delta_{\text{box}}$  is obvious, since the constraints  $x \geq 0$  and  $Ax \leq \mathbb{1}$  together imply  $x_i \leq \frac{1}{\|A_{:i}\|_\infty}$  for every  $i \in [n]$ . □

This bounding-box constraint allows us to focus only on searching  $x$  in  $\Delta_{\text{box}}$ .

**Our regularized objective** We now introduce the smoothed objective  $f_\mu(x)$  that we minimize over  $\Delta_{\text{box}}$  in order to approximately solve packing LP. At a high level, this objective  $f_\mu(x)$  turns each row of the hard, non-smooth LP constraint  $Ax \leq \mathbb{1}$  into an exponential penalty function so that we only need to require  $x \in \Delta_{\text{box}}$  throughout the algorithm.

Formally, the packing LP can be written as the following minimization LP by introducing the Lagrangian variable  $y \in \mathbb{R}^m$ :

$$\min_{x \in \Delta_{\text{box}}} \left\{ -\mathbb{1}^T x + \max_{y \geq 0} \{y^T Ax - \mathbb{1}^T y\} \right\}. \tag{2.2}$$

The problem can be now smoothed by introducing a concave regularizer over  $y \geq 0$ . We take this regularizer to be the generalized entropy  $H(y) = -\sum_{j=1}^m y_j \log y_j + y_j$  over the first orthant  $y \geq 0$ , and minimize the following smoothed objective  $f_\mu(x)$  over  $x \in \Delta_{\text{box}}$ :

$$f_\mu(x) \stackrel{\text{def}}{=} -\mathbb{1}^T x + \max_{y \geq 0} \{y^T Ax - \mathbb{1}^T y + \boxed{\mu \cdot H(y)}\}. \tag{2.3}$$

Above,  $\mu > 0$  is some smoothing parameter to be chosen later. By explicitly computing the maximization over  $y \geq 0$ ,  $f_\mu(x)$  can be rewritten as

---

<sup>6</sup> If  $\min_{i \in [n]} \{\|A_{:i}\|_\infty\} = 0$  then the packing LP is unbounded so we are done. Otherwise, if  $\min_{i \in [n]} \{\|A_{:i}\|_\infty\} = v > 0$  we scale all entries of  $A$  by  $1/v$ , and scale  $\text{OPT}$  by  $v$ .

**Fact 2.2**  $f_\mu(x) = \mu \sum_{j=1}^m e^{\frac{1}{\mu}((Ax)_j-1)} - \mathbb{1}^T x$ .

We study the *minimization* problem on  $f_\mu(x)$  over  $x \in \Delta_{\text{box}}$ . Intuitively  $f_\mu(x)$  captures the original packing LP (1.1) as follows. Firstly, since we want to maximize  $\mathbb{1}^T x$ , the negative term  $-\mathbb{1}^T x$  shows up in  $f_\mu(x)$ . Secondly, if a packing constraint  $j \in [m]$  is violated by  $\varepsilon$ , that is,  $(Ax)_j \geq 1 + \varepsilon$ , the exponential penalty in  $f_\mu(x)$  introduces a penalty at least  $\mu e^{\varepsilon/\mu}$ ; this will be a large penalty if  $\mu \leq O(\varepsilon/\log(n/\varepsilon))$ .

*Remark 2.3* The use of exponential function at least traces back to [31] in 1991 (implicitly) and to [20] in 1994 (explicitly). The way most previous results minimize  $f_\mu(x)$  is by taking a logarithm  $g(x) = \log(\sum_{j=1}^m e^{((Ax)_j-1)/\mu})$ , explicitly or implicitly arguing that  $g(x)$  is Lipschitz smooth (i.e.,  $\|\nabla^2 f(x)\|$  is bounded), and then taking gradient descent.<sup>7</sup> Unfortunately, the Lipschitz smoothness parameter of  $g(x)$  depends on the width of the LP, and thus first-order iterative approaches based on directly minimizing  $g(x)$  are mostly width-dependent [7, 26, 29, 31]. One can also reduce the width parameter of  $g(x)$  which yields super linear-time solvers [13, 15, 27].

In this paper, we directly perform gradient descent and mirror descent on  $f_\mu(x)$ —without taking the logarithm. Note that traditional accelerated gradient methods [28, 29] should not be applied directly to minimize  $f_\mu$  because it is not Lipschitz smooth.<sup>8</sup>

Our  $f_\mu(x)$  incurs a regularization error. The next proposition bounds this error following a similar treatment in [4].

**Proposition 2.4** Let  $\mu = \frac{\varepsilon}{4 \log(nm/\varepsilon)}$  and recall  $x^*$  is an optimal solution for packing LP.

- (a)  $f_\mu(u^*) \leq -(1 - \varepsilon)\text{OPT}$  for  $u^* \stackrel{\text{def}}{=} (1 - \varepsilon/2)x^* \in \Delta_{\text{box}}$ .
- (b)  $f_\mu(x) \geq -(1 + \varepsilon)\text{OPT}$  for every  $x \in \Delta_{\text{box}}$ .
- (c) If  $x \in \Delta_{\text{box}}$  satisfies  $f_\mu(x) \leq -(1 - \theta)\text{OPT}$  for some  $\theta \in [0, 1]$ , then  $\frac{1}{1+\varepsilon}x$  is a  $\frac{1-\theta}{1+\varepsilon}$ -approximate solution to the packing LP.

*Remark 2.5* Our box constraint  $x \in \Delta_{\text{box}}$  is almost redundant for minimizing  $f_\mu(x)$ : whenever  $x \geq 0$  and  $f_\mu(x) \leq 0$ , one should automatically have  $x_i \leq \frac{1+\varepsilon}{\|A_{:,i}\|_\infty}$ . However, this constraint shall be used to make sure that our updates are always inside  $\Delta_{\text{box}}$ .

*Proof of Proposition 2.4* (a) We have  $\mathbb{1}^T u^* = (1 - \varepsilon/2)\text{OPT}$  by the definition of OPT. from the feasibility  $Ax^* \leq \mathbb{1}$  in the packing LP, we have  $Au^* - \mathbb{1} \leq -\varepsilon/2 \cdot \mathbb{1}$ , and can compute  $f_\mu(u^*)$  as follows:

<sup>7</sup> Note that some of the previous results (such as [7, 31]) appear to directly minimize  $\sum_{j=1}^m e^{((Ax)_j-1)/\mu}$  as opposed to its logarithm  $g(x)$ . However, their per-iteration objective decrease is multiplicative, meaning it is essentially equivalent to performing a single gradient-descent step on  $g(x)$  with additive objective decrease.

<sup>8</sup> The exact same  $f_\mu(x)$  also appeared in our previous work [4], albeit without this smoothing interpretation and without the constraint  $x \in \Delta_{\text{box}}$ . The techniques in [4] only leads to  $\varepsilon^{-2}$  convergence (see Table 1).



$$\begin{aligned}
 f_\mu(u^*) &= \mu \sum_j e^{\frac{1}{\mu}((Au^*)_j-1)} - \mathbb{1}^T u^* \\
 &\leq \mu \sum_j e^{\frac{-\varepsilon/2}{\mu}} - (1 - \varepsilon/2)\text{OPT} \\
 &\leq \frac{\mu m}{(nm)^2} - (1 - \varepsilon/2)\text{OPT} \leq -(1 - \varepsilon)\text{OPT}.
 \end{aligned}$$

(b) Suppose towards contradiction that  $f_\mu(x) < -(1 + \varepsilon)\text{OPT}$ . Since  $f_\mu(x) > -\mathbb{1}^T x$ , it must satisfy that  $\mathbb{1}^T x > (1 + \varepsilon)\text{OPT}$ . Suppose that  $\mathbb{1}^T x = (1 + v)\text{OPT}$  for some  $v > \varepsilon$ . By the definition of **OPT**, we must have that  $Ax < (1 + v)\mathbb{1}$  is broken, and therefore there exists some  $j \in [m]$  satisfying that  $(Ax)_j \geq 1 + v$ . In such a case, the objective

$$\begin{aligned}
 f_\mu(x) &\geq \mu e^{v/\mu} - (1 + v)\text{OPT} = \frac{\varepsilon}{4 \log(nm/\varepsilon)} \left( \left( \frac{nm}{\varepsilon} \right)^4 \right)^{v/\varepsilon} - (1 + v)\text{OPT} \\
 &\geq \left( \left( \left( \frac{nm}{\varepsilon} \right)^2 \right)^{v/\varepsilon} - (1 + v) \right) \text{OPT} > 0
 \end{aligned}$$

giving a contradiction to the assumption that  $f_\mu(x) < 0$ .

(c) Note that  $x$  satisfies  $f_\mu(x) \leq -(1 - \delta)\text{OPT} \leq 0$ , and we first show  $Ax \leq (1 + \varepsilon)\mathbb{1}$ . Let us assume that  $v = \max_j((Ax)_j - 1) \geq 0$  because otherwise we will have  $Ax \leq \mathbb{1}$ . Under this definition, we have  $Ax \leq (1 + v)\mathbb{1}$  and therefore  $\mathbb{1}^T x \leq (1 + v)\text{OPT}$  by the definition of **OPT**. We compute  $f_\mu(x)$  as follows.

$$\begin{aligned}
 f_\mu(x) &\geq \mu e^{\frac{v}{\mu}} - (1 + v)\text{OPT} \geq \mu \left( \left( \frac{nm}{\varepsilon} \right)^4 \right)^{v/\varepsilon} - (1 + v)n \\
 &= \frac{\varepsilon}{4 \log(nm/\varepsilon)} \left( \left( \frac{nm}{\varepsilon} \right)^4 \right)^{v/\varepsilon} - (1 + v)n.
 \end{aligned}$$

The above quantity is positive whenever  $v \geq \varepsilon$ , and therefore, to satisfy  $f_\mu(x) \leq 0$  we must have  $v \leq \varepsilon$ , which is equivalent to  $Ax \leq (1 + \varepsilon)\mathbb{1}$ . Next, because  $-\mathbb{1}^T x \leq f_\mu(x) \leq -(1 - \delta)\text{OPT}$ , we know  $\mathbb{1}^T x \geq (1 - \delta)\text{OPT}$ . Letting  $x' = \frac{1}{1+\varepsilon}x$ , we both have that  $x'$  is feasible (i.e.,  $Ax' \leq \mathbb{1}$ ), and  $x'$  has an objective  $\mathbb{1}^T x'$  at least as large as  $\frac{1-\delta}{1+\varepsilon}\text{OPT}$ . □

**Some non-standard smoothness properties** The gradient and Hessian of  $f_\mu(x)$  can be written in the following closed forms:

**Fact 2.6**  $\nabla f_\mu(x) = A^T p(x) - \mathbb{1}$  and  $\nabla^2 f_\mu(x) = \frac{1}{\mu} A^T \text{diag}\{p(x)\}A$ , where  $p_j(x) \stackrel{\text{def}}{=} e^{\frac{1}{\mu}((Ax)_j-1)}$ .

By staring at these closed forms, we note that  $f_\mu(x)$  is not Lipschitz-smooth: for instance, each  $\nabla_{ii}^2 f_\mu(x)$  can go to infinity so the spectral norm of  $\nabla^2 f_\mu(x)$  is

unbounded. However, the non-negativity of  $A$  guarantees that whenever  $\nabla_{ii}^2 f_\mu(x)$  is large for some coordinate  $i$ , the corresponding entry of the gradient  $\nabla_i f_\mu(x)$  must also be large. This still allows us to take a larger step in direction  $\mathbf{e}_i$  than traditionally allowed by coordinate descent.

The above intuition is formalized in the next lemma, whose proof is by simple manipulation of Hessian. The first half of the lemma is the same as the traditional coordinate Lipschitz-smoothness property, but holds only conditionally; the second half is a salient characteristic of this work and requires the non-negativity of  $A$ . These smoothness properties will be crucial in applying gradient descent arguments in Sect. 3.3, and are the main motivation for us to adopt the  $\|\cdot\|_A$  norm for our proposed algorithms.

**Lemma 2.7** *Let  $L \stackrel{\text{def}}{=} \frac{4}{\mu}$ . Then, for every  $x \geq 0$ , every  $i \in [n]$ , and every  $\lambda \in \left[-\frac{1}{L\|A_{:i}\|_\infty}, \frac{1}{L\|A_{:i}\|_\infty}\right]$ :*

- (a) *If  $|\nabla_i f_\mu(x)| \leq 1$ , then  $|\nabla_i f_\mu(x + \lambda \mathbf{e}_i) - \nabla_i f_\mu(x)| \leq L\|A_{:i}\|_\infty \cdot |\lambda|$ .*
- (b) *If  $\nabla_i f_\mu(x) \geq 1$ , then  $\nabla_i f_\mu(x + \lambda \mathbf{e}_i) \geq \left(1 - \frac{\|A_{:i}\|_\infty L}{2}\right) \nabla_i f_\mu(x)$ .*

*Proof of Lemma 2.7* Using the fact that  $\nabla_i f_\mu(x) > -1$  for all  $x$ , we have:

$$\begin{aligned} \left| \log \frac{\nabla_i f_\mu(x + \lambda \mathbf{e}_i) + 1}{\nabla_i f_\mu(x) + 1} \right| &\stackrel{\textcircled{1}}{=} \left| \int_0^\lambda \frac{\nabla_{ii}^2 f_\mu(x + v \mathbf{e}_i)}{\nabla_i f_\mu(x + v \mathbf{e}_i) + 1} dv \right| \\ &\stackrel{\textcircled{2}}{=} \frac{1}{\mu} \left| \int_0^\lambda \frac{(A^T \text{diag}\{p(x + v \mathbf{e}_i)\}A)_{ii}}{(A^T p(x + v \mathbf{e}_i))_i} dv \right| \\ &\stackrel{\textcircled{3}}{\leq} \frac{\|A_{:i}\|_\infty}{\mu} |\lambda| \stackrel{\textcircled{4}}{=} \frac{\|A_{:i}\|_\infty L}{4} |\lambda|. \end{aligned}$$

Above,  $\textcircled{1}$  holds because  $\int_0^\lambda g'(v)dv = g(\lambda) - g(0)$  where  $g(v) = \log(\nabla_i f_\mu(x + v \mathbf{e}_i) + 1)$ ;  $\textcircled{2}$  holds according to Fact 2.6;  $\textcircled{3}$  is because the numerator is  $\sum_j A_{j,i}^2 p_j$  while the denominator is  $\sum_j A_{j,i} p_j$ ;  $\textcircled{4}$  holds because  $L = \frac{4}{\mu}$ . This immediately implies

$$e^{-\frac{\|A_{:i}\|_\infty L}{4} |\lambda|} \leq \frac{\nabla_i f_\mu(x + \lambda \mathbf{e}_i) + 1}{\nabla_i f_\mu(x) + 1} \leq e^{\frac{\|A_{:i}\|_\infty L}{4} |\lambda|}.$$

Our assumption on  $\lambda$  implies  $\frac{\|A_{:i}\|_\infty L}{4} |\lambda| \leq \frac{1}{4}$ , so that we can use the approximation  $x \leq e^x - 1 \leq 1.2x$  over  $x \in [-\frac{1}{4}, \frac{1}{4}]$ . This yields the simpler bound:

$$-\frac{\|A_{:i}\|_\infty L}{4} |\lambda| \leq \frac{\nabla_i f_\mu(x + \lambda \mathbf{e}_i) - \nabla_i f_\mu(x)}{\nabla_i f_\mu(x) + 1} \leq 1.2 \frac{\|A_{:i}\|_\infty L}{4} |\lambda|.$$

- (a) Assuming that  $\nabla_i f_\mu(x) \in (-1, 1]$ , we have:

$$\left| \nabla_i f_\mu(x + \lambda \mathbf{e}_i) - \nabla_i f_\mu(x) \right| \leq 2.4 \cdot \frac{\|A_{:i}\|_\infty L}{4} |\lambda| \leq \|A_{:i}\|_\infty L |\lambda|.$$

(b) Assuming  $\nabla_i f_\mu(x) \geq 1$ , we have

$$\begin{aligned} \nabla_i f_\mu(x + \lambda \mathbf{e}_i) &\geq \nabla_i f_\mu(x) - \frac{\|A_{:i}\|_\infty L}{4} |\lambda| (\nabla_i f_\mu(x) + 1) \\ &\geq \left(1 - \frac{\|A_{:i}\|_\infty L}{2} |\lambda|\right) \nabla_i f_\mu(x). \end{aligned}$$

□

**Initialization** Iterative methods require a starting point, and we use the following one

**Fact 2.8** Let  $x_i^{\text{start}} \stackrel{\text{def}}{=} \frac{1-\varepsilon/2}{n\|A_{:i}\|_\infty}$  for each  $i \in [n]$ . Then,  $x^{\text{start}} \in \Delta_{\text{box}}$  and  $f_\mu(x^{\text{start}}) \leq -\frac{1-\varepsilon}{n}$ .

*Proof* Using the fact that  $Ax^{\text{start}} - \mathbb{1} \leq -\varepsilon/2 \cdot \mathbb{1}$ , we compute  $f_\mu(x^{\text{start}})$  as follows:

$$\begin{aligned} f_\mu(x^{\text{start}}) &= \mu \sum_j e^{\frac{1}{\mu}((Ax^{\text{start}})_j - 1)} - \mathbb{1}^T x^{\text{start}} \leq \mu \sum_j e^{-\frac{\varepsilon/2}{\mu}} - \frac{1 - \varepsilon/2}{n} \\ &\leq \frac{\mu m}{(nm)^2} - \frac{1 - \varepsilon/2}{n} \leq -\frac{1 - \varepsilon}{n}. \end{aligned}$$

Above, we have used  $\mathbb{1}^T x^{\text{start}} \geq x_i^{\text{start}} = \frac{1-\varepsilon/2}{n}$ , where  $i$  is the column s.t.  $\|A_{:i}\|_\infty = 1$ . □

### 3 Our packing LP solver

Recall traditional (accelerated or not) gradient descent [28,29] or coordinate descent [6,17,30] should not be applied directly to minimize  $f_\mu$ , because  $f_\mu$  is not Lipschitz-smooth.

Our proposed algorithm *PacLPSolver* starts with some initial vector  $\mathbf{x}_0 = \mathbf{y}_0 = x^{\text{start}}$  (see Fact 2.8) and  $\mathbf{z}_0 = 0$ , and is divided into  $T$  iterations. In each iteration  $k$ , it computes a weighted midpoint  $\mathbf{x}_k \leftarrow \tau \mathbf{z}_{k-1} + (1 - \tau) \mathbf{y}_{k-1}$  for some parameter  $\tau \in (0, 1)$ . This step is analogous to that in traditional accelerated coordinate descent [6, 17, 30]. We then compute  $\mathbf{y}_k$  and  $\mathbf{z}_k$  as follows.

We select  $i \in [n]$  uniformly at random. Let  $\xi_k^{(i)} = (0, \dots, 0, \mathbb{T}^p(v), 0, \dots, 0)$  be the vector that is only non-zero at coordinate  $i$ , where  $v = \nabla_i f_\mu(\mathbf{x}_k) \in [-1, \infty)$ , and  $\mathbb{T}^p(v)$  is the thresholding function  $\mathbb{T}^p(v) \stackrel{\text{def}}{=} \min\{v, 1\}$ . We refer to  $\xi_k^{(i)}$  as the *truncated gradient*.<sup>9</sup> Next,

- Perform a *mirror (descent) step*  $\mathbf{z}_k \leftarrow \mathbf{z}_k^{(i)} \stackrel{\text{def}}{=} \operatorname{argmin}_{\mathbf{z} \in \Delta_{\text{box}}} \left\{ \frac{1}{2} \|\mathbf{z} - \mathbf{z}_{k-1}\|_A^2 + \langle n\alpha_k \xi_k^{(i)}, \mathbf{z} \rangle \right\}$  for some parameter  $\alpha_k \ll 1/n$  to be chosen later.

<sup>9</sup> A similar gradient truncation was developed in our prior work [4], but for a different purpose (to ensure parallelism) and not applied to coordinate gradient. The truncation idea of this paper also inspired later works in matrix scaling [2] and in SDP [1].

**Algorithm 1** *PacLP Solver* ( $A, x^{\text{start}}, \varepsilon$ )

**Input:**  $A \in \mathbb{R}_{\geq 0}^{m \times n}, x^{\text{start}} \in \Delta_{\text{box}}, \varepsilon \in (0, 1/30)$ .

**Output:**  $x \in \Delta_{\text{box}}$ .

- 1:  $\mu \leftarrow \frac{\varepsilon}{4 \log(nm/\varepsilon)}, L \leftarrow \frac{4}{\mu}, \tau \leftarrow \frac{1}{3 \cdot nL}$  and  $\alpha_0 \leftarrow \frac{1}{nL}$ . ▷ recall  $\Delta_{\text{box}} \stackrel{\text{def}}{=} \{x \in \mathbb{R}^n : x_i \in [0, \frac{1}{\|A_{:i}\|_\infty}]\}$
- 2:  $T \leftarrow \lceil 3nL \log(1/\varepsilon) \rceil = O(n \cdot \frac{\log(nm/\varepsilon) \cdot \log(1/\varepsilon)}{\varepsilon})$ . ▷ parameters
- 3:  $x_0 = y_0 \leftarrow x^{\text{start}}, z_0 \leftarrow 0$ . ▷ number of iterations
- 4: **for**  $k \leftarrow 1$  **to**  $T$  **do**
- 5:    $\alpha_k \leftarrow \frac{1}{1-\tau} \alpha_{k-1}$
- 6:    $x_k \leftarrow \tau z_{k-1} + (1-\tau)y_{k-1}$ .
- 7:   Randomly select  $i \in [n]$  uniformly at random.
- 8:   Define vector  $\xi_k^{(i)}$  to be all-zero except at coordinate  $i$ , where  $\xi_{k,i}^{(i)} = \min\{1, \nabla_i f_\mu(x_k)\}$ .
- 9:    $z_k \leftarrow z_k^{(i)} \stackrel{\text{def}}{=} \operatorname{argmin}_{z \in \Delta_{\text{box}}} \left\{ \frac{1}{2} \|z - z_{k-1}\|_A^2 + \langle n\alpha_k \xi_k^{(i)}, z \rangle \right\}$ . ▷ See Proposition 3.2
- 10:    $y_k \leftarrow y_k^{(i)} \stackrel{\text{def}}{=} x_k + \frac{1}{n\alpha_k L} (z_k^{(i)} - z_{k-1})$ .
- 11: **end for**
- 12: **return**  $y_T$ .

– Perform a *gradient (descent) step*  $y_k \leftarrow y_k^{(i)} \stackrel{\text{def}}{=} x_k + \frac{1}{n\alpha_k L} (z_k^{(i)} - z_{k-1})$ .

This finishes the description of our *PacLP Solver*.

*Remark 3.1* We use the superscript  $(i)$  on  $\xi_k^{(i)}, y_k^{(i)}$  and  $z_k^{(i)}$  to emphasize that the value depends on the choice of  $i$ . We use generic parameters  $\tau, \alpha_k, T$  in the above description and their precise values are presented in Algorithm 1.

Our update on  $y_k$  is a “gradient descent step” because we shall prove that it strictly decreases the objective (i.e.,  $f_\mu(x_k) - f_\mu(y_k^{(i)}) \geq 0$ ). Our update on  $z_k$  is a “mirror descent step” because we shall apply standard mirror descent analysis [12] to it. We explicitly describe how to implement the mirror step (its proof is straightforward by computing the gradient):

**Proposition 3.2** *If  $\Delta_{\text{box}} = \{x \in \mathbb{R}^n : x_i \in [0, \frac{C}{\|A_{:i}\|_\infty}]\}$  for some constant  $C > 0$ , the minimizer  $z = \operatorname{argmin}_{z \in \Delta_{\text{box}}} \left\{ \frac{1}{2} \|z - z_{k-1}\|_A^2 + \langle \delta e_i, z \rangle \right\}$  for any  $\delta \in \mathbb{R}$  and basis vector  $e_i$  can be computed as follows:*

1.  $z \leftarrow z_{k-1}$ .
2.  $z_i \leftarrow z_i - \delta / \|A_{:i}\|_\infty$ .
3. *If  $z_i < 0$ , then  $z_i \leftarrow 0$ ; if  $z_i > C / \|A_{:i}\|_\infty$ ,  $z_i \leftarrow C / \|A_{:i}\|_\infty$ .*
4. *Return  $z$ .*

We also point out that

**Lemma 3.3** *Each iteration of *PacLP Solver* can be implemented to run in expected  $O(N/n)$  time. The total expected running time is  $O(TN/n)$ .*

Lemma 3.3 is not hard to prove, but anyways included in “Appendix E.5”. It follows from standard implementation tricks which compute  $x_k$  and  $y_k$  only *implicitly*: that is to express  $x_k$  and  $y_k$  as linear combinations of two less-frequently-updated vectors.

### 3.1 Convergence statement

In this section, we focus on proving the following main theorem.

**Theorem 3.4** *PacLP Solver*  $(A, x^{\text{start}}, \varepsilon)$  outputs some  $y_T$  satisfying

$$\mathbf{E}[f_\mu(y_T)] \leq -(1 - 3\varepsilon)\text{OPT}.$$

It is straightforward to use Markov’s bound to turn Theorem 3.4 into a probabilistic one

**Corollary 3.5** With probability at least  $2/3$ , the output  $y_T = \text{PacLP Solver}(A, x^{\text{start}}, \varepsilon)$  satisfies that  $\frac{y_T}{1+\varepsilon}$  is a  $(1 - O(\varepsilon))$  approximate solution to the packing LP program. The expected running time is  $O(\frac{\log(nm/\varepsilon)\log(1/\varepsilon)}{\varepsilon} N)$ .

*Proof of Corollary 3.5* Since for every  $x \in \Delta_{\text{box}}$  it satisfies  $f_\mu(x) \geq -(1 + \varepsilon)\text{OPT}$  according to Proposition 2.4.b, we obtain that  $f_\mu(y_T) + (1 + \varepsilon)\text{OPT}$  is a random variable that is non-negative, whose expectation  $\mathbf{E}[f_\mu(y_T) + (1 + \varepsilon)\text{OPT}] \leq 4\varepsilon$  according to Theorem 3.4. By Markov bound, with at least probability  $2/3$ , we obtain some  $y_T$  satisfying  $f_\mu(y_T) \leq -(1 - 11\varepsilon)\text{OPT}$ , which yields a  $(1 - O(\varepsilon))$  approximate solution to packing LP according to Proposition 2.4.c. The running time follows from Lemma 3.3. □

Before we prove Theorem 3.4 in subsequent subsections, let us first point out that our iterates  $x_k, y_k, z_k$  never leave the bounding box  $\Delta_{\text{box}}$ :

**Lemma 3.6** We have  $x_k, y_k, z_k \in \Delta_{\text{box}}$  for all  $k = 0, 1, \dots, T$ .

(The proof of Lemma 3.6 is included in “Appendix A.1”, and the main technique already appeared in randomized coordinate descent [17].)

### 3.2 Step 1: mirror descent guarantee

Following almost classical analysis of mirror descent (cf. textbook [12]), our update  $z_k^{(i)} = \text{argmin}_{z \in \Delta_{\text{box}}} \{ \frac{1}{2} \|z - z_{k-1}\|_A^2 + \langle n\alpha_k \xi_k^{(i)}, z \rangle \}$  satisfies

**Lemma 3.7** (mirror descent) For every  $u \in \Delta_{\text{box}}$ , it satisfies

$$\langle n\alpha_k \xi_k^{(i)}, z_{k-1} - u \rangle \leq n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, x_k - y_k^{(i)} \rangle + \frac{1}{2} \|z_{k-1} - u\|_A^2 - \frac{1}{2} \|z_k^{(i)} - u\|_A^2.$$

*Proof* Denoting by  $V_a(b) = \frac{1}{2} \|b - a\|_A^2$  as a function of  $b \in \Delta_{\text{box}}$  parameterized at  $a \in \Delta_{\text{box}}$ , we have  $\nabla_i V_a(b) = \|A_{\cdot i}\|_\infty \cdot (a_i - b_i)$ . In the optimization language,  $V_a(b)$  is the Bregman divergence of the  $\|\cdot\|_A^2$  regularizer [12]. We derive the following sequence of inequalities:

$$\begin{aligned} \langle n\alpha_k \xi_k^{(i)}, z_{k-1} - u \rangle &= \langle n\alpha_k \xi_k^{(i)}, z_{k-1} - z_k^{(i)} \rangle + \langle n\alpha_k \xi_k^{(i)}, z_k^{(i)} - u \rangle \\ &\stackrel{\textcircled{1}}{\leq} \langle n\alpha_k \xi_k^{(i)}, z_{k-1} - z_k^{(i)} \rangle + \langle -\nabla V_{z_{k-1}}(z_k^{(i)}), z_k^{(i)} - u \rangle \end{aligned}$$

$$\begin{aligned}
 &\stackrel{\textcircled{2}}{=} \langle n\alpha_k \xi_k^{(i)}, \mathbf{z}_{k-1} - \mathbf{z}_k^{(i)} \rangle - \frac{1}{2} \|\mathbf{z}_{k-1} - \mathbf{z}_k^{(i)}\|_A^2 + \frac{1}{2} \|\mathbf{z}_{k-1} - u\|_A^2 \\
 &\quad - \frac{1}{2} \|\mathbf{z}_k^{(i)} - u\|_A^2 \\
 &\stackrel{\textcircled{3}}{=} n^2 \alpha_k^2 L \left( \langle \xi_k^{(i)}, \mathbf{x}_k - \mathbf{y}_k \rangle - \frac{L}{2} \|\mathbf{x}_k - \mathbf{y}_k\|_A^2 \right) + \frac{1}{2} \|\mathbf{z}_{k-1} - u\|_A^2 \\
 &\quad - \frac{1}{2} \|\mathbf{z}_k^{(i)} - u\|_A^2 \\
 &\leq n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, \mathbf{x}_k - \mathbf{y}_k \rangle + \frac{1}{2} \|\mathbf{z}_{k-1} - u\|_A^2 - \frac{1}{2} \|\mathbf{z}_k^{(i)} - u\|_A^2.
 \end{aligned}$$

Above, ① is due to the minimality of  $\mathbf{z}_k^{(i)} = \operatorname{argmin}_{z \in \Delta_{\text{box}}} \{V_{\mathbf{z}_{k-1}}(z) + \langle n\alpha_k \xi_k^{(i)}, z \rangle\}$ , which implies that  $\langle \nabla V_{\mathbf{z}_{k-1}}(\mathbf{z}_k^{(i)}) + n\alpha_k \xi_k^{(i)}, u - \mathbf{z}_k^{(i)} \rangle \geq 0$  for all  $u \in \Delta_{\text{box}}$ . Equality ② can be checked for every coordinate  $\ell \in [n]$  as follows:

$$\begin{aligned}
 -\nabla_\ell V_{\mathbf{z}_{k-1}}(\mathbf{z}_k^{(i)}) \cdot (\mathbf{z}_{k,\ell}^{(i)} - u_\ell) &= \|A_{:i}\|_\infty (\mathbf{z}_{k-1,\ell} - \mathbf{z}_{k,\ell}^{(i)}) \cdot (\mathbf{z}_{k,\ell}^{(i)} - u_\ell) \\
 &= \|A_{:i}\|_\infty \left( -\frac{1}{2} (\mathbf{z}_{k-1,\ell} - \mathbf{z}_{k,\ell}^{(i)})^2 + \frac{1}{2} (u_\ell - \mathbf{z}_{k-1,\ell})^2 \right. \\
 &\quad \left. - \frac{1}{2} (\mathbf{z}_{k,\ell}^{(i)} - u_\ell)^2 \right).
 \end{aligned}$$

③ is by our choice of  $\mathbf{y}_k$  which satisfies that  $\mathbf{z}_{k-1} - \mathbf{z}_k^{(i)} = n\alpha_k L(\mathbf{x}_k - \mathbf{y}_k^{(i)})$ . □

In addition, as a simple corollary of Proposition 3.2, we have the following fact

**Fact 3.8**  $|\mathbf{z}_{k,i}^{(i)} - \mathbf{z}_{k-1,i}| \leq \frac{n\alpha_k |\xi_{k,i}^{(i)}|}{\|A_{:i}\|_\infty}$  and  $|\mathbf{y}_{k,i}^{(i)} - \mathbf{x}_{k,i}| = \frac{1}{n\alpha_k L} |\mathbf{z}_{k,i}^{(i)} - \mathbf{z}_{k-1,i}| \leq \frac{|\xi_{k,i}^{(i)}|}{L\|A_{:i}\|_\infty} \leq \frac{1}{L\|A_{:i}\|_\infty}$ . If  $\xi_{k,i}^{(i)} \geq 0$ , then  $\mathbf{z}_{k,i}^{(i)} \leq \mathbf{z}_{k-1,i}$  and  $\mathbf{y}_{k,i}^{(i)} \leq \mathbf{x}_{k,i}$ ; if  $\xi_{k,i}^{(i)} \leq 0$ , then  $\mathbf{z}_{k,i}^{(i)} \geq \mathbf{z}_{k-1,i}$  and  $\mathbf{y}_{k,i}^{(i)} \geq \mathbf{x}_{k,i}$ .

### 3.3 Step 2: gradient descent guarantee

We call our update  $\mathbf{y}_k^{(i)} \leftarrow \mathbf{x}_k + \frac{1}{n\alpha_k L} (\mathbf{z}_k^{(i)} - \mathbf{z}_{k-1})$  a gradient descent step, because the following lemma guarantees  $f_\mu(\mathbf{y}_k^{(i)}) \leq f_\mu(\mathbf{x}_k)$ , that is, the objective only decreases; moreover, the objective decreases at least by  $\frac{1}{2} \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle$ .

**Lemma 3.9** (gradient descent) *We have  $f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)}) \geq \frac{1}{2} \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \geq 0$ .*

This Lemma 3.9, which is characteristic of the PC-LP setting, is strong in the following sense. Even though the update  $\mathbf{y}_k^{(i)}$  only depends on the truncated gradient  $\xi_k^{(i)}$ , the progress we make is a function of the true gradient  $\nabla_i f_\mu(\mathbf{x}_k)$ , including the large component that was discarded. This is possible because the smoothness guarantee of Lemma 2.7.b allows us to take a long coordinate step even though  $f_\mu(x)$  is not Lipschitz-smooth.

*Proof of Lemma 3.9* Using Fact 3.8, we write  $\mathbf{y}_k^{(i)} = \mathbf{x}_k + s\lambda\mathbf{e}_i$  for some  $s \in \{-1, +1\}$  and step length  $\lambda \in [0, \frac{1}{L\|A_{:i}\|_\infty}]$ . We first focus on the case  $\nabla_i f_\mu(\mathbf{x}_k) \in [-1, 1]$  so  $\xi_{k,i}^{(i)} = \nabla_i f_\mu(\mathbf{x}_k)$ .

$$\begin{aligned} f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)}) &= f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{x}_k + s\lambda\mathbf{e}_i) = -s \int_0^\lambda (\nabla_i f_\mu(\mathbf{x}_k + s\chi\mathbf{e}_i)) d\chi \\ &\stackrel{\textcircled{1}}{\geq} s \int_0^\lambda (-\nabla_i f_\mu(\mathbf{x}_k) - L\|A_{:i}\|_\infty \cdot s\chi) d\chi = -\nabla_i f_\mu(\mathbf{x}_k) \cdot s\lambda - \frac{L\|A_{:i}\|_\infty}{2} \cdot \lambda^2 \\ &\stackrel{\textcircled{2}}{\geq} -\nabla_i f_\mu(\mathbf{x}_k) \cdot s\lambda - \frac{L\|A_{:i}\|_\infty}{2} \cdot \lambda \cdot \frac{|\xi_{k,i}^{(k)}|}{L\|A_{:i}\|_\infty} \stackrel{\textcircled{3}}{=} -\frac{1}{2} \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{y}_k^{(i)} - \mathbf{x}_k \rangle. \end{aligned}$$

Above,  $\textcircled{1}$  uses Lemma 2.7.a,  $\textcircled{2}$  uses Fact 3.8, and  $\textcircled{3}$  uses  $|\xi_{k,i}^{(k)}| = -s\nabla_i f_\mu(\mathbf{x}_k)$  (see also Fact 3.8). Next, we turn to the case of  $\nabla_i f_\mu(\mathbf{x}_k) > 1$ . In this case, we have  $s = -1$  and

$$\begin{aligned} f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)}) &= f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{x}_k - \lambda\mathbf{e}_i) = \int_0^\lambda \nabla_i f_\mu(\mathbf{x}_k - \chi\mathbf{e}_i) d\chi \\ &\stackrel{\textcircled{1}}{\geq} \int_0^\lambda \left(1 - \frac{\|A_{:i}\|_\infty L}{2} \chi\right) \nabla_i f_\mu(x) d\chi \\ &\stackrel{\textcircled{2}}{\geq} \int_0^\lambda \frac{1}{2} \nabla_i f_\mu(x) d\chi = \frac{1}{2} \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle. \end{aligned}$$

Above,  $\textcircled{1}$  uses Lemma 2.7.b and  $\textcircled{2}$  uses  $\chi \leq \lambda \leq \frac{1}{L\|A_{:i}\|_\infty}$ . Finally, we have  $\langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \geq 0$  because  $\nabla_i f_\mu(\mathbf{x}_k)$  and  $\mathbf{x}_{k,i} - \mathbf{y}_{k,i}^{(i)}$  have the same sign, and  $\mathbf{x}_{k,\ell} = \mathbf{y}_{k,\ell}^{(i)}$  for  $\ell \neq i$ .

### 3.4 Step 3: putting all together

We denote by  $\eta_k^{(i)} \in \mathbb{R}_{\geq 0}^n$  the vector that is only non-zero at coordinate  $i$ , and satisfies  $\eta_{k,i}^{(i)} = \nabla_i f_\mu(\mathbf{x}_k) - \xi_{k,i}^{(i)} \in [0, \infty)$ . In other words, the full gradient

$$\nabla f_\mu(\mathbf{x}_k) = \mathbf{E}_i[(0, \dots, n\nabla_i f_\mu(\mathbf{x}_k), \dots, 0)] = \mathbf{E}_i[n\eta_k^{(i)} + n\xi_k^{(i)}]$$

can be (in expectation) decomposed into the a large non-negative component  $\eta_k^{(i)} \in [0, \infty)^n$  and a truncated component  $\xi_k^{(i)} \in [-1, 1]^n$ . Recall that  $\eta_k^{(i)}$  did not contribute to the descent steps (see Line 9 of *PacLP Solver*). Now, for any  $u \in \Delta_{\text{box}}$ , we can use a basic convexity argument and the mirror descent lemma to compute that

$$\begin{aligned} \alpha_k(f_\mu(\mathbf{x}_k) - f_\mu(u)) &\leq \langle \alpha_k \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - u \rangle \\ &= \langle \alpha_k \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{z}_{k-1} \rangle + \langle \alpha_k \nabla f_\mu(\mathbf{x}_k), \mathbf{z}_{k-1} - u \rangle \end{aligned}$$

$$\begin{aligned}
 &= \langle \alpha_k \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{z}_{k-1} \rangle + \mathbf{E}_i \left[ \langle n\alpha_k \eta_k^{(i)}, \mathbf{z}_{k-1} - u \rangle + \langle n\alpha_k \xi_k^{(i)}, \mathbf{z}_{k-1} - u \rangle \right] \\
 &\stackrel{\textcircled{1}}{=} \frac{(1 - \tau)\alpha_k}{\tau} \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{y}_{k-1} - \mathbf{x}_k \rangle \\
 &\quad + \mathbf{E}_i \left[ \langle n\alpha_k \eta_k^{(i)}, \mathbf{z}_{k-1} - u \rangle + \langle n\alpha_k \xi_k^{(i)}, \mathbf{z}_{k-1} - u \rangle \right] \tag{3.1}
 \end{aligned}$$

$$\begin{aligned}
 &\stackrel{\textcircled{2}}{\leq} \frac{(1 - \tau)\alpha_k}{\tau} (f_\mu(\mathbf{y}_{k-1}) - f_\mu(\mathbf{x}_k)) \\
 &\quad + \mathbf{E}_i \left[ \langle n\alpha_k \eta_k^{(i)}, \mathbf{z}_{k-1} - u \rangle + n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \right] \\
 &\quad + \frac{1}{2} \|\mathbf{z}_{k-1} - u\|_A^2 - \frac{1}{2} \|\mathbf{z}_k^{(i)} - u\|_A^2 \tag{3.2}
 \end{aligned}$$

Above,  $\textcircled{1}$  is because  $\mathbf{x}_k = \tau \mathbf{z}_{k-1} + (1 - \tau) \mathbf{y}_{k-1}$ , which implies that  $\tau(\mathbf{x}_k - \mathbf{z}_{k-1}) = (1 - \tau)(\mathbf{y}_{k-1} - \mathbf{x}_k)$ .  $\textcircled{2}$  uses convexity and Lemma 3.7. This above computation is motivated by [5], and as we shall see below, it allows one to linearly couple gradient and mirror steps.

Intuitively, the first term in the box of (3.2) is the loss introduced by the large gradient  $\eta_k^{(i)}$ . This part was truncated so did not contribute to the mirror step. The second term in the box is the loss introduced by mirror descent on the small gradient  $\xi_k^{(i)}$  in Lemma 3.7.

Now comes an important observation. As shown by Lemma 3.10 below, the performance of the gradient step—that is, the objective decrease of  $f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)})$ —is at least proportional to the total loss incurred in the box. Intuitively, this means that the progress in the gradient step is so large that it outweighs not only the loss from mirror descent (as is typical in accelerated gradient analyses [5,29]) but also the loss term introduced by  $\eta_k^{(i)}$ .

**Lemma 3.10** (gradient descent total guarantee) *For every  $u \geq 0$ ,*

$$\langle n\alpha_k \eta_k^{(i)}, \mathbf{z}_{k-1} - u \rangle + n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \leq 3n\alpha_k L \cdot (f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)})).$$

The proof of Lemma 3.10 is a careful case analysis and several simple applications of Lemma 3.9. We remark that to properly upper bound  $\langle n\alpha_k \eta_k^{(i)}, \mathbf{z}_{k-1} - u \rangle$ , one needs to have some good upper bound the coordinates of  $\mathbf{z}_{k-1}$ . This is exactly the place we need our redundant constraint which ensures  $\mathbf{z}_{k-1,i} \leq \frac{1}{\|A_i\|_\infty}$  (see Remark 2.5).

*Proof of Lemma 3.10* There are three possibilities:

- If  $\eta_{k,i}^{(i)} = 0$ , then we must have  $\xi_{k,i}^{(i)} = \nabla_i f_\mu(\mathbf{x}_k) \in [-1, 1]$ . Lemma 3.9 implies

$$\begin{aligned}
 &\langle n\alpha_k \eta_k^{(i)}, \mathbf{z}_{k-1} - u \rangle + n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \\
 &\quad = n^2 \alpha_k^2 L \cdot \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \leq 2n^2 \alpha_k^2 L \cdot (f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)}))
 \end{aligned}$$



- If  $\eta_{k,i}^{(i)} > 0$  and  $z_{k,i}^{(i)} > 0$ , then we precisely have  $z_{k,i}^{(i)} = z_{k-1,i} - \frac{n\alpha_k}{\|A_{:i}\|_\infty}$  (see Proposition 3.2), and accordingly  $y_{k,i}^{(i)} = x_{k,i} - \frac{1}{L\|A_{:i}\|_\infty} < x_{k,i}$ . In this case,

$$\begin{aligned} & \langle n\alpha_k \eta_k^{(i)}, z_{k-1} - u \rangle + n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, x_k - y_k^{(i)} \rangle \\ & \stackrel{\textcircled{1}}{\leq} n\alpha_k \cdot \nabla_i f_\mu(x_k) \cdot \frac{1}{\|A_{:i}\|_\infty} + n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, x_k - y_k^{(i)} \rangle \\ & \stackrel{\textcircled{2}}{<} n\alpha_k \cdot \nabla_i f_\mu(x_k) \cdot \frac{1}{\|A_{:i}\|_\infty} + n^2 \alpha_k^2 L \cdot \langle \nabla f_\mu(x_k), x_k - y_k^{(i)} \rangle \\ & \stackrel{\textcircled{3}}{=} n\alpha_k L \cdot \langle \nabla f_\mu(x_k), x_k - y_k^{(i)} \rangle + n^2 \alpha_k^2 L \cdot \langle \nabla f_\mu(x_k), x_k - y_k^{(i)} \rangle \\ & \stackrel{\textcircled{4}}{\leq} (2n\alpha_k L + 2n^2 \alpha_k^2 L) \cdot (f_\mu(x_k) - f_\mu(y_k^{(i)})). \end{aligned}$$

- Above,  $\textcircled{1}$  follows from the fact that  $z_{k-1} \in \Delta_{\text{box}}$  and therefore  $z_{k-1,i} \leq \frac{1}{\|A_{:i}\|_\infty}$  by the definition of  $\Delta_{\text{box}}$ , and  $u \geq 0$ ;  $\textcircled{2}$  follows from the fact that  $x_k$  and  $y_k^{(i)}$  are only different at coordinate  $i$ , and  $\xi_{k,i}^{(i)} = 1 < \nabla_i f_\mu(x_k)$  (since  $\eta_{k,i}^{(i)} > 0$ );  $\textcircled{3}$  follows from the fact that  $y_k^{(i)} = x_k - \frac{e_i}{L\|A_{:i}\|_\infty}$ ; and  $\textcircled{4}$  uses Lemma 3.9.
- If  $\eta_{k,i}^{(i)} > 0$  and  $z_{k,i}^{(i)} = 0$ , then we have

$$\begin{aligned} & \langle n\alpha_k \eta_k^{(i)}, z_{k-1} - u \rangle + n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, x_k - y_k^{(i)} \rangle \\ & \stackrel{\textcircled{1}}{\leq} (n\alpha_k \nabla_i f_\mu(x_k) \cdot z_{k-1,i}) + n^2 \alpha_k^2 L \cdot \langle \nabla f_\mu(x_k), x_k - y_k^{(i)} \rangle \\ & \stackrel{\textcircled{2}}{=} \langle n\alpha_k \nabla f_\mu(x_k), z_{k-1} - z_k^{(i)} \rangle + n^2 \alpha_k^2 L \cdot \langle \nabla f_\mu(x_k), x_k - y_k^{(i)} \rangle \\ & \stackrel{\textcircled{3}}{=} n^2 \alpha_k^2 L \cdot \langle \nabla f_\mu(x_k), x_k - y_k^{(i)} \rangle + n^2 \alpha_k^2 L \cdot \langle \nabla f_\mu(x_k), x_k - y_k^{(i)} \rangle \\ & \stackrel{\textcircled{4}}{\leq} 4n^2 \alpha_k^2 L \cdot (f_\mu(x_k) - f_\mu(y_k^{(i)})). \end{aligned}$$

- Above,  $\textcircled{1}$  is because  $u \geq 0$ ,  $\nabla_i f_\mu(x_k) = \eta_{k,i}^{(i)} + 1 > \eta_{k,i}^{(i)}$  and  $\nabla_i f_\mu(x_k) > \xi_{k,i}^{(i)}$ ;  $\textcircled{2}$  uses the assumption that  $z_{k,i}^{(i)} = 0$  and the fact that  $z_{k-1,\ell} = z_{k,\ell}^{(i)}$  for every  $\ell \neq i$ ;  $\textcircled{3}$  is from our choice of  $y_k$  which satisfies that  $z_{k-1} - z_k^{(i)} = n\alpha_k L(x_k - y_k^{(i)})$ ; and  $\textcircled{4}$  uses Lemma 3.9.

Combining the three cases, and using the fact that  $f_\mu(x_k) - f_\mu(y_k^{(i)}) \geq 0$ , we conclude that

$$\begin{aligned} & \langle n\alpha_k \eta_k^{(i)}, z_{k-1} - u \rangle + n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, x_k - y_k^{(i)} \rangle \\ & \leq (2n\alpha_k L + 4n^2 \alpha_k^2 L) \cdot (f_\mu(x_k) - f_\mu(y_k^{(i)})) \\ & \leq 3n\alpha_k L \cdot (f_\mu(x_k) - f_\mu(y_k^{(i)})). \end{aligned}$$

Above, the last inequality uses our choice of  $\alpha_k$ , which implies  $n\alpha_k \leq n\alpha_T = \frac{1}{\varepsilon L} \leq \frac{1}{4}$ .

Plugging Lemma 3.10 back to (3.2), we have

$$\begin{aligned}
 \alpha_k(f_\mu(\mathbf{x}_k) - f_\mu(u)) &\leq (\alpha_k \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - u) \\
 &\stackrel{\textcircled{1}}{\leq} \frac{(1 - \tau)\alpha_k}{\tau} (f_\mu(\mathbf{y}_{k-1}) - f_\mu(\mathbf{x}_k)) \\
 &\quad + \mathbf{E}_i \left[ 3n\alpha_k L \cdot (f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)})) + \frac{1}{2} \|\mathbf{z}_{k-1} - u\|_A^2 - \frac{1}{2} \|\mathbf{z}_k - u\|_A^2 \right] \\
 &\stackrel{\textcircled{2}}{\leq} \alpha_k f_\mu(\mathbf{x}_k) + (3n\alpha_k L - \alpha_k) f_\mu(\mathbf{y}_{k-1}) \\
 &\quad + \mathbf{E}_i \left[ -3n\alpha_k L \cdot f_\mu(\mathbf{y}_k^{(i)}) + \frac{1}{2} \|\mathbf{z}_{k-1} - u\|_A^2 - \frac{1}{2} \|\mathbf{z}_k - u\|_A^2 \right]. \tag{3.3}
 \end{aligned}$$

Above,  $\textcircled{1}$  uses Lemma 3.10; and  $\textcircled{2}$  is because we have chosen  $\tau$  to satisfy  $\frac{1}{\tau} = 3nL$ .

Next, recall that we have picked  $\alpha_k$  so that  $(3nL - 1)\alpha_k = 3nL \cdot \alpha_{k-1}$  in Algorithm 1. Telescoping (3.3) for  $k = 1, \dots, T$  and choosing  $u^* = (1 - \varepsilon/2)x^*$ , we have

$$\begin{aligned}
 -\sum_{k=1}^T \alpha_k f_\mu(u^*) &\leq 3f_\mu(\mathbf{y}_0) - 3n\alpha_T L \cdot \mathbf{E}[f_\mu(\mathbf{y}_T)] + \|\mathbf{z}_0 - u^*\|_A^2 \\
 &\leq -3n\alpha_T L \cdot \mathbf{E}[f_\mu(\mathbf{y}_T)] + \text{OPT}.
 \end{aligned}$$

Here, the second inequality is due to  $f_\mu(\mathbf{y}_0) = f_\mu(x^{\text{start}}) \leq 0$  from Fact 2.8, and the fact that

$$\|\mathbf{z}_0 - u^*\|_A^2 = \|u^*\|_A^2 = \sum_{i=1}^n (u_i^*)^2 \cdot \|A_{:i}\|_\infty \leq \sum_{i=1}^n (x_i^*)^2 \cdot \|A_{:i}\|_\infty \leq \sum_{i=1}^n x_i^* = \text{OPT}.$$

Finally, using the fact that  $\sum_{k=1}^T \alpha_k = \alpha_T \cdot \sum_{k=0}^{T-1} (1 - \frac{1}{3nL})^k = 3n\alpha_T L (1 - (1 - \frac{1}{3nL})^T)$ , we rearrange and obtain that

$$\begin{aligned}
 \mathbf{E}[f_\mu(\mathbf{y}_T)] &\leq \frac{\sum_k \alpha_k}{3n\alpha_T L} f_\mu(u^*) + \frac{1}{3n\alpha_T L} \text{OPT} = (1 - (1 - \frac{1}{3nL})^T) f_\mu(u^*) \\
 &\quad + \frac{1}{3n\alpha_T L} \text{OPT}.
 \end{aligned}$$

We choose  $T = \lceil 3nL \log(1/\varepsilon) \rceil$  so that  $\frac{1}{n\alpha_T L} = (1 - \frac{1}{3nL})^T \leq \varepsilon$ . Combining this with the fact that  $f_\mu(u^*) \leq -(1 - \varepsilon)\text{OPT} < 0$  (see Proposition 2.4.a), we obtain

$$\mathbf{E}[f_\mu(\mathbf{y}_T)] \leq (1 - \varepsilon) f_\mu(u^*) + \varepsilon/3 \cdot \text{OPT} < -(1 - 3\varepsilon)\text{OPT}.$$

Therefore, we have finished proving Theorem 3.4. □

### 4 Relaxation of the covering linear program

Since *PacLP Solver* gives only an approximate packing LP solution, we cannot infer from it a dual covering LP solution. Therefore, we have to work on a relaxed version of covering LP directly. For input matrix  $A \in \mathbb{R}_{\geq 0}^{m \times n}$ , we rewrite the covering LP problem (1.2) as follows in order to be notationally close to packing LP:

$$\min_{x \geq 0} \{\mathbb{1}^T x : Ax \geq \mathbb{1}\}. \tag{4.1}$$

We denote by OPT the optimal value to this LP, and by  $x^*$  any of its optimal solutions. We say that  $x$  is a  $(1 + \varepsilon)$ -approximation for the covering LP if  $Ax \geq \mathbb{1}$  and  $\mathbb{1}^T x \leq (1 + \varepsilon)\text{OPT}$ .

Again, we use indices  $i \in [n]$  for the columns of  $A$ , and indices  $j \in [m]$  for the rows of  $A$ . We denote by  $A_{:i}$  the  $i$ -th column vector of  $A$ , and  $A_j$  the  $j$ -th row vector of  $A$ . We assume without loss of generality by simple scaling that <sup>10</sup>

$$\min_{j \in [m]} \{\|A_j\|_\infty\} = 1. \tag{4.2}$$

**Proposition 4.1** *The normalization (4.2) ensures  $\text{OPT} \in [1, m]$ .*

*Proof* Suppose that  $j^*$  is the row that achieves the smallest infinite norm  $\|A_j\|_\infty$  over all rows  $j \in [m]$ . Then, for any solution  $x \in \mathbb{R}_{\geq 0}^n$  satisfying  $\langle A_{j^*}, x \rangle \geq 1$ , we must have  $\mathbb{1}^T x \geq 1/\|A_{j^*}\|_\infty = 1$  using (4.2). On the other hand, we can construct a feasible solution  $x$  as follows. Initialize  $x = 0$ , and then for each row  $j$ , let us find the coordinate  $i$  that maximizes the value of  $A_{ij}$  among all columns  $i$ . Then, we increase  $x_i$  by  $1/A_{ij} = 1/\|A_j\|_\infty$ . After we have exhausted all the  $m$  rows, we arrive at some  $x \geq 0$  satisfying  $Ax \geq \mathbb{1}$  as well as  $\mathbb{1}^T x = \sum_j 1/\|A_j\|_\infty \leq m$ .  $\square$

In our covering LP solvers, we assume that an initial solution, achieving a constant approximation, is available to the algorithm. Such a solution can be obtained for instance by the covering LP solver from Young [36] with constant  $\epsilon$  in time  $O(N \log N)$ .

**Definition 4.2** Let  $x^\ddagger$  be a given 2-approximate solution to the covering problem given and let  $\text{OPT}' \stackrel{\text{def}}{=} \mathbb{1}^T x^\ddagger \in [\text{OPT}, 2\text{OPT}]$ . Without loss of generality, assume  $\text{OPT}' \geq 2$ .

We now introduce the smoothed objective  $f_\mu(x)$  we are going to minimize in order to solve covering LP. Symmetric to the case in the packing solver, this smoothed objective  $f_\mu(x)$  turns each row of the LP constraint  $Ax \geq \mathbb{1}$  into an exponential penalty function.

---

<sup>10</sup> If  $\min_{j \in [m]} \{\|A_j\|_\infty\} = 0$  then the covering LP is infeasible so we are done. Otherwise, if  $\min_{j \in [m]} \{\|A_j\|_\infty\} = v > 0$  we scale all entries of  $A$  by  $1/v$ , and scale OPT by  $v$ .

**Definition 4.3** Letting  $\mu \stackrel{\text{def}}{=} \frac{\varepsilon}{4 \log(nm/\varepsilon)}$ , we define the smoothed objective  $f_\mu(x)$  as

$$f_\mu(x) \stackrel{\text{def}}{=} \mu \sum_{j=1}^m e^{\frac{1}{\mu}(1-(Ax)_j)} + \mathbb{1}^T x.$$

**Fact 4.4**  $\nabla f_\mu(x) = \mathbb{1} - A^T p(x)$  and  $\nabla^2 f_\mu(x) = \frac{1}{\mu} A^T \text{diag}\{p(x)\}A$ , where  $p_j(x) \stackrel{\text{def}}{=} e^{\frac{1}{\mu}(1-(Ax)_j)}$ .

We present some properties about  $f_\mu(x)$ . They together imply that the minimum of  $f_\mu(x)$  is around  $\text{OPT}$ , and if one approximately finds the minimum of  $f_\mu(x)$  up to an additive error  $O(\varepsilon \text{OPT})$ , this corresponds to a  $(1 + O(\varepsilon))$ -approximate solution to the covering LP (4.1). The proofs are analogous to Sect. 2, and included in ‘‘Appendix B.2’’ for completeness’ sake.

**Proposition 4.5** (a)  $f_\mu(u^*) \leq (1 + \varepsilon)\text{OPT}$  for  $u^* \stackrel{\text{def}}{=} (1 + \varepsilon/2)x^*$ .

(b)  $f_\mu(x) \geq (1 - \varepsilon)\text{OPT}$  for every  $x \geq 0$ .

(c) For any  $x \geq 0$  satisfying  $f_\mu(x) \leq 2\text{OPT}$ , we must have  $Ax \geq (1 - \varepsilon)\mathbb{1}$ .

(d) If  $x \geq 0$  satisfies  $f_\mu(x) \leq (1 + \delta)\text{OPT}$  for some  $\delta \in [0, 1]$ , then  $\frac{1}{1-\varepsilon}x$  is a  $\frac{1+\delta}{1-\varepsilon}$ -approximate solution to the covering LP.

## 5 Our covering LP solver in the well-conditioned case

Recall in packing LPs, since it satisfies  $0 \leq x_i^* \leq \frac{1}{\|A_{:i}\|_\infty}$  (see Fact 2.1), we can minimize  $f_\mu$  over a bounding box  $\Delta_{\text{box}}$ . Unfortunately, it no longer satisfies  $x_i^* \leq \frac{1}{\|A_{:i}\|_\infty}$  in covering LPs, so one cannot directly turn *PacLP Solver* into its symmetric version to solve covering LP.

In this section, we show that this symmetric covering LP solver still solves all *well-behaved* covering LP instances. Specifically, we say the covering LP is well-behaved if:<sup>11</sup>

**Assumption 5.1** There exists some optimal covering LP solution  $x^*$  satisfying  $x_i^* \leq \frac{9}{\|A_{:i}\|_\infty}$ ; and the initial point  $x^\sharp$  satisfies  $x_i^\sharp \leq \frac{9}{\|A_{:i}\|_\infty}$ .

For instance, well-behaved instances naturally arise from those where the constraints  $Ax \geq \mathbb{1}$  are non-redundant. If the optimal covering LP solution  $x^*$  and the initial point  $x^\sharp$  satisfy  $\mathbb{1} \leq Ax^* \leq 9 \cdot \mathbb{1}$  and  $\mathbb{1} \leq Ax^\sharp \leq 9 \cdot \mathbb{1}$ , then Assumption 5.1 is satisfied.

Well-behaved covering LP problems immediately satisfy the following:

**Fact 5.2** Define  $\Delta_{\text{box}} \stackrel{\text{def}}{=} \{x \in \mathbb{R}^n : x_i \in [0, \frac{10}{\|A_{:i}\|_\infty}]\}$ . Under Assumption 5.1, we have  $u^* \stackrel{\text{def}}{=} (1 + \varepsilon/2)x^* \in \Delta_{\text{box}}$  and  $x^{\text{start}} \stackrel{\text{def}}{=} (1 + \varepsilon/2) \cdot x^\sharp \in \Delta_{\text{box}}$ . Also, it satisfies  $f_\mu(x^{\text{start}}) \leq 3\text{OPT}$ .

<sup>11</sup> The constant 9 in this section can be replaced with any other constant greater than 1.

**Algorithm 2** *CovLP Solver*<sup>wb</sup>( $A, x^{\text{start}}, \varepsilon$ )

**Input:**  $A \in \mathbb{R}_{>0}^{m \times n}, x^{\text{start}} \in \Delta_{\text{box}}, \varepsilon \in (0, 1/30)$ .

**Output:**  $x \in \Delta_{\text{box}}$ .

▷ recall  $\Delta_{\text{box}} \stackrel{\text{def}}{=} \{x \in \mathbb{R}^n : x_i \in [0, \frac{10}{\|A_{\cdot i}\|_\infty}]\}$

1:  $\mu \leftarrow \frac{\varepsilon}{4 \log(nm/\varepsilon)}, L \leftarrow \frac{4}{\mu}, \tau \leftarrow \frac{1}{21 \cdot nL}$  and  $\alpha_0 \leftarrow \frac{1}{nL}$ .

▷ parameters

2:  $T \leftarrow \lceil 21nL \log(1/\varepsilon) \rceil = O(n \cdot \frac{\log(nm/\varepsilon) \cdot \log(1/\varepsilon)}{\varepsilon})$ .

▷ number of iterations

3:  $x_0 = y_0 \leftarrow x^{\text{start}}, z_0 \leftarrow 0$ .

4: **for**  $k \leftarrow 1$  **to**  $T$  **do**

5:  $\alpha_k \leftarrow \frac{1}{1-\tau} \alpha_{k-1}$

6:  $x_k \leftarrow \tau z_{k-1} + (1 - \tau)y_{k-1}$ .

7: Randomly select  $i \in [n]$  uniformly at random.

8: Define vector  $\xi_k^{(i)}$  be all-zero except at coordinate  $i$ , where  $\xi_{k,i}^{(i)} = \max\{-1, \nabla_i f_\mu(x_k)\}$ .

9:  $z_k \leftarrow z_k^{(i)} \stackrel{\text{def}}{=} \operatorname{argmin}_{z \in \Delta_{\text{box}}} \{ \frac{1}{2} \|z - z_{k-1}\|_A^2 + \langle n\alpha_k \xi_k^{(i)}, z \rangle \}$ .

▷ See Proposition 3.2

10:  $y_k \leftarrow y_k^{(i)} \stackrel{\text{def}}{=} x_k + \frac{1}{n\alpha_k L} (z_k^{(i)} - z_{k-1})$ .

11: **end for**

12: **return**  $y_T$ .

*Proof* The claims  $u^*, x^{\text{start}} \in \Delta_{\text{box}}$  are trivial after noticing  $\varepsilon \leq 1/30$ . Using the fact that  $Ax^{\text{start}} - \mathbb{1} \geq (1 + \varepsilon/2)Ax^{\sharp} - \mathbb{1} \geq \varepsilon/2 \cdot \mathbb{1}$ , we compute  $f_\mu(x^{\text{start}})$  as follows:

$$\begin{aligned} f_\mu(x^{\text{start}}) &= \mu \sum_j e^{\frac{1}{\mu}(1 - (Ax^{\text{start}})_j)} + \mathbb{1}^T x^{\text{start}} \leq \mu \sum_j e^{\frac{-\varepsilon/2}{\mu}} + 2\text{OPT} \\ &\leq \frac{\mu m}{(nm)^2} + 2\text{OPT} < 3\text{OPT}. \end{aligned}$$

□

We now describe *CovLP Solver*<sup>wb</sup> (which is a symmetric variant of *PacLP Solver*) that solves well-behaved covering LP problems, see Algorithm 2. It starts with the initial vector  $x_0 = y_0 = x^{\text{start}}$  and  $z_0 = 0$ . Then, *CovLP Solver*<sup>wb</sup> is divided into  $T$  iterations. In each iteration  $k$ , it computes a weighted midpoint  $x_k \leftarrow \tau z_{k-1} + (1 - \tau)y_{k-1}$  for some parameter  $\tau \in (0, 1)$ , and then proceeds to compute  $y_k$  and  $z_k$  as follows.

We select  $i \in [n]$  uniformly at random. Let  $\xi_k^{(i)} = (0, \dots, 0, -\mathbb{T}^P(v), 0, \dots, 0)$  be the vector that is only non-zero at coordinate  $i$ , where  $v = -\nabla_i f_\mu(x_k) \in [-1, \infty)$ , and recall  $\mathbb{T}^P(v) \stackrel{\text{def}}{=} \min\{v, 1\}$ . We refer to  $\xi_k^{(i)}$  as the *truncated gradient*. Next,

- Perform a *mirror (descent) step*  $z_k \leftarrow z_k^{(i)} \stackrel{\text{def}}{=} \operatorname{argmin}_{z \in \Delta_{\text{box}}} \{ \frac{1}{2} \|z - z_{k-1}\|_A^2 + \langle n\alpha_k \xi_k^{(i)}, z \rangle \}$  for some parameter  $\alpha_k \ll 1/n$  to be chosen later.
- Perform a *gradient (descent) step*  $y_k \leftarrow y_k^{(i)} \stackrel{\text{def}}{=} x_k + \frac{1}{n\alpha_k L} (z_k^{(i)} - z_{k-1})$ .

This finishes the description of *CovLP Solver*<sup>wb</sup>. It is not surprising to deduce the following theorem similar to Theorem 3.4. We include its proof in ‘‘Appendix C.3’’ for completeness.

**Theorem 5.3** Under the well-behavior assumption 5.1 in the covering LP problem,  $CovLP\text{Solver}^{wb}(A, x^{\text{start}}, \varepsilon)$  outputs some  $y_T$  satisfying

$$\mathbf{E}[f_\mu(y_T)] \leq (1 + 4.6\varepsilon)\text{OPT}.$$

Again, using the same proof as Corollary 3.5, one can apply Markov’s bound to turn Theorem 5.3 into a probabilistic statement:

**Corollary 5.4** Under the well-behavior assumption 5.1 in the covering LP problem, with probability at least  $2/3$ ,  $y_T = CovLP\text{Solver}^{wb}(A, x^{\text{start}}, \varepsilon)$  satisfies that  $\frac{y_T}{1-\varepsilon}$  is a  $(1 + O(\varepsilon))$  approximate solution to covering LP. The expected running time is

$$O\left(\frac{\log(nm/\varepsilon) \log(1/\varepsilon)}{\varepsilon} N\right).$$

**Removing the well-behavior assumption** In subsequent work, after the conference presentation of this paper, Wang, Mahoney and Rao [34] showed the following theorem.

**Theorem 5.5** ([34]) Any covering LP with constraint matrix  $A \in \mathbb{R}_{\geq 0}^{m \times n}$  of sparsity  $N$  can be converted into an equivalent but well-behaved covering LP with matrix  $\tilde{A} \in \mathbb{R}^{m \times n \cdot O(\log(mn/\varepsilon))}$  and sparsity  $N \cdot O(\log(mn/\varepsilon))$ . The conversion takes time  $N \cdot O(\log(mn/\varepsilon))$ .

As a result, we can apply our covering solver  $CovLP\text{Solver}^{wb}$  to this modified LP and apply our Theorem 5.3 to solve any covering LP in expected time  $O\left(\frac{\log^2(nm/\varepsilon) \log(1/\varepsilon)}{\varepsilon} N\right)$ .

## 6 Our covering LP solver in the general case

In this section, we remove the well-behavior assumption and propose a different algorithm  $CovLP\text{Solver}$  to solve all covering LP instances. This algorithm introduces a factor  $1/\sqrt{\varepsilon}$  loss in the running time, but is a *direct* covering LP solver without using any reduction.

The main difference to  $PacLP\text{Solver}$  and  $CovLP\text{Solver}^{wb}$  is that, this time we abandon the box constraint and study the minimization of  $f_\mu(x)$  over a simplex

$$x \in \Delta_{\text{simplex}} \stackrel{\text{def}}{=} \{x \in \mathbb{R}^n : x_i \geq 0 \wedge \mathbb{1}^T x \leq 2\text{OPT}'\}.$$

Again, this constraint  $\mathbb{1}^T x \leq 2\text{OPT}'$  is redundant just like the old  $\Delta_{\text{box}}$  constraint for packing LP (recall Remark 2.5); however, it shall be used to make sure that our updates are always inside  $\Delta_{\text{simplex}}$ . It is a simple fact that

**Fact 6.1**  $u^* \stackrel{\text{def}}{=} (1 + \varepsilon/2)x^* \in \Delta_{\text{simplex}}$ .

Recall that the initial vector  $x^\sharp$  is defined in Definition 4.2, and  $\text{OPT}'$  is a crude approximation to  $\text{OPT}$ , satisfying  $\text{OPT}' \stackrel{\text{def}}{=} \mathbb{1}^T x^\sharp \in [\text{OPT}, 2\text{OPT}]$ . We choose different starting vector  $x^{\text{start}}$  from Sect. 5:

**Algorithm 3** *CovLP Solver* ( $A, x^{\text{start}}, \varepsilon$ )

**Input:**  $A \in \mathbb{R}_{\geq 0}^{m \times n}, x^{\text{start}} \in \Delta_{\text{simplex}}, \varepsilon \in (0, 1/30]$ .

**Output:**  $x \in \Delta_{\text{simplex}}$ .

- 1:  $\mu \leftarrow \frac{\varepsilon}{4 \log(nm/\varepsilon)}, \beta \leftarrow \sqrt{\varepsilon}, \tau \leftarrow \frac{\mu\beta}{12n}$ . ▷ parameters
- 2:  $T \leftarrow \lceil \frac{1}{\tau} \log(1/\varepsilon) \rceil = O(\frac{\log(nm/\varepsilon) \log(1/\varepsilon)}{\varepsilon^{1.5}} n)$ . ▷ number of iterations
- 3:  $\alpha_0 \leftarrow (1 - \tau)^T \frac{\varepsilon}{12n\beta}$  and  $\gamma \leftarrow \frac{\varepsilon}{6\beta}$ . ▷ so that  $\alpha_T = \frac{\varepsilon}{12n\beta}$  and  $\gamma = 2\alpha_T n$
- 4:  $x_0 = y_0 = z_0 \leftarrow x^{\text{start}}$ .
- 5: **for**  $k \leftarrow 1$  **to**  $T$  **do**
- 6:    $\alpha_k \leftarrow \frac{1}{1-\tau} \alpha_{k-1}$ .
- 7:    $x_k \leftarrow \tau z_{k-1} + (1 - \tau)y_{k-1}$ .
- 8:   Randomly select  $i$  uniformly at random from  $[n]$ .
- 9:   Define vector  $\xi_k^{(i)}$  to be all-zero except at coordinate  $i$ , where  $\xi_{k,i}^{(i)} = \max\{-\beta, \nabla_i f_\mu(x_k)\}$ .
- 10:    $z_k \leftarrow z_k^{(i)} \stackrel{\text{def}}{=} \operatorname{argmin}_{z \in \Delta_{\text{simplex}}} \{V_{z_{k-1}}(z) + ((1 + \gamma)n\alpha_k \xi_k^{(i)}, z)\}$ . ▷ See Proposition 6.4
- 11:   **if**  $\nabla_i f_\mu(x_k) < -\beta$  **then**
- 12:     Denote by  $\pi$  the permutation that satisfies  $A_{\pi(1),i} \leq \dots \leq A_{\pi(m),i}$ .
- 13:     Pick  $j^* \in [m]$  such that  $\begin{cases} \sum_{j < j^*} A_{\pi(j),i} \cdot p_{\pi(j)}(x_k) < 1 + \beta \\ \sum_{j \leq j^*} A_{\pi(j),i} \cdot p_{\pi(j)}(x_k) \geq 1 + \beta \end{cases}$  ▷  $j^* \in [m]$  always exists, see
- (6.1)
- 14:      $y_k \leftarrow y_k^{(i)} \stackrel{\text{def}}{=} x_k + \delta \cdot e_i$  where  $\delta = \frac{\mu\beta}{2A_{\pi(j^*),i}}$ .
- 15:   **else**
- 16:      $y_k \leftarrow y_k^{(i)} \stackrel{\text{def}}{=} x_k$ .
- 17:   **end if**
- 18: **end for**
- 19: **return**  $y_T$ .

**Proposition 6.2** *Letting  $x^{\text{start}} \stackrel{\text{def}}{=} (1 + \varepsilon/2) \cdot x^\sharp + (\frac{1}{n}, \dots, \frac{1}{n})$ , we have  $x^{\text{start}} \in \Delta_{\text{simplex}}$  and  $f_\mu(x^{\text{start}}) \leq 4\text{OPT}$ .*

*Proof* Using  $Ax^{\text{start}} - \mathbf{1} \geq (1 + \varepsilon/2)Ax^\sharp - \mathbf{1} \geq \varepsilon/2 \cdot \mathbf{1}$ , we compute  $f_\mu(x^{\text{start}})$  as follows:

$$\begin{aligned} f_\mu(x^{\text{start}}) &= \mu \sum_j e^{\frac{1}{\mu}(1 - (Ax^{\text{start}})_j)} + \mathbf{1}^T x^{\text{start}} \leq \mu \sum_j e^{\frac{-\varepsilon/2}{\mu}} + 2\text{OPT} + 1 \\ &\leq \frac{\mu m}{(nm)^2} + 3\text{OPT} < 4\text{OPT}. \end{aligned}$$

Also, we have  $\mathbf{1}^T x^{\text{start}} \leq (1 + \varepsilon/2)\text{OPT}' + 1 \leq 2\text{OPT}'$ . (Recall  $\text{OPT}' \geq 2$  in Definition 4.2.) □

Our proposed algorithm *CovLP Solver* starts with the initial vector  $x_0 = y_0 = z_0 = x^{\text{start}}$  and is divided into  $T$  iterations. In each iteration  $k$ , as usual, it computes a weighted midpoint  $x_k \leftarrow \tau z_{k-1} + (1 - \tau)y_{k-1}$  for some parameter  $\tau \in (0, 1)$ , and then computes  $y_k$  and  $z_k$  as follows.

We select  $i \in [n]$  uniformly at random, and let  $\xi_k^{(i)} = (0, \dots, 0, \mathbb{T}^c(v), 0, \dots, 0)$  be the vector that is only non-zero at coordinate  $i$ , where  $v = \nabla_i f_\mu(x_k) \in (-\infty, 1]$  and  $\mathbb{T}^c(v) \stackrel{\text{def}}{=} \max\{-\beta, v\}$  is the new thresholding function for some parameter  $\beta \stackrel{\text{def}}{=} \sqrt{\varepsilon}$ . Then,

- Perform a *mirror (descent) step*  $\mathbf{z}_k \leftarrow \mathbf{z}_k^{(i)} \stackrel{\text{def}}{=} \operatorname{argmin}_{z \in \Delta_{\text{simplex}}} \{V_{\mathbf{z}_{k-1}}(z) + \langle (1 + \gamma)n\alpha_k \xi_k^{(i)}, z \rangle\}$  for some positive parameters  $\gamma \ll 1$  and  $\alpha_k \ll 1/n$ , where

$$V_x(\mathbf{y}) \stackrel{\text{def}}{=} \sum_{i=1}^n y_i \log \frac{y_i}{x_i} + x_i - y_i$$

is the so-called Bregman divergence of the generalized entropy function (c.f. [12]).

- If  $\nabla_i f_\mu(x_k) < -\beta$ , perform a *gradient (descent) step*  $\mathbf{y}_k \leftarrow \mathbf{y}_k^{(i)} \stackrel{\text{def}}{=} \mathbf{x}_k + \delta \mathbf{e}_i$  for some  $\delta > 0$ . In practice, one can line-search over  $\delta$ , but we choose an explicit  $\delta$  as follows.
  - Denote by  $\pi$  the permutation that satisfies  $A_{\pi(1),i} \leq \dots \leq A_{\pi(m),i}$
  - Pick  $j^* \in [m]$  s.t.  $\begin{cases} \sum_{j < j^*} A_{\pi(j),i} \cdot p_{\pi(j)}(\mathbf{x}_k) < 1 + \beta \\ \sum_{j \leq j^*} A_{\pi(j),i} \cdot p_{\pi(j)}(\mathbf{x}_k) \geq 1 + \beta \end{cases}$ . Such  $j^*$  exists because

$$\sum_{j=1}^m A_{\pi(j),i} \cdot p_{\pi(j)}(\mathbf{x}_k) = \sum_{j=1}^m A_{ji} \cdot p_j(\mathbf{x}_k) = 1 - \nabla_i f_\mu(\mathbf{x}_k) \geq 1 + \beta. \tag{6.1}$$

- Set  $\delta = \frac{\mu\beta}{2A_{\pi(j^*),i}}$ .

This finishes the description of our *CovLP Solver*.

*Remark 6.3* We use the superscript  $(i)$  on  $\xi_k^{(i)}$ ,  $\mathbf{y}_k^{(i)}$  and  $\mathbf{z}_k^{(i)}$  to emphasize that the value depends on the choice of  $i$ . We have used generic parameters  $\tau, \alpha_k, T$  in the above description and their precise values are presented in Algorithm 3.

Our update on  $\mathbf{y}_k$  is a ‘‘gradient descent step’’ because we shall prove that it strictly decreases the objective (i.e.,  $f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)}) \geq 0$ ). Our update on  $\mathbf{z}_k$  is a ‘‘mirror descent step’’ because we shall apply standard mirror descent analysis [12] to it. We explicitly describe how to implement this mirror step: (proved in ‘‘Appendix D.4’’)

**Proposition 6.4** *If  $\mathbf{z}_{k-1} \in \Delta_{\text{simplex}}$  and  $\mathbf{z}_{k-1} > 0$ , the minimizer  $z = \operatorname{argmin}_{z \in \Delta_{\text{simplex}}} \{V_{\mathbf{z}_{k-1}}(z) + \langle \delta \mathbf{e}_i, z \rangle\}$  for any scalar  $\delta \in \mathbb{R}$  and basis vector  $\mathbf{e}_i$  can be computed as follows:*

1.  $z \leftarrow \mathbf{z}_{k-1}$ .
2.  $z_i \leftarrow z_i \cdot e^{-\delta}$ .
3. If  $\mathbf{1}^T z > 2\text{OPT}'$ ,  $z \leftarrow \frac{2\text{OPT}'}{\mathbf{1}^T z} z$ .
4. Return  $z$ .

We also point out that

**Lemma 6.5** *Each iteration of *CovLP Solver* can be implemented to run in expected  $O(N/n)$  time. The total expected running time is  $O(TN/n)$ .*

The proof of Lemma 6.5 is analogous to its packing counterpart, and included in Sect. F.6.



### 6.1 Main proof ideas and ingredients

In *CovLPSolver*, we pick a random coordinate  $i \in [n]$  at each iteration, and decompose  $\nabla_i f(x_k) = \xi + \eta$ , where  $\eta \in (-\infty, 0]$  is the (negative) large gradient component and  $\xi \in [-\sqrt{\varepsilon}, 1]$  is the truncated gradient component. In other words, we truncate the gradient  $\nabla_i f(x_k)$  at a negative threshold  $-\beta = -\sqrt{\varepsilon}$ , rather than at  $-1$  as in *CovLPSolver<sup>wb</sup>*.

The reason for this new threshold  $-\sqrt{\varepsilon}$  can be understood as follows. In *PacLPSolver* (and symmetrically in *CovLPSolver<sup>wb</sup>*), we used Lemma 3.10 to show that our gradient descent step  $y_k$  decreases the objective by an amount that both includes the  $\xi$  and  $\eta$  components. Unfortunately, for covering LP, this decrease amount is only proportional to  $\eta$  but not to  $\xi$  (compare Lemma 3.10 with Lemma 6.14 later). This forces us to treat the small gradient  $\xi$  separately using mirror descent, but not gradient descent.

If  $\xi$  were in  $[-1, 1]$ , classical theory of mirror descent [12] (or multiplicative weight update [7]) would imply that the mirror step  $z_k$  converges at a rate  $\propto \varepsilon^{-2}$ . This is too slow. Instead, since truncated  $\xi$  into a smaller interval  $[-\sqrt{\varepsilon}, 1]$ , using a *negative-width technique* (see Sect. 6.5), we improve this mirror-descent convergence rate from  $\varepsilon^{-2}$  to  $\varepsilon^{-1.5}$ .

On the other hand, due to this truncation at  $-\sqrt{\varepsilon}$  instead of  $-1$ , our gradient step on  $y_k$  also converges slower, at a rate  $1/\varepsilon^{1.5}$  instead of  $1/\varepsilon$ . This is why  $\beta = \sqrt{\varepsilon}$  is the best truncation threshold, as it balances gradient and mirror descent.

Another ingredient behind our proof is a new distance bound that is uncommon in first-order analysis. Recall that, given convex function  $g(x)$ , traditional analysis applies convexity argument  $g(x) - g(x^*) \leq \langle \nabla g(x), x - x^* \rangle$  to bound the objective distance to optimum. If  $g(x) = e^{-x}$  is univariate,  $x = -1$ , and  $x^* = -100$ , this bound becomes  $e^{-1} \approx e^{-1} - e^{-100} \leq e^{-1} \cdot 99$ , which is very loose. In our analysis, we replace this convexity argument with a more benign bound, specifically designed for covering LP (see Lemma 6.10).

### 6.2 Convergence statement

The main convergence theorem of this section is as follows:

**Theorem 6.6** *CovLPSolver*( $A, x^{\text{start}}, \varepsilon$ ) outputs some  $y_T$  satisfying

$$\mathbf{E}[f_\mu(y_T)] \leq (1 + 9\varepsilon)\text{OPT}.$$

Again, using the same proof as Corollary 3.5, one can apply Markov’s bound to turn Theorem 5.3 into a probabilistic statement:

**Corollary 6.7** With probability at least  $2/3$ ,  $y_T = \text{CovLPSolver}(A, x^{\text{start}}, \varepsilon)$  satisfies that  $\frac{y_T}{1-\varepsilon}$  is a  $(1 + O(\varepsilon))$  approximate solution to the covering LP program. The expected running time is  $O(\frac{\log(nm/\varepsilon) \log(1/\varepsilon)}{\varepsilon^{1.5}} N)$ .

Before delving into the proof of Theorem 6.6, we make the following observations:

**Fact 6.8** For every  $k \in \{0, 1, \dots, T\}$ , it satisfies  $\mathbf{x}_k, \mathbf{y}_k \geq 0, \mathbf{z}_k > 0$ , and  $\mathbf{z}_k \in \Delta_{\text{simplex}}$ .

*Proof* Since the  $x^{\text{start}}$  satisfies  $\mathbb{1}^T x^{\text{start}} \leq 2\text{OPT}'$  by Proposition 6.2, we have  $\mathbf{z}_0 = x^{\text{start}} \in \Delta_{\text{simplex}}$ . Also, the mirror descent step (see Proposition 6.4) ensures  $\mathbf{z}_{k,i} > 0$  for all rounds  $k$  and coordinates  $i$ , as well as  $\mathbf{z}_k \in \Delta_{\text{simplex}}$  for all rounds  $k$ . However, we  $\mathbf{x}_k$  and  $\mathbf{y}_k$  may not necessarily lie inside  $\Delta_{\text{simplex}}$ , but will always stay non-negative. □

We prove Theorem 6.6 in the subsequent subsections.

### 6.3 Step 1: distance adjustment

Using convexity, one can obtain

$$f_\mu(\mathbf{x}_k) - f_\mu(u) \leq \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - u \rangle \quad \text{for every } u \in \Delta_{\text{simplex}}. \tag{6.2}$$

Note that inequality (6.2) can be very loose for exponential functions. For instance, if  $f_\mu(x)$  were as simple as  $e^x$ , then the convexity inequality  $e^b - e^a \leq e^b \cdot (b - a)$  says

- when  $b = 2$  and  $a = -10$ , we have  $e^2 - e^{-10} \leq 12e^2$ ;
- when  $b = 2$  and  $a = -100$ , we have  $e^2 - e^{-100} \leq 102e^2$ .

Although  $e^{-100} \approx e^{-10}$ , the two upper bounds are off from each other by a factor of 10.

In this section, we strengthen (6.2) in the special case of  $u = u^* \stackrel{\text{def}}{=} (1 + \varepsilon/2)x^*$ . For analysis purpose, let  $\tilde{A}$  be the *adjusted matrix* of  $A$  described as follows.

**Definition 6.9** (adjusted matrix  $\tilde{A}$ ) For each row  $j \in [m]$ , if  $(Au^*)_j \leq 2$  then we keep it and let  $\tilde{A}_j \stackrel{\text{def}}{=} A_j$ . Otherwise,—that is, if  $(Au^*)_j > 2$ —we define  $\tilde{A}_j \stackrel{\text{def}}{=} \frac{2}{(Au^*)_j} \cdot A_j$ : to be the same  $j$ -th row  $A_j$ , but scaled down by a factor of  $\frac{2}{(Au^*)_j}$ . It is clear from this definition that

$$A_{ji} \geq \tilde{A}_{ji} \quad \text{for all } (i, j) \in [n] \times [m] \quad \text{and} \quad (1 + \varepsilon)\mathbb{1} \leq \tilde{A}u^* \leq 2\mathbb{1}.$$

**Lemma 6.10** (distance adjustment)

$$\begin{aligned} f_\mu(\mathbf{x}_k) - f_\mu(u^*) &\leq \langle \mathbb{1} - A^T p(\mathbf{x}_k), \mathbf{x}_k - u^* \rangle + \langle \tilde{A}^T p(\mathbf{x}_k) - A^T p(\mathbf{x}_k), u^* \rangle + \varepsilon\text{OPT} \\ &= \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - u^* \rangle + \langle \tilde{A}^T p(\mathbf{x}_k) - A^T p(\mathbf{x}_k), u^* \rangle + \varepsilon\text{OPT} \end{aligned}$$

At high level, ignoring the negligible term  $\varepsilon\text{OPT}$ , Lemma 6.10 strengthens the classical bound due to the extra term of  $\langle \tilde{A}^T p(\mathbf{x}_k) - A^T p(\mathbf{x}_k), u^* \rangle$ . This extra term is always non-positive since  $\tilde{A} \leq A$  coordinate-wise, but may be very negative in certain cases.

*Proof of Lemma 6.10*

$$\begin{aligned}
 f_\mu(\mathbf{x}_k) - f_\mu(u^*) &= \mu \sum_{j=1}^m \left( e^{\frac{1}{\mu}(1-(A\mathbf{x}_k)_j)} - e^{\frac{1}{\mu}(1-(Au^*)_j)} \right) + \langle \mathbb{1}, \mathbf{x}_k - u^* \rangle \\
 &\stackrel{\textcircled{1}}{\leq} \mu \sum_{j=1}^m \left( e^{\frac{1}{\mu}(1-(A\mathbf{x}_k)_j)} - e^{\frac{1}{\mu}(1-(\tilde{A}u^*)_j)} \right) + \langle \mathbb{1}, \mathbf{x}_k - u^* \rangle + \mu \cdot m \cdot e^{-1/\mu} \\
 &\stackrel{\textcircled{2}}{\leq} \sum_{j=1}^m e^{\frac{1}{\mu}(1-(A\mathbf{x}_k)_j)} \cdot ((\tilde{A}u^*)_j - (A\mathbf{x}_k)_j) + \langle \mathbb{1}, \mathbf{x}_k - u^* \rangle + \varepsilon \text{OPT} \\
 &= \sum_{j=1}^m p_j(\mathbf{x}_k) \cdot ((\tilde{A}u^*)_j - (A\mathbf{x}_k)_j) + \langle \mathbb{1}, \mathbf{x}_k - u^* \rangle + \varepsilon \text{OPT} \\
 &= \sum_{j=1}^m p_j(\mathbf{x}_k) \cdot ((Au^*)_j - (A\mathbf{x}_k)_j) + \langle \mathbb{1}, \mathbf{x}_k - u^* \rangle \\
 &\quad + \sum_{j=1}^m p_j(\mathbf{x}_k) \cdot ((\tilde{A}u^*)_j - (Au^*)_j) + \varepsilon \text{OPT} \\
 &= \langle -A^T p(\mathbf{x}_k), \mathbf{x}_k - u^* \rangle + \langle \mathbb{1}, \mathbf{x}_k - u^* \rangle + \langle \tilde{A}^T p(\mathbf{x}_k) - A^T p(\mathbf{x}_k), u^* \rangle + \varepsilon \text{OPT}.
 \end{aligned}$$

Above,  $\textcircled{1}$  is because if  $(Au^*)_j \neq (\tilde{A}u^*)_j$  for some  $j$ , then it must satisfy that  $(\tilde{A}u^*)_j = 2$ , and therefore  $-e^{\frac{1}{\mu}(1-(Au^*)_j)} \leq -e^{\frac{1}{\mu}(1-(\tilde{A}u^*)_j)} + e^{-1/\mu}$ .  $\textcircled{2}$  uses the convexity inequality of  $e^b - e^a \leq e^b \cdot (b - a)$ , and the fact that  $\mu m e^{-1/\mu} \ll \varepsilon \text{OPT}$ .

**6.4 Step 2: gradient truncation**

For analysis purpose, let us separate the indices  $i \in [n]$  into large and small ones.

**Definition 6.11** We make the following definitions.

- Let  $B_k \stackrel{\text{def}}{=} \{i \in [n] : \nabla_i f_\mu(\mathbf{x}_k) < -\beta\}$  and  $[n] \setminus B_k$  be the set of *large and small indices*.
- Let  $\xi_k \in [-\beta, 1]^n$  be the *truncated gradient* so that  $\xi_{k,i} = \mathbb{T}^c(\nabla_i f_\mu(x_k))$  for each  $i \in [n]$ .
- Let  $\eta_k \in (-\infty, 0]^n$  be the *large gradient* so that  $\nabla f_\mu(\mathbf{x}_k) = \xi_k + \eta_k$ . It is clear that

$$\eta_{k,i} = 0 \text{ for every } i \notin B, \text{ and } \eta_{k,i} = (1+\beta) - (A^T p(\mathbf{x}_k))_i \text{ for every } i \in B.$$

- Let  $\tilde{\eta}_k \in (-\infty, \infty)^n$  be the *adjusted large gradient* so that

$$\tilde{\eta}_{k,i} = 0 \text{ for every } i \notin B, \text{ and } \tilde{\eta}_{k,i} = (1+\beta) - (\tilde{A}^T p(\mathbf{x}_k))_i \text{ for every } i \in B.$$

We denote by  $\eta_k^{(i)} = (0, \dots, 0, \eta_{k,i}, 0, \dots, 0)$ , the vector that is zero at all coordinates other than  $i$ , and similarly  $\xi_k^{(i)} = (0, \dots, \xi_{k,i}, \dots, 0)$  and  $\tilde{\eta}_k^{(i)} = (0, \dots, \tilde{\eta}_{k,i}, \dots, 0)$ . We emphasize that  $\eta_k^{(i)} \neq \eta_k, \tilde{\eta}_k^{(i)} \neq \tilde{\eta}_k$ , and  $\xi_k^{(i)} \neq \xi_k$ .

The following key lemma is very analogous to (3.1) in the packing LP analysis.

**Lemma 6.12** (*distance upper bound*)

$$f_\mu(\mathbf{x}_k) - f_\mu(u^*) \leq \frac{(1-\tau)}{\tau} (f_\mu(\mathbf{y}_{k-1}) - f_\mu(\mathbf{x}_k)) + \mathbf{E}_i \left[ \langle n\xi_k^{(i)}, \mathbf{z}_{k-1} - u^* \rangle \right] + \mathbf{E}_i \left[ \langle n\tilde{\eta}_k^{(i)}, -u^* \rangle \right] + \varepsilon \text{OPT}.$$

Note that if one uses  $\eta_k^{(i)}$  instead of  $\tilde{\eta}_k^{(i)}$ , then Lemma 6.12 becomes trivial to prove just like (3.1). The reason we can have the stronger term  $\tilde{\eta}_k^{(i)}$  is precisely due to the distance adjustment Lemma 6.10.

*Proof of Lemma 6.12* We derive the following sequence of inequalities:

$$\begin{aligned} & (f_\mu(\mathbf{x}_k) - f_\mu(u^*)) - \varepsilon \text{OPT} \\ & \stackrel{\textcircled{1}}{\leq} \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - u^* \rangle + \langle \tilde{A}^T p(\mathbf{x}_k) - A^T p(\mathbf{x}_k), u^* \rangle \\ & = \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{z}_{k-1} \rangle + \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{z}_{k-1} - u^* \rangle + \langle \tilde{A}^T p(\mathbf{x}_k) - A^T p(\mathbf{x}_k), u^* \rangle \\ & \stackrel{\textcircled{2}}{=} \frac{(1-\tau)}{\tau} \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{y}_{k-1} - \mathbf{x}_k \rangle + \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{z}_{k-1} - u^* \rangle \\ & \quad + \langle \tilde{A}^T p(\mathbf{x}_k) - A^T p(\mathbf{x}_k), u^* \rangle \\ & \stackrel{\textcircled{3}}{\leq} \frac{(1-\tau)}{\tau} (f_\mu(\mathbf{y}_{k-1}) - f_\mu(\mathbf{x}_k)) + \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{z}_{k-1} - u^* \rangle \\ & \quad + \langle \tilde{A}^T p(\mathbf{x}_k) - A^T p(\mathbf{x}_k), u^* \rangle \\ & = \frac{(1-\tau)}{\tau} (f_\mu(\mathbf{y}_{k-1}) - f_\mu(\mathbf{x}_k)) + \langle \xi_k + \eta_k, \mathbf{z}_{k-1} - u^* \rangle \\ & \quad + \langle \tilde{A}^T p(\mathbf{x}_k) - A^T p(\mathbf{x}_k), u^* \rangle \\ & \stackrel{\textcircled{4}}{\leq} \frac{(1-\tau)}{\tau} (f_\mu(\mathbf{y}_{k-1}) - f_\mu(\mathbf{x}_k)) + \langle \xi_k, \mathbf{z}_{k-1} - u^* \rangle \\ & \quad + \langle \tilde{A}^T p(\mathbf{x}_k) - A^T p(\mathbf{x}_k) - \eta_k, u^* \rangle \\ & \stackrel{\textcircled{5}}{\leq} \frac{(1-\tau)}{\tau} (f_\mu(\mathbf{y}_{k-1}) - f_\mu(\mathbf{x}_k)) + \langle \xi_k, \mathbf{z}_{k-1} - u^* \rangle + \langle -\tilde{\eta}_k, u^* \rangle \\ & = \frac{(1-\tau)}{\tau} (f_\mu(\mathbf{y}_{k-1}) - f_\mu(\mathbf{x}_k)) + \mathbf{E}_i \left[ \langle n\xi_k^{(i)}, \mathbf{z}_{k-1} - u^* \rangle + \langle -n\tilde{\eta}_k^{(i)}, u^* \rangle \right]. \end{aligned}$$

Above, ① is due to Lemma 6.10. ② is because  $\mathbf{x}_k = \tau \mathbf{z}_{k-1} + (1-\tau)\mathbf{y}_{k-1}$ , which implies that  $\tau(\mathbf{x}_k - \mathbf{z}_{k-1}) = (1-\tau)(\mathbf{y}_{k-1} - \mathbf{x}_k)$ . ③ is by the convexity of  $f_\mu(\cdot)$ . ④ is because  $\langle \eta_k, \mathbf{z}_{k-1} \rangle \leq 0$ , since  $\eta_k \leq 0$  while  $\mathbf{z}_{k-1} \geq 0$ . ⑤ needs some careful justification: for every  $i \notin B_k$ , we have  $(\tilde{A}^T p(\mathbf{x}_k) - A^T p(\mathbf{x}_k))_i - \eta_{k,i} \leq 0 - 0 = -\tilde{\eta}_{k,i}$ ; for every  $i \in B_k$ , we have

$$\begin{aligned}
 (\tilde{A}^T p(\mathbf{x}_k) - A^T p(\mathbf{x}_k))_i - \eta_{k,i} &= (\tilde{A}^T p(\mathbf{x}_k) - A^T p(\mathbf{x}_k))_i - ((1 + \beta) - (A^T p(\mathbf{x}_k))_i) \\
 &= -((1 + \beta) - (\tilde{A}^T p(\mathbf{x}_k))_i) = -\tilde{\eta}_{k,i},
 \end{aligned}$$

where the two equalities follow from the definitions of  $\eta_{k,i}$  and  $\tilde{\eta}_{k,i}$ .

### 6.5 Step 3: mirror descent guarantee

Our update  $\mathbf{z}_k^{(i)} \stackrel{\text{def}}{=} \operatorname{argmin}_{z \in \Delta_{\text{simplex}}} \{V_{\mathbf{z}_{k-1}}(z) + ((1 + \gamma)n\alpha_k \xi_k^{(i)}, z)\}$  is, by its definition, a mirror descent step [12]. We begin by explaining an attempt that is too weak for obtaining the  $\varepsilon^{-1.5}$  convergence rate.

Using the classical theory, it is not hard to repeat the proof of Lemma 3.7—although changing the distance function from  $\|\cdot\|_A^2$  to  $V_x(y)$ —and obtain that, as long as  $\xi_{k,i}$  is in  $[-1, +1]$  for each coordinate  $i$ , for every  $u \in \Delta_{\text{simplex}}$ ,

$$\mathbf{E}_i[\alpha_k \langle n\xi_k^{(i)}, \mathbf{z}_{k-1} - u \rangle] \leq V_{\mathbf{z}_{k-1}}(u) - \mathbf{E}_i[V_{\mathbf{z}_k^{(i)}}(u)] + O(\alpha_k^2 n) \text{OPT}.$$

This inequality only yields a slower  $\varepsilon^{-2}$  convergence rate, and  $\pm 1$  is also known as the *width parameter* from the multiplicative-weight-update language [7].

In our lemma below, we make use of the fact  $\xi_{k,i}$  is in  $[-\beta, +1] \subseteq [-1, +1]$ . In essence, this allows us to replace the  $O(\alpha_k^2 n)$  factor with a better  $O(\alpha_k^2 \beta n)$  factor. We call it the *negative-width technique*.<sup>12</sup> Formally,

**Lemma 6.13** (mirror descent) *Denoting by  $\gamma \stackrel{\text{def}}{=} 2\alpha_T n$ , we have*

$$\mathbf{E}_i[\alpha_k \langle n\xi_k^{(i)}, \mathbf{z}_{k-1} - u^* \rangle] \leq V_{\mathbf{z}_{k-1}}\left(\frac{u^*}{1 + \gamma}\right) - \mathbf{E}_i\left[V_{\mathbf{z}_k^{(i)}}\left(\frac{u^*}{1 + \gamma}\right)\right] + 12\text{OPT} \cdot \gamma\alpha_k\beta.$$

The proof is somewhat technical and included in “Appendix D.4”.

### 6.6 Step 4: gradient descent guarantee

We show our gradient step never increases the objective for all choices of  $i$ . In addition, it decreases the objective by an amount proportional to the adjusted large gradient  $\tilde{\eta}_k^{(i)}$ .

**Lemma 6.14** (gradient descent) *For every  $i \in [n]$ , we have*

- (a)  $f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)}) \geq 0$ , and
- (b)  $f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)}) \geq \frac{\mu\beta}{12} \cdot \langle -\tilde{\eta}_k^{(i)}, u^* \rangle$ .

<sup>12</sup> This negative width technique is related to [7, Definition 3.2], where the authors analyze the multiplicative weight update method in a special case when the oracle returns loss values only in  $[-\ell, +\rho]$ , for some  $\ell \ll \rho$ . This technique is also a sub-case of a more general theory of mirror descent, known as the local-norm convergence, that we have summarized in a separate and later paper [3].

The proof of Lemma 6.14 is quite technical and can be found in ‘‘Appendix D.4’’.

At high level, one would hope to prove that the gradient step decreases the objective by an amount proportional to the large gradient  $\eta_k^{(i)}$ , rather than the adjusted large gradient  $\tilde{\eta}_k^{(i)}$ . If that were true, the entire proof structure of our covering LP convergence would become much closer to that of packing LP, and there would be absolutely no need for the introduction of the distance adjustment in Sect. 6.3, as well as the definitions of  $\tilde{A}$  and  $\tilde{\eta}$ .

Unfortunately, if one replaces  $\tilde{\eta}$  with  $\eta$  in the above lemma, the inequality is *false*. The reason behind it is very similar to what we have summarized in Sect. 6.3, and related to the fact that the exponential penalty function is not Lipschitz smooth.

### 6.7 Step 5: putting all together

Combining Lemma 6.12, Lemma 6.13, and Lemma 6.14, we obtain that

$$\begin{aligned} & \alpha_k (f_\mu(\mathbf{x}_k) - f_\mu(u^*)) - \alpha_k \varepsilon \text{OPT} \\ & \leq \frac{(1 - \tau)\alpha_k}{\tau} (f_\mu(\mathbf{y}_{k-1}) - f_\mu(\mathbf{x}_k)) + \mathbf{E}_i \left[ \alpha_k \langle n \xi_k^{(i)}, \mathbf{z}_{k-1} - u^* \rangle \right] \\ & \quad + \mathbf{E}_i \left[ \alpha_k \langle n \tilde{\eta}_k^{(i)}, -u^* \rangle \right] \\ & \leq \frac{(1 - \tau)\alpha_k}{\tau} (f_\mu(\mathbf{y}_{k-1}) - f_\mu(\mathbf{x}_k)) + V_{\mathbf{z}_{k-1}} \left( \frac{u^*}{1 + \gamma} \right) - \mathbf{E}_i \left[ V_{\mathbf{z}_k^{(i)}} \left( \frac{u^*}{1 + \gamma} \right) \right] \\ & \quad + 12\text{OPT} \cdot \gamma \alpha_k \beta + \mathbf{E}_i \left[ \frac{12\alpha_k n}{\mu\beta} (f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)})) \right] \end{aligned}$$

*Remark 6.15* Above, the quantity ‘‘ $12\text{OPT} \cdot \gamma \alpha_k \beta$ ’’ is the loss term introduced by the mirror descent. Unlike the packing LP case—see (3.2)—this loss term is not dominated by the gradient step. (If one could do so, this would give *CovLPSolver* an  $\varepsilon^{-1}$  convergence rate.)

The quantity ‘‘ $\alpha_k \langle n \xi_k^{(i)}, \mathbf{z}_{k-1} - u^* \rangle$ ’’ is the loss introduced by the (adjusted) large gradient  $\tilde{\eta}$ , and is dominated by our gradient step progress owing to Lemma 6.14. This is similar to the packing LP case—see Lemma 3.10.

From here, let us use the special choice of  $\tau = \frac{\mu\beta}{12n}$ . We obtain that

$$\begin{aligned} & -\alpha_k (f_\mu(u^*) + \varepsilon \text{OPT}) \\ & \leq 12\gamma \alpha_k \beta \text{OPT} + \frac{(1 - \tau)\alpha_k}{\tau} f_\mu(\mathbf{y}_{k-1}) + V_{\mathbf{z}_{k-1}} \left( \frac{u^*}{1 + \gamma} \right) \\ & \quad - \mathbf{E}_i \left[ \frac{\alpha_k}{\tau} f_\mu(\mathbf{y}_k^{(i)}) + V_{\mathbf{z}_k^{(i)}} \left( \frac{u^*}{1 + \gamma} \right) \right]. \end{aligned}$$

Using the choice  $\alpha_k = \frac{\alpha_{k-1}}{1 - \tau}$  and telescoping the above inequality for  $k = 1, \dots, T$ , we have

$$\begin{aligned}
 -\left(\sum_{k=1}^T \alpha_k\right)(f_\mu(u^*) + \varepsilon\text{OPT}) &\leq \left(\sum_{k=1}^T \alpha_k\right) \cdot 12\gamma\beta\text{OPT} + \frac{\alpha_0}{\tau} f_\mu(y_0) \\
 &\quad + V_{z_0}\left(\frac{u^*}{1+\gamma}\right) - \frac{\alpha_T}{\tau} \mathbf{E}[f_\mu(y_T)].
 \end{aligned}$$

We compute that  $\sum_{k=1}^T \alpha_k = \alpha_T \cdot \sum_{k=0}^{T-1} (1-\tau)^k = \alpha_T \cdot \frac{1-(1-\tau)^T}{\tau} < \frac{\alpha_T}{\tau}$ , and recall that  $\gamma = 2\alpha_T n$ . Therefore, we rearrange and get

$$\begin{aligned}
 \frac{\alpha_T}{\tau} \mathbf{E}[f_\mu(y_T)] &\leq \frac{\alpha_T}{\tau} (f_\mu(u^*) + \varepsilon\text{OPT}) + \frac{\alpha_T}{\tau} \cdot 12\gamma\beta\text{OPT} + \frac{\alpha_0}{\tau} f_\mu(y_0) \\
 &\quad + V_{z_0}\left(\frac{u^*}{1+\gamma}\right), \\
 \implies \mathbf{E}[f_\mu(y_T)] &\leq f_\mu(u^*) + \varepsilon\text{OPT} + 24\alpha_T n\beta\text{OPT} + (1-\tau)^T f_\mu(y_0) \\
 &\quad + \frac{\tau}{\alpha_T} V_{z_0}\left(\frac{u^*}{1+\gamma}\right). \tag{6.3}
 \end{aligned}$$

From this point, we need to use our special choice of the initial point  $x_0 = y_0 = z_0 = x^{\text{start}}$  (see Proposition 6.2), which implies that  $f_\mu(y_0) \leq 4\text{OPT}$  and  $\mathbb{1}^T x^{\text{start}} \leq 4\text{OPT}$ . We also have

$$\begin{aligned}
 V_{z_0}\left(\frac{u^*}{1+\gamma}\right) &= V_{x^{\text{start}}}\left(\frac{u^*}{1+\gamma}\right) = \sum_{i=1}^n \frac{u_i^*}{1+\gamma} \log \frac{u_i^*}{(1+\gamma)x_i^{\text{start}}} + x_i^{\text{start}} - \frac{u_i^*}{1+\gamma} \\
 &\stackrel{\textcircled{1}}{\leq} \sum_{i=1}^n u_i^* \log(u_i^* \cdot n) + 4\text{OPT} \stackrel{\textcircled{2}}{\leq} (2 \log(nm) + 4) \cdot \text{OPT}.
 \end{aligned}$$

Above, inequality  $\textcircled{1}$  follows because  $x_i^{\text{start}} \geq 1/n$  for all  $i \in [n]$  according to the definition in Proposition 6.2; inequality  $\textcircled{2}$  follows because each  $u_i^* \leq (1 + \varepsilon/2)x_i^* \leq (1 + \varepsilon/2)\text{OPT} \leq (1 + \varepsilon/2)m$  and  $\mathbb{1}^T u_i^* = (1 + \varepsilon/2)\text{OPT}$ , as well as the fact that  $\varepsilon$  is sufficiently small.

Finally, we choose  $\beta = \sqrt{\varepsilon}$ ,  $T = \lceil \frac{1}{\tau} \log(1/\varepsilon) \rceil$ , and  $\alpha_0$  such that  $\alpha_T = \frac{\varepsilon}{12n\beta}$ . Substituting into (6.3) all of these parameters, along with the aforementioned inequalities  $f_\mu(y_0) \leq 4\text{OPT}$  and  $V_{z_0}\left(\frac{u^*}{1+\gamma}\right) \leq (2 \log(nm) + 4) \cdot \text{OPT}$ , as well as  $f_\mu(u^*) \leq (1 + \varepsilon)\text{OPT}$  from Proposition 4.5.a, we obtain that

$$\begin{aligned}
 \mathbf{E}[f_\mu(y_T)] &\leq (1 + \varepsilon)\text{OPT} + \varepsilon\text{OPT} + 2\varepsilon\text{OPT} + \varepsilon f_\mu(y_0) \\
 &\quad + \frac{\mu\beta/12n}{\varepsilon/12n\beta} (2 \log(nm) + 4)\text{OPT} \\
 &\leq (1 + 9\varepsilon)\text{OPT}.
 \end{aligned}$$

This finishes the proof of Theorem 6.6. □

## Appendix

### A.1 Proof of Lemma 3.6

**Lemma 3.6** *We have  $\mathbf{x}_k, \mathbf{y}_k, \mathbf{z}_k \in \Delta_{\text{box}}$  for all  $k = 0, 1, \dots, T$ .*

*Proof* This is true at the beginning as  $\mathbf{x}_0 = \mathbf{y}_0 = x^{\text{start}} \in \Delta_{\text{box}}$  (see Fact 2.8) and  $\mathbf{z}_0 = 0 \in \Delta_{\text{box}}$ . In fact, it suffices for us to show that for every  $k \geq 1$ ,  $\mathbf{y}_k = \sum_{l=0}^k \gamma_k^l \mathbf{z}_l$  for some scalars  $\gamma_k^l$  satisfying  $\sum_l \gamma_k^l = 1$  and  $\gamma_k^l \geq 0$  for each  $l = 0, \dots, k$ . If this is true, we can prove the lemma by induction: at each iteration  $k \geq 1$ ,

1.  $\mathbf{x}_k = \tau \mathbf{z}_{k-1} + (1 - \tau) \mathbf{y}_{k-1}$  must be in  $\Delta_{\text{box}}$  because  $\mathbf{y}_{k-1}$  and  $\mathbf{z}_{k-1}$  are and  $\tau \in [0, 1]$ ,
2.  $\mathbf{z}_k$  is in  $\Delta_{\text{box}}$  by the definition that  $\mathbf{z}_k = \operatorname{argmin}_{\mathbf{z} \in \Delta_{\text{box}}} \{\dots\}$ , and
3.  $\mathbf{y}_k$  is also in  $\Delta_{\text{box}}$  because  $\mathbf{y}_k = \sum_{l=0}^k \gamma_k^l \mathbf{z}_l$  is a convex combination of the  $\mathbf{z}_l$ 's and  $\Delta_{\text{box}}$  is convex.

For the rest of the proof, we show that  $\mathbf{y}_k = \sum_{l=0}^k \gamma_k^l \mathbf{z}_l$  for every  $k \geq 1$  with coefficients<sup>13</sup>

$$\gamma_k^l = \begin{cases} (1 - \tau) \gamma_{k-1}^l, & l = 0, \dots, k-2; \\ \left( \frac{1}{n\alpha_{k-1}L} - \frac{1}{n\alpha_k L} \right) + \tau \left( 1 - \frac{1}{n\alpha_{k-1}L} \right), & l = k-1; \\ \frac{1}{n\alpha_k L}, & l = k. \end{cases}$$

This is true at the base case  $k = 1$  because  $\mathbf{y}_1 = \mathbf{x}_1 + \frac{1}{n\alpha_1 L} (\mathbf{z}_1 - \mathbf{z}_0) = \frac{1}{n\alpha_1 L} \mathbf{z}_1 + \left( 1 - \frac{1}{n\alpha_1 L} \right) \mathbf{z}_0$ . For the general  $k \geq 2$ , we have

$$\begin{aligned} \mathbf{y}_k &= \mathbf{x}_k + \frac{1}{n\alpha_k L} (\mathbf{z}_k - \mathbf{z}_{k-1}) \\ &= \tau \mathbf{z}_{k-1} + (1 - \tau) \mathbf{y}_{k-1} + \frac{1}{n\alpha_k L} (\mathbf{z}_k - \mathbf{z}_{k-1}) \\ &= \tau \mathbf{z}_{k-1} + (1 - \tau) \left( \sum_{l=0}^{k-2} \gamma_{k-1}^l \mathbf{z}_l + \frac{1}{n\alpha_{k-1} L} \mathbf{z}_{k-1} \right) + \frac{1}{n\alpha_k L} (\mathbf{z}_k - \mathbf{z}_{k-1}) \\ &= \left( \sum_{l=0}^{k-2} (1 - \tau) \gamma_{k-1}^l \mathbf{z}_l \right) + \left( \left( \frac{1}{n\alpha_{k-1} L} - \frac{1}{n\alpha_k L} \right) + \tau \left( 1 - \frac{1}{n\alpha_{k-1} L} \right) \right) \mathbf{z}_{k-1} \\ &\quad + \frac{1}{n\alpha_k L} \mathbf{z}_k. \end{aligned}$$

Therefore, we obtain  $\mathbf{y}_k = \sum_{l=0}^k \gamma_k^l \mathbf{z}_l$  as desired.

<sup>13</sup> We wish to point out that this proof coincides with a lemma from the accelerated coordinate descent theory of Fercoq and Richtárik [17]. Their paper is about optimizing an objective function that is Lipschitz smooth, and thus irrelevant to our work.



It is now easy to check that under our definition of  $\alpha_k$  (which satisfies  $\alpha_k \geq \alpha_{k-1}$  and  $\alpha_k \geq \alpha_0 = \frac{1}{nL}$ , we must have  $\gamma_k^l \geq 0$  for all  $k$  and  $l$ . Also,

$$\begin{aligned} \sum_l \gamma_k^l &= \sum_{l=0}^{k-2} (1 - \tau)\gamma_{k-1}^l + \left( \left( \frac{1}{n\alpha_{k-1}L} - \frac{1}{n\alpha_kL} \right) + \tau \left( 1 - \frac{1}{n\alpha_{k-1}L} \right) \right) + \frac{1}{n\alpha_kL} \\ &= (1 - \tau) \left( 1 - \frac{1}{n\alpha_{k-1}L} \right) + \left( \left( \frac{1}{n\alpha_{k-1}L} - \frac{1}{n\alpha_kL} \right) + \tau \left( 1 - \frac{1}{n\alpha_{k-1}L} \right) \right) \\ &\quad + \frac{1}{n\alpha_kL} = 1. \end{aligned}$$

□

### B.2 Proof of Proposition 4.5

**Proposition 4.5** (a)  $f_\mu(u^*) \leq (1 + \varepsilon)\text{OPT}$  for  $u^* \stackrel{\text{def}}{=} (1 + \varepsilon/2)x^*$ .

(b)  $f_\mu(x) \geq (1 - \varepsilon)\text{OPT}$  for every  $x \geq 0$ .

(c) For any  $x \geq 0$  satisfying  $f_\mu(x) \leq 2\text{OPT}$ , we must have  $Ax \geq (1 - \varepsilon)\mathbb{1}$ .

(d) If  $x \geq 0$  satisfies  $f_\mu(x) \leq (1 + \delta)\text{OPT}$  for some  $\delta \in [0, 1]$ , then  $\frac{1+\delta}{1-\varepsilon}x$  is a  $\frac{1+\delta}{1-\varepsilon}$ -approximate solution to the covering LP.

*Proof* (a) We have  $\mathbb{1}^T u^* = (1 + \varepsilon/2)\text{OPT}$  by the definition of  $\text{OPT}$ . Also, from the feasibility constraint  $Ax^* \geq \mathbb{1}$  in the covering LP, we have  $Au^* - \mathbb{1} \geq \varepsilon/2 \cdot \mathbb{1}$ , and can compute  $f_\mu(u^*)$  as follows:

$$\begin{aligned} f_\mu(u^*) &= \mu \sum_j e^{\frac{1}{\mu}(1 - (Au^*)_j)} + \mathbb{1}^T u^* \leq \mu \sum_j e^{\frac{-\varepsilon/2}{\mu}} + (1 + \varepsilon/2)\text{OPT} \\ &\leq \frac{\mu m}{(nm)^2} + (1 + \varepsilon/2)\text{OPT} \leq (1 + \varepsilon)\text{OPT}. \end{aligned}$$

(b) Suppose towards contradiction that  $f_\mu(x) < (1 - \varepsilon)\text{OPT}$ . Since  $f_\mu(x) < \text{OPT} \leq m$ , we must have that for every  $j \in [m]$ , it satisfies that  $e^{\frac{1}{\mu}(1 - (Ax)_j)} \leq f_\mu(x)/\mu \leq m/\mu$ . This further implies  $(Ax)_j \geq 1 - \varepsilon$  by the definition of  $\mu$ . In other words,  $Ax \geq (1 - \varepsilon)\mathbb{1}$ . By the definition of  $\text{OPT}$ , we must then have  $\mathbb{1}^T x \geq (1 - \varepsilon)\text{OPT}$ , finishing the proof that  $f_\mu(x) \geq \mathbb{1}^T x \geq (1 - \varepsilon)\text{OPT}$ , giving a contradiction.

(c) To show  $Ax \geq (1 - \varepsilon)\mathbb{1}$ , we can assume that  $v = \max_j (1 - (Ax)_j) > \varepsilon$  because otherwise we are done. Under this definition, we have

$$f_\mu(x) \geq \mu e^{\frac{v}{\mu}} = \mu \left( \left( \frac{nm}{\varepsilon} \right)^4 \right)^{v/\varepsilon} \geq \frac{\varepsilon}{4 \log(nm/\varepsilon)} \left( \frac{nm}{\varepsilon} \right)^4 \gg 2\text{OPT},$$

contradicting to our assumption that  $f_\mu(x) \leq 2\text{OPT}$ . Therefore, we must have  $v \leq \varepsilon$ , that is,  $Ax \geq (1 - \varepsilon)\mathbb{1}$ .

(d) For any  $x$  satisfying  $f_\mu(x) \leq (1 + \theta)\text{OPT} \leq 2\text{OPT}$ , owing to Proposition 4.5c, we first have that  $x$  is approximately feasible, i.e.,  $Ax \geq (1 - \varepsilon)\mathbb{1}$ . Next, because

$\mathbb{1}^T x \leq f_\mu(x) \leq (1 + \theta)\text{OPT}$ , we know that  $x$  yields an objective  $\mathbb{1}^T x \leq (1 + \theta)\text{OPT}$ . Letting  $x' = \frac{1}{1-\varepsilon}x$ , we both have that  $x'$  is feasible (i.e.,  $Ax' \geq \mathbb{1}$ ), and  $x'$  has an objective  $\mathbb{1}^T x'$  at most  $\frac{1+\delta}{1-\varepsilon}\text{OPT}$ . □

### C.3 Missing proofs for Sect. 5

In this section we prove Theorem 5.3. Because the proof structure is almost identical to that of Theorem 3.4, we spend most of the discussions only pointing out the difference rather than repeating the proofs. The following three lemmas are completely identical to the ones in the packing LP case, so we restate them below:

**Lemma C.1** (cf. Lemma 3.3) *Each iteration of  $\text{CovLP Solver}^{\text{wb}}$  can be implemented to run in expected  $O(N/n)$  time.*

**Lemma C.2** (cf. Lemma 3.6) *We have  $\mathbf{x}_k, \mathbf{y}_k, \mathbf{z}_k \in \Delta_{\text{box}}$  for all  $k = 0, 1, \dots, T$ .*

**Lemma C.3** (cf. Lemma 3.7) *For every  $u \in \Delta_{\text{box}}$ , it satisfies  $\langle n\alpha_k \xi_k^{(i)}, \mathbf{z}_{k-1} - u \rangle \leq n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle + \frac{1}{2} \|\mathbf{z}_{k-1} - u\|_A^2 - \frac{1}{2} \|\mathbf{z}_k^{(i)} - u\|_A^2$ .*

For the gradient descent guarantee of Sect. 3.3, one can first note that Lemma 2.7 remains true: this can be verified by replacing  $\nabla_i f_\mu(x) + 1$  in its proof with  $1 - \nabla_i f_\mu(x)$ . For this reason, Lemma 3.9 (which is built on Lemma 2.7) also remains true. We state it below:

**Lemma C.4** (cf. Lemma 3.9) *We have  $f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)}) \geq \frac{1}{2} \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \geq 0$ .*

**Putting all together** Denote by  $\eta_k^{(i)} \in \mathbb{R}_{\leq 0}^n$  the vector that is only non-zero at coordinate  $i$ , and satisfies  $\eta_{k,i}^{(i)} = \nabla_i f_\mu(\mathbf{x}_k) - \xi_{k,i}^{(i)} \in (-\infty, 0]$ . In other words, the full gradient

$$\nabla f_\mu(\mathbf{x}_k) = \mathbf{E}_i[(0, \dots, n\nabla_i f_\mu(\mathbf{x}_k), \dots, 0)] = \mathbf{E}_i[n\eta_k^{(i)} + n\xi_k^{(i)}]$$

can be (in expectation) decomposed into the a large but non-positive component  $\eta_k^{(i)} \in (-\infty, 0]^n$  and a small component  $\xi_k^{(i)} \in [-1, 1]^n$ . Similar as Sect. 3.4, for any  $u \in \Delta_{\text{box}}$ , we can use a basic convexity argument and the mirror descent lemma to compute that

$$\begin{aligned} \alpha_k(f_\mu(\mathbf{x}_k) - f_\mu(u)) &\leq \langle \alpha_k \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - u \rangle \\ &= \langle \alpha_k \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{z}_{k-1} \rangle + \langle \alpha_k \nabla f_\mu(\mathbf{x}_k), \mathbf{z}_{k-1} - u \rangle \\ &= \langle \alpha_k \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{z}_{k-1} \rangle + \mathbf{E}_i \left[ \langle n\alpha_k \eta_k^{(i)}, \mathbf{z}_{k-1} - u \rangle + \langle n\alpha_k \xi_k^{(i)}, \mathbf{z}_{k-1} - u \rangle \right] \end{aligned}$$

$$\begin{aligned} &\stackrel{\textcircled{1}}{=} \frac{(1 - \tau)\alpha_k}{\tau} \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{y}_{k-1} - \mathbf{x}_k \rangle \\ &\quad + \mathbf{E}_i \left[ \langle n\alpha_k \eta_k^{(i)}, \mathbf{z}_{k-1} - u \rangle + \langle n\alpha_k \xi_k^{(i)}, \mathbf{z}_{k-1} - u \rangle \right] \end{aligned} \tag{C.1}$$

$$\begin{aligned} &\stackrel{\textcircled{2}}{\leq} \frac{(1 - \tau)\alpha_k}{\tau} (f_\mu(\mathbf{y}_{k-1}) - f_\mu(\mathbf{x}_k)) \\ &\quad + \mathbf{E}_i \left[ \langle n\alpha_k \eta_k^{(i)}, \mathbf{z}_{k-1} - u \rangle + n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \right] \\ &\quad + \frac{1}{2} \|\mathbf{z}_{k-1} - u\|_A^2 - \frac{1}{2} \|\mathbf{z}_k^{(i)} - u\|_A^2 \end{aligned} \tag{C.2}$$

Above,  $\textcircled{1}$  is because  $\mathbf{x}_k = \tau \mathbf{z}_{k-1} + (1 - \tau) \mathbf{y}_{k-1}$ , which implies that  $\tau(\mathbf{x}_k - \mathbf{z}_{k-1}) = (1 - \tau)(\mathbf{y}_{k-1} - \mathbf{x}_k)$ .  $\textcircled{2}$  uses convexity and Lemma C.3. We can establish the following lemma to upper bound the boxed term in (C.2). Its proof is in the same spirit to that of Lemma 3.10, and is the *only* place that we require all vectors to reside in  $\Delta_{\text{box}}$ .

**Lemma C.5** (cf. Lemma 3.10) *For every  $u \in \Delta_{\text{box}}$ ,*

$$\langle n\alpha_k \eta_k^{(i)}, \mathbf{z}_{k-1} - u \rangle + n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \leq 21n\alpha_k L \cdot (f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)})).$$

*Proof of Lemma C.5* Now there are three possibilities:

– If  $\eta_{k,i}^{(i)} = 0$ , then we must have  $\xi_{k,i}^{(i)} = \nabla_i f_\mu(\mathbf{x}_k) \in [-1, 1]$ . Lemma C.4 implies

$$\begin{aligned} &\langle n\alpha_k \eta_k^{(i)}, \mathbf{z}_{k-1} - u \rangle + n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \\ &\quad = n^2 \alpha_k^2 L \cdot \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \leq 2n^2 \alpha_k^2 L \cdot (f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)})) \end{aligned}$$

– If  $\eta_{k,i}^{(i)} < 0$  and  $\mathbf{z}_{k,i}^{(i)} < \frac{10}{\|A_{:,i}\|_\infty}$  (thus  $\mathbf{z}_k^{(i)}$  is not on the boundary of  $\Delta_{\text{box}}$ ), then we precisely have  $\mathbf{z}_{k,i}^{(i)} = \mathbf{z}_{k-1,i} + \frac{n\alpha_k}{\|A_{:,i}\|_\infty}$ , and accordingly  $\mathbf{y}_{k,i}^{(i)} = \mathbf{x}_{k,i} + \frac{1}{L\|A_{:,i}\|_\infty} > \mathbf{x}_{k,i}$ . In this case,

$$\begin{aligned} &\langle n\alpha_k \eta_k^{(i)}, \mathbf{z}_{k-1} - u \rangle + n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \\ &\quad \stackrel{\textcircled{1}}{\leq} n\alpha_k \cdot \nabla_i f_\mu(\mathbf{x}_k) \cdot \frac{-10}{\|A_{:,i}\|_\infty} + n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \\ &\quad \stackrel{\textcircled{2}}{<} n\alpha_k \cdot \nabla_i f_\mu(\mathbf{x}_k) \cdot \frac{-10}{\|A_{:,i}\|_\infty} + n^2 \alpha_k^2 L \cdot \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \\ &\quad \stackrel{\textcircled{3}}{=} 10n\alpha_k L \cdot \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle + n^2 \alpha_k^2 L \cdot \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \\ &\quad \stackrel{\textcircled{4}}{\leq} (20n\alpha_k L + 2n^2 \alpha_k^2 L) \cdot (f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)})). \end{aligned}$$

Above,  $\textcircled{1}$  follows from the fact that  $\mathbf{z}_{k-1}, u \in \Delta_{\text{box}}$  and therefore  $z_{k-1,i} \geq 0$  and  $u_i \leq \frac{10}{\|A_{:,i}\|_\infty}$  by the definition of  $\Delta_{\text{box}}$ , and  $u \geq 0$ ;  $\textcircled{2}$  follows from the fact that  $\mathbf{x}_k$  and  $\mathbf{y}_k^{(i)}$  are only different at coordinate  $i$ , and  $\xi_{k,i}^{(i)} = -1 > \nabla_i f_\mu(\mathbf{x}_k)$  (since  $\eta_{k,i}^{(i)} < 0$ );  $\textcircled{3}$  follows from the fact that  $\mathbf{y}_k^{(i)} = \mathbf{x}_k + \frac{\mathbf{e}_i}{L\|A_{:,i}\|_\infty}$ ; and  $\textcircled{4}$  uses Lemma C.4.

– If  $\eta_{k,i}^{(i)} < 0$  and  $\mathbf{z}_{k,i}^{(i)} = \frac{10}{\|A_{:,i}\|_\infty}$ , then we have

$$\begin{aligned} & \langle n\alpha_k \eta_k^{(i)}, \mathbf{z}_{k-1} - u \rangle + n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \\ & \stackrel{\textcircled{1}}{\leq} \langle n\alpha_k \eta_k^{(i)}, \mathbf{z}_{k-1} - \mathbf{z}_k^{(i)} \rangle + n^2 \alpha_k^2 L \cdot \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \\ & \stackrel{\textcircled{2}}{\leq} \langle n\alpha_k \nabla f_\mu(\mathbf{x}_k), \mathbf{z}_{k-1} - \mathbf{z}_k^{(i)} \rangle + n^2 \alpha_k^2 L \cdot \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \\ & \stackrel{\textcircled{3}}{=} n^2 \alpha_k^2 L \cdot \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle + n^2 \alpha_k^2 L \cdot \langle \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \\ & \stackrel{\textcircled{4}}{\leq} 4n^2 \alpha_k^2 L \cdot (f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)})). \end{aligned}$$

Above,  $\textcircled{1}$  is because  $u_i \leq \frac{10}{\|A_{:,i}\|_\infty} = \mathbf{z}_{k,i}^{(i)}$  and  $\eta_{k,i}^{(i)} < 0$ , together with  $\nabla_i f_\mu(\mathbf{x}_k) < \xi_{k,i}^{(i)}$  and  $\mathbf{x}_{k,i} \leq \mathbf{y}_{k,i}^{(i)}$ ;  $\textcircled{2}$  uses  $\nabla_i f_\mu(\mathbf{x}_k) = \eta_{k,i}^{(i)} - 1 < \eta_{k,i}^{(i)}$  and  $\mathbf{z}_{k,i}^{(i)} \geq \mathbf{z}_{k-1,i}$ ;  $\textcircled{3}$  is from our choice of  $\mathbf{y}_k$  which satisfies that  $\mathbf{z}_{k-1} - \mathbf{z}_k^{(i)} = n\alpha_k L(\mathbf{x}_k - \mathbf{y}_k^{(i)})$ ; and  $\textcircled{4}$  uses Lemma C.4.

Combining the three cases, and using the fact that  $f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)}) \geq 0$ , we conclude that

$$\begin{aligned} & \langle n\alpha_k \eta_k^{(i)}, \mathbf{z}_{k-1} - u \rangle + n^2 \alpha_k^2 L \cdot \langle \xi_k^{(i)}, \mathbf{x}_k - \mathbf{y}_k^{(i)} \rangle \\ & \leq (20n\alpha_k L + 4n^2 \alpha_k^2 L) \cdot (f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)})) \\ & \leq 21n\alpha_k L \cdot (f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)})). \end{aligned}$$

Above, the last inequality uses our choice of  $\alpha_k$ , which implies  $n\alpha_k \leq n\alpha_T = \frac{1}{\varepsilon L} \leq \frac{1}{4}$ .

Plugging Lemma C.5 back to (C.2), we have

$$\begin{aligned} \alpha_k (f_\mu(\mathbf{x}_k) - f_\mu(u)) & \leq \langle \alpha_k \nabla f_\mu(\mathbf{x}_k), \mathbf{x}_k - u \rangle \\ & \stackrel{\textcircled{1}}{\leq} \frac{(1 - \tau)\alpha_k}{\tau} (f_\mu(\mathbf{y}_{k-1}) - f_\mu(\mathbf{x}_k)) \\ & \quad + \mathbf{E}_i \left[ 21n\alpha_k L \cdot (f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)})) + \frac{1}{2} \|\mathbf{z}_{k-1} - u\|_A^2 - \frac{1}{2} \|\mathbf{z}_k - u\|_A^2 \right] \\ & \stackrel{\textcircled{2}}{\leq} \alpha_k f_\mu(\mathbf{x}_k) + (21n\alpha_k L - \alpha_k) f_\mu(\mathbf{y}_{k-1}) \\ & \quad + \mathbf{E}_i \left[ -21n\alpha_k L \cdot f_\mu(\mathbf{y}_k^{(i)}) + \frac{1}{2} \|\mathbf{z}_{k-1} - u\|_A^2 - \frac{1}{2} \|\mathbf{z}_k - u\|_A^2 \right]. \tag{C.3} \end{aligned}$$

Above,  $\textcircled{1}$  uses Lemma C.5; and  $\textcircled{2}$  is because we have chosen  $\tau$  to satisfy  $\frac{1}{\tau} = 21nL$ .

Next, recall that we have picked  $\alpha_k$  so that  $(21nL - 1)\alpha_k = 21nL \cdot \alpha_{k-1}$  in *CovLPSolver*<sup>wb</sup>. Telescoping (C.3) for  $k = 1, \dots, T$  and choosing  $u^* = (1 + \varepsilon/2)x^*$ , we have

$$\begin{aligned}
 -\sum_{k=1}^T \alpha_k f_\mu(u^*) &\leq 21f_\mu(y_0) - 21n\alpha_T L \\
 \cdot \mathbf{E}[f_\mu(y_T)] + \|z_0 - u^*\|_A^2 &\leq -21n\alpha_T L \cdot \mathbf{E}[f_\mu(y_T)] + 75\text{OPT}.
 \end{aligned}$$

Here, the second inequality is due to  $f_\mu(y_0) = f_\mu(x^{\text{start}}) \leq 3\text{OPT}$  from Fact 5.2, and the fact that

$$\begin{aligned}
 \|z_0 - u^*\|_A^2 &= \|u^*\|_A^2 = \sum_{i=1}^n (u_i^*)^2 \cdot \|A_{:i}\|_\infty \leq (1 + \varepsilon/2)^2 \sum_{i=1}^n (x_i^*)^2 \cdot \|A_{:i}\|_\infty \\
 &\leq 10(1 + \varepsilon/2)^2 \sum_{i=1}^n x_i^* < 12\text{OPT}.
 \end{aligned}$$

Finally, using the fact that  $\sum_{k=1}^T \alpha_k = \alpha_T \cdot \sum_{k=0}^{T-1} (1 - \frac{1}{21nL})^k = 21n\alpha_T L(1 - (1 - \frac{1}{21nL})^T)$ , we rearrange and obtain that

$$\begin{aligned}
 \mathbf{E}[f_\mu(y_T)] &\leq \frac{\sum_k \alpha_k}{21n\alpha_T L} f_\mu(u^*) + \frac{75}{21n\alpha_T L} \text{OPT} = (1 - (1 - \frac{1}{21nL})^T) f_\mu(u^*) \\
 &\quad + \frac{75}{21n\alpha_T L} \text{OPT}.
 \end{aligned}$$

We choose  $T = \lceil 21nL \log(1/\varepsilon) \rceil$  so that  $\frac{1}{n\alpha_T L} = (1 - \frac{1}{21nL})^T \leq \varepsilon$ . Combining this with the fact that  $f_\mu(u^*) \leq (1 + \varepsilon)\text{OPT}$  (see Proposition 4.5a), we obtain

$$\mathbf{E}[f_\mu(y_T)] \leq (1 + \varepsilon)\text{OPT} + 3.6\varepsilon \cdot \text{OPT} < (1 + 4.6\varepsilon)\text{OPT}.$$

Therefore, we have finished proving Theorem 5.3. □

### D.4 Missing proofs for Sect. 6

**Proposition 6.4** *If  $z_{k-1} \in \Delta_{\text{simplex}}$  and  $z_{k-1} > 0$ , the minimizer  $z = \operatorname{argmin}_{z \in \Delta_{\text{simplex}}} \{V_{z_{k-1}}(z) + \langle \delta e_i, z \rangle\}$  for any scalar  $\delta \in \mathbb{R}$  and basis vector  $e_i$  can be computed as follows:*

1.  $z \leftarrow z_{k-1}$ .
2.  $z_i \leftarrow z_i \cdot e^{-\delta}$ .
3. If  $\mathbf{1}^T z > 2\text{OPT}'$ ,  $z \leftarrow \frac{2\text{OPT}'}{\mathbf{1}^T z} z$ .
4. Return  $z$ .

*Proof* Let us denote by  $z$  the returned value of the described procedure, and  $g(u) \stackrel{\text{def}}{=} V_{z_{k-1}}(u) + \langle \delta e_i, u \rangle$ . Since  $\Delta_{\text{simplex}}$  is a convex body and  $g(\cdot)$  is convex, to show  $z = \operatorname{argmin}_{z \in \Delta_{\text{simplex}}} \{g(u)\}$ , it suffices for us to prove that for every  $u \in \Delta_{\text{simplex}}$ ,

$\langle \nabla g(z), u - z \rangle \geq 0$ . Since the gradient  $\nabla g(z)$  can be written explicitly, this is equivalent to

$$\delta(u_i - z_i) + \sum_{\ell=1}^n \log \frac{z_\ell}{z_{k-1,\ell}} \cdot (u_\ell - z_\ell) \geq 0.$$

If the re-scaling in step 3 is not executed, then we have  $z_\ell = z_{k-1,\ell}$  for every  $\ell \neq i$ , and  $z_i = z_{k-1,i} \cdot e^{-\delta}$ ; thus, the left-hand side is zero so the above inequality is true for every  $u \in \Delta_{\text{simplex}}$ .

Otherwise, we have  $\mathbb{1}^T z = 2\text{OPT}'$  and there exists some constant  $Z > 1$  such that,  $z_\ell = z_{k-1,\ell}/Z$  for every  $\ell \neq i$ , and  $z_i = z_{k-1,i} \cdot e^{-\delta}/Z$ . In such a case, the left-hand side equals to

$$(u_i - z_i) \cdot (\delta - \delta) + \sum_{\ell=1}^n -\log Z \cdot (u_\ell - z_\ell).$$

It is clear at this moment that since  $\log Z > 0$  and  $\mathbb{1}^T u \leq 2\text{OPT}' = \mathbb{1}^T z$ , the above quantity is always non-negative, finishing the proof. □

**Lemma 6.13** Denoting by  $\gamma \stackrel{\text{def}}{=} 2\alpha_T n$ , we have

$$\mathbf{E}_i \left[ \alpha_k \langle n\xi_k^{(i)}, z_{k-1} - u^* \rangle \right] \leq V_{z_{k-1}} \left( \frac{u^*}{1 + \gamma} \right) - \mathbf{E}_i \left[ V_{z_k^{(i)}} \left( \frac{u^*}{1 + \gamma} \right) \right] + 12\text{OPT} \cdot \gamma \alpha_k \beta.$$

*Proof* Define  $w(x) \stackrel{\text{def}}{=} \sum_i x_i \log(x_i) - x_i$  and accordingly,  $V_x(y) = w(y) - \langle \nabla w(x), y - x \rangle - w(x) = \sum_i y_i \log \frac{y_i}{x_i} + x_i - y_i$ . We first compute using the classical analysis of mirror descent step as follows:

$$\begin{aligned} & \gamma \alpha_k \langle n\xi_k^{(i)}, z_{k-1} \rangle + \alpha_k \langle n\xi_k^{(i)}, z_{k-1} - u^* \rangle \\ &= (1 + \gamma) \alpha_k \left\langle n\xi_k^{(i)}, z_k - \frac{u^*}{1 + \gamma} \right\rangle + (1 + \gamma) \alpha_k \langle n\xi_k^{(i)}, z_{k-1} - z_k^{(i)} \rangle \\ &\stackrel{\text{①}}{\leq} \left\langle \nabla w(z_{k-1}) - \nabla w(z_k^{(i)}), z_k^{(i)} - \frac{u^*}{1 + \gamma} \right\rangle + (1 + \gamma) \alpha_k \langle n\xi_k^{(i)}, z_{k-1} - z_k^{(i)} \rangle \\ &= \left( w \left( \frac{u^*}{1 + \gamma} \right) - w(z_{k-1}) - \left\langle \nabla w(z_{k-1}), \frac{u^*}{1 + \gamma} - z_{k-1} \right\rangle \right) \\ &\quad - \left( w \left( \frac{u^*}{1 + \gamma} \right) - w(z_k^{(i)}) - \left\langle \nabla w(z_k^{(i)}), \frac{u^*}{1 + \gamma} - z_k^{(i)} \right\rangle \right) \\ &\quad + \left( w(z_{k-1}) - w(z_k^{(i)}) - \langle \nabla w(z_{k-1}), z_{k-1} - z_k^{(i)} \rangle \right) \\ &\quad + (1 + \gamma) \alpha_k \langle n\xi_k^{(i)}, z_{k-1} - z_k^{(i)} \rangle \\ &= V_{z_{k-1}} \left( \frac{u^*}{1 + \gamma} \right) - V_{z_k^{(i)}} \left( \frac{u^*}{1 + \gamma} \right) + \boxed{(1 + \gamma) \alpha_k \langle n\xi_k^{(i)}, z_{k-1} - z_k^{(i)} \rangle - V_{z_{k-1}}(z_k^{(i)})}. \end{aligned} \tag{D.1}$$

Above, ① is because  $\mathbf{z}_k^{(i)} = \operatorname{argmin}_{z \in \Delta_{\text{simplex}}} \{V_{\mathbf{z}_{k-1}}(z) + \langle (1 + \gamma)\alpha_k n \xi_k^{(i)}, z \rangle\}$ , which is equivalent to saying

$$\begin{aligned} \forall u \in \Delta_{\text{simplex}}, \quad & \langle \nabla V_{\mathbf{z}_{k-1}}(\mathbf{z}_k^{(i)}) + (1 + \gamma)\alpha_k n \xi_k^{(i)}, u - \mathbf{z}_k^{(i)} \rangle \geq 0 \\ \iff \forall u \in \Delta_{\text{simplex}}, \quad & \langle \nabla w(\mathbf{z}_k^{(i)}) - \nabla w(\mathbf{z}_{k-1}) + (1 + \gamma)\alpha_k n \xi_k^{(i)}, u - \mathbf{z}_k^{(i)} \rangle \geq 0. \end{aligned}$$

In particular, we have  $\mathbb{1}^T \frac{u^*}{1+\gamma} = \mathbb{1}^T \frac{(1+\varepsilon/2)x^*}{1+\gamma} < 2\text{OPT} \leq 2\text{OPT}'$  and therefore  $\frac{u^*}{1+\gamma} \in \Delta_{\text{simplex}}$ . Substituting  $u = \frac{u^*}{1+\gamma}$  into the above inequality we get ①.

Next, we upper bound the term in the box:

$$\begin{aligned} & (1 + \gamma)\alpha_k \langle n \xi_k^{(i)}, \mathbf{z}_{k-1} - \mathbf{z}_k^{(i)} \rangle - V_{\mathbf{z}_{k-1}}(\mathbf{z}_k^{(i)}) \\ & \stackrel{\textcircled{1}}{\leq} (1 + \gamma)\alpha_k n \xi_{k,i} \cdot (\mathbf{z}_{k-1,i} - \mathbf{z}_{k,i}^{(i)}) - \left( \mathbf{z}_{k,i}^{(i)} \log \frac{\mathbf{z}_{k,i}^{(i)}}{\mathbf{z}_{k-1,i}} + \mathbf{z}_{k-1,i} - \mathbf{z}_{k,i}^{(i)} \right) \\ & \stackrel{\textcircled{2}}{\leq} (1 + \gamma)\alpha_k n \xi_{k,i} \cdot (\mathbf{z}_{k-1,i} - \mathbf{z}_{k,i}^{(i)}) - \frac{|\mathbf{z}_{k,i}^{(i)} - \mathbf{z}_{k-1,i}|^2}{2 \max\{\mathbf{z}_{k,i}^{(i)}, \mathbf{z}_{k-1,i}\}} \\ & \stackrel{\textcircled{3}}{\leq} (1 + \gamma)\alpha_k n \xi_{k,i} \cdot (\mathbf{z}_{k-1,i} - \mathbf{z}_{k,i}^{(i)}) - \frac{|\mathbf{z}_{k,i}^{(i)} - \mathbf{z}_{k-1,i}|^2}{4\mathbf{z}_{k-1,i}} \\ & \stackrel{\textcircled{4}}{\leq} (1 + \gamma)^2 \mathbf{z}_{k-1,i} \cdot (\alpha_k n \xi_{k,i})^2 \stackrel{\textcircled{5}}{\leq} 2\mathbf{z}_{k-1,i} \cdot (\alpha_k n \xi_{k,i})^2 \stackrel{\textcircled{6}}{\leq} \mathbf{z}_{k-1,i} \cdot \gamma \alpha_k n |\xi_{k,i}| \\ & \stackrel{\textcircled{7}}{\leq} \mathbf{z}_{k-1,i} \cdot \gamma \alpha_k n \xi_{k,i} + 2\mathbf{z}_{k-1,i} \cdot \gamma \alpha_k n \beta = \gamma \alpha_k \langle n \xi_k^{(i)}, \mathbf{z}_{k-1} \rangle + 2\mathbf{z}_{k-1,i} \cdot \gamma \alpha_k n \beta. \end{aligned} \tag{D.2}$$

Above, ① uses the facts (i)  $a \log \frac{a}{b} + b - a \geq 0$  for any  $a, b > 0$ , (ii)  $\mathbf{z}_{k-1,i} - \mathbf{z}_k^{(i)}$  and  $\xi_{k,i}$  have the same sign, and (iii)  $\xi_{k,i'} = 0$  for every  $i' \neq i$ ; ② uses the inequality that for every  $a, b > 0$ , we have  $a \log \frac{a}{b} + b - a \geq \frac{(a-b)^2}{2 \max\{a,b\}}$ . ③ uses the fact that  $\mathbf{z}_{k,i}^{(i)} \leq 2\mathbf{z}_{k-1,i}$ .<sup>14</sup> ④ uses Cauchy-Shwarz:  $ab - b^2/4 \leq a^2$ . ⑤ uses  $(1 + \gamma)^2 < 2$ . ⑥ uses  $|\xi_{k,i}| \leq 1$  and  $\gamma = 2\alpha_T n \geq 2\alpha_k n$ . ⑦ uses  $\xi_{k,i} \geq -\beta$ .

Next, we combine (D.1) and (D.2) to conclude that

$$\alpha_k \langle n \xi_k^{(i)}, \mathbf{z}_{k-1} - u^* \rangle \leq V_{\mathbf{z}_{k-1}}\left(\frac{u^*}{1+\gamma}\right) - V_{\mathbf{z}_k^{(i)}}\left(\frac{u^*}{1+\gamma}\right) + 2\mathbf{z}_{k-1,i} \cdot \gamma \alpha_k n \beta.$$

Taking expectation on both sides with respect to  $i$ , and using the property that  $\mathbb{1}^T \mathbf{z}_{k-1} \leq 3\text{OPT}' \leq 6\text{OPT}$ , we obtain that

<sup>14</sup> This is because, our parameter choices ensure that  $(1 + \gamma)\alpha_k n < 1/2\beta$ , which further means  $-(1 + \gamma)\alpha_k n \xi_{k,i}^{(i)} \leq 1/2$ . As a result, we must have  $\mathbf{z}_{k,i}^{(i)} \leq \mathbf{z}_{k-1,i} \cdot e^{0.5} < 2\mathbf{z}_{k-1,i}$  (see the explicit definition of the mirror step at Proposition 6.4).

$$\mathbf{E}_i[\alpha_k \langle n \xi_k^{(i)}, \mathbf{z}_{k-1} - u^* \rangle] \leq V_{\mathbf{z}_{k-1}}\left(\frac{u^*}{1 + \gamma}\right) - \mathbf{E}_i\left[V_{\mathbf{z}_k^{(i)}}\left(\frac{u^*}{1 + \gamma}\right)\right] + 12\text{OPT} \cdot \gamma \alpha_k \beta.$$

□

**Lemma 6.14** *For every  $i \in [n]$ , we have*

- (a)  $f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)}) \geq 0$ , and
- (b)  $f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)}) \geq \frac{\mu\beta}{12} \cdot \langle -\tilde{\eta}_k^{(i)}, u^* \rangle$ .

*Proof of Lemma 6.14 part (a)* Since if  $i \notin B_k$  is not a large index we have  $\mathbf{y}_k^{(i)} = \mathbf{x}_k$  and the claim is trivial, we focus on  $i \in B_k$  in the remaining proof. Recall that  $\mathbf{y}_k^{(i)} = \mathbf{x}_k + \delta \mathbf{e}_i$  for some  $\delta > 0$  defined in Algorithm 3, so we have

$$f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)}) = \int_{\tau=0}^\delta \langle -\nabla f_\mu(\mathbf{x}_k + \tau \mathbf{e}_i), \mathbf{e}_i \rangle d\tau = \int_{\tau=0}^\delta (\langle A_{:i}, p(\mathbf{x}_k + \tau \mathbf{e}_i) \rangle - 1) d\tau.$$

It is clear that  $\langle A_{:i}, p(\mathbf{x}_k + \tau \mathbf{e}_i) \rangle$  decreases as  $\tau$  increases, and therefore it suffices to prove that  $\langle A_{:i}, p(\mathbf{x}_k + \delta \mathbf{e}_i) \rangle \geq 1$ .

Suppose that the rows of  $A_{:i}$  are sorted (for the simplicity of notation) by the increasing order of  $A_{j,i}$ . Now, by the definition of the algorithm (recall (6.1)), there exists some  $j^* \in [m]$  satisfying that

$$\sum_{j < j^*} A_{j,i} \cdot p_j(\mathbf{x}_k) < 1 + \beta \quad \text{and} \quad \sum_{j \leq j^*} A_{j,i} \cdot p_j(\mathbf{x}_k) \geq 1 + \beta.$$

Next, by our choice of  $\delta$  which satisfies  $\delta = \frac{\mu\beta}{2A_{j^*,i}} \leq \frac{\mu\beta}{2A_{j,i}}$  for every  $j \leq j^*$ , we have for every  $j \leq j^*$ :

$$p_j(\mathbf{x}_k + \delta \mathbf{e}_i) = p_j(\mathbf{x}_k) \cdot e^{-\frac{A_{j,i}\delta}{\mu}} \geq p_j(\mathbf{x}_k) \cdot e^{-\beta/2} \geq p_j(\mathbf{x}_k) \cdot (1 - \beta/2),$$

and as a result,

$$\begin{aligned} \langle A_{:i}, p(\mathbf{x}_k + \delta \mathbf{e}_i) \rangle &\geq \sum_{j \leq j^*} A_{j,i} \cdot p_j(\mathbf{x}_k + \delta \mathbf{e}_i) \geq (1 - \beta/2) \sum_{j \leq j^*} A_{j,i} \cdot p_j(\mathbf{x}_k) \\ &\geq (1 - \beta/2)(1 + \beta) \geq 1. \end{aligned}$$

□

*Proof of Lemma 6.14 part (b)* Owing to part (a), for every coordinate  $i$  such that  $\tilde{\eta}_{k,i} \geq 0$ , we automatically have  $f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)}) \geq 0$  so the lemma is obvious. Therefore, let us focus only on coordinates  $i$  such that  $\tilde{\eta}_{k,i} < 0$ ; these are necessarily large indices  $i \in B$ . Recall from Definition 6.11 that  $\tilde{\eta}_{k,i} = (1 + \beta) - (\tilde{A}^T p(\mathbf{x}_k))_i$ , so we have

$$\sum_{j=1}^m \tilde{A}_{j,i} \cdot p_j(\mathbf{x}_k) - (1 + \beta) > 0.$$



For the simplicity of description, suppose again that each  $i$ -th column is sorted in non-decreasing order, that is,  $A_{1,i} \leq \dots \leq A_{m,i}$ . The definition of  $j^*$  can be simplified as

$$\sum_{j < j^*} A_{j,i} \cdot p_j(\mathbf{x}_k) < 1 + \beta \quad \text{and} \quad \sum_{j \leq j^*} A_{j,i} \cdot p_j(\mathbf{x}_k) \geq 1 + \beta.$$

Let  $j^b \in [m]$  be the row such that

$$\sum_{j < j^b} \tilde{A}_{j,i} \cdot p_j(\mathbf{x}_k) < 1 + \beta \quad \text{and} \quad \sum_{j \leq j^b} \tilde{A}_{j,i} \cdot p_j(\mathbf{x}_k) \geq 1 + \beta.$$

Note that such a  $j^b$  must exist because  $\sum_{j=1}^m \tilde{A}_{j,i} \cdot p_j > 1 + \beta$ . It is clear that  $j^b \geq j^*$ , owing to the definition that  $\tilde{A}_{ji} \leq A_{ji}$  for all  $i \in [n], j \in [m]$ . Defining  $\delta^b = \frac{\mu\beta}{2A_{j^b,i}} \leq \delta$ , the objective decrease is lower bounded as

$$\begin{aligned} f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)}) &= \int_{\tau=0}^{\delta} \langle -\nabla f_\mu(\mathbf{x}_k + \tau \mathbf{e}_i), \mathbf{e}_i \rangle d\tau = \int_{\tau=0}^{\delta} (\langle A_{:,i}, p(\mathbf{x}_k + \tau \mathbf{e}_i) \rangle - 1) d\tau \\ &\geq \int_{\tau=0}^{\delta^b} (\langle A_{:,i}, p(\mathbf{x}_k + \tau \mathbf{e}_i) \rangle - 1) d\tau \\ &= \underbrace{\int_{\tau=0}^{\delta^b} \left( -1 + \sum_{j \leq j^b} A_{j,i} \cdot p_j(\mathbf{x}_k + \tau \mathbf{e}_i) \right) d\tau}_I \\ &\quad + \underbrace{\sum_{j > j^b} \int_{\tau=0}^{\delta^b} A_{j,i} \cdot p_j(\mathbf{x}_k + \tau \mathbf{e}_i) d\tau}_{I'} \end{aligned}$$

where the inequality is because  $\delta^b \leq \delta$  and  $\langle A_{:,i}, p(\mathbf{x}_k + \tau \mathbf{e}_i) \rangle \geq 1$  for all  $\tau \leq \delta$  (see the proof of part (a)).

**Part I** To lower bound  $I$ , we use the monotonicity of  $p_j(\cdot)$  and obtain that

$$\begin{aligned} I &= \int_{\tau=0}^{\delta^b} \left( -1 + \sum_{j \leq j^b} A_{j,i} \cdot p_j(\mathbf{x}_k + \tau \mathbf{e}_i) \right) d\tau \geq \delta^b \\ &\quad \cdot \left( -1 + \sum_{j \leq j^b} A_{j,i} \cdot p_j(\mathbf{x}_k + \delta^b \mathbf{e}_i) \right). \end{aligned}$$

However, our choice of  $\delta^b = \frac{\mu\beta}{2A_{j^b,i}} \leq \frac{\mu\beta}{2A_{j,i}}$  for all  $j \leq j^b$  ensures that

$$\begin{aligned} \sum_{j \leq j^b} A_{j,i} \cdot p_j(\mathbf{x}_k + \delta^b \mathbf{e}_i) &\geq \sum_{j \leq j^b} A_{j,i} \cdot p_j(\mathbf{x}_k) \cdot e^{\frac{-A_{j,i} \cdot \delta^b}{\mu}} \\ &\geq \sum_{j \leq j^b} A_{j,i} \cdot p_j(\mathbf{x}_k) \cdot (1 - \beta/2). \end{aligned}$$

Therefore, we obtain that

$$I \geq \delta^b \left( -1 + (1 - \beta/2) \sum_{j \leq j^b} A_{j,i} \cdot p_j(\mathbf{x}_k) \right) \geq \frac{\delta^b}{3} \left( -1 + \sum_{j \leq j^b} A_{j,i} \cdot p_j(\mathbf{x}_k) \right),$$

where the inequality is because  $(\frac{2}{3} - \frac{\beta}{2}) \sum_{j \leq j^b} A_{j,i} \cdot p_j(\mathbf{x}_k) \geq \frac{4-3\beta}{6} \cdot (1 + \beta) \geq \frac{2}{3}$  whenever  $\beta \leq \frac{1}{3}$  (or equivalently, whenever  $\varepsilon \leq 1/9$ ).

Now, suppose that  $\sum_{j \leq j^b} \tilde{A}_{j,i} \cdot p_j(\mathbf{x}_k) - (1 + \beta) = b \cdot \tilde{A}_{j^b,i} \cdot p_{j^b}(\mathbf{x}_k)$  for some  $b \in [0, 1]$ . Note that we can do so by the very definition of  $j^b$ . Then, we must have

$$\begin{aligned} -1 + \sum_{j \leq j^b} A_{j,i} \cdot p_j(\mathbf{x}_k) &\geq -1 + \sum_{j < j^b} \tilde{A}_{j,i} \cdot p_j(\mathbf{x}_k) + A_{j^b,i} \cdot p_{j^b}(\mathbf{x}_k) \\ &= -1 + (1 + \beta) - (1 - b)\tilde{A}_{j^b,i} \cdot p_{j^b}(\mathbf{x}_k) + A_{j^b,i} \cdot p_{j^b}(\mathbf{x}_k) \\ &\geq \beta + b \cdot A_{j^b,i} \cdot p_{j^b}(\mathbf{x}_k). \end{aligned}$$

Therefore, we conclude that

$$\begin{aligned} I &\geq \frac{\delta^b}{3} \left( -1 + \sum_{j \leq j^b} A_{j,i} \cdot p_j(\mathbf{x}_k) \right) > \frac{\delta^b}{3} \cdot b \cdot A_{j^b,i} \cdot p_{j^b}(\mathbf{x}_k) \\ &= \frac{\mu\beta}{6\tilde{A}_{j^b,i}} \cdot b \cdot \tilde{A}_{j^b,i} \cdot p_{j^b}(\mathbf{x}_k) \\ &= \frac{\mu\beta}{6\tilde{A}_{j^b,i}} \cdot \left( -(1 + \beta) + \sum_{j \leq j^b} \tilde{A}_{j,i} \cdot p_j(\mathbf{x}_k) \right) \\ &\geq \frac{\mu\beta}{12} \cdot u_i^* \cdot \left( -(1 + \beta) + \sum_{j \leq j^b} \tilde{A}_{j,i} \cdot p_j(\mathbf{x}_k) \right). \end{aligned}$$

Above, the last inequality is because  $u_i^* \cdot \tilde{A}_{j^b,i} \leq \langle \tilde{A}_{j^b,\cdot}, u^* \rangle \leq 2$  by our definition of  $\tilde{A}$ .

**Part I'** To lower bound  $I'$ , consider every  $j > j^b$  and the integral

$$\int_{\tau=0}^{\delta^b} A_{j,i} \cdot p_j(\mathbf{x}_k + \tau \mathbf{e}_i) d\tau.$$

Note that whenever  $\tau \leq \frac{\mu\beta}{2A_{j,i}} \leq \frac{\mu\beta}{2\tilde{A}_{j,i}} = \delta^b$ , we have that  $p_j(\mathbf{x}_k + \tau \mathbf{e}_i) \geq p_j(\mathbf{x}_k) \cdot e^{-\beta/2} \geq \frac{1}{2} p_j(\mathbf{x}_k)$ . Therefore,

$$\int_{\tau=0}^{\delta^b} A_{j,i} \cdot p_j(\mathbf{x}_k + \tau \mathbf{e}_i) d\tau \geq \int_{\tau=0}^{\frac{\mu\beta}{2\tilde{A}_{j,i}}} A_{j,i} \cdot p_j(\mathbf{x}_k + \tau \mathbf{e}_i) d\tau \geq \frac{\mu\beta}{2A_{j,i}} \cdot A_{j,i} \cdot \frac{1}{2} p_j(\mathbf{x}_k).$$

This implies a lower bound on  $I'$ :

$$I' \geq \sum_{j>j^b} \frac{\mu\beta}{4A_{j,i}} \cdot A_{j,i} \cdot p_j(\mathbf{x}_k) \geq \frac{\mu\beta}{8} \cdot \sum_{j>j^b} u_i^* \cdot \tilde{A}_{j,i} \cdot p_j(\mathbf{x}_k),$$

where again in the last inequality we have used  $u_i^* \cdot \tilde{A}_{j^b,i} \leq \langle \tilde{A}_{j^b}, u^* \rangle \leq 2$  by our definition of  $\tilde{A}$ .

**Together** Combining the lower bounds on  $I$  and  $I'$ , we obtain

$$\begin{aligned} f_\mu(\mathbf{x}_k) - f_\mu(\mathbf{y}_k^{(i)}) &\geq I + I' \geq \frac{\mu\beta}{12} \cdot u_i^* \cdot \left( -(1 + \beta) + \sum_{j=1}^m \tilde{A}_{j,i} \cdot p_j(\mathbf{x}_k) \right) \\ &= \frac{\mu\beta}{12} \cdot \langle -\tilde{\eta}_k^{(i)}, u^* \rangle. \end{aligned}$$

□

### E.5 Proof of Lemma 3.3: Efficient Implementation of *PacLP Solver*

In this section, we illustrate how to implement each iteration of *PacLP Solver* to run in an expected  $O(N/n)$  time. We maintain the following quantities

$$\mathbf{z}_k \in \mathbb{R}_{\geq 0}^n, \quad \mathbf{az}_k \in \mathbb{R}_{\geq 0}^m, \quad \mathbf{y}'_k \in \mathbb{R}^n, \quad \mathbf{ay}'_k \in \mathbb{R}^m, \quad B_{k,1}, B_{k,2} \in \mathbb{R}_+$$

throughout the algorithm, so as to ensure the following invariants are always satisfied

$$A\mathbf{z}_k = \mathbf{az}_k, \tag{E.1}$$

$$\mathbf{y}_k = B_{k,1} \cdot \mathbf{z}_k + B_{k,2} \cdot \mathbf{y}'_k, \quad A\mathbf{y}'_k = \mathbf{ay}'_k. \tag{E.2}$$

It is clear that when  $k = 0$ , letting  $\mathbf{az}_k = A\mathbf{z}_0$ ,  $\mathbf{y}'_k = \mathbf{y}_0$ ,  $\mathbf{ay}'_k = A\mathbf{y}_0$ ,  $B_{k,1} = 0$ , and  $B_{k,2} = 1$ , we can ensure that all the invariants are satisfied initially. We denote  $\|A_{\cdot i}\|_0$  the number of nonzeros elements in vector  $A_{\cdot i}$ . In each iteration  $k = 1, 2, \dots, T$ :

- The step  $\mathbf{x}_k = \tau\mathbf{z}_{k-1} + (1 - \tau)\mathbf{y}_{k-1}$  does not need to be implemented.
- The value  $\nabla_i f(\mathbf{x}_k)$  requires the knowledge of  $p_j(\mathbf{x}_k) = e^{\frac{1}{\mu} \langle A\mathbf{x}_k, \mathbf{j} \rangle}$  for each  $j$  such that  $A_{ij} \neq 0$ . Accordingly, for each  $j$ , we need to know the value

$$\begin{aligned} (A\mathbf{x}_k)_j &= \tau(A\mathbf{z}_{k-1})_j + (1 - \tau)(A\mathbf{y}_{k-1})_j \\ &= (\tau + (1 - \tau)B_{k-1,1})\mathbf{az}_{k-1,j} + (1 - \tau)B_{k-1,2}\mathbf{ay}'_{k-1,j}. \end{aligned}$$

This can be computed in  $O(1)$  time for each  $j$ , and  $O(\|A_{\cdot i}\|_0)$  time in total.

- Recall the step  $\mathbf{z}_k \leftarrow \operatorname{argmin}_{\mathbf{z} \in \Delta_{\text{box}}} \left\{ \frac{1}{2} \|\mathbf{z} - \mathbf{z}_{k-1}\|_A^2 + \langle n\alpha_k \xi_k^{(i)}, \mathbf{z} \rangle \right\}$  can be written as  $\mathbf{z}_k = \mathbf{z}_{k-1} + \delta \mathbf{e}_i$  for some  $\delta \in \mathbb{R}$  that can be computed in  $O(1)$  time (see Proposition 3.2). Observe also  $\mathbf{z}_k = \mathbf{z}_{k-1} + \delta \mathbf{e}_i$  yields  $\mathbf{y}_k = \tau\mathbf{z}_{k-1} + (1 - \tau)\mathbf{y}_{k-1} +$

$\frac{\delta \mathbf{e}_i}{n\alpha_k L}$  due to Line 6 and Line 10 of Algorithm 1. Therefore, we perform two explicit updates on  $\mathbf{z}_k$  and  $\mathbf{az}_k$  as

$$\mathbf{z}_k \leftarrow \mathbf{z}_{k-1} + \delta \mathbf{e}_i, \quad \mathbf{az}_k \leftarrow \mathbf{az}_{k-1} + \delta A_{:i}$$

and two implicit updates on  $\mathbf{y}_k$  as

$$B_{k,1} = \tau + (1 - \tau)B_{k-1,1}, \quad B_{k,2} = (1 - \tau)B_{k-1,2}, \\ \mathbf{y}'_k \leftarrow \mathbf{y}'_{k-1} + \delta \mathbf{e}_i \cdot \left( -\frac{B_{k,1}}{B_{k,2}} + \frac{1}{n\alpha_k L} \frac{1}{B_{k,2}} \right), \quad \mathbf{ay}'_k \leftarrow \mathbf{ay}'_{k-1} + \delta A_{:i} \cdot \left( -\frac{B_{k,1}}{B_{k,2}} + \frac{1}{n\alpha_k L} \frac{1}{B_{k,2}} \right)$$

It is not hard to verify that after these updates,  $\mathbf{Ay}'_k = \mathbf{ay}'_k$  and we have

$$\begin{aligned} B_{k,1} \cdot \mathbf{z}_k + B_{k,2} \cdot \mathbf{y}'_k &= B_{k,1} \cdot (\mathbf{z}_{k-1} + \delta \mathbf{e}_i) \\ &+ B_{k,2} \cdot \left( \mathbf{y}'_{k-1} + \delta \mathbf{e}_i \cdot \left( -\frac{B_{k,1}}{B_{k,2}} + \frac{1}{n\alpha_k L} \frac{1}{B_{k,2}} \right) \right) \\ &= B_{k,1} \cdot \mathbf{z}_{k-1} + B_{k,2} \cdot \left( \mathbf{y}'_{k-1} + \delta \mathbf{e}_i \cdot \left( \frac{1}{n\alpha_k L} \frac{1}{B_{k,2}} \right) \right) \\ &= B_{k,1} \cdot \mathbf{z}_{k-1} + B_{k,2} \cdot \mathbf{y}'_{k-1} + \frac{\delta \mathbf{e}_i}{n\alpha_k L} \\ &= (\tau + (1 - \tau)B_{k-1,1}) \cdot \mathbf{z}_{k-1} + ((1 - \tau)B_{k-1,2}) \cdot \mathbf{y}'_{k-1} + \frac{\delta \mathbf{e}_i}{n\alpha_k L} \\ &= \tau \mathbf{z}_{k-1} + (1 - \tau)\mathbf{y}_{k-1} + \frac{\delta \mathbf{e}_i}{n\alpha_k L} = \mathbf{y}_k, \end{aligned}$$

so the invariant  $\mathbf{y}_k = B_{k,1} \cdot \mathbf{z}_k + B_{k,2} \cdot \mathbf{y}'_k$  also holds. In sum, after performing updates on  $\mathbf{AZ}_k$  and  $\mathbf{ay}'_k$  in time  $O(\|A_{:i}\|_0)$ , we can ensure that the invariants in (E.1) and (E.2) are satisfied at iteration  $k$ .

In sum, we only need  $O(\|A_{:i}\|_0)$  time to perform the updates in *PacLP Solver* for an iteration  $k$  if the coordinate  $i$  is selected. Therefore, each iteration of *PacLP Solver* can be implemented to run in an expected  $O(\mathbf{E}_i[\|A_{:i}\|_0]) = O(N/n)$  time.

## F.6 Proof of Lemma 6.5: Efficient Implementation of *CovLP Solver*

In this section we illustrate how to implement each iteration of *CovLP Solver* to run in an expected  $O(N/n)$  time. We maintain the following quantities

$$\mathbf{z}'_k \in \mathbb{R}_+^n, \quad \mathbf{sz}_k \in \mathbb{R}_+, \quad \mathbf{sumz}_k \in \mathbb{R}_+, \quad \mathbf{az}'_k \in \mathbb{R}_{\geq 0}^m, \quad \mathbf{y}'_k \in \mathbb{R}^n, \\ \mathbf{ay}'_k \in \mathbb{R}^m, \quad B_{k,1}, B_{k,2} \in \mathbb{R}_+$$

throughout the algorithm, so as to maintain the following invariants

$$\mathbf{z}_k = \mathbf{z}'_k / \mathbf{sz}_k, \quad \mathbf{sumz}_k = \mathbf{1}^T \mathbf{z}'_k, \quad \mathbf{AZ}_k = \mathbf{az}'_k / \mathbf{sz}_k, \quad (\text{F.1})$$

$$\mathbf{y}_k = B_{k,1} \cdot \mathbf{z}'_k + B_{k,2} \cdot \mathbf{y}'_k, \quad \mathbf{Ay}_k = \mathbf{ay}'_k. \quad (\text{F.2})$$

It is clear that when  $k = 0$ , letting  $\mathbf{z}'_k = \mathbf{z}_0$ ,  $\mathbf{sz}_k = 1$ ,  $\text{sumz}_k = \mathbb{1}^T \mathbf{z}_0$ ,  $\mathbf{az}'_k = A\mathbf{z}_0$ ,  $\mathbf{y}'_k = \mathbf{y}_0$ ,  $\mathbf{ay}'_k = A\mathbf{y}_0$ ,  $B_{k,1} = 0$ , and  $B_{k,2} = 1$ , we can ensure that all the invariants are satisfied initially.

We denote by  $\|A_{:i}\|_0$  the number of nonzero elements in vector  $A_{:i}$ . In each iteration  $k = 1, 2, \dots, T$ :

- The step  $\mathbf{x}_k = \tau \mathbf{z}_{k-1} + (1 - \tau) \mathbf{y}_{k-1}$  does not need to be implemented.
- The value  $p_j(\mathbf{x}_k) = e^{\frac{1}{\mu}(1 - (A\mathbf{x}_k)_j)}$  for each  $j$  only requires the knowledge of

$$\begin{aligned} (A\mathbf{x}_k)_j &= \tau(A\mathbf{z}_{k-1})_j + (1 - \tau)(A\mathbf{y}_{k-1})_j \\ &= (\tau + (1 - \tau)B_{k-1,1}) \frac{\mathbf{az}'_{k-1,j}}{\mathbf{sz}_{k-1}} + (1 - \tau)B_{k-1,2} \mathbf{ay}'_{k-1,j}. \end{aligned}$$

This can be computed in  $O(1)$  time.

- The value  $\nabla_i f(\mathbf{x}_k)$  requires the knowledge of  $p_j(\mathbf{x}_k)$  for each  $j \in [m]$  such that  $A_{ij} \neq 0$ . Since we have  $\|A_{:i}\|_0$  such  $j$ 's, we can compute  $\nabla_i f(\mathbf{x}_k)$  in  $O(\|A_{:i}\|_0)$  time.
- Letting  $\delta = (1 + \gamma)n\alpha_k \xi_{k,i}^{(i)}$ , recall that the mirror step  $\mathbf{z}_k \leftarrow \text{argmin}_{\mathbf{z} \in \Delta_{\text{simplex}}} \{V_{\mathbf{z}_{k-1}}(\mathbf{z}) + \langle \delta \mathbf{e}_i, \mathbf{z} \rangle\}$  has a very simple form (see Proposition 6.4): first multiply the  $i$ -th coordinate of  $\mathbf{z}_{k-1}$  by  $e^{-\delta}$  and then, if the sum of all coordinates have exceeded  $2\text{OPT}'$ , scale everything down so as to sum up to  $2\text{OPT}'$ . This can be implemented as follows: setting  $\delta_1 = \mathbf{z}'_{k-1,i}(e^{-\delta} - 1)$ ,

$$\begin{aligned} \mathbf{z}'_k &\leftarrow \mathbf{z}'_{k-1} + \delta_1 \mathbf{e}_i, \mathbf{az}'_k \leftarrow \mathbf{az}'_{k-1} + \delta_1 A_{:i}, \\ \text{sumz}_k &\leftarrow \text{sumz}_{k-1} + \delta_1, \mathbf{sz}_k \leftarrow \mathbf{sz}_{k-1} \cdot \max \left\{ 1, \frac{\text{sumz}_k}{\mathbf{sz}_{k-1} \cdot 2\text{OPT}'} \right\}. \end{aligned}$$

These updates can be implemented to run in  $O(\|A_{:i}\|_0)$  time, and they together ensure that the invariants in (F.1) are satisfied at iteration  $k$ .

- Recall that the gradient step is of the form  $\mathbf{y}_k \leftarrow \mathbf{x}_k + \delta_2 \cdot \mathbf{e}_i$  for some value  $\delta_2 \geq 0$ . This value  $\delta_2$  can be computed in  $O(\|A_{:i}\|_0)$  time, since each  $p_j(\mathbf{x}_k)$  can be computed in  $O(1)$  time, and we can sort the rows of each column of  $A$  by preprocessing.

Since  $\mathbf{y}_k = \mathbf{x}_k + \delta_2 \cdot \mathbf{e}_i = \tau \mathbf{z}_{k-1} + (1 - \tau) \mathbf{y}_{k-1} + \delta_2 \mathbf{e}_i$ , we can implement this update by letting

$$\begin{aligned} B_{k,1} &= \frac{\tau}{\mathbf{sz}_{k-1}} + (1 - \tau)B_{k-1,1}, B_{k,2} = (1 - \tau)B_{k-1,2} \\ \mathbf{y}'_k &\leftarrow \mathbf{y}'_{k-1} + \mathbf{e}_i \cdot \left( -\frac{B_{k,1}\delta_1}{B_{k,2}} + \frac{\delta_2}{B_{k,2}} \right), \mathbf{ay}'_k \leftarrow \mathbf{ay}'_{k-1} + A_{:i} \cdot \left( -\frac{B_{k,1}\delta_1}{B_{k,2}} + \frac{\delta_2}{B_{k,2}} \right) \end{aligned}$$

It is not hard to verify that after these updates,  $\mathbf{ay}'_k = A\mathbf{y}'_k$  and we have

$$\begin{aligned} B_{k,1} \cdot \mathbf{z}'_k + B_{k,2} \cdot \mathbf{y}'_k &= B_{k,1} \cdot (\mathbf{z}'_{k-1} + \delta_1 \mathbf{e}_i) \\ &\quad + B_{k,2} \cdot \left( \mathbf{y}'_{k-1} + \mathbf{e}_i \cdot \left( -\frac{B_{k,1}\delta_1}{B_{k,2}} + \frac{\delta_2}{B_{k,2}} \right) \right) \\ &= B_{k,1} \cdot \mathbf{z}'_{k-1} + B_{k,2} \cdot (\mathbf{y}'_{k-1} + \delta_2 \mathbf{e}_i / B_{k,2}) \end{aligned}$$

$$\begin{aligned}
&= B_{k,1} \cdot \mathbf{z}'_{k-1} + B_{k,2} \cdot \mathbf{y}'_{k-1} + \delta_2 \mathbf{e}_i \\
&= \left( \frac{\tau}{\mathbf{S}\mathbf{z}_{k-1}} + (1 - \tau)B_{k-1,1} \right) \cdot \mathbf{z}'_{k-1} + \left( (1 - \tau)B_{k-1,2} \right) \cdot \mathbf{y}'_{k-1} + \delta_2 \mathbf{e}_i \\
&= \tau \mathbf{z}_{k-1} + (1 - \tau)\mathbf{y}_{k-1} + \delta_2 \mathbf{e}_i = \mathbf{y}_k,
\end{aligned}$$

so that the invariant  $\mathbf{y}_k = B_{k,1} \cdot \mathbf{z}'_k + B_{k,2} \cdot \mathbf{y}'_k$  is also satisfied. In sum, after running time  $O(\|A_{:i}\|_0)$ , we can ensure that the invariants in (F2) are satisfied at iteration  $k$ .

In sum, we only need  $O(\|A_{:i}\|_0)$  time to perform the updates in *CovLP Solver* for an iteration  $k$  if the coordinate  $i$  is selected. Therefore, each iteration of *CovLP Solver* can be implemented to run in an expected  $O(\mathbf{E}_i[\|A_{:i}\|_0]) = O(N/n)$  time.

## References

- Allen-Zhu, Z., Lee, Y.T., Orecchia, L.: Using optimization to obtain a width-independent, parallel, simpler, and faster positive SDP solver. In: SODA (2016)
- Allen-Zhu, Z., Li, Y., Oliveira, R., Wigderson, A.: Much faster algorithms for matrix scaling. In: FOCS, 2017. [arXiv:1704.02315](https://arxiv.org/abs/1704.02315)
- Allen-Zhu, Z., Liao, Z., Orecchia, L.: Spectral sparsification and regret minimization beyond multiplicative updates. In: STOC (2015)
- Allen-Zhu, Z., Orecchia, L.: Using optimization to break the epsilon barrier: a faster and simpler width-independent algorithm for solving positive linear programs in parallel. In: SODA (2015)
- Allen-Zhu, Z., Orecchia, L.: Linear coupling: an ultimate unification of gradient and mirror descent. In: ITCS (2017)
- Allen-Zhu, Z., Qu, Z., Richtárik, P., Yuan, Y.: Even faster accelerated coordinate descent using non-uniform sampling. In: ICML (2016)
- Arora, S., Hazan, E., Kale, S.: The multiplicative weights update method: a meta-algorithm and applications. *Theory Comput.* **8**, 121–164 (2012)
- Awerbuch, B., Khandekar, R.: Stateless distributed gradient descent for positive linear programs. In: STOC (2008)
- Awerbuch, B., Khandekar, R., Rao, S.: Distributed algorithms for multicommodity flow problems via approximate steepest descent framework. *ACM Trans. Algorithms* **9**(1), 1–14 (2012)
- Bartal, Y., Byers, J.W., Raz, D.: Global optimization using local information with applications to flow control. In: Proceedings 38th Annual Symposium on Foundations of Computer Science, pp. 303–312. IEEE Computer Society (1997)
- Bartal, Y., Byers, J.W., Raz, D.: Fast, distributed approximation algorithms for positive linear programming with applications to flow control. *SIAM J. Comput.* **33**(6), 1261–1279 (2004)
- Ben-Tal, A., Arkadi, N.: Lectures on modern convex optimization. *Soc. Ind. Appl. Math.* 315–341 (2013)
- Bienstock, D., Iyengar, G.: Faster approximation algorithms for packing and covering problems. Technical report, Columbia University, September 2004. Preliminary version published in STOC '04
- Byers, J., Nasser, G.: Utility-based decision-making in wireless sensor networks. In: 2000 first annual workshop on mobile and ad hoc networking and computing, 2000. MobiHOC, pp. 143–144. IEEE (2000)
- Chudak, F.A., Eleutério, V.: Improved approximation schemes for linear programming relaxations of combinatorial optimization problems. In: Proceedings of the 11th International IPCO Conference on Integer Programming and Combinatorial Optimization, pp. 81–96 (2005)
- Duan, R., Pettie, S.: Linear-time approximation for maximum weight matching. *J. ACM* **61**(1), 1–23 (2014)
- Fercoq, O., Richtárik, P.: Accelerated, parallel and proximal coordinate descent. *SIAM J. Optim.* **25**(4), 1997–2023 (2015)
- Flischer, L.K.: Approximating fractional multicommodity flow independent of the number of commodities. *SIAM J. Discrete Math.* **13**(4), 505–520 (2000)

19. Garg, N., Könemann, J.: Faster and simpler algorithms for multicommodity flow and other fractional packing problems. *SIAM J. Comput.* **37**(2), 630–652 (2007)
20. Grigoriadis, M.D., Khachiyan, L.G.: Fast approximation schemes for convex programs with many blocks and coupling constraints. *SIAM J. Optim.* **4**(1), 86–107 (1994)
21. Jain, R., Ji, Z., Upadhyay, S., Watrous, J.: QIP = PSPACE. *J. ACM JACM* **58**(6), 30 (2011)
22. Klein, P., Young, N.E.: On the number of iterations for Dantzig–Wolfe optimization and packing-covering approximation algorithms. *SIAM J. Comput.* **44**(4), 1154–1172 (2015)
23. Koufogiannakis, C., Young, N.E.: A nearly linear-time PTAS for explicit fractional packing and covering linear programs. *Algorithmica* **70**, 494–506 (2013). **(Previously appeared in FOCS '07)**
24. Luby, M., Nisan, N.: A parallel approximation algorithm for positive linear programming. In: *STOC*, pp. 448–457. ACM Press, New York (1993)
25. Madry, A.: Faster approximation schemes for fractional multicommodity flow problems via dynamic graph algorithms. In: *STOC*. ACM Press, New York (2010)
26. Nemirovski, A.: Prox-method with rate of convergence  $O(1/t)$  for variational inequalities with Lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM J. Optim.* **15**(1), 229–251 (2004)
27. Nesterov, Y.: Rounding of convex sets and efficient gradient methods for linear programming problems. *Optim. Methods Softw.* **23**(1), 109–128 (2008)
28. Nesterov, Y.: A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ . *Sov. Math. Dokl.* **269**, 543–547 (1983)
29. Nesterov, Y.: Smooth minimization of non-smooth functions. *Math. Program.* **103**(1), 127–152 (2005)
30. Nesterov, Y.: Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM J. Optim.* **22**(2), 341–362 (2012)
31. Plotkin, S.A., Shmoys, D.B., Tardos, É.: Fast approximation algorithms for fractional packing and covering problems. *Math. Oper. Res.* **20**(2), 257–301 (1995). **(conference version published in FOCS 1991)**
32. Trevisan, L.: Parallel approximation algorithms by positive linear programming. *Algorithmica* **21**(1), 72–88 (1998)
33. Wang, D., Mahoney, M., Mohan, N., Rao, S.: Faster parallel solver for positive linear programs via dynamically-bucketed selective coordinate descent. ArXiv e-prints [arXiv:1511.06468](https://arxiv.org/abs/1511.06468) (2015)
34. Wang, D., Rao, S., Mahoney, M.W.: Unified acceleration method for packing and covering problems via diameter reduction. In: *ICALP* (2016)
35. Young, N.E.: Sequential and parallel algorithms for mixed packing and covering. In: *42nd Annual IEEE Symposium on Foundations of Computer Science (FOCS'01)*, pp. 538–546. IEEE Computer Society (2001)
36. Young, N.E.: Nearly linear-time approximation schemes for mixed packing/covering and facility-location linear programs. ArXiv e-prints [arXiv:1407.3015](https://arxiv.org/abs/1407.3015) (2014)
37. Zurel, E., Nisan, N.: An efficient approximate allocation algorithm for combinatorial auctions. In: *Proceedings of the 3rd ACM Conference on Electronic Commerce*, pp. 125–136. ACM (2001)