
FASTER ALGORITHMS FOR STRUCTURED LINEAR AND KERNEL SUPPORT VECTOR MACHINES

Anonymous authors

Paper under double-blind review

ABSTRACT

Quadratic programming is a ubiquitous prototype in convex programming. Many machine learning problems can be formulated as quadratic programming, including the famous Support Vector Machines (SVMs). Linear and kernel SVMs have been among the most popular models in machine learning over the past three decades, prior to the deep learning era.

Generally, a quadratic program has an input size of $\Theta(n^2)$, where n is the number of variables. Assuming the Strong Exponential Time Hypothesis (SETH), it is known that no $O(n^{2-o(1)})$ time algorithm exists when the quadratic objective matrix is positive semidefinite (Backurs, Indyk, and Schmidt, NeurIPS'17). However, problems such as SVMs usually admit much smaller input sizes: one is given n data points, each of dimension d , and d is oftentimes much smaller than n . Furthermore, the SVM program has only $O(1)$ equality linear constraints. This suggests that faster algorithms are feasible, provided the program exhibits certain structures.

In this work, we design the first nearly-linear time algorithm for solving quadratic programs whenever the quadratic objective admits a low-rank factorization, and the number of linear constraints is small. Consequently, we obtain results for SVMs:

- For linear SVM when the input data is d -dimensional, our algorithm runs in time $\tilde{O}(nd^{(\omega+1)/2} \log(1/\epsilon))$ where $\omega \approx 2.37$ is the fast matrix multiplication exponent;
- For Gaussian kernel SVM, when the data dimension $d = O(\log n)$ and the squared dataset radius is sub-logarithmic in n , our algorithm runs in time $O(n^{1+o(1)} \log(1/\epsilon))$. We also prove that when the squared dataset radius is at least $\Omega(\log^2 n)$, then $\Omega(n^{2-o(1)})$ time is required. This improves upon the prior best lower bound in both the dimension d and the squared dataset radius.

1 INTRODUCTION

Quadratic programming (QP) represents a class of convex optimization problems that optimize a quadratic objective over the intersection of an affine subspace and the non-negative orthant¹. QPs naturally extend linear programming by incorporating a quadratic objective, and they find extensive applications in operational research, theoretical computer science, and machine learning (Kozlov et al., 1979; Wright, 1999; Gould & Toint, 2000; Gould et al., 2001; Propato & Uber, 2004; Cornuejols & Tütüncü, 2006). The quadratic objective introduces challenges: QPs with a general (not necessarily positive semidefinite) symmetric quadratic objective matrix are NP-hard to solve (Sahni, 1974; Pardalos & Vavasis, 1991). When the quadratic objective matrix is positive semidefinite, the problem becomes weakly polynomial-time solvable, as it can be reduced to convex empirical risk minimization (Lee et al., 2019) (refer to Section C for further discussion).

Formally, the QP problem is defined as follows:

Definition 1.1 (Quadratic Programming). *Given an $n \times n$ symmetric, positive semidefinite objective matrix Q , a vector $c \in \mathbb{R}^n$, and a polytope described by a pair $(A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m)$, the linearly*

¹There are classes of QPs with quadratic constraints as well. However, in this paper, we focus on cases where the constraints are linear.

constrained quadratic programming (LCQP) or simply quadratic programming (QP) problem seeks to solve the following optimization problem:

$$\begin{aligned} \min_{x \in \mathbb{R}^n} \quad & \frac{1}{2} x^\top Q x + c^\top x \\ \text{s.t.} \quad & A x = b \\ & x \geq 0. \end{aligned} \quad (1)$$

A classic application of QP is the Support Vector Machine (SVM) problem (Boser et al., 1992; Cortes & Vapnik, 1995). In SVMs, a dataset $x_1, \dots, x_n \in \mathbb{R}^d$ is provided, along with corresponding labels $y_1, \dots, y_n \in \{\pm 1\}$. The objective is to identify a hyperplane that separates the two groups of points with opposite labels, while maintaining a large margin from both. Remarkably, this popular machine learning problem can be formulated as a QP and subsequently solved using specialized QP solvers (Muller et al., 2001). Thus, advancements in QP algorithms could potentially lead to runtime improvements for SVMs.

Despite its practical and theoretical significance, algorithmic quadratic programming has garnered relatively less attention compared to its close relatives in convex programming, such as linear programming (Cohen et al., 2019; Jiang et al., 2021; Brand, 2020; Song & Yu, 2021), convex empirical risk minimization (Lee et al., 2019; Qin et al., 2023), and semidefinite programming (Jiang et al., 2020; Huang et al., 2022; Gu & Song, 2022). In this work, we take a pioneering step in developing fast and robust interior point-type algorithms for solving QPs. We particularly focus on improving the runtime for high-precision hard- and soft-margin SVMs. For the purposes of this discussion, we will concentrate on hard-margin SVMs, with the understanding that our results naturally extend to soft-margin variants. We begin by introducing the hard-margin linear SVMs:

Definition 1.2 (Linear SVM). *Given a dataset $X \in \mathbb{R}^{n \times d}$ and a collection of labels y_1, \dots, y_n each in ± 1 , the linear SVM problem requires solving the following quadratic program:*

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^n} \quad & \mathbf{1}_n^\top \alpha - \frac{1}{2} \alpha^\top (y y^\top \circ X X^\top) \alpha, \\ \text{s.t.} \quad & \alpha^\top y = 0, \\ & \alpha \geq 0. \end{aligned} \quad (2)$$

where \circ denotes the Hadamard product.

It should be noted that this formulation is actually the *dual* of the SVM optimization problem. The primal program seeks a vector $w \in \mathbb{R}^d$ such that

$$\begin{aligned} \min_{w \in \mathbb{R}^d} \quad & \frac{1}{2} \|w\|_2^2, \\ \text{s.t.} \quad & y_i (w^\top x_i - b) \geq 1, \quad \forall i \in [n], \end{aligned}$$

where $b \in \mathbb{R}$ is the bias term. Given the solution $\alpha \in \mathbb{R}^n$, one can conveniently convert it to a primal solution: $w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$. At first glance, one might be inclined to solve the primal problem directly, especially in cases where $d \ll n$, as it presents a lower-dimensional optimization problem compared to the dual. The dual formulation becomes particularly advantageous when solving the kernel SVM, which maps features to a high or potentially infinite-dimensional space.

Definition 1.3 (Kernel SVM). *Given a dataset $X \in \mathbb{R}^{n \times d}$ and a positive definite kernel function $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, let $K \in \mathbb{R}^{n \times n}$ denote the kernel matrix, where $K_{i,j} = K(x_i, x_j)$. With a collection of labels y_1, \dots, y_n each in $\{\pm 1\}$, the kernel SVM problem requires solving the following quadratic program:*

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^n} \quad & \mathbf{1}_n^\top \alpha - \frac{1}{2} \alpha^\top (y y^\top \circ K) \alpha, \\ \text{s.t.} \quad & \alpha^\top y = 0, \\ & \alpha \geq 0. \end{aligned} \quad (3)$$

The positive definite kernel function K corresponds to a feature mapping, implying that $K(x_i, x_j) = \phi(x_i)^\top \phi(x_j)$ for some $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^s$. Thus, solving the primal SVM can be viewed as solving the

108 optimization problem on the transformed dataset. However, the primal program’s dimension depends
 109 on the (transformed) data’s dimension s , which can be infinite. Conversely, the dual program, with
 110 dimension n , is typically easier to solve. Throughout this paper, when discussing the SVM program,
 111 we implicitly refer to the dual quadratic program, not the primal.

112 One key aspect of the SVM program is its minimal equality constraints. Specifically, for both linear
 113 and kernel SVMs, there is only a single equality constraint of the form $\alpha^\top y = 0$. This constraint
 114 arises naturally from the bias term in the primal SVM formulation and its Lagrangian. The limited
 115 number of constraints enables us to design QP solvers with favorable dependence on the number of
 116 data points n , albeit with a higher dependence on the number of constraints m , thus offering effective
 117 end-to-end guarantees for SVMs.

118 Previous efforts to solve the SVM program typically involve breaking down the large QP into smaller,
 119 constant-sized QPs. These algorithms, while easy to implement and well-suited to modern hardware
 120 architectures, oftentimes lack tight theoretical analysis and the estimation of iteration count is usually
 121 pessimistic (Chang & Lin, 2011). Theoretically, Joachims (2006) systematically analyzed this class
 122 of algorithms, demonstrating that to achieve an ϵ -approximation solution, $\tilde{O}(\epsilon^{-2}B \cdot \text{nnz}(A))$ time is
 123 sufficient, where B is the squared-radius of the dataset and $\text{nnz}(\cdot)$ denotes the number of nonzero
 124 entries. **This is subsequently improved in Shalev-Shwartz et al. (2011) with a subgradient-based**
 125 **method that runs in $\tilde{O}(\epsilon^{-1}d)$ time.** Unfortunately, the polynomial dependence on the precision ϵ^{-1}
 126 makes them hard to be adapted for even moderately small ϵ . For example, when ϵ is set to be 10^{-3} to
 127 account for the usual machine precision errors, these algorithms would require at least 10^3 iterations
 128 to converge.

129 To develop a high-precision algorithm with $\text{poly} \log(\epsilon^{-1})$ dependence instead of $\text{poly}(\epsilon^{-1})$, we
 130 focus on second-order methods for QPs. A variety of approaches have been explored in previous
 131 works, including the interior point method (Karmarkar, 1984), active set methods (Murty, 1988),
 132 augmented Lagrangian techniques (Delbos & Gilbert, 2003), conjugate gradient, gradient projection,
 133 and extensions of the simplex algorithm (Dantzig, 1955; Wolfe, 1959; Murty, 1988). Our interest is
 134 particularly piqued by the interior point method (IPM). Recent advances in the robust IPM framework
 135 have led to significant successes for convex programming problems (Cohen et al., 2019; Lee et al.,
 136 2019; Brand, 2020; Jiang et al., 2020; Brand et al., 2020; Jiang et al., 2021; Song & Yu, 2021; Jiang
 137 et al., 2022; Huang et al., 2022; Gu & Song, 2022; Qin et al., 2023). These successes are a result of
 138 combining robust analysis of IPM with dedicated data structure design.

139 Applying IPM to solve QPs with a constant number of constraints is not entirely novel; existing
 140 work (Ferris & Munson, 2002) has already adapted IPM to solve the linear SVM problem. However,
 141 the runtime of their algorithm is sub-optimal. Each iteration of their algorithm requires multiplying a
 142 $d \times n$ matrix with an $n \times d$ matrix in $O(nd^\omega)$ time, where $\omega \approx 2.37$ is the fast matrix multiplication
 143 exponent (Duan et al., 2023; Williams et al., 2024; Gall, 2024). Moreover, the IPM requires
 144 $O(\sqrt{n} \log(1/\epsilon))$ iterations to converge. This ends up with an overall runtime $O(n^{1.5}d^{\omega-1} \log(1/\epsilon))$,
 145 which is super-linear in the dataset size even when the dimension d is small. In practical scenarios
 146 where n is usually large, the $n^{1.5}$ dependence becomes prohibitive. Therefore, it is crucial to develop
 147 an algorithm with almost- or nearly-linear dependence on n and logarithmic dependence on ϵ^{-1} .

148 For linear SVM, we propose a nearly-linear time algorithm with high-precision guarantees, applicable
 149 when the dimension of the dataset is smaller than the number of points:

150 **Theorem 1.4** (Low-rank QP and Linear SVM, informal version of Theorem E.1). *Given a quadratic*
 151 *program as defined in Definition 1.1, and assuming a low-rank factorization of the quadratic objective*
 152 *matrix $Q = UV^\top$, where $U, V \in \mathbb{R}^{n \times k}$, there exists an algorithm that can solve the program (1) up*
 153 *to ϵ -error² in $\tilde{O}(n(k+m)^{(\omega+1)/2} \log(n/\epsilon))$ time.*

154 *Specifically, for linear SVM (as per Definition 1.2) with $d \leq n$, one can solve program (2) up to*
 155 *ϵ -error in $\tilde{O}(nd^{(\omega+1)/2} \log(n/\epsilon))$ time.*

157 While a nearly-linear time algorithm for linear SVMs is appealing, most applications look at kernel
 158 SVMs as they provide more expressive power to the linear classifier. This poses significant challenge
 159 in algorithm design, as forming the kernel matrix exactly would require $\Omega(n^2)$ time. Moreover, the
 160

161 ²We say an algorithm that solves the program up to ϵ -error if it returns an approximate solution vector $\tilde{\alpha}$
 whose objective value is at most ϵ more than the optimal objective value.

kernel matrix could be full-rank without any structural assumptions, rendering our low-rank QP solver inapplicable. In fact, it has been shown that for data dimension $d = \omega(\log n)$, no algorithm can approximately solve kernel SVM within an error $\exp(-\omega(\log^2 n))$ in time $O(n^{2-o(1)})$, assuming the famous Strong Exponential Time Hypothesis (SETH)³ (Backurs et al., 2017).

Conversely, a long line of works aim to speed up computation with the kernel matrix faster than quadratic, especially when the kernel has certain smooth and Lipschitz properties (Alman et al., 2020; Aggarwal & Alman, 2022; Bakshi et al., 2023; Charikar et al., 2024). For instance, when kernel functions are sufficiently smooth, efficient approximation using low-degree polynomials is feasible, leading to an approximate low-rank factorization. A prime example is the Gaussian RBF kernel, where Aggarwal & Alman (2022) showed that for dimension $d = \Theta(\log n)$ and squared dataset radius (defined as $\max_{i,j \in [n]} \|x_i - x_j\|_2^2$) $B = o(\log n)$, there exists low-rank matrices $U, V \in \mathbb{R}^{n \times n^{o(1)}}$ such that for any vector $x \in \mathbb{R}^n$, $\|(K - UV^T)x\|_\infty \leq \epsilon \|x\|_1$. They subsequently develop an algorithm to solve the Batch Gaussian KDE problem in $O(n^{1+o(1)})$ time.

Based on this dichotomy in fast kernel matrix algebra, we establish two results: 1) Solving Gaussian kernel SVM in $O(n^{1+o(1)} \log(1/\epsilon))$ time is feasible when $B = o(\frac{\log n}{\log \log n})$, and 2) Assuming SETH, no sub-quadratic time algorithm exists for $B = \Omega(\log^2 n)$ in SVMs without bias and $B = \Omega(\log^6 n)$ in SVMs with bias. This improves the lower bound established by Backurs et al. (2017) in terms of dimension d .

Theorem 1.5 (Gaussian Kernel SVM, informal version of Theorem G.7 and G.12). *Given a dataset $X \in \mathbb{R}^{n \times d}$ with dimension d and squared radius denoted by B , let $K(x_i, x_j) = \exp(-\|x_i - x_j\|_2^2)$ be the Gaussian kernel function. Then, for the kernel SVM problem defined in Definition 1.3,*

- *If $d = O(\log n)$, $B = o(\frac{\log n}{\log \log n})$, there exists an algorithm that solves the Gaussian kernel SVM up to ϵ -error in time $O(n^{1+o(1)} \log(1/\epsilon))$;*
- *If $d = \Omega(\log n)$, $B = \Omega(\log^2 n)$, then assuming SETH, any algorithm that solves the Gaussian kernel SVM without a bias term up to $\exp(-\omega(\log^2 n))$ error would require $\Omega(n^{2-o(1)})$ time;*
- *If $d = \Omega(\log n)$, $B = \Omega(\log^6 n)$, then assuming SETH, any algorithm that solves the Gaussian kernel SVM with a bias term up to $\exp(-\omega(\log^2 n))$ error would require $\Omega(n^{2-o(1)})$ time.*

To our knowledge, this is the first almost-linear time algorithm for Gaussian kernel SVM even when $d = \log n$ and the radius is small. Our algorithm effectively utilizes the rank- $n^{o(1)}$ factorization of the Gaussian kernel matrix alongside our low-rank QP solver.

1.1 RELATED WORK

Support Vector Machines. SVM, one of the most prominent machine learning models before the rise of deep learning, has a rich literature dedicated to its algorithmic speedup. For linear SVM, Joachims (2006) offers a first-order algorithm that solves its QP in nearly-linear time, but with a runtime dependence of ϵ^{-2} , limiting its use in high precision settings. **This runtime is later significantly improved by Shalev-Shwartz et al. (2011) to $\tilde{O}(\epsilon^{-1}d)$ via a stochastic subgradient descent algorithm.** For SVM classification, existing algorithms such as SVM-Light (Joachims, 1999), SMO (Platt, 1998), LIBSVM (Chang & Lin, 2011), and SVM-Torch (Collobert & Bengio, 2001) perform well in high-dimensional data settings. However, their runtime scales super-linearly with n , making them less viable for large datasets. Previous investigations into solving linear SVM via interior point methods (Ferris & Munson, 2002) have been somewhat basic, leading to an overall runtime of $O(n^{1.5}d^{\omega-1} \log(1/\epsilon))$. For a more comprehensive survey on efficient algorithms for SVM, refer to Cervantes et al. (2020). On the hardness side, Backurs et al. (2017) provides an efficient reduction from the Bichromatic Closest Pair problem to Gaussian kernel SVM, establishing an almost-quadratic lower bound assuming SETH.

³SETH is a standard complexity theoretical assumption (Impagliazzo et al., 1998; Impagliazzo & Paturi, 2001). Informally, it states that for a Conjunctive Normal Form (CNF) formula with m clauses and n variables, there is no algorithm for checking its feasibility in time less than $O(c^n \cdot \text{poly}(m))$ for $c < 2$.

Interior Point Method. The interior point method, a well-established approach for solving convex programs under constraints, was first proposed by Karmarkar (1984) as a (weakly) polynomial-time algorithm for linear programs, later improved by Vaidya (1989) in terms of runtime. Recent work by Cohen et al. (2019) has shown how to solve linear programs with interior point methods in the current matrix multiplication time, utilizing a robust IPM framework. Subsequent studies (Lee et al., 2019; Brand, 2020; Jiang et al., 2021; Song & Yu, 2021; Huang et al., 2022; Jiang et al., 2022; Qin et al., 2023) have further refined their algorithm or applied it to different optimization problems.

Kernel Matrix Algebra. Kernel methods, fundamental in machine learning, enable feature mappings to potentially infinite dimensions for n data points in d dimensions. The kernel matrix, a crucial component of kernel methods, often has a prohibitive quadratic size for explicit formation. Recent active research focuses on computing and approximating kernel matrices and related tasks in sub-quadratic time, such as kernel matrix-vector multiplication, spectral sparsification, and Laplacian system solving. The study by Alman et al. (2020) introduces a comprehensive toolkit for solving these problems in almost-linear time for small dimensions, leveraging techniques like polynomial methods and ultra Johnson-Lindenstrauss transforms. Alternatively, Backurs et al. (2021); Bakshi et al. (2023) reduce various kernel matrix algebra tasks to kernel density estimation (KDE), which recent advancements in KDE data structures (Charikar & Siminelakis, 2017; Backurs et al., 2018; Charikar et al., 2020) have made more efficient. A recent contribution by Aggarwal & Alman (2022) provides a tighter characterization of the low-degree polynomial approximation for the e^{-x} function, leading to more efficient algorithms for the Batch Gaussian KDE problem.

2 TECHNIQUE OVERVIEW

In this section, we provide an overview of the techniques employed in our development of two nearly-linear time algorithms for structured QPs. In Section 2.1, we detail the robust IPM framework, which forms the foundation of our algorithms. Subsequent section, namely Section 2.2 and 2.3, delves into dedicated data structures designed for efficiently solving low-treewidth and low-rank QPs, respectively. Finally, in Section 2.4, we discuss the adaptation of these advanced QP solvers for both linear and kernel SVMs.

Due to the heavily-technical nature, we recommend that in the first read, the audience can skip Section 2.2 and 2.3.

2.1 GENERAL STRATEGY

Our algorithm is built upon the robust IPM framework, an efficient variant of the primal-dual central path method (Renegar, 1988). This framework maintains a primal-dual solution pair $(x, s) \in \mathbb{R}^n \times \mathbb{R}^n$. To understand the central path for QPs, we first consider the central path equations for linear programming (see Cohen et al. (2019); Lee et al. (2019) for reference):

$$\begin{aligned} s/t + \nabla\phi(x) &= \mu, \\ Ax &= b, \\ A^\top y + s &= c, \end{aligned}$$

where x is the primal variable, s is the slack variable, y is the dual variable, $\phi(x)$ is a self-concordant barrier function, and μ denotes the error. The central path is defined by the trajectory of (x, s) as t approaches 0.

In quadratic programming, we modify these equations:

$$\begin{aligned} s/t + \nabla\phi(x) &= \mu, \\ Ax &= b, \\ -Qx + A^\top y + s &= c, \end{aligned}$$

where Q is the positive semidefinite objective matrix. The key difference in the central path equations for LP and QP is the inclusion of the $-Qx$ term in the third equation, significantly affecting algorithm design.

Fundamentally, IPM is a Newton’s method in which we update the variables x, y and s through the second-order information from the self-concordant barrier function. We derive the update rules for QP (detailed derivation in Section F.2):

$$\begin{aligned} \delta_x &= tM^{-1/2}(I - P)M^{-1/2}\delta_\mu, \\ \delta_y &= -t(AM^{-1}A^\top)^{-1}AM^{-1}\delta_\mu, \\ \delta_s &= t\delta_\mu - t^2HM^{-1/2}(I - P)M^{-1/2}\delta_\mu, \\ \text{where } H &= \nabla^2\phi(x), \quad M = Q + tH, \\ P &= M^{-1/2}A^\top(AM^{-1}A^\top)^{-1}AM^{-1/2}, \end{aligned}$$

where $\delta_x, \delta_y, \delta_s$, and δ_μ are the incremental steps for x, y, s , and μ , respectively.

The robust IPM approximates these updates rather than computing them exactly. It maintains an approximate primal-dual solution pair $(\bar{x}, \bar{s}) \in \mathbb{R}^n \times \mathbb{R}^n$ and computes the steps using this approximation. Provided the approximation is sufficiently accurate, it can be shown (see Section F for more details) that the algorithm converges efficiently to the optimal solution along the robust central path.

Therefore, the critical challenge lies in efficiently maintaining (\bar{x}, \bar{s}) , an approximation to (x, s) , when (x, s) evolves following the robust central path steps. The primary difficulty is that explicitly managing the primal-dual solution pair (x, s) is inefficient due to potential dense changes. Such changes can lead to dense updates in H , slowing down the computation of steps. The innovative aspect of robust IPM is recognizing that (x, s) are only required at the algorithm’s conclusion, not during its execution. Instead, we can identify entries with significant changes and update the approximation (\bar{x}, \bar{s}) correspondingly. With IPM’s lazy updates, only a nearly-linear number of entries are adjusted throughout the algorithm:

$$\sum_{t=1}^T \|\bar{x}^{(t)} - \bar{x}^{(t-1)}\|_0 + \|\bar{s}^{(t)} - \bar{s}^{(t-1)}\|_0 = \tilde{O}(n \log(1/\epsilon))$$

where $T = \tilde{O}(\sqrt{n} \log(1/\epsilon))$ is the number of iterations for IPM convergence. This indicates that, on average, each entry of \bar{x} and \bar{s} is updated $\log(1/\epsilon)$ times, facilitating rapid updates to these quantities and, consequently, to H .

In the special case where $Q = 0$, the path reverts to the LP case, with $M = tH$ being a diagonal matrix, allowing for efficient computation and updates of M^{-1} . This simplifies maintaining $AM^{-1}A^\top$, as updates to M^{-1} correspond to row and column scaling of A . However, in the QP scenario, where M is symmetric positive semidefinite, maintaining the term $AM^{-1}A^\top$ becomes more complex. Nevertheless, when the number of constraints is small, as in SVMs, this issue is less problematic. Yet, even with this simplification, the challenge is far from trivial, given the presence of terms like $M^{-1/2}$ in the robust central path steps. While the matrix Woodbury identity could be considered, it falls short when maintaining a square root term. Despite these hurdles, we construct efficient data structures for $M^{-1/2}$ maintenance when Q possesses succinct representations, such as low-rank.

Before diving into the particular techniques for low-rank QPs, we start by exploring the *low-treewidth QPs*, which could be viewed as a structured sparsity condition. It provides valuable insights for the low-rank scenario.

2.2 LOW-TREewidth SETTING: HOW TO LEVERAGE SPARSITY

Treewidth is parameter for graphs that captures the sparsity pattern. Given a graph $G = (V, E)$ with n vertices and m edges, a *tree decomposition* of G arranges its vertices into bags, which collectively form a tree structure. For any two bags X_i, X_j , if a vertex v is present in both, it must also be included in all bags along the path between X_i and X_j . Additionally, each pair of adjacent vertices in the graph must be present together in at least one bag. The treewidth τ is defined as the maximum size of a bag minus one. Intuitively, a graph G with a small treewidth τ implies a structure akin to a tree. For a formal definition, see Definition A.1. When relating this combinatorial structure to linear algebra, we could treat the quadratic objective matrix Q as a generalized adjacency matrix, where we put a vertex v_i on i -th row of Q , and put an edge $\{v_i, v_j\}$ whenever the entry $Q_{i,j}$ is

324 nonzero. The low-treewidth structure of the graph corresponds to a sparsity pattern that allows one to
 325 compute a column-sparse Cholesky factorization of Q . Since $M = Q + tH$ and H is diagonal, we
 326 can decompose $M = LL^\top$ into sparse Cholesky factors⁴.
 327

328 Under any coordinate update to \bar{x} , M is updated on only one diagonal entry, enabling efficient updates
 329 to L . The remaining task is to use this Cholesky decomposition to maintain the central path step.

330 By expanding the central path equations and substituting $M = LL^\top$, we derive

$$\begin{aligned}
 331 \delta_x &= tM^{-1/2}(I - P)M^{-1/2}\delta_\mu \\
 332 &= tM^{-1}\delta_\mu - tM^{-1}A^\top(AM^{-1}A^\top)^{-1}AM^{-1}\delta_\mu \\
 333 &= tL^{-\top}L^{-1}\delta_\mu - tL^{-\top}L^{-1}A^\top(AL^{-\top}L^{-1}A^\top)^{-1}AL^{-\top}L^{-1}\delta_\mu, \\
 334 \delta_s &= t\delta_\mu - t^2HM^{-1/2}(I - P)M^{-1/2}\delta_\mu \\
 335 &= t\delta_\mu - t^2L^{-\top}L^{-1}\delta_\mu + t^2L^{-\top}L^{-1}A^\top(AL^{-\top}L^{-1}A^\top)^{-1}AL^{-\top}L^{-1}\delta_\mu.
 \end{aligned}$$

336 Updates to the diagonal of M do not change L 's nonzero pattern, allowing for efficient utilization
 337 of the sparse factor and maintenance of $L^{-1}A^\top \in \mathbb{R}^{n \times m}$ and $L^{-1}\delta_\mu \in \mathbb{R}^n$. Terms like
 338 $(AL^{-\top}L^{-1}A^\top)^{-1}AL^{-\top}L^{-1}\delta_\mu \in \mathbb{R}^m$ can also be explicitly maintained.
 339

340 With this approach, we propose the following implicit representation for maintaining (x, s) :

$$341 x = \hat{x} + H^{-1/2}\mathcal{W}^\top(h\beta_x - \tilde{h}\tilde{\beta}_x + \epsilon_x), \quad (4)$$

$$342 s = \hat{s} + H^{1/2}c_s\beta_{c_s} - H^{1/2}\mathcal{W}^\top(h\beta_s - \tilde{h}\tilde{\beta}_s + \epsilon_s), \quad (5)$$

343 where $\hat{x}, \hat{s} \in \mathbb{R}^n$, $\mathcal{W} = L^{-1}H^{1/2} \in \mathbb{R}^{n \times n}$, $h = L^{-1}\delta_\mu \in \mathbb{R}^n$, $c_s = H^{-1/2}\delta_\mu \in \mathbb{R}^n$, $\beta_x, \beta_s, \beta_{c_s} \in$
 344 \mathbb{R} , $\tilde{h} = L^{-1}A^\top \in \mathbb{R}^{n \times m}$, $\tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$, $\epsilon_x, \epsilon_s \in \mathbb{R}^n$. All quantities except for \mathcal{W} can be explicitly
 345 maintained. For linear programming, the implicit representation is as follows:
 346

$$\begin{aligned}
 347 x &= \hat{x} + H^{-1/2}\beta_x c_x - H^{-1/2}\mathcal{W}^\top(\beta_x h + \epsilon_x) \\
 348 s &= \hat{s} + H^{1/2}\mathcal{W}^\top(\beta_s h + \epsilon_s),
 \end{aligned}$$

349 with $\mathcal{W} = L^{-1}AH^{-1/2}$ maintained implicitly and the other terms explicitly.
 350

351 The representation in (4) and (5) enables us to maintain the central path step using a combination of
 352 ‘‘coefficients’’ $h + \tilde{h}\tilde{\beta}_x$ and ‘‘basis’’ \mathcal{W}^\top . We need to detect entries of \bar{x} that deviate significantly from x
 353 and capture these changes with $\|H^{1/2}(\bar{x} - x)\|_2$. We maintain this vector using $x_0 + \mathcal{W}^\top(h + \tilde{h}\tilde{\beta}_x)$.
 354 Here, \mathcal{W}^\top acts as a wavelet basis and the vector $h + \tilde{h}\tilde{\beta}_x$ as its multiscale coefficients. While
 355 computing and maintaining $\mathcal{W}^\top h$ seems challenging, leveraging column-sparsity of L^{-1} is possible
 356 through contraction with a vector v :
 357

$$\begin{aligned}
 358 v^\top \mathcal{W}^\top &= (\mathcal{W}v)^\top \\
 359 &= (L^{-1}H^{1/2}v)^\top.
 \end{aligned}$$

360 By applying the Johnson-Lindenstrauss transform (JL) in place of v , we can quickly approximate
 361 $\|\mathcal{W}^\top h\|_2$ by maintaining $\Phi\mathcal{W}^\top$ for a JL matrix Φ . Similarly, we handle $\mathcal{W}^\top\tilde{h}\tilde{\beta}_x$ by explicitly
 362 computing $A^\top\tilde{\beta}_x$ and using the sparsity of L^{-1} for $\tilde{h}\tilde{\beta}_x$.
 363

364 We focus on entries significantly deviating from x_0 , the heavy entries of $\mathcal{W}^\top(h + \tilde{h}\tilde{\beta}_x)$. Here, the
 365 treewidth- τ decomposition enables quick computation of an elimination tree based on L^{-1} 's sparsity,
 366 facilitating efficient estimation of $\|(\mathcal{W}^\top(h + \tilde{h}\tilde{\beta}_x))_{\chi(v)}\|_2$ for any subtree $\chi(v)$ ⁵. With an elimination
 367 tree of height $\tilde{O}(\tau)$, we can employ heavy-light decomposition (Sleator & Tarjan, 1981) for an
 368 $O(\log n)$ -height tree.
 369

370 Using these data structures, convergence is established using the robust IPM framework (Ye, 2020;
 371 Lee & Vempala, 2021). While the framework is generally applicable to QPs, computing an initial
 372 point remains a challenge. We propose a simpler objective $x_0 = \arg \min_{x \in \mathbb{R}^n} \sum_{i=1}^n \phi_i(x_i)$ with ϕ_i
 373 as the log-barrier function, resembling the initial point reduction in Lee et al. (2019). This initial
 374 point enables us to solve an augmented quadratic program that increases dimension by 1.
 375

376 ⁴Note that adding a non-negative diagonal matrix to Q does not change its sparsity pattern, hence M also
 377 retains the treewidth τ .

⁵Given any tree node v , we use $\chi(v)$ to denote the subtree rooted at v .

2.3 LOW-RANK SETTING: HOW TO UTILIZE SMALL FACTORIZATION

The low-treewidth structure can be considered a form of sparsity, allowing for a sparse factorization $M = LL^\top$. Another significant structure arises when the matrix Q admits a low-rank factorization. Let $Q = UV^\top$ where $U, V \in \mathbb{R}^{n \times k}$ and $k \ll n$, then $M = Q + tH = UV^\top + tH$. Although Q has a low-rank structure, M may not be low-rank due to the diagonal matrix being dense. However, in the central path equations, we need only handle M^{-1} , which can be efficiently maintained using the matrix Woodbury identity:

$$M^{-1} = t^{-1}H^{-1} - t^{-2}H^{-1}U(I + t^{-1}V^\top H^{-1}U)^{-1}V^\top H^{-1},$$

Given that H is diagonal, the complex term $(I + t^{-1}V^\top H^{-1}U)^{-1}$ can be quickly updated under sparse changes to H^{-1} by simply scaling rows of U and V . With only a nearly-linear number of updates to H^{-1} , the total update time across $\tilde{O}(\sqrt{n} \log(1/\epsilon))$ iterations is bounded by $\tilde{O}(nk^{\omega-1} + k^\omega)$. We modify the (x, s) implicit representation as follows:

$$x = \hat{x} + H^{-1/2}h\beta_x + H^{-1/2}\hat{h}\hat{\beta}_x + H^{-1/2}\tilde{h}\tilde{\beta}_x, \quad (6)$$

$$s = \hat{s} + H^{1/2}h\beta_s + H^{1/2}\hat{h}\hat{\beta}_s + H^{1/2}\tilde{h}\tilde{\beta}_s, \quad (7)$$

where $\bar{x}, \bar{s} \in \mathbb{R}^n$, $h = H^{-1/2}\bar{\delta}_\mu \in \mathbb{R}^n$, $\hat{h} = H^{-1/2}U \in \mathbb{R}^{n \times k}$, and $\tilde{h} = H^{-1/2}A^\top \in \mathbb{R}^{n \times m}$, with $\tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$ and $\beta_x, \beta_s \in \mathbb{R}$. The nontrivial terms to maintain are \hat{h} and \tilde{h} , but both can be managed straightforwardly: updates to $H^{-1/2}$ correspond to scaling rows of U and A^\top , and can be performed in total $\tilde{O}(nk)$ and $\tilde{O}(nm)$ time, respectively. The key observation is that we *never explicitly form* $M^{-1/2}$, hence matrix Woodbury identity suffices for fast updates.

The remaining task is to design a data structure for detecting heavy entries. Instead of starting with an elimination tree and re-balancing it through heavy-light decomposition, we construct a balanced tree on n nodes, hierarchically dividing length- n vectors by their indices. Sampling is then performed by traversing down to the tree's leaves. While a heavy-hitter data structure could lead to improvements in poly-logarithmic and sub-logarithmic factors, we primarily focus on polynomial dependencies on various parameters and leave this enhancement for future exploration.

2.4 GAUSSIAN KERNEL SVM: ALGORITHM AND HARDNESS

Our specialized QP solvers provide fast implementations for linear SVMs when the data dimension d is much smaller than n . However, for kernel SVM, forming the kernel matrix exactly would take $\Theta(n^2)$ time. Fortunately, advancements in kernel matrix algebra (Alman et al., 2020; Bakshi et al., 2021; Aggarwal & Alman, 2022; Bakshi et al., 2023) have enabled sub-quadratic algorithms when the data dimension d is small or the kernel matrix has a relatively large minimum entry. Both Alman et al. (2020) and Bakshi et al. (2023) introduce algorithms for spectral sparsification, generating an approximate matrix $\tilde{K} \in \mathbb{R}^{n \times n}$ such that $(1 - \epsilon) \cdot K \preceq \tilde{K} \preceq (1 + \epsilon) \cdot K$, with \tilde{K} having only $O(\epsilon^{-2}n \log n)$ nonzero entries. Alman et al. (2020) achieves this in $O(n^{1+o(1)})$ time for multiplicatively Lipschitz kernels when $d = O(\log n)$, while Bakshi et al. (2023) overcomes limitations for Gaussian kernels by basing their algorithm on KDE and the magnitude of the minimum entry of the kernel matrix, parameterized by τ . Their algorithm for Gaussian kernels runs in time $\tilde{O}(nd/\tau^{3.173+o(1)})$. Unfortunately, spectral sparsifiers do not aid our primitives since a sparsifier only reduces the number of nonzero entries, but not the rank of the kernel matrix.

Besides spectral sparsification, Alman et al. (2020); Aggarwal & Alman (2022) also demonstrate that with $d = d = O(\log n)$ and suitable kernels, there exists an $O(n^{1+o(1)})$ time algorithm to multiply the kernel matrix with an arbitrary vector $v \in \mathbb{R}^n$. This operation is crucial in Batch KDE as shown in Aggarwal & Alman (2022). Moreover, Aggarwal & Alman (2022) establishes an almost-quadratic lower bound for this operation when the squared dataset radius $B = \omega(\log n)$, assuming SETH. These results rely on computing a rank- $n^{o(1)}$ factorization for the Gaussian kernel matrix. The function e^{-x} can be approximated by a low-degree polynomial of degree

$$q := \Theta(\max\{\sqrt{B \log(1/\epsilon)}, \frac{\log(1/\epsilon)}{\log(\log(1/\epsilon)/B)}\})$$

for $x \in [0, B]$. Using this polynomial, one can create matrices U, V with $\text{rank} \binom{2d+2q}{2q} = n^{o(1)}$ in time $O(n^{1+o(1)})$. Given this factorization, multiplying it with a vector v as $U(V^\top v)$ takes $O(n^{1+o(1)})$ time. Let $\tilde{K} = UV^\top$ where $\tilde{K}_{i,j} = f(\|x_i - x_j\|_2^2)$, we have for any $(i, j) \in [n] \times [n]$,

$$|f(\|x_i - x_j\|_2^2) - \exp(-\|x_i - x_j\|_2^2)| \leq \epsilon,$$

and for any row $i \in [n]$,

$$\begin{aligned} |(\tilde{K}v)_i - (Kv)_i| &= \left| \sum_{j=1}^n v_j (f(\|x_i - x_j\|_2^2) - \exp(-\|x_i - x_j\|_2^2)) \right| \\ &\leq (\max_{j \in [n]} |f(\|x_i - x_j\|_2^2) - \exp(-\|x_i - x_j\|_2^2)|) \|v\|_1 \\ &\leq \epsilon \|v\|_1, \end{aligned}$$

using Hölder's inequality. This provides an ℓ_∞ -guarantee of the error vector $(\tilde{K} - K)v$, useful for Batch Gaussian KDE. Transforming this ℓ_∞ -guarantee into a spectral approximator yields

$$(1 - \epsilon n) \cdot K \preceq \tilde{K} \preceq (1 + \epsilon n) \cdot K.$$

Setting $\epsilon = 1/n^2$, the low-rank factorization offers an adequate spectral approximation to the exact kernel matrix K .

Given $\tilde{K} = UV^\top$ for $U, V \in \mathbb{R}^{n \times n^{o(1)}}$, we can solve program (3) with \tilde{K} using our low-rank QP algorithm in time $O(n^{1+o(1)} \log(1/\epsilon))$.⁶ This is the first almost-linear time algorithm for Gaussian kernel SVM, even in low-precision settings, as prior works either lack machinery to approximately form the kernel matrix efficiently, or do not possess faster convex optimization solvers for solving a structured quadratic program associated with a kernel SVM.

The requirements $d = O(\log n)$ and $B = o(\frac{\log n}{\log \log n})$ may seem restrictive, but they are necessary, as no sub-quadratic time algorithm exists for Gaussian kernel SVM without bias when $d = \Omega(\log n)$ and $B = \Omega(\log^2 n)$, and with bias when $B = \Omega(\log^6 n)$, assuming SETH. This is based on a reduction from Bichromatic Closet Pair to Gaussian kernel SVM, as established by Backurs et al. (2017). Our assumptions on d and B are therefore justified for seeking almost-linear time algorithms.

We note that in other variants of definitions for Gaussian kernels, one requires an additional parameter called the *kernel width*, and the kernel function is defined as $\exp(-\frac{\|x_i - x_j\|_2^2}{2\sigma^2})$. In commonly used heuristics (Ramdas et al., 2015), $\sigma = O(\sqrt{d})$, hence we could without loss of generality assuming $\sigma = 1$ by requiring the squared radius to be B/d .

3 CONCLUSION

On the algorithmic front, we introduce the first nearly-linear time algorithms for low-rank convex quadratic programming, leading to nearly-linear time algorithms for linear SVMs. For Gaussian kernel SVMs, we utilize a low-rank approximation from Aggarwal & Alman (2022) when $d = O(\log n)$ and the squared dataset radius is small, enabling an almost-linear time algorithm. On the hardness aspect, we establish that when $d = \Omega(\log n)$, if the squared dataset radius is sufficiently large ($\Omega(\log^2 n)$ without bias and $\Omega(\log^6 n)$ with bias), then assuming SETH, no sub-quadratic algorithm exists. As our work is theoretical in nature, we do not foresee any potential negative societal impact. Several open problems arise from our work:

Better dependence on k for low-rank QPs. Our low-rank QP solver exhibits a dependence of $k^{(\omega+1)/2}$ on the rank k . Given the precomputed factorization, can we improve the exponents on k ? Ideally, an algorithm with nearly-linear dependence on k would align more closely with input size.

⁶Additional requirement: $B = o(\frac{\log n}{\log \log n})$. See Section G for further discussion.

486 **Better dependence on m for general QPs.** Focusing on SVMs with a few equality constraints,
487 our QP solvers do not exhibit strong dependence on the number of equality constraints m . Without
488 structural assumptions on the constraint matrix A , this is expected. However, many QPs, particularly
489 in graph contexts, involve large m . Is there a pathway to an algorithm with better dependence on m ?
490 More broadly, can we achieve a result akin to that of Lee & Sidford (2019), where the number of
491 iterations depends on the square root of the rank of A , with minimal per iteration cost?
492

493 **Stronger lower bound in terms of B for Gaussian kernel SVMs.** We establish hardness results
494 for Gaussian kernel SVM when $B = \Omega(\log^2 n)$ without bias and $B = \Omega(\log^6 n)$ with bias. This
495 contrasts with our algorithm, which requires B to have sub-logarithmic dependence on n . For
496 Batch Gaussian KDE, Aggarwal & Alman (2022) demonstrated that fast algorithms are feasible for
497 $B = o(\log n)$, with no sub-quadratic time algorithms for $B = \omega(\log n)$ assuming SETH. Can a
498 stronger lower bound be shown for SVM programs with a bias term, reflecting a more natural setting?
499

500 REFERENCES

- 501 Amol Aggarwal and Josh Alman. Optimal-degree polynomial approximations for exponentials
502 and gaussian kernel density estimation. In *Proceedings of the 37th Computational Complexity
503 Conference, CCC '22*, Dagstuhl, DEU, 2022. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik.
504
- 505 Josh Alman, Timothy Chu, Aaron Schild, and Zhao Song. Algorithms and hardness for linear algebra
506 on geometric graphs. In *2020 IEEE 61st Annual Symposium on Foundations of Computer Science
507 (FOCS)*, 2020.
508
- 509 Arturs Backurs, Piotr Indyk, and Ludwig Schmidt. On the fine-grained complexity of empirical risk
510 minimization: Kernel methods and neural networks. *Advances in Neural Information Processing
511 Systems*, 30, 2017.
- 512 Arturs Backurs, Moses Charikar, Piotr Indyk, and Paris Siminelakis. Efficient density evaluation
513 for smooth kernels. In *2018 IEEE 59th Annual Symposium on Foundations of Computer Science
514 (FOCS)*, 2018.
515
- 516 Arturs Backurs, Piotr Indyk, Cameron Musco, and Tal Wagner. Faster kernel matrix algebra via
517 density estimation. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International
518 Conference on Machine Learning*, Proceedings of Machine Learning Research. PMLR, 2021.
519
- 520 Ainesh Bakshi, Piotr Indyk, Praneeth Kacham, Sandeep Silwal, and Samson Zhou. Sub-quadratic
521 algorithms for kernel matrices via kernel density estimation. In *International Conference on
522 Learning Representation, ICLR'23*, 2023.
- 523 Aaron Bernstein, Maximilian Probst Gutenberg, and Thatchaphol Saranurak. Deterministic decre-
524 mental sssp and approximate min-cost flow in almost-linear time. In *2021 IEEE 62nd Annual
525 Symposium on Foundations of Computer Science (FOCS)*, pp. 1000–1008. IEEE, 2022.
526
- 527 Bernhard E. Boser, Isabelle M. Guyon, and Vladimir N. Vapnik. A training algorithm for optimal
528 margin classifiers. In *Proceedings of the Fifth Annual Workshop on Computational Learning
529 Theory, COLT '92*, 1992.
530
- 531 Jan van den Brand. A deterministic linear program solver in current matrix multiplication time. In
532 *Proceedings of the Thirty-First Annual ACM-SIAM Symposium on Discrete Algorithms, SODA
533 '20*, 2020.
- 534 Jan van den Brand, Yin Tat Lee, Aaron Sidford, and Zhao Song. Solving tall dense linear programs
535 in nearly linear time. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of
536 Computing*, pp. 775–788, 2020.
537
- 538 Jair Cervantes, Farid Garcia Lamont, Lisbeth Rodriguez Mazahua, and Asdrubal Lopez. A com-
539 prehensive survey on support vector machine classification: Applications, challenges and trends.
Neurocomputing, 408:189–215, 2020.

-
- 540 Chih-Chung Chang and Chih-Jen Lin. Training ν -support vector classifiers: Theory and algorithms.
541 *Neural Comput.*, sep 2001.
- 542
- 543 Chih-Chung Chang and Chih-Jen Lin. Libsvm: A library for support vector machines. *ACM Trans.*
544 *Intell. Syst. Technol.*, may 2011.
- 545
- 546 Moses Charikar and Paris Siminelakis. Hashing-based-estimators for kernel density in high dimen-
547 sions. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*,
548 2017.
- 549 Moses Charikar, Michael Kapralov, Navid Nouri, and Paris Siminelakis. Kernel density estimation
550 through density constrained near neighbor search. In *2020 IEEE 61st Annual Symposium on*
551 *Foundations of Computer Science (FOCS)*. IEEE Computer Society, 2020.
- 552
- 553 Moses Charikar, Michael Kapralov, and Erik Waingarten. *A Quasi-Monte Carlo Data Structure for*
554 *Smooth Kernel Evaluations*, pp. 5118–5144. 2024.
- 555
- 556 Michael B Cohen, Yin Tat Lee, and Zhao Song. Solving linear programs in the current matrix
557 multiplication time. *STOC*, 2019.
- 558
- 559 Ronan Collobert and Samy Bengio. Svmtorch: Support vector machines for large-scale regression
560 problems. *Journal of machine learning research*, 1(Feb):143–160, 2001.
- 561
- 562 Gerard Cornuejols and Reha Tütüncü. *Optimization methods in finance*, volume 5. Cambridge
563 University Press, 2006.
- 564
- 565 Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Mach. Learn.*, 1995.
- 566
- 567 George B Dantzig. Linear programming under uncertainty. *Management science*, 1(3-4):197–206,
568 1955.
- 569
- 570 Timothy A Davis and William W Hager. Modifying a sparse cholesky factorization. *SIAM Journal*
571 *on Matrix Analysis and Applications*, 20(3):606–627, 1999.
- 572
- 573 Frédéric Delbos and Jean Charles Gilbert. *Global linear convergence of an augmented Lagrangian*
574 *algorithm for solving convex quadratic optimization problems*. PhD thesis, INRIA, 2003.
- 575
- 576 James R Driscoll, Neil Sarnak, Daniel D Sleator, and Robert E Tarjan. Making data structures
577 persistent. *Journal of computer and system sciences*, 38(1):86–124, 1989.
- 578
- 579 Ran Duan, Hongxun Wu, and Renfei Zhou. Faster matrix multiplication via asymmetric hashing. In
580 *FOCS*, 2023.
- 581
- 582 Michael C Ferris and Todd S Munson. Interior-point methods for massive support vector machines.
583 *SIAM Journal on Optimization*, 13(3):783–804, 2002.
- 584
- 585 Francois Le Gall. Faster rectangular matrix multiplication by combination loss analysis. In *Pro-*
586 *ceedings of the Thirty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA'24*,
587 2024.
- 588
- 589 Nicholas IM Gould and Philippe L Toint. A quadratic programming bibliography. *Numerical Analysis*
590 *Group Internal Report*, 1:32, 2000.
- 591
- 592 Nicholas IM Gould, Mary E Hribar, and Jorge Nocedal. On the solution of equality constrained
593 quadratic programming problems arising in optimization. *SIAM Journal on Scientific Computing*,
23(4):1376–1395, 2001.
- 594
- 595 Yuzhou Gu and Zhao Song. A faster small treewidth sdp solver. *arXiv preprint arXiv:2211.06033*,
2022.
- 596
- 597 Baihe Huang, Shunhua Jiang, Zhao Song, Runzhou Tao, and Ruizhe Zhang. Solving sdp faster: A
598 robust ipm framework and efficient implementation. In *FOCS*, 2022.

-
- 594 R. Impagliazzo, R. Paturi, and F. Zane. Which problems have strongly exponential complexity? In
595 *Proceedings 39th Annual Symposium on Foundations of Computer Science*, 1998.
596
- 597 Russell Impagliazzo and Ramamohan Paturi. On the complexity of k-sat. *J. Comput. Syst. Sci.*, 62:
598 367–375, mar 2001.
- 599 Haotian Jiang, Tarun Kathuria, Yin Tat Lee, Swati Padmanabhan, and Zhao Song. A faster interior
600 point method for semidefinite programming. In *2020 IEEE 61st annual symposium on foundations*
601 *of computer science (FOCS)*, pp. 910–918. IEEE, 2020.
602
- 603 Shunhua Jiang, Zhao Song, Omri Weinstein, and Hengjie Zhang. A faster algorithm for solving
604 general lps. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*,
605 pp. 823–832, 2021.
- 606 Shunhua Jiang, Bento Natura, and Omri Weinstein. A faster interior-point method for sum-of-
607 squares optimization. In *49th EATCS International Conference on Automata, Languages, and*
608 *Programming, LIPIcs*. Leibniz Int. Proc. Inform., 2022.
609
- 610 Thorsten Joachims. Making large-scale support vector machine learning practical. In *Advances in*
611 *kernel methods: support vector learning*, pp. 169. MIT press, 1999.
- 612 Thorsten Joachims. Training linear svms in linear time. In *Proceedings of the 12th ACM SIGKDD*
613 *international conference on Knowledge discovery and data mining*, pp. 217–226, 2006.
614
- 615 William B Johnson and Joram Lindenstrauss. Extensions of lipschitz mappings into a hilbert space.
616 *Contemporary mathematics*, 26(189-206):1, 1984.
- 617 Narendra Karmarkar. A new polynomial-time algorithm for linear programming. In *Proceedings of*
618 *the sixteenth annual ACM symposium on Theory of computing*, pp. 302–311, 1984.
619
- 620 Mikhail K Kozlov, Sergei Pavlovich Tarasov, and Leonid Genrikhovich Khachiyan. Polynomial
621 solvability of convex quadratic programming. In *Doklady Akademii Nauk*, pp. 1049–1051. Russian
622 Academy of Sciences, 1979.
- 623 Yin Tat Lee and Aaron Sidford. Solving linear programs with sqrt(rank) linear system solves. *arXiv*
624 *preprint arXiv:1910.08033*, 2019.
625
- 626 Yin Tat Lee and Santosh S. Vempala. Tutorial on the robust interior point method, 2021.
627
- 628 Yin Tat Lee, Zhao Song, and Qiuyi Zhang. Solving empirical risk minimization in the current matrix
629 multiplication time. In *Conference on Learning Theory*, pp. 2140–2157. PMLR, 2019.
- 630 K.-R. Muller, S. Mika, G. Ratsch, K. Tsuda, and B. Scholkopf. An introduction to kernel-based
631 learning algorithms. *IEEE Transactions on Neural Networks*, 12, 2001.
632
- 633 Katta G Murty. *Linear complementarity, linear and nonlinear programming*, volume 3. Citeseer,
634 1988.
- 635 Yurii Nesterov. Introductory lectures on convex programming volume i: Basic course. *Lecture notes*,
636 3(4):5, 1998.
637
- 638 Panos M Pardalos and Stephen A Vavasis. Quadratic programming with one negative eigenvalue is
639 np-hard. *Journal of Global optimization*, 1(1):15–22, 1991.
640
- 641 John Platt. Sequential minimal optimization: A fast algorithm for training support vector machines.
642 *MSR-TR-98-14*, 1998.
- 643 Marco Propato and James G Uber. Booster system design using mixed-integer quadratic programming.
644 *Journal of Water Resources Planning and Management*, 130(4):348–352, 2004.
645
- 646 Lianke Qin, Zhao Song, Lichen Zhang, and Danyang Zhuo. An online and unified algorithm
647 for projection matrix vector multiplication with application to empirical risk minimization. In
AISTATS, 2023.

648 Aaditya Ramdas, Sashank J. Reddi, Barnabás Póczos, Aarti Singh, and Larry Wasserman. On the
649 decreasing power of kernel and distance based nonparametric hypothesis tests in high dimensions.
650 In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, AAAI'15, pp.
651 3571–3577. AAAI Press, 2015. ISBN 0262511290.

652 James Renegar. A polynomial-time algorithm, based on newton’s method, for linear programming.
653 *Math. Program.*, 40(1–3):59–93, jan 1988.

654
655 Aviad Rubinfeld. Hardness of approximate nearest neighbor search. In *Proceedings of the 50th*
656 *Annual ACM SIGACT Symposium on Theory of Computing*, STOC 2018, pp. 1260–1268, New
657 York, NY, USA, 2018. Association for Computing Machinery.

658
659 Sartaj Sahni. Computationally related problems. *SIAM Journal on computing*, 3(4):262–279, 1974.

660
661 Robert Schreiber. A new implementation of sparse gaussian elimination. *ACM Transactions on*
662 *Mathematical Software (TOMS)*, 8(3):256–276, 1982.

663
664 Shai Shalev-Shwartz, Yoram Singer, Nathan Srebro, and Andrew Cotter. Pegasos: primal estimated
665 sub-gradient solver for svm. *Math. Program.*, 127, 2011.

666
667 Daniel D Sleator and Robert Endre Tarjan. A data structure for dynamic trees. In *Proceedings of the*
668 *thirteenth annual ACM symposium on Theory of computing*, pp. 114–122, 1981.

669
670 Zhao Song and Zheng Yu. Oblivious sketching-based central path method for linear programming.
671 In *International Conference on Machine Learning*, pp. 9835–9847. PMLR, 2021.

672
673 Pravin M Vaidya. Speeding-up linear programming using fast matrix multiplication. In *30th Annual*
674 *Symposium on Foundations of Computer Science*, pp. 332–337. IEEE, 1989.

675
676 Virginia Vassilevska Williams, Yinzhan Xu, Zixuan Xu, and Renfei Zhou. New bounds for matrix
677 multiplication: from alpha to omega. In *Proceedings of the Thirty-Fifth Annual ACM-SIAM*
678 *Symposium on Discrete Algorithms*, SODA’24, 2024.

679
680 Philip Wolfe. The simplex method for quadratic programming. *Econometrica: Journal of the*
681 *Econometric Society*, pp. 382–398, 1959.

682
683 Stephen J Wright. Continuous optimization (nonlinear and linear programming). *Foundations of*
684 *Computer-Aided Process Design*, 1999.

685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
Guanghao Ye. Fast algorithm for solving structured convex programs. *The University of Washington*,
Master Thesis, 2020.

Roadmap. In Section A, we present some basic definitions and tools that will be used in the remainder of the paper. In Section B, we introduce a few more SVM formulations, including classification and distribution estimation. In Section C, we show convex quadratic programming can be reduced to convex empirical risk minimization, and therefore can be solved in the current matrix multiplication time owing to Lee et al. (2019). In Section D and E, we prove results on low-treewidth and low-rank QPs, respectively. In Section F, we present a robust IPM framework for QPs, generalize beyond LPs and convex ERM with linear objective. In Section G, we present our algorithms for Gaussian kernel SVMs, with complementary lower bound.

A PRELIMINARY

A.1 NOTATIONS

For a positive integer n , we use $[n]$ to denote the set $\{1, 2, \dots, n\}$. For a matrix A , we use A^\top to denote its transpose. For a matrix A , we define $\|A\|_{p \rightarrow q} := \sup_x \|Ax\|_q / \|x\|_p$. When $p = q = 2$, we recover the spectral norm.

We define the entrywise ℓ_p -norm of a matrix A as $\|A\|_p := (\sum_{i,j} |A_{i,j}|^p)^{1/p}$.

For any function $f : \mathbb{N} \rightarrow \mathbb{N}$ and $n \in \mathbb{N}$, we use $\tilde{O}(f(n))$ to denote $O(f(n) \text{ poly log } f(n))$. We use $\mathbb{1}\{E\}$ to denote the indicator for event E , i.e., if E happens, $\mathbb{1}\{E\} = 1$ and otherwise it's 0.

A.2 TREEWIDTH

Treewidth captures the sparsity and tree-like structures of a graph.

Definition A.1 (Tree Decomposition and Treewidth). *Let $G = (V, E)$ be a graph, a tree decomposition of G is a tree T with b vertices, and b sets $J_1, \dots, J_b \subseteq V$ (called bags), satisfying the following properties:*

- For every edge $(u, v) \in E$, there exists $j \in [b]$ such that $u, v \in J_j$;
- For every vertex $v \in V$, $\{j \in [b] : v \in J_j\}$ is a non-empty subtree of T .

The treewidth of G is defined as the minimum value of $\max\{|J_j| : j \in [b]\} - 1$ over all tree decompositions.

A near-optimal tree decomposition of a graph can be computed in almost linear time.

Theorem A.2 (Bernstein et al. (2022)). *Given a graph G , there is an $O(m^{1+o(1)})$ time algorithm that produces a tree decomposition of G of maximum bag size $O(\tau \log^3 n)$, where τ is the actual (unknown) treewidth of G .*

Therefore, when $\tau = m^{\Theta(1)}$, we can compute an $\tilde{O}(\tau)$ -size tree decomposition in time $O(m\tau^{o(1)})$, which is negligible in the final running time of Theorem D.1.

A.3 SPARSE CHOLESKY DECOMPOSITION

In this section we state a few results on sparse Cholesky decomposition. Fast sparse Cholesky decomposition algorithms are based on the concept of elimination tree, introduced in Schreiber (1982).

Definition A.3 (Elimination tree). *Let G be an undirected graph on n vertices. An elimination tree \mathcal{T} is a rooted tree on $V(G)$ together with an ordering π of $V(G)$ such that for any vertex v , its parent is the smallest (under π) element u such that there exists a path P from v to u , such that $\pi(w) \leq \pi(v)$ for all $w \in P - u$.*

The following lemma relates the elimination tree and the structure of Cholesky factors.

Lemma A.4 (Schreiber (1982)). *Let M be a PSD matrix and \mathcal{T} be an elimination tree of the adjacency graph of M (i.e., $(i, j) \in E(G)$ iff $M_{i,j} \neq 0$) together with an elimination ordering π .*

Let P be the permutation matrix $P_{i,v} = \mathbb{1}\{v = \pi(i)\}$. Then the Cholesky factor L of PMP^\top (i.e., $PMP^\top = LL^\top$) satisfies $L_{i,j} \neq 0$ only if $\pi(i)$ is an ancestor of $\pi(j)$.

The following result is the current best result for computing a sparse Cholesky decomposition.

Lemma A.5 ((Gu & Song, 2022, Lemma 8.4)). *Let $M \in \mathbb{R}^{n \times n}$ be a PSD matrix whose adjacency graph has treewidth τ . Then we can compute the Cholesky factorization $M = LL^\top$ in $\tilde{O}(n\tau^{\omega-1})$ time.*

The following result is the current best result for updating a sparse Cholesky decomposition.

Lemma A.6 (Davis & Hager (1999)). *Let $M \in \mathbb{R}^{n \times n}$ be a PSD matrix whose adjacency graph has treewidth τ . Assume that we are given the Cholesky factorization $M = LL^\top$. Let $w \in \mathbb{R}^n$ be a vector such that $M + ww^\top$ has the same adjacency graph as M . Then we can compute $\Delta_L \in \mathbb{R}^{n \times n}$ such that $L + \Delta_L$ is the Cholesky factor of $M + ww^\top$ in $O(\tau^2)$ time.*

Throughout our algorithm, we need to compute matrix-vector multiplications involving Cholesky factors. We use the following results from Gu & Song (2022).

Lemma A.7 ((Gu & Song, 2022, Lemma 4.7)). *Let $M \in \mathbb{R}^{n \times n}$ be a PSD matrix whose adjacency graph has treewidth τ . Assume that we are given the Cholesky factorization $M = LL^\top$. Then we have the following running time for matrix-vector multiplications.*

- (i) For $v \in \mathbb{R}^n$, computing Lv , $L^\top v$, $L^{-1}v$, $L^{-\top}v$ takes $O(n\tau)$ time.
- (ii) For $v \in \mathbb{R}^n$, computing Lv takes $O(\|v\|_0\tau)$ time.
- (iii) For $v \in \mathbb{R}^n$, computing $L^{-1}v$ takes $O(\|L^{-1}v\|_0\tau)$ time.
- (iv) For $v \in \mathbb{R}^n$, if v is supported on a path in the elimination tree, then computing $L^{-1}v$ takes $O(\tau^2)$ time.
- (v) For $v \in \mathbb{R}^n$, computing $\mathcal{W}^\top v$ takes $O(n\tau)$ time, where $\mathcal{W} = L^{-1}H^{1/2}$ with $H \in \mathbb{R}^{n \times n}$ is a non-negative diagonal matrix.

Lemma A.8 ((Gu & Song, 2022, Lemma 4.8)). *Let $M \in \mathbb{R}^{n \times n}$ be a PSD matrix whose adjacency graph has treewidth τ . Assume that we are given the Cholesky factorization $M = LL^\top$. Then we have the following running time for matrix-vector multiplications, when we only need result for a subset of coordinates.*

- (i) Let S be a path in the elimination tree whose one endpoint is the root. For $v \in \mathbb{R}^n$, computing $(L^{-\top}v)_S$ takes $O(\tau^2)$ time.
- (ii) For $v \in \mathbb{R}^n$, for $i \in [n]$, computing $(\mathcal{W}^\top v)_i$ takes $O(\tau^2)$ time, where $\mathcal{W} = L^{-1}H^{1/2}$ with $H \in \mathbb{R}^{n \times n}$ be a non-negative diagonal matrix.

A.4 JOHNSON-LINDENTRAUSS LEMMA

We recall the Johnson-Lindenstrauss lemma, a powerful algorithmic primitive that reduces dimension while preserving ℓ_2 norms.

Lemma A.9 (Johnson & Lindenstrauss (1984)). *Let $\epsilon \in (0, 1)$ be the precision parameter. Let $\delta \in (0, 1)$ be the failure probability. Let $A \in \mathbb{R}^{m \times n}$ be a real matrix. Let $r = \epsilon^{-2} \log(mn/\delta)$. For $R \in \mathbb{R}^{r \times n}$ whose entries are i.i.d $\mathcal{N}(0, \frac{1}{r})$, the following holds with probability at least $1 - \delta$:*

$$(1 - \epsilon)\|a_i\|_2 \leq \|Ra_i\|_2 \leq (1 + \epsilon)\|a_i\|_2, \forall i \in [m],$$

where for a matrix A , a_i^\top denotes the i -th row of matrix $A \in \mathbb{R}^{m \times n}$.

A.5 HEAVY-LIGHT DECOMPOSITION

Heavy-light decomposition is useful when one wants to re-balance a binary tree with height $O(\log n)$.

Lemma A.10 (Sleator & Tarjan (1981)). *Given a rooted tree \mathcal{T} with n vertices, we can construct in $O(n)$ time an ordering π of the vertices such that (1) every path in \mathcal{T} can be decomposed into $O(\log n)$ contiguous subsequences under π , and (2) every subtree in \mathcal{T} is a single contiguous subsequence under π .*

B SVM FORMULATIONS

In this section, we review a list of formulations of SVM. These formulations have been implemented in the LIBSVM library [Chang & Lin \(2011\)](#).

Throughout this section, we use $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^s$ to denote the feature mapping, K to denote the associated kernel function and $K \in \mathbb{R}^{n \times n}$ to denote the kernel matrix. For linear SVM, ϕ is just the identity mapping. We will focus on the dual quadratic program formulation as usual. We will also assume for each problem, a dataset $X \in \mathbb{R}^{n \times d}$ is given together with binary labels $y \in \mathbb{R}^n$. Let $Q := (yy^\top) \circ K$.

B.1 C-SUPPORT VECTOR CLASSIFICATION

This formulation is also referred as the *soft-margin SVM*. It can be viewed as imposing a regularization on the primal program to allow mis-classification.

Definition B.1 (*C-Support Vector Classification*). *Given a parameter $C > 0$, the C-support vector classification (C-SVC) is defined as*

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^n} \quad & \mathbf{1}_n^\top \alpha - \frac{1}{2} \alpha^\top Q \alpha \\ \text{s.t.} \quad & \alpha^\top y = 0, \\ & 0 \leq \alpha \leq C \cdot \mathbf{1}_n. \end{aligned}$$

B.2 ν -SUPPORT VECTOR CLASSIFICATION

The C-SVC (Definition B.1) penalizes large values of α by limiting the magnitude of it. The ν -SVC (Definition B.2) turns $\mathbf{1}_n^\top \alpha$ from an objective into a constraint on ℓ_1 norm.

Definition B.2 (*ν -Support Vector Classification*). *Given a parameter $\nu > 0$, the ν -support vector classification (ν -SVC) is defined as*

$$\begin{aligned} \min_{\alpha \in \mathbb{R}^n} \quad & \frac{1}{2} \alpha^\top Q \alpha \\ \text{s.t.} \quad & \alpha^\top y = 0, \\ & \mathbf{1}_n^\top \alpha = \nu, \\ & 0 \leq \alpha \leq \frac{1}{n} \cdot \mathbf{1}_n. \end{aligned}$$

One can interpret this formulation as to find a vector that lives in the orthogonal complement of y that is non-negative, each entry is at most $\frac{1}{n}$ and its ℓ_1 norm is ν . Clearly, we must have $\nu \in (0, 1]$. More specifically, let k_+ be the number of positive labels and k_- be the number of negative labels. It is shown by [Chang & Lin \(2001\)](#) that the above problem is feasible if and only if

$$\nu \leq \frac{2 \min\{k_-, k_+\}}{n}.$$

B.3 DISTRIBUTION ESTIMATION

SVM is widely-used for predicting binary labels. It can also be used to estimate the support of a high-dimensional distribution. The formulation is similar to ν -SVC, except the PSD matrix Q is *label-less*.

Definition B.3 (*Distribution Estimation*). *Given a parameter $\nu > 0$, the ν -distribution estimation problem is defined as*

$$\begin{aligned} \min_{\alpha \in \mathbb{R}^n} \quad & \frac{1}{2} \alpha^\top K \alpha \\ \text{s.t.} \quad & 0 \leq \alpha \leq \frac{1}{n} \cdot \mathbf{1}_n, \\ & \mathbf{1}_n^\top \alpha = \nu. \end{aligned}$$

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

B.4 ϵ -SUPPORT VECTOR REGRESSION

In addition to classification, support vector framework can also be adapted for regression.

Definition B.4 (ϵ -Support Vector Regression). *Given parameters $\epsilon, C > 0$, the ϵ -support vector regression (ϵ -SVR) is defined as*

$$\begin{aligned} \min_{\alpha, \alpha^* \in \mathbb{R}^n} \quad & \frac{1}{2}(\alpha - \alpha^*)^\top K(\alpha - \alpha^*) + \epsilon \mathbf{1}_n^\top (\alpha + \alpha^*) + y^\top (\alpha - \alpha^*) \\ \text{s.t.} \quad & \mathbf{1}_n^\top (\alpha - \alpha^*) = 0, \\ & 0 \leq \alpha \leq C \cdot \mathbf{1}_n, \\ & 0 \leq \alpha^* \leq C \cdot \mathbf{1}_n. \end{aligned}$$

B.5 ν -SUPPORT VECTOR REGRESSION

One can similar adapt the parameter ν to control the ℓ_1 norm of the regression.

Definition B.5 (ν -Support Vector Regression). *Given parameters $\nu, C > 0$, the ν -support vector regression (ν -SVR) is defined as*

$$\begin{aligned} \min_{\alpha, \alpha^* \in \mathbb{R}^n} \quad & \frac{1}{2}(\alpha - \alpha^*)^\top K(\alpha - \alpha^*) + y^\top (\alpha - \alpha^*) \\ \text{s.t.} \quad & \mathbf{1}_n^\top (\alpha - \alpha^*) = 0, \\ & \mathbf{1}_n^\top (\alpha + \alpha^*) = C\nu, \\ & 0 \leq \alpha \leq \frac{C}{n} \cdot \mathbf{1}_n, \\ & 0 \leq \alpha^* \leq \frac{C}{n} \cdot \mathbf{1}_n. \end{aligned}$$

B.6 ONE EQUALITY CONSTRAINT

We classify C -SVC (Definition B.1), ϵ -SVR (Definition B.4) and ν -distribution estimation (Definition B.3) into the following generic form:

$$\begin{aligned} \min_{\alpha \in \mathbb{R}^n} \quad & \frac{1}{2} \alpha^\top Q \alpha + p^\top \alpha \\ \text{s.t.} \quad & \alpha^\top y = \Delta \\ & 0 \leq \alpha \leq C \cdot \mathbf{1}_n. \end{aligned}$$

Note that C -SVC (Definition B.1) and distribution estimation (Definition B.3) are readily in this form. For ϵ -SVR (Definition B.4), we need to perform a simple transformation:

Set $\hat{\alpha} = \begin{bmatrix} \alpha \\ \alpha^* \end{bmatrix} \in \mathbb{R}^{2n}$, then it can be written as

$$\begin{aligned} \min_{\hat{\alpha} \in \mathbb{R}^{2n}} \quad & \frac{1}{2} \hat{\alpha}^\top \begin{bmatrix} Q & -Q \\ -Q & Q \end{bmatrix} \hat{\alpha} + \begin{bmatrix} \epsilon \mathbf{1}_n + y \\ \epsilon \mathbf{1}_n - y \end{bmatrix}^\top \hat{\alpha} \\ \text{s.t.} \quad & \begin{bmatrix} \mathbf{1}_n \\ -\mathbf{1}_n \end{bmatrix}^\top \hat{\alpha} = 0 \\ & 0 \leq \hat{\alpha} \leq C \cdot \mathbf{1}_{2n}. \end{aligned}$$

B.7 TWO EQUALITY CONSTRAINTS

Both ν -SVC (Definition B.2) and ν -SVR (Definition B.5) require one extra constraint. They can be formulated as follows:

$$\min_{\alpha \in \mathbb{R}^n} \quad \frac{1}{2} \alpha^\top Q \alpha + p^\top \alpha$$

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

$$\begin{aligned} \text{s.t. } \mathbf{1}_n^\top \alpha &= \Delta_1, \\ y^\top \alpha &= \Delta_2, \\ 0 &\leq \alpha \leq C \cdot \mathbf{1}_n. \end{aligned}$$

For ν -SVR (Definition B.5), one can leverage a similar transformation as ϵ -SVR (Definition B.4).

Remark B.6. All variants of SVM-related quadratic programs can all be solved using our QP solvers for three cases:

- Linear SVM with $n \gg d$, we can solve it in $\tilde{O}(nd^{(\omega+1)/2} \log(1/\epsilon))$ time;
- Linear SVM with a small treewidth decomposition with width τ on XX^\top , we can solve it in $\tilde{O}(n\tau^{(\omega+1)/2} \log(1/\epsilon))$ time;
- Gaussian kernel SVM with $d = \Theta(\log n)$ and $B = o(\frac{\log n}{\log \log n})$, we can solve it in $O(n^{1+o(1)} \log(1/\epsilon))$ time.

Even though our solvers have relatively bad dependence on the number of equality constraints, for all these SVM formulations, at most 2 equality constraints are presented and thus can be solved very fast.

C ALGORITHMS FOR GENERAL QP

In this section, we discuss algorithms for general (convex) quadratic programming. We show that they can be solved in the current matrix multiplication time via reduction to linear programming with convex constraints Lee et al. (2019).

C.1 LCQP IN THE CURRENT MATRIX MULTIPLICATION TIME

Proposition C.1. There is an algorithm that solves LCQP (Definition 1.1) up to ϵ error in $\tilde{O}((n^\omega + n^{2.5-\alpha/2} + n^{2+1/6}) \log(1/\epsilon))$ time, where $\omega \leq 2.373$ is the matrix multiplication constant and $\alpha \geq 0.32$ is the dual matrix multiplication constant.

Proof. Let $Q = PDP^\top$ be an eigen-decomposition of Q where D is diagonal and P is orthogonal. Let $\tilde{x} := P^{-1}x$. Then it suffices to solve

$$\begin{aligned} \min \quad & \frac{1}{2} \tilde{x}^\top D \tilde{x} + c^\top P \tilde{x} \\ \text{s.t.} \quad & AP \tilde{x} = b \\ & P \tilde{x} \geq 0. \end{aligned}$$

By adding n non-negative variables and n constraints $x = P \tilde{x}$ we can make all constraints equality constraints. There are n non-negative variables and n unconstrained variables. If we want to ensure all variables are non-negative, we need to split every coordinate of \tilde{x} into two. In this way the coefficient matrix Q will be block diagonal with block size 2.

We perform the above reduction, and assume that we have a program of form (1) where Q is diagonal. Let $q_i := Q_{i,i}$ be the i -th element on the diagonal. Then the QP is equivalent to the following program

$$\begin{aligned} \min \quad & c^\top x + q^\top t \\ \text{s.t.} \quad & Ax = b \\ & t_i \geq \frac{1}{2} x_i^2 \quad \forall i \in [n] \\ & x \geq 0 \end{aligned}$$

Note that the set $\{(x_i, t_i) \in \mathbb{R}^2 : t_i \geq \frac{1}{2} x_i^2\}$ is a convex set. So we can apply Lee et al. (2019) here with n variables, each in the convex set $\{(a, b) \in \mathbb{R}^2 : a \geq 0, b \geq \frac{1}{2} a^2\}$. \square

972 C.2 ALGORITHM FOR QCQP

973
974 Our algorithm for LCQP in the previous section can be generalized to quadratically constrained
975 quadratic programs (QCQP). QCQP is defined as follows.

976 **Definition C.2 (QCQP).** Let $Q_0, \dots, Q_m \in \mathbb{R}^{n \times n}$ be PSD matrices. Let $q_0, \dots, q_m \in \mathbb{R}^n$. Let
977 $r \in \mathbb{R}^m$. Let $A \in \mathbb{R}^{d \times n}$, $b \in \mathbb{R}^d$. The quadratically constrained quadratic programming (QCQP)
978 problem asks to solve the following program.

$$979 \begin{aligned} 980 \min_{x \in \mathbb{R}^n} \quad & \frac{1}{2} x^\top Q_0 x + q_0^\top x \\ 981 \text{s.t.} \quad & \frac{1}{2} x^\top Q_i x + q_i^\top x + r_i \leq 0 \quad \forall i \in [m] \\ 982 & Ax = b \\ 983 & x \geq 0 \end{aligned}$$

984
985
986 **Proposition C.3.** There is an algorithm that solves QCQP (Definition C.2) up to ϵ error in
987 $\tilde{O}((mn)^\omega + (mn)^{2.5-\alpha/2} + (mn)^{2+1/6}) \log(1/\epsilon)$ time, where $\omega \leq 2.373$ is the matrix multi-
988 plication constant and $\alpha \geq 0.32$ is the dual matrix multiplication constant.
989

990 *Proof.* We first rewrite the program as following.

$$991 \begin{aligned} 992 \min \quad & -r_0 \\ 993 \text{s.t.} \quad & \frac{1}{2} x^\top Q_i x + q_i^\top x + r_i \leq 0 \quad \forall 0 \leq i \leq m \\ 994 & Ax = b \\ 995 & x \geq 0 \end{aligned}$$

996
997
998 Write $Q_i = P_i D_i P_i^\top$ be an eigen-decomposition of Q_i where D_i is diagonal and P_i is orthogonal.
999 Let $\tilde{x}_i \in \mathbb{R}^n$ be defined as $\tilde{x}_i := P_i^{-1} x$. Then we can rewrite the program as

$$1000 \begin{aligned} 1001 \min \quad & -r_0 \\ 1002 \text{s.t.} \quad & \frac{1}{2} \tilde{x}_i^\top D_i \tilde{x}_i + q_i^\top P_i \tilde{x}_i + r_i \leq 0 \quad \forall 0 \leq i \leq m \\ 1003 & Ax = b \\ 1004 & \tilde{x}_i = P_i^{-1} x \\ 1005 & x \geq 0 \end{aligned}$$

1006
1007
1008 For $0 \leq i \leq m$ and $j \in [n]$, define variable $t_{i,j} \in \mathbb{R}$. Then we can rewrite the program as

$$1009 \begin{aligned} 1010 \min \quad & -r_0 \\ 1011 \text{s.t.} \quad & \sum_{j \in [n]} D_{i,(j,j)} t_{i,j} + q_i^\top P_i \tilde{x}_i + r_i \leq 0 \quad \forall 0 \leq i \leq m \\ 1012 & Ax = b \\ 1013 & \tilde{x}_i = P_i^{-1} x \\ 1014 & t_{i,j} \geq \tilde{x}_{i,j}^2 \\ 1015 & x \geq 0 \end{aligned}$$

1016
1017
1018
1019 We can consider $(\tilde{x}_{i,j}, t_{i,j})_{0 \leq i \leq m, j \in [n]}$ as $(m+1)n$ variables in the convex set $\{(a, b) : b \geq \frac{1}{2}a^2\}$.
1020 Then we can apply Lee et al. (2019). \square
1021

1022 D ALGORITHM FOR LOW-TREEWIDTH QP

1023
1024 In this section we present a nearly-linear time algorithm for solving low-treewidth QP with small
1025 number of linear constraints. We briefly describe the outline of this section.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

- In Section D.1, we present the main statement of Section D.
- In Section D.2, we present the main data structure CENTRALPATHMAINTENANCE.
- In Section D.3, we present several data structures used in CENTRALPATHMAINTENANCE, including EXACTDS (Section D.3.1), APPROXDS (Section D.3.2), BATCHSKETCH (Section D.3.3), VECTORSKETCH (Section D.3.4), BALANCEDSKETCH (Section D.3.5).
- In Section D.4, we prove correctness and running time of CENTRALPATHMAINTENANCE data structure.
- In Section D.5, we prove the main result (Theorem D.1).

D.1 MAIN STATEMENT

We consider programs of the form (16), i.e.,

$$\begin{aligned} \min_{x \in \mathbb{R}^n} \quad & \frac{1}{2} x^\top Q x + c^\top x \\ \text{s.t.} \quad & A x = b \\ & x_i \in \mathcal{K}_i \quad \forall i \in [n] \end{aligned}$$

where $Q \in \mathcal{S}^{n_{\text{tot}}}$, $c \in \mathbb{R}^{n_{\text{tot}}}$, $A \in \mathbb{R}^{m \times n_{\text{tot}}}$, $b \in \mathbb{R}^m$, $\mathcal{K}_i \subset \mathbb{R}^{n_i}$ is a convex set. For simplicity, we assume that $n_i = O(1)$ for all $i \in [n]$.

Theorem D.1. *Consider the convex program (16). Let $\phi_i : \mathcal{K}_i \rightarrow \mathbb{R}$ be a ν_i -self-concordant barrier for all $i \in [n]$. Suppose the program satisfies the following properties:*

- *Inner radius r : There exists $z \in \mathbb{R}^{n_{\text{tot}}}$ such that $Az = b$ and $B(z, r) \in \mathcal{K}$.*
- *Outer radius R : $\mathcal{K} \subseteq B(0, R)$ where $0 \in \mathbb{R}^{n_{\text{tot}}}$.*
- *Lipschitz constant L : $\|Q\|_{2 \rightarrow 2} \leq L$, $\|c\|_2 \leq L$.*
- *Treewidth τ : Treewidth (Definition A.1) of the adjacency graph of Q is at most τ .*

Let $(w_i)_{i \in [n]} \in \mathbb{R}_{\geq 1}^n$ and $\kappa = \sum_{i \in [n]} w_i \nu_i$. Given any $0 < \epsilon \leq \frac{1}{2}$, we can find an approximate solution $x \in \mathcal{K}$ satisfying

$$\begin{aligned} \frac{1}{2} x^\top Q x + c^\top x &\leq \min_{Ax=b, x \in \mathcal{K}} \left(\frac{1}{2} x^\top Q x + c^\top x \right) + \epsilon LR(R+1), \\ \|Ax - b\|_1 &\leq 3\epsilon(R\|A\|_1 + \|b\|_1), \end{aligned}$$

in expected time

$$\tilde{O}((\sqrt{\kappa} n^{-1/2} + \log(R/(r\epsilon))) \cdot n(\tau^2 m + \tau m^2)^{1/2} (\tau^{\omega-1} + \tau m + m^{\omega-1})^{1/2}).$$

When $\max_{i \in [n]} \nu_i = \tilde{O}(1)$, $w_i = 1$, $m = \tilde{O}(\tau^{\omega-2})$, the running time simplifies to

$$\tilde{O}(n\tau^{(\omega+1)/2} m^{1/2} \log(R/(r\epsilon))).$$

D.2 ALGORITHM STRUCTURE AND CENTRAL PATH MAINTENANCE

Our algorithm is based on the robust Interior Point Method (robust IPM). Details of the robust IPM will be given in Section F. During the algorithm, we maintain a primal-dual solution pair $(x, s) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$ on the robust central path. In addition, we maintain a sparsely-changing approximation $(\bar{x}, \bar{s}) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$ to (x, s) . In each iteration, we implicitly perform update

$$\begin{aligned} x &\leftarrow x + \bar{t} B_{w, \bar{x}, \bar{t}}^{-1/2} (I - P_{w, \bar{x}, \bar{t}}) B_{w, \bar{x}, \bar{t}}^{-1/2} \delta_\mu \\ s &\leftarrow s + \bar{t} \delta_\mu - \bar{t}^2 H_{w, \bar{x}} B_{w, \bar{x}, \bar{t}}^{-1/2} (I - P_{w, \bar{x}, \bar{t}}) B_{w, \bar{x}, \bar{t}}^{-1/2} \delta_\mu \end{aligned}$$

where

$$H_{w, \bar{x}} = \nabla^2 \phi_w(\bar{x}) \quad (\text{see Eq. (24)})$$

$$B_{w,\bar{x},\bar{t}} = Q + \bar{t}H_{w,\bar{x}} \quad (\text{see Eq. (25)})$$

$$P_{w,\bar{x},\bar{t}} = B_{w,\bar{x},\bar{t}}^{-1/2}A^\top(AB_{w,\bar{x},\bar{t}}^{-1}A^\top)^{-1}AB_{w,\bar{x},\bar{t}}^{-1/2} \quad (\text{see Eq. (26)})$$

and explicitly maintain (\bar{x}, \bar{s}) such that they remain close to (x, s) in ℓ_∞ -distance.

This task is handled by the CENTRALPATHMAINTENANCE data structure, which is our main data structure. The robust IPM algorithm (Algorithm 19, 20) directly calls it in every iteration.

The CENTRALPATHMAINTENANCE data structure (Algorithm 1) has two main sub data structures, EXACTDS (Algorithm 2, 3) and APPROXDS (Algorithm 4, 5). EXACTDS is used to maintain (x, s) , and APPROXDS is used to maintain (\bar{x}, \bar{s}) .

Algorithm 1 Our main data structure for low-treewidth QP solver.

```

1: data structure CENTRALPATHMAINTENANCE ▷ Theorem D.2
2: private : member
3:   EXACTDS exact ▷ Algorithm 2, 3
4:   APPROXDS approx ▷ Algorithm 4
5:    $\ell \in \mathbb{N}$ 
6: end members
7: procedure INITIALIZE( $x, s \in \mathbb{R}^{n_{\text{tot}}}, t \in \mathbb{R}_+, \bar{\epsilon} \in (0, 1)$ ) ▷ Algorithm 2
8:   exact.INITIALIZE( $x, s, x, s, t$ )
9:    $\ell \leftarrow 0$ 
10:   $w \leftarrow \nu_{\max}, N \leftarrow \sqrt{\kappa} \log n \log \frac{n\kappa R}{\bar{\epsilon}r}$ 
11:   $q \leftarrow n^{1/2}(\tau^2 m + \tau m^2)^{-1/2}(\tau^{\omega-1} + \tau m + m^{\omega-1})^{1/2}$ 
12:   $\epsilon_{\text{apx},x} \leftarrow \bar{\epsilon}, \zeta_x \leftarrow 2\alpha, \delta_{\text{apx}} \leftarrow \frac{1}{N}$ 
13:   $\epsilon_{\text{apx},s} \leftarrow \bar{\epsilon} \cdot \bar{t}, \zeta_s \leftarrow 3\alpha \bar{t}$ 
14:
15:      approx.INITIALIZE( $x, s, h, \tilde{h}, \epsilon_x, \epsilon_s, H_{w,\bar{x}}^{1/2}\hat{x}, H_{w,\bar{x}}^{-1/2}\hat{s}, c_s, \beta_x, \beta_s, \beta_{c_s},$ 
16:           $\tilde{\beta}_x, \tilde{\beta}_s, q, \&\text{exact}, \epsilon_{\text{apx},x}, \epsilon_{\text{apx},s}, \delta_{\text{apx}}$ )
17:
18:      ▷ Algorithm 4. Parameters from  $x$  to  $\tilde{\beta}_s$  come from exact. &exact is pointer to exact
19: end procedure
20: procedure MULTIPLYANDMOVE( $t \in \mathbb{R}_+$ )
21:    $\ell \leftarrow \ell + 1$ 
22:   if  $|\bar{t} - t| > \bar{t} \cdot \epsilon_t$  or  $\ell > q$  then
23:      $x, s \leftarrow \text{exact.OUTPUT}()$  ▷ Algorithm 2
24:     INITIALIZE( $x, s, t, \bar{\epsilon}$ )
25:   end if
26:    $\beta_x, \beta_s, \beta_{c_s}, \tilde{\beta}_x, \tilde{\beta}_s \leftarrow \text{exact.MOVE}()$  ▷ Algorithm 2
27:    $\delta_{\bar{x}}, \delta_{\bar{s}} \leftarrow \text{approx.MOVEANDQUERY}(\beta_x, \beta_s, \beta_{c_s}, \tilde{\beta}_x, \tilde{\beta}_s)$  ▷ Algorithm 4
28:    $\delta_h, \delta_{\tilde{h}}, \delta_{\epsilon_x}, \delta_{\epsilon_s}, \delta_{H_{w,\bar{x}}^{1/2}\hat{x}}, \delta_{H_{w,\bar{x}}^{-1/2}\hat{s}}, \delta_{c_s} \leftarrow \text{exact.UPDATE}(\delta_{\bar{x}}, \delta_{\bar{s}})$  ▷ Algorithm 3
29:   approx.UPDATE( $\delta_{\bar{x}}, \delta_h, \delta_{\tilde{h}}, \delta_{\epsilon_x}, \delta_{\epsilon_s}, \delta_{H_{w,\bar{x}}^{1/2}\hat{x}}, \delta_{H_{w,\bar{x}}^{-1/2}\hat{s}}, \delta_{c_s}$ ) ▷ Algorithm 4
30: end procedure
31: procedure OUTPUT()
32:   return exact.OUTPUT() ▷ Algorithm 2
33: end procedure
34: end data structure

```

Theorem D.2. Data structure CENTRALPATHMAINTENANCE (Algorithm 1) implicitly maintains the central path primal-dual solution pair $(x, s) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$ and explicitly maintains its approximation $(\bar{x}, \bar{s}) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$ using the following functions:

- INITIALIZE($x \in \mathbb{R}^{n_{\text{tot}}}, s \in \mathbb{R}^{n_{\text{tot}}}, t_0 \in \mathbb{R}_{>0}, \epsilon \in (0, 1)$): Initializes the data structure with initial primal-dual solution pair $(x, s) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$, initial central path timestamp $t_0 \in \mathbb{R}_{>0}$ in $\tilde{O}(n(\tau^{\omega-1} + \tau m + m^{\omega-1}))$ time.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

- MULTIPLYANDMOVE($t \in \mathbb{R}_{>0}$): *It implicitly maintains*

$$x \leftarrow x + \bar{t} B_{w,\bar{x},\bar{t}}^{-1/2} (I - P_{w,\bar{x},\bar{t}}) B_{w,\bar{x},\bar{t}}^{-1/2} \delta_\mu(\bar{x}, \bar{s}, \bar{t})$$

$$s \leftarrow s + \bar{t} \delta_\mu - \bar{t}^2 H_{w,\bar{x}} B_{w,\bar{x},\bar{t}}^{-1/2} (I - P_{w,\bar{x},\bar{t}}) B_{w,\bar{x},\bar{t}}^{-1/2} \delta_\mu(\bar{x}, \bar{s}, \bar{t})$$

where $H_{w,\bar{x}}$, $B_{w,\bar{x},\bar{t}}$, $P_{w,\bar{x},\bar{t}}$ are defined in Eq. (24)(25)(26) respectively, and \bar{t} is some timestamp satisfying $|\bar{t} - t| \leq \epsilon_t \cdot \bar{t}$.

It also explicitly maintains $(\bar{x}, \bar{s}) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$ such that $\|\bar{x}_i - x_i\|_{\bar{x}_i} \leq \bar{\epsilon}$ and $\|\bar{s}_i - s_i\|_{\bar{x}_i}^* \leq \bar{t} \bar{\epsilon} w_i$ for all $i \in [n]$ with probability at least 0.9.

Assuming the function is called at most N times and t decreases from t_{max} to t_{min} , the total running time is

$$\tilde{O}((Nn^{-1/2} + \log(t_{\text{max}}/t_{\text{min}})) \cdot n(\tau^2 m + \tau m^2)^{1/2} (\tau^{\omega-1} + \tau m + m^{\omega-1})^{1/2}).$$

- OUTPUT: *Computes $(x, s) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$ exactly and outputs them in $\tilde{O}(n\tau m)$ time.*

D.3 DATA STRUCTURES USED IN CENTRALPATHMAINTENANCE

In this section we present several data structures used in CENTRALPATHMAINTENANCE, including:

- EXACTDS (Section D.3.1): This data structure maintains an implicit representation of the primal-dual solution pair (x, s) . This is directly used by CENTRALPATHMAINTENANCE.
- APPROXDS (Section D.3.2): This data structure explicitly maintains an approximation (\bar{x}, \bar{s}) of (x, s) . This data structure is directly used by CENTRALPATHMAINTENANCE.
- BATCHSKETCH (Section D.3.3): This data structure maintains a sketch of (x, s) . This data structure is used by APPROXDS.
- VECTORSKETCH (Section D.3.4): This data structure maintains a sketch of sparsely-changing vectors. This data structure is used by BATCHSKETCH.
- BALANCEDSKETCH (Section D.3.5): This data structure maintains a sketch of vectors of form $\mathcal{W}^\top v$, where v is sparsely-changing. This data structure is used by BATCHSKETCH.

Notation: In this section, for simplicity, we write $B_{\bar{x}}$ for $B_{w,\bar{x},\bar{t}}$, and $L_{\bar{x}}$ for the Cholesky factor of $B_{\bar{x}}$, i.e., $B_{\bar{x}} = L_{\bar{x}} L_{\bar{x}}^\top$.

D.3.1 EXACTDS

In this section we present the data structure EXACTDS. It maintains an implicit representation of the primal-dual solution pair (x, s) by maintaining several sparsely-changing vectors (see Eq. (8)(9)). This data structure has a similar spirit as EXACTDS in Gu & Song (2022), but we have a different representation from the previous works because we are working with quadratic programming rather than linear programming.

Theorem D.3. *Data structure EXACTDS (Algorithm 2, 3) implicitly maintains the primal-dual pair $(x, s) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$, computable via the expression*

$$x = \hat{x} + H_{w,\bar{x}}^{-1/2} \mathcal{W}^\top (h\beta_x - \tilde{h}\tilde{\beta}_x + \epsilon_x), \quad (8)$$

$$s = \hat{s} + H_{w,\bar{x}}^{1/2} c_s \beta_{c_s} - H_{w,\bar{x}}^{1/2} \mathcal{W}^\top (h\beta_s - \tilde{h}\tilde{\beta}_s + \epsilon_s), \quad (9)$$

where $\hat{x}, \hat{s} \in \mathbb{R}^{n_{\text{tot}}}$, $\mathcal{W} = L_{\bar{x}}^{-1} H_{w,\bar{x}}^{1/2} \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$, $h = L_{\bar{x}}^{-1} \delta_\mu \in \mathbb{R}^{n_{\text{tot}}}$, $c_s = H_{w,\bar{x}}^{-1/2} \bar{\delta}_\mu \in \mathbb{R}^{n_{\text{tot}}}$, $\beta_x, \beta_s, \beta_{c_s} \in \mathbb{R}$, $\tilde{h} = L_{\bar{x}}^{-1} A^\top \in \mathbb{R}^{n_{\text{tot}}} \times m$, $\tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$, $\epsilon_x, \epsilon_s \in \mathbb{R}^{n_{\text{tot}}}$.

The data structure supports the following functions:

- INITIALIZE($x, s, \bar{x}, \bar{s} \in \mathbb{R}^{n_{\text{tot}}}, \bar{t} \in \mathbb{R}_{>0}$): *Initializes the data structure in $\tilde{O}(n\tau^{\omega-1} + n\tau m + nm^{\omega-1})$ time, with initial value of the primal-dual pair (x, s) , its initial approximation (\bar{x}, \bar{s}) , and initial approximate timestamp \bar{t} .*

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

- MOVE(): *Performs robust central path step*

$$x \leftarrow x + \bar{t}B_{\bar{x}}^{-1}\delta_{\mu} - \bar{t}B_{\bar{x}}^{-1}A^{\top}(AB_{\bar{x}}^{-1}A^{\top})^{-1}AB_{\bar{x}}^{-1}\delta_{\mu}, \quad (10)$$

$$s \leftarrow s + \bar{t}\delta_{\mu} - \bar{t}^2B_{\bar{x}}^{-1}\delta_{\mu} + \bar{t}^2B_{\bar{x}}^{-1}A^{\top}(AB_{\bar{x}}^{-1}A^{\top})^{-1}AB_{\bar{x}}^{-1}\delta_{\mu} \quad (11)$$

in $O(m^{\omega})$ time by updating its implicit representation.

- UPDATE($\delta_{\bar{x}}, \delta_{\bar{s}} \in \mathbb{R}^{n_{\text{tot}}}$): *Updates the approximation pair (\bar{x}, \bar{s}) to $(\bar{x}^{\text{new}} = \bar{x} + \delta_{\bar{x}} \in \mathbb{R}^{n_{\text{tot}}}, \bar{s}^{\text{new}} = \bar{s} + \delta_{\bar{s}} \in \mathbb{R}^{n_{\text{tot}}})$ in $\tilde{O}((\tau^2m + \tau m^2)(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0))$ time, and output the changes in variables $\delta_{H_{w,\bar{x}}^{1/2}\hat{x}}, \delta_h, \delta_{\beta_x}, \delta_{\tilde{h}}, \delta_{\tilde{\beta}_x}, \delta_{\epsilon_x}, \delta_{H_{w,\bar{x}}^{-1/2}\hat{s}}, \delta_{\beta_s}, \delta_{\tilde{\beta}_s}, \delta_{\epsilon_s}$.*

Furthermore, $h, \epsilon_x, \epsilon_s$ change in $O(\tau(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0))$ coordinates, \tilde{h} changes in $\tilde{O}(\tau m(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0))$ coordinates, and $H_{\bar{x}}^{1/2}\hat{x}, H_{\bar{x}}^{-1/2}\hat{s}, c_s$ change in $O(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0)$ coordinates.

- OUTPUT(): *Output x and s in $\tilde{O}(n\tau m)$ time.*
- QUERY $x(i \in [n])$: *Output x_i in $\tilde{O}(\tau^2m)$ time. This function is used by APPROXDS.*
- QUERY $s(i \in [n])$: *Output s_i in $\tilde{O}(\tau^2m)$ time. This function is used by APPROXDS.*

Proof of Theorem D.3. By combining Lemma D.4 and D.5. □

Lemma D.4. EXACTDS correctly maintains an implicit representation of (x, s) , i.e., invariant

$$\begin{aligned} x &= \hat{x} + H_{w,\bar{x}}^{-1/2}\mathcal{W}^{\top}(h\beta_x - \tilde{h}\tilde{\beta}_x + \epsilon_x), \\ s &= \hat{s} + H_{w,\bar{x}}^{1/2}c_s\beta_{c_s} - H_{w,\bar{x}}^{1/2}\mathcal{W}^{\top}(h\beta_s - \tilde{h}\tilde{\beta}_s + \epsilon_s), \\ h &= L_{\bar{x}}^{-1}\bar{\delta}_{\mu}, \quad c_s = H_{w,\bar{x}}^{-1/2}\bar{\delta}_{\mu}, \quad \tilde{h} = L_{\bar{x}}^{-1}A^{\top}, \\ \tilde{u} &= \tilde{h}^{\top}\tilde{h}, \quad u = \tilde{h}^{\top}h, \\ \bar{\alpha} &= \sum_{i \in [n]} w_i^{-1} \cosh^2\left(\frac{\lambda}{w_i} \gamma_i(\bar{x}, \bar{s}, \bar{t})\right), \\ \bar{\delta}_{\mu} &= \bar{\alpha}^{1/2}\delta_{\mu}(\bar{x}, \bar{s}, \bar{t}) \end{aligned}$$

always holds after every external call, and return values of the queries are correct.

Proof. INITIALIZE: By checking the definitions we see that all invariants are satisfied after INITIALIZE.

MOVE: By comparing the implicit representation (8)(9) and the robust central path step (10)(11), we see that MOVE updates (x, s) correctly.

UPDATE: We would like to prove that UPDATE correctly updates the values of $h, c_s, \tilde{h}, \tilde{u}, u, \bar{\alpha}, \bar{\delta}_{\mu}$, while preserving the values of (x, s) .

First note that $H_{w,\bar{x}}, B_{\bar{x}}, L_{\bar{x}}$ are updated correctly. The remaining updates are separated into two steps: UPDATE h and UPDATE $\bar{\delta}_{\mu}$.

Step UPDATE h : The first few lines of UPDATE h updates $\bar{\alpha}$ and $\bar{\delta}_{\mu}$ correctly.

We define $H_{w,\bar{x}}^{\text{new}} := H_{w,\bar{x}} + \Delta_{H_{w,\bar{x}}}, B_{\bar{x}}^{\text{new}} := B_{\bar{x}} + \Delta_{B_{\bar{x}}}, L_{\bar{x}}^{\text{new}} := L_{\bar{x}} + \Delta_{L_{\bar{x}}}, \bar{\delta}_{\mu}^{\text{new}} := \bar{\delta}_{\mu} + \delta_{\bar{\delta}_{\mu}}$, and so on. Immediately after Algorithm 3, Line 26, we have

$$\begin{aligned} h + \delta_h &= L_{\bar{x}}^{-1}\bar{\delta}_{\mu} + L_{\bar{x}}^{-1}\delta_{\bar{\delta}_{\mu}} - (L_{\bar{x}} + \Delta_{L_{\bar{x}}})^{-1}\Delta_{L_{\bar{x}}}(L_{\bar{x}}^{-1}\bar{\delta}_{\mu} + L_{\bar{x}}^{-1}\delta_{\bar{\delta}_{\mu}}) \\ &= (L_{\bar{x}}^{-1} - (L_{\bar{x}} + \Delta_{L_{\bar{x}}})^{-1}\Delta_{L_{\bar{x}}}L_{\bar{x}}^{-1})\bar{\delta}_{\mu}^{\text{new}} \\ &= L_{\bar{x}}^{\text{new}}\bar{\delta}_{\mu}^{\text{new}}, \\ c_s + \delta_{c_s} &= H_{w,\bar{x}}^{-1/2}\bar{\delta}_{\mu} + \Delta_{H_{w,\bar{x}}^{-1/2}}(\bar{\delta}_{\mu} + \delta_{\bar{\delta}_{\mu}}) + H_{w,\bar{x}}^{-1/2}\delta_{\bar{\delta}_{\mu}} \end{aligned}$$

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

Algorithm 2 The EXACTDS data structure used in Algorithm 1.

▷ Theorem D.3

```

1: data structure EXACTDS
2: members
3:    $\bar{x}, \bar{s} \in \mathbb{R}^{n_{\text{tot}}}, \bar{t} \in \mathbb{R}_+, H_{w,\bar{x}}, B_{\bar{x}}, L_{\bar{x}} \in \mathbb{R}^{n_{\text{tot}} \times n_{\text{tot}}}$ 
4:    $\hat{x}, \hat{s}, \hat{h}, \epsilon_x, \epsilon_s, c_s \in \mathbb{R}^{n_{\text{tot}}}, \tilde{h} \in \mathbb{R}^{n_{\text{tot}} \times m}, \beta_x, \beta_s, \beta_{c_s} \in \mathbb{R}, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$ 
5:    $\tilde{u} \in \mathbb{R}^{m \times m}, u \in \mathbb{R}^m, \bar{\alpha} \in \mathbb{R}, \bar{\delta}_\mu \in \mathbb{R}^n$ 
6:    $k \in \mathbb{N}$ 
7: end members
8: procedure INITIALIZE( $x, s, \bar{x}, \bar{s} \in \mathbb{R}^{n_{\text{tot}}}, \bar{t} \in \mathbb{R}_+$ )
9:    $\bar{x} \leftarrow x, \bar{s} \leftarrow s, \bar{t} \leftarrow t$ 
10:   $\hat{x} \leftarrow x, \hat{s} \leftarrow s, \epsilon_x \leftarrow 0, \epsilon_s \leftarrow 0, \beta_x \leftarrow 0, \beta_s \leftarrow 0, \tilde{\beta}_x \leftarrow 0, \tilde{\beta}_s \leftarrow 0, \beta_{c_s} \leftarrow 0$ 
11:   $H_{w,\bar{x}} \leftarrow \nabla^2 \phi_w(\bar{x}), B_{\bar{x}} \leftarrow Q + \bar{t} H_{w,\bar{x}}$ 
12:  Compute lower Cholesky factor  $L_{\bar{x}}$  where  $L_{\bar{x}} L_{\bar{x}}^\top = B_{\bar{x}}$ 
13:  INITIALIZE $h(\bar{x}, \bar{s}, H_{w,\bar{x}}, L_{\bar{x}})$ 
14: end procedure
15: procedure INITIALIZE $h(\bar{x}, \bar{s} \in \mathbb{R}^{n_{\text{tot}}}, H_{w,\bar{x}}, L_{\bar{x}} \in \mathbb{R}^{n_{\text{tot}} \times n_{\text{tot}}})$ 
16:   for  $i \in [n]$  do
17:      $(\bar{\delta}_\mu)_i \leftarrow -\frac{\alpha \sinh(\frac{\lambda}{w_i} \gamma_i(\bar{x}, \bar{s}, \bar{t}))}{\gamma_i(\bar{x}, \bar{s}, \bar{t})} \cdot \mu_i(\bar{x}, \bar{s}, \bar{t})$ 
18:      $\bar{\alpha} \leftarrow \bar{\alpha} + w_i^{-1} \cosh^2(\frac{\lambda}{w_i} \gamma_i(\bar{x}, \bar{s}, \bar{t}))$ 
19:   end for
20:    $h \leftarrow L_{\bar{x}}^{-1} \bar{\delta}_\mu, \tilde{h} \leftarrow L_{\bar{x}}^{-1} A^\top, c_s \leftarrow H_{w,\bar{x}}^{-1/2} \bar{\delta}_\mu$ 
21:    $\tilde{u} \leftarrow \tilde{h}^\top \tilde{h}, u \leftarrow \tilde{h}^\top h$ 
22: end procedure
23: procedure MOVE()
24:    $\beta_x \leftarrow \beta_x + \bar{t} \cdot (\bar{\alpha})^{-1/2}$ 
25:    $\tilde{\beta}_x \leftarrow \tilde{\beta}_x + \bar{t} \cdot (\bar{\alpha})^{-1/2} \cdot \tilde{u}^{-1} u$ 
26:    $\beta_{c_s} \leftarrow \beta_{c_s} + \bar{t} \cdot (\bar{\alpha})^{-1/2}$ 
27:    $\beta_s \leftarrow \beta_s + \bar{t}^2 \cdot (\bar{\alpha})^{-1/2}$ 
28:    $\tilde{\beta}_s \leftarrow \tilde{\beta}_s + \bar{t}^2 \cdot (\bar{\alpha})^{-1/2} \cdot \tilde{u}^{-1} u$ 
29:   return  $\beta_x, \beta_s, \beta_{c_s}, \tilde{\beta}_x, \tilde{\beta}_s$ 
30: end procedure
31: procedure OUTPUT()
32:   return  $\hat{x} + H_{w,\bar{x}}^{-1/2} \mathcal{W}^\top (h \beta_x - \tilde{h} \tilde{\beta}_x + \epsilon_x), \hat{s} + H_{w,\bar{x}}^{1/2} c_s \beta_{c_s} - H_{w,\bar{x}}^{1/2} \mathcal{W}^\top (h \beta_s - \tilde{h} \tilde{\beta}_s + \epsilon_s)$ 
33: end procedure
34: procedure QUERY $x(i \in [n])$ 
35:   return  $\hat{x}_i + H_{w,\bar{x},(i,i)}^{-1/2} (\mathcal{W}^\top (h \beta_x - \tilde{h} \tilde{\beta}_x + \epsilon_x))_i$ 
36: end procedure
37: procedure QUERY $s(i \in [n])$ 
38:   return  $\hat{s}_i + H_{w,\bar{x},(i,i)}^{1/2} c_{s,i} \beta_{c_s} + H_{w,\bar{x},(i,i)}^{1/2} (\mathcal{W}^\top (h \beta_s - \tilde{h} \tilde{\beta}_s + \epsilon_s))_i$ 
39: end procedure
40: end data structure

```

$$\begin{aligned}
&= (H_{w,\bar{x}}^{\text{new}})^{-1/2} \bar{\delta}_\mu^{\text{new}}, \\
\tilde{h} + \delta_{\tilde{h}} &= L_{\bar{x}}^{-1} A^\top - (L_{\bar{x}} + \Delta_{L_{\bar{x}}})^{-1} \Delta_{L_{\bar{x}}} A^\top \\
&= (L_{\bar{x}}^{-1} - (L_{\bar{x}} + \Delta_{L_{\bar{x}}})^{-1} \Delta_{L_{\bar{x}}} L_{\bar{x}}^{-1}) A^\top \\
&= L_{\bar{x}}^{\text{new}} A^\top.
\end{aligned}$$

So h, c_s, \tilde{h} are updated correctly. Also

$$\begin{aligned}
\tilde{u} + \delta_{\tilde{u}} &= \tilde{h}^\top \tilde{h} + \delta_{\tilde{h}}^\top (\tilde{h} + \delta_{\tilde{h}}) + \tilde{h}^\top \delta_{\tilde{h}} = (\tilde{h} + \delta_{\tilde{h}})^\top (\tilde{h} + \delta_{\tilde{h}}), \\
u + \delta_u &= \tilde{h}^\top h + \delta_{\tilde{h}}^\top (h + \delta_h) + \tilde{h}^\top \delta_h = (\tilde{h} + \delta_{\tilde{h}})^\top (h + \delta_h).
\end{aligned}$$

1296

Algorithm 3 Algorithm 2 continued.

1297

1: **data structure** EXACTDS

▷ Theorem D.3

1298

2: **procedure** UPDATE($\delta_{\bar{x}}, \delta_{\bar{s}} \in \mathbb{R}^{n_{\text{tot}}}$)

1299

3: $\Delta_{H_{w,\bar{x}}} \leftarrow \nabla^2 \phi_w(\bar{x} + \delta_{\bar{x}}) - H_{w,\bar{x}}$ ▷ $\Delta_{H_{w,\bar{x}}}$ is non-zero only for diagonal blocks (i, i) for which $\delta_{\bar{x},i} \neq 0$

1300

4: Compute $\Delta_{L_{\bar{x}}}$ where $(L_{\bar{x}} + \Delta_{L_{\bar{x}}})(L_{\bar{x}} + \Delta_{L_{\bar{x}}})^\top = B_{\bar{x}} + \bar{t} \Delta_{H_{w,\bar{x}}}$

1301

5: UPDATE $h(\delta_{\bar{x}}, \delta_{\bar{s}}, \Delta_{H_{w,\bar{x}}}, \Delta_{L_{\bar{x}}})$

1302

6: UPDATE $\mathcal{W}(\Delta_{H_{w,\bar{x}}}, \Delta_{L_{\bar{x}}})$

1303

7: $\bar{x} \leftarrow \bar{x} + \delta_{\bar{x}}, \bar{s} \leftarrow \bar{s} + \delta_{\bar{s}}$

1304

8: $H_{w,\bar{x}} \leftarrow H_{w,\bar{x}} + \Delta_{H_{w,\bar{x}}}, B_{\bar{x}} \leftarrow B_{\bar{x}} + \bar{t} \Delta_{H_{w,\bar{x}}}, L_{\bar{x}} \leftarrow L_{\bar{x}} + \Delta_{L_{\bar{x}}}$

1305

9: **return** $\delta_h, \delta_{\tilde{h}}, \delta_{\epsilon_x}, \delta_{\epsilon_s}, \delta_{H_{w,\bar{x}}^{1/2} \hat{x}}, \delta_{H_{w,\bar{x}}^{-1/2} \tilde{s}}, \delta_{c_s}$

1306

10: **end procedure**

1307

11: **procedure** UPDATE $h(\delta_{\bar{x}}, \delta_{\bar{s}} \in \mathbb{R}^{n_{\text{tot}}}, \Delta_{H_{w,\bar{x}}}, \Delta_{L_{\bar{x}}} \in \mathbb{R}^{n_{\text{tot}} \times n_{\text{tot}}})$

1308

12: $S \leftarrow \{i \in [n] \mid \delta_{\bar{x},i} \neq 0 \text{ or } \delta_{\bar{s},i} \neq 0\}$

1309

13: $\delta_{\bar{\delta}_\mu} \leftarrow 0$

1310

14: **for** $i \in S$ **do**

1311

15: Let $\gamma_i = \gamma_i(\bar{x}, \bar{s}, \bar{t}), \gamma_i^{\text{new}} = \gamma_i(\bar{x} + \delta_{\bar{x}}, \bar{s} + \delta_{\bar{s}}, \bar{t}), \mu_i^{\text{new}} = \mu_i(\bar{x} + \delta_{\bar{x}}, \bar{s} + \delta_{\bar{s}}, \bar{t})$

1312

16: $\bar{\alpha} \leftarrow \bar{\alpha} - w_i^{-1} \cosh^2(\frac{\lambda}{w_i} \gamma_i) + w_i^{-1} \cosh^2(\frac{\lambda}{w_i} \gamma_i^{\text{new}})$

1313

17: $\delta_{\bar{\delta}_\mu, i} \leftarrow -\alpha \sinh(\frac{\lambda}{w_i} \gamma_i^{\text{new}}) \cdot \frac{1}{\gamma_i^{\text{new}}} \cdot \mu_i^{\text{new}} - \bar{\delta}_\mu, i$

1314

18: **end for**

1315

19: $\delta_{\tilde{h}} \leftarrow L_{\bar{x}}^{-1} \delta_{\bar{\delta}_\mu} - (L_{\bar{x}} + \Delta_{L_{\bar{x}}})^{-1} \Delta_{L_{\bar{x}}}(h + L_{\bar{x}}^{-1} \delta_{\bar{\delta}_\mu})$

1316

20: $\delta_{c_s} \leftarrow \Delta_{H_{w,\bar{x}}^{-1/2}}(\bar{\delta}_\mu + \delta_{\bar{\delta}_\mu}) + H_{w,\bar{x}}^{-1/2} \delta_{\bar{\delta}_\mu}$

1317

21: $\tilde{\delta}_{\tilde{h}} \leftarrow -(L_{\bar{x}} + \Delta_{L_{\bar{x}}})^{-1} \Delta_{L_{\bar{x}}} \tilde{h}$

1318

22: $\tilde{\delta}_{\tilde{s}} \leftarrow -\delta_{\tilde{\delta}_\mu} \beta_{c_s}$

1319

23: $\delta_{\epsilon_x} \leftarrow -\delta_{\tilde{h}} \beta_x + \tilde{\delta}_{\tilde{h}} \tilde{\beta}_x$

1320

24: $\delta_{\epsilon_s} \leftarrow -\delta_{\tilde{h}} \beta_s + \tilde{\delta}_{\tilde{h}} \tilde{\beta}_s$

1321

25: $\tilde{\delta}_{\tilde{u}} \leftarrow \delta_{\tilde{h}}^\top (\tilde{h} + \delta_{\tilde{h}}) + \tilde{h}^\top \tilde{\delta}_{\tilde{h}}$

1322

26: $\delta_u \leftarrow \delta_{\tilde{h}}^\top (h + \delta_{\tilde{h}}) + \tilde{h}^\top \delta_{\tilde{h}}$

1323

27: $\tilde{\delta}_\mu \leftarrow \tilde{\delta}_\mu + \delta_{\tilde{\delta}_\mu}, h \leftarrow h + \delta_{\tilde{h}}, \tilde{h} \leftarrow \tilde{h} + \delta_{\tilde{h}}, \epsilon_x \leftarrow \epsilon_x + \delta_{\epsilon_x}, \epsilon_s \leftarrow \epsilon_s + \delta_{\epsilon_s}, \tilde{u} \leftarrow \tilde{u} + \delta_{\tilde{u}},$

1324

 $u \leftarrow u + \delta_u$

1325

28: **end procedure**

1326

29: **procedure** UPDATE $\mathcal{W}(\Delta_{H_{w,\bar{x}}}, \Delta_{L_{\bar{x}}} \in \mathbb{R}^{n_{\text{tot}}})$

1327

30: $\delta_{\epsilon_x} \leftarrow \Delta_{L_{\bar{x}}}^\top L_{\bar{x}}^{-\top} (h \beta_x - \tilde{h} \tilde{\beta}_x + \epsilon_x)$

1328

31: $\delta_{\epsilon_s} \leftarrow \Delta_{L_{\bar{x}}}^\top L_{\bar{x}}^{-\top} (h \beta_s - \tilde{h} \tilde{\beta}_s + \epsilon_s)$

1329

32: $\epsilon_x \leftarrow \epsilon_x + \delta_{\epsilon_x}, \epsilon_s \leftarrow \epsilon_s + \delta_{\epsilon_s}$

1330

33: **end procedure**

1331

34: **end data structure**

1332

1333

1334

So \tilde{u} and u are maintained correctly. Furthermore, immediately after Algorithm 3, Line 26, we have

1335

$$(\hat{x} + L_{\bar{x}}^{-\top} (h^{\text{new}} \beta_x - \tilde{h}^{\text{new}} \tilde{\beta}_x + \epsilon_x^{\text{new}})) - (\hat{x} + L_{\bar{x}}^{-\top} (h \beta_x - \tilde{h} \tilde{\beta}_x + \epsilon_x))$$

1336

$$= L_{\bar{x}}^{-\top} (\delta_{\tilde{h}} \beta_x - \delta_{\tilde{h}} \tilde{\beta}_s + \delta_{\epsilon_x})$$

1337

$$= 0.$$

1338

1339

Therefore, after UPDATE h finishes, we have

1340

$$x = \hat{x} + L_{\bar{x}}^{-\top} (h \beta_x - \tilde{h} \tilde{\beta}_x + \epsilon_x).$$

1341

1342

For s , we have

1343

$$(\tilde{s}^{\text{new}} + (H_{w,\bar{x}}^{\text{new}})^{1/2} c_s^{\text{new}} \beta_{c_s} - L_{\bar{x}}^{-\top} (h^{\text{new}} \beta_s - \tilde{h}^{\text{new}} \tilde{\beta}_s + \epsilon_s^{\text{new}}))$$

1344

$$- (\tilde{s} + H_{w,\bar{x}}^{1/2} c_s \beta_{c_s} - L_{\bar{x}}^{-\top} (h \beta_s - \tilde{h} \tilde{\beta}_s + \epsilon_s))$$

1345

1346

$$\begin{aligned}
&= \delta_{\bar{s}} + \delta_{\bar{s}} \beta_{c_s} - L_{\bar{x}}^{-\top} (\delta_h \beta_s - \delta_{\tilde{h}} \tilde{\beta}_s + \delta_{\epsilon_s}) \\
&= 0.
\end{aligned}$$

Therefore, after UPDATE h finishes, we have

$$s = \hat{s} + (H_{w,\bar{x}}^{\text{new}})^{1/2} c_s \beta_{c_s} - L_{\bar{x}}^{-\top} (h \beta_s - \tilde{h} \tilde{\beta}_s + \epsilon_s).$$

So x and s are both updated correctly. This proves the correctness of UPDATE h .

Step UPDATE \mathcal{W} : Define $\epsilon_x^{\text{new}} := \epsilon_x + \delta_{\epsilon_x}$, $\epsilon_s^{\text{new}} := \epsilon_s + \delta_{\epsilon_s}$. Immediately after Algorithm 3, Line 31, we have

$$\begin{aligned}
&(\hat{x} + (L_{\bar{x}}^{\text{new}})^{-\top} (h \beta_x - \tilde{h} \tilde{\beta}_x + \epsilon_x^{\text{new}})) - (\hat{x} + L_{\bar{x}}^{-\top} (h \beta_x - \tilde{h} \tilde{\beta}_x + \epsilon_x)) \\
&= ((L_{\bar{x}}^{\text{new}})^{-\top} - L_{\bar{x}}^{-\top}) (h \beta_x - \tilde{h} \tilde{\beta}_x + \epsilon_x) + (L_{\bar{x}}^{\text{new}})^{-\top} \delta_{\epsilon_x} \\
&= 0, \\
&(\hat{s} + (H_{w,\bar{x}}^{\text{new}})^{1/2} c_s \beta_{c_s} - (L_{\bar{x}}^{\text{new}})^{-\top} (h \beta_s - \tilde{h} \tilde{\beta}_s + \epsilon_s^{\text{new}})) \\
&\quad - (\hat{s} + (H_{w,\bar{x}}^{\text{new}})^{1/2} c_s \beta_{c_s} - L_{\bar{x}}^{-\top} (h \beta_s - \tilde{h} \tilde{\beta}_s + \epsilon_s)) \\
&= (-(L_{\bar{x}}^{\text{new}})^{-\top} + L_{\bar{x}}^{-\top}) (h \beta_s - \tilde{h} \tilde{\beta}_s + \epsilon_s) - (L_{\bar{x}}^{\text{new}})^{-\top} \delta_{\epsilon_s} \\
&= 0.
\end{aligned}$$

Therefore, after UPDATE \mathcal{W} finishes, we have

$$\begin{aligned}
x &= \hat{x} + (L_{\bar{x}}^{\text{new}})^{-\top} (h \beta_x - \tilde{h} \tilde{\beta}_x + \epsilon_x), \\
s &= \hat{s} + (H_{w,\bar{x}}^{\text{new}})^{1/2} c_s \beta_{c_s} - (L_{\bar{x}}^{\text{new}})^{-\top} (h \beta_s - \tilde{h} \tilde{\beta}_s + \epsilon_s).
\end{aligned}$$

So x and s are both updated correctly. This proves the correctness of UPDATE \mathcal{W} . \square

Lemma D.5. *We bound the running time of EXACTDS as following.*

- (i) EXACTDS.INITIALIZE (Algorithm 2) runs in $\tilde{O}(n\tau^{\omega-1} + n\tau m + nm^{\omega-1})$ time.
- (ii) EXACTDS.MOVE (Algorithm 2) runs in $\tilde{O}(m^{\omega})$ time.
- (iii) EXACTDS.OUTPUT (Algorithm 2) runs in $\tilde{O}(n\tau m)$ time and correctly outputs (x, s) .
- (iv) EXACTDS.QUERY x and EXACTDS.QUERY s (Algorithm 2) runs in $\tilde{O}(\tau^2 m)$ time and returns the correct answer.
- (v) EXACTDS.UPDATE (Algorithm 2) runs in $\tilde{O}((\tau^2 m + \tau m^2)(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0))$ time. Furthermore, $h, \epsilon_x, \epsilon_s$ change in $O(\tau(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0))$ coordinates, \tilde{h} changes in $\tilde{O}(\tau m(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0))$ coordinates, and $H_{\bar{x}}^{1/2} \hat{x}, H_{\bar{x}}^{-1/2} \hat{s}, c_s$ change in $O(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0)$ coordinates.

Proof. (i) Computing $L_{\bar{x}}$ takes $\tilde{O}(n\tau^{\omega-1})$ time by Lemma A.5. Computing h and \tilde{h} takes $\tilde{O}(n\tau m)$ by Lemma A.7(i).⁷ Computing \tilde{u} and u takes $\mathcal{T}_{\text{mat}}(m, n, m) = \tilde{O}(nm^{\omega-1})$ time. All other operations are cheap.

(ii) Computing \tilde{u}^{-1} takes $\tilde{O}(m^{\omega})$ time. All other operations take $O(m^2)$ time.

(iii) Running time is by Lemma A.7(v). Correctness is by Lemma D.4.

(iv) Running time is by Lemma A.8(ii). Correctness is by Lemma D.4.

(v) Computing $\Delta_{L_{\bar{x}}}$ takes $\tilde{O}(\tau^2 \|\delta_{\bar{x}}\|_0)$ time by Lemma A.6. It is easy to see that $\text{nnz}(\Delta_{H_{w,\bar{x}}}) = O(\|\delta_{\bar{x}}\|_0)$ and $\text{nnz}(\Delta_{L_{\bar{x}}}) = \tilde{O}(\tau^2 \|\delta_{\bar{x}}\|_0)$. It remains to analyze UPDATE h and UPDATE \mathcal{W} . For simplicity, we write $k = \delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0$ in this proof only.

⁷Here we compute \tilde{h} by computing $\tilde{h}_{*,i} = L_{\bar{x}}^{-1}(A_{i,*})^{\top}$ for $i \in [m]$ independently. Using fast rectangular matrix multiplication is possible to improve this running time and other similar places. We keep the current bounds for simplicity.

- 1404 • UPDATE h : Updating $\bar{\alpha}$ and computing $\delta_{\bar{\delta}_\mu}$ takes $O(k)$ time. Also, $\|\delta_{\bar{\delta}_\mu}\|_0 = O(k)$.
 1405
 1406 Computing δ_h takes $\tilde{O}(\tau^2 k)$ time by Lemma A.7(i). Also, δ_h is supported on $O(k)$
 1407 paths in the elimination tree, so $\|\delta_h\|_0 = \tilde{O}(\tau k)$. Similarly we see that computing $\delta_{\tilde{h}}$
 1408 take $\tilde{O}(\tau^2 mk)$ time and $\text{nnz}(\delta_{\tilde{h}}) = \tilde{O}(\tau mk)$.
 1409 Computing δ_{c_s} and $\delta_{\tilde{s}}$ takes $O(\tau^2 k)$ time and $\|\delta_{c_s}\|_0, \|\delta_{\tilde{s}}\|_0 = O(k)$.
 1410 Computing δ_{ϵ_x} and δ_{ϵ_s} takes $O(\tau mk)$ time after computing δ_h and $\delta_{\tilde{h}}$. Furthermore,
 1411 $\|\delta_{\epsilon_x}\|_0, \|\delta_{\epsilon_s}\|_0 = O(\tau k)$.
 1412 Computing $\delta_{\tilde{u}}$ takes $\mathcal{T}_{\text{mat}}(m, \tau k, m) = \tilde{O}(\tau m^2 k)$ time. Computing δ_u takes $\tilde{O}(\tau mk)$
 1413 time.
 1414 • UPDATE W : To compute δ_{ϵ_x} and δ_{ϵ_s} , we first compute $L_{\bar{x}}^{-\top}(h\beta_x - \tilde{h}\tilde{\beta}_x + \epsilon_x)$ and
 1415 $L_{\bar{x}}^{-\top}(h\beta_s - \tilde{h}\tilde{\beta}_s + \epsilon_s)$, where $S \subseteq [n_{\text{tot}}]$ is the row support of $\Delta_{L_{\bar{x}}}$, which can
 1416 be decomposed into at most $O(\|\delta_{\bar{x}}\|_0)$ paths. This takes $\tilde{O}(\tau^2 m \|\delta_{\bar{x}}\|_0)$ time by
 1417 Lemma A.8(i) (the extra m factor is due to \tilde{h}).
 1418
 1419 Combining everything we finish the proof of running time of EXACTDS.UPDATE. \square
 1420

1421 D.3.2 APPROXDS

1422 In this section we present the data structure APPROXDS. Given BATCHSKETCH, a data structure
 1423 maintaining a sketch of the primal-dual pair $(x, s) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$, APPROXDS maintains a sparsely-
 1424 changing ℓ_∞ -approximation of (x, s) . This data structure is a slight variation of APPROXDS in Gu &
 1425 Song (2022).
 1426

1427 **Algorithm 4** The APPROXDS data structure used in Algorithm 1.

1428 1: **data structure** APPROXDS ▷ Theorem D.6
 1429 2: **private : members**
 1430 3: $\epsilon_{\text{apx},x}, \epsilon_{\text{apx},s} \in \mathbb{R}$
 1431 4: $\ell \in \mathbb{N}$
 1432 5: BATCHSKETCH bs ▷ This maintains a sketch of $H_{w,\bar{x}}^{1/2}x$ and $H_{w,\bar{x}}^{-1/2}s$. See Algorithm 6, 7, 8.
 1433 6: EXACTDS* exact ▷ This is a pointer to the EXACTDS (Algorithm 2, 3) we maintain in parallel to
 1434 APPROXDS.
 1435 7: $\tilde{x}, \tilde{s} \in \mathbb{R}^{n_{\text{tot}}}$ ▷ (\tilde{x}, \tilde{s}) is a sparsely-changing approximation of (x, s) . They have the same value as
 1436 (\bar{x}, \bar{s}) , but for these local variables we use (\tilde{x}, \tilde{s}) to avoid confusion.
 1437 8: **end members**
 1438 9: **procedure** INITIALIZE($x, s \in \mathbb{R}^{n_{\text{tot}}}, h \in \mathbb{R}^{n_{\text{tot}}}, \tilde{h} \in \mathbb{R}^{n_{\text{tot}} \times m}, \epsilon_x, \epsilon_s, H_{w,\bar{x}}^{1/2}\hat{x}, H_{w,\bar{x}}^{-1/2}\hat{s}, c_s \in$
 1439 $\mathbb{R}^{n_{\text{tot}}}, \beta_x, \beta_s, \beta_{c_s} \in \mathbb{R}, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m, q \in \mathbb{N}, \text{EXACTDS* exact}, \epsilon_{\text{apx},x}, \epsilon_{\text{apx},s}, \delta_{\text{apx}} \in \mathbb{R}$)
 1440 10: $\ell \leftarrow 0, q \leftarrow q$
 1441 11: $\epsilon_{\text{apx},x} \leftarrow \epsilon_{\text{apx},x}, \epsilon_{\text{apx},s} \leftarrow \epsilon_{\text{apx},s}$
 1442 12: bs.INITIALIZE($x, h, \tilde{h}, \epsilon_x, \epsilon_s, H_{w,\bar{x}}^{1/2}\hat{x}, H_{w,\bar{x}}^{-1/2}\hat{s}, c_s, \beta_x, \beta_s, \beta_{c_s}, \tilde{\beta}_x, \tilde{\beta}_s, \delta_{\text{apx}}/q$) ▷ Algorithm 6
 1443 13: $\tilde{x} \leftarrow x, \tilde{s} \leftarrow s$
 1444 14: exact \leftarrow exact
 1445 15: **end procedure**
 1446 16: **procedure** UPDATE($\delta_{\bar{x}} \in \mathbb{R}^{n_{\text{tot}}}, \delta_h \in \mathbb{R}^{n_{\text{tot}}}, \delta_{\tilde{h}} \in \mathbb{R}^{n_{\text{tot}} \times m}, \delta_{\epsilon_x}, \delta_{\epsilon_s}, \delta_{H_{w,\bar{x}}^{1/2}\hat{x}}, \delta_{H_{w,\bar{x}}^{-1/2}\hat{s}}, \delta_{c_s} \in \mathbb{R}^{n_{\text{tot}}}$)
 1447 17: bs.UPDATE($\delta_{\bar{x}}, \delta_h, \delta_{\tilde{h}}, \delta_{\epsilon_x}, \delta_{\epsilon_s}, \delta_{H_{w,\bar{x}}^{1/2}\hat{x}}, \delta_{H_{w,\bar{x}}^{-1/2}\hat{s}}, \delta_{c_s}$) ▷ Algorithm 7
 1448 18: $\ell \leftarrow \ell + 1$
 1449 19: **end procedure**
 1450 20: **procedure** MOVEANDQUERY($\beta_x, \beta_s, \beta_{c_s} \in \mathbb{R}, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$)
 1451 21: bs.MOVE($\beta_x, \beta_s, \beta_{c_s}, \tilde{\beta}_x, \tilde{\beta}_s$) ▷ Algorithm 7. Do not update ℓ yet
 1452 22: $\delta_{\bar{x}} \leftarrow \text{QUERY}x(\epsilon_{\text{apx},x}/(2 \log q + 1))$ ▷ Algorithm 5
 1453 23: $\delta_{\bar{s}} \leftarrow \text{QUERY}s(\epsilon_{\text{apx},s}/(2 \log q + 1))$ ▷ Algorithm 5
 1454 24: $\tilde{x} \leftarrow \tilde{x} + \delta_{\bar{x}}, \tilde{s} \leftarrow \tilde{s} + \delta_{\bar{s}}$
 1455 25: **return** $(\delta_{\bar{x}}, \delta_{\bar{s}})$
 1456 26: **end procedure**
 1457 27: **end data structure**

Theorem D.6. Given parameters $\epsilon_{\text{apx},x}, \epsilon_{\text{apx},s} \in (0, 1), \delta_{\text{apx}} \in (0, 1), \zeta_x, \zeta_s \in \mathbb{R}$ such that

$$\|H_{w,\bar{x}^{(\ell)}}^{1/2}x^{(\ell)} - H_{w,\bar{x}^{(\ell)}}^{1/2}x^{(\ell+1)}\|_2 \leq \zeta_x, \quad \|H_{w,\bar{x}^{(\ell)}}^{-1/2}s^{(\ell)} - H_{w,\bar{x}^{(\ell)}}^{-1/2}s^{(\ell+1)}\|_2 \leq \zeta_s$$

1458 **Algorithm 5** APPROXDS Algorithm 4 continued.

1459 1: **data structure** APPROXDS ▷ Theorem D.6

1460 2: **private:**

1461 3: **procedure** QUERY $x(\epsilon \in \mathbb{R})$

1462 4: $\mathcal{I} \leftarrow 0$

1463 5: **for** $j = 0 \rightarrow \lfloor \log_2 \ell \rfloor$ **do**

1464 6: **if** $\ell \bmod 2^j = 0$ **then**

1465 7: $\mathcal{I} \leftarrow \mathcal{I} \cup \text{bs.QUERY}_x(\ell - 2^j + 1, \epsilon)$ ▷ Algorithm 8

1466 8: **end if**

1467 9: **end for**

1468 10: $\delta_{\tilde{x}} \leftarrow 0$

1469 11: **for all** $i \in \mathcal{I}$ **do**

1470 12: $x_i \leftarrow \text{exact.QUERY}_x(i)$ ▷ Algorithm 2

1471 13: **if** $\|\tilde{x}_i - x_i\|_{\tilde{x}_i} > \epsilon$ **then**

1472 14: $\delta_{\tilde{x},i} \leftarrow x_i - \tilde{x}_i$

1473 15: **end if**

1474 16: **end for**

1475 17: **return** $\delta_{\tilde{x}}$

1476 18: **end procedure**

1477 19: **procedure** QUERY $s(\epsilon \in \mathbb{R})$

1478 20: Same as QUERY x , except for replacing x, \tilde{x}, \dots with s, \tilde{s}, \dots , and replacing “ $\|\tilde{x}_i - x_i\|_{\tilde{x}_i}$ ”
in Line 13 with “ $\|\tilde{s}_i - s_i\|_{\tilde{s}_i}^*$ ”.

1479 21: **end procedure**

1480 22: **end data structure**

1481

1482 *for all* $\ell \in \{0, \dots, q-1\}$, *data structure* APPROXDS (Algorithm 4 and Algorithm 5) *supports the*
1483 *following operations:*

- 1485 • INITIALIZE($x, s \in \mathbb{R}^{n_{\text{tot}}}, h \in \mathbb{R}^{n_{\text{tot}}}, \tilde{h} \in \mathbb{R}^{n_{\text{tot}} \times m}, \epsilon_x, \epsilon_s, H_{w,\tilde{x}}^{1/2} \hat{x}, H_{w,\tilde{x}}^{-1/2} \hat{s}, c_s \in$
1486 $\mathbb{R}^{n_{\text{tot}}}, \beta_x, \beta_s, \beta_{c_s} \in \mathbb{R}, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m, q \in \mathbb{N}, \text{EXACTDS}^* \text{ exact}, \epsilon_{\text{apx},x}, \epsilon_{\text{apx},s}, \delta_{\text{apx}} \in \mathbb{R}$):
1487 Initialize the data structure in $\tilde{O}(n\tau^{\omega-1} + n\tau m)$ time.
- 1488
- 1489 • MOVEANDQUERY($\beta_x, \beta_s, \beta_{c_s} \in \mathbb{R}, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$): Update values of $\beta_x, \beta_s, \beta_{c_s}, \tilde{\beta}_x, \tilde{\beta}_s$
1490 by calling BATCHSKETCH.MOVE. This effectively moves $(x^{(\ell)}, s^{(\ell)})$ to $(x^{(\ell+1)}, s^{(\ell+1)})$
1491 while keeping $\tilde{x}^{(\ell)}$ unchanged.

1492 Then return two sets $L_x^{(\ell)}, L_s^{(\ell)} \subset [n]$ where

1493

$$1494 L_x^{(\ell)} \supseteq \{i \in [n] : \|H_{w,\tilde{x}^{(\ell)}}^{1/2} x_i^{(\ell)} - H_{w,\tilde{x}^{(\ell)}}^{1/2} x_i^{(\ell+1)}\|_2 \geq \epsilon_{\text{apx},x}\},$$

1495

$$1496 L_s^{(\ell)} \supseteq \{i \in [n] : \|H_{w,\tilde{x}^{(\ell)}}^{-1/2} s_i^{(\ell)} - H_{w,\tilde{x}^{(\ell)}}^{-1/2} s_i^{(\ell+1)}\|_2 \geq \epsilon_{\text{apx},s}\},$$

1497 satisfying

1498

$$1499 \sum_{0 \leq \ell \leq q-1} |L_x^{(\ell)}| = \tilde{O}(\epsilon_{\text{apx},x}^{-2} \zeta_x^2 q^2),$$

1500

$$1501 \sum_{0 \leq \ell \leq q-1} |L_s^{(\ell)}| = \tilde{O}(\epsilon_{\text{apx},s}^{-2} \zeta_s^2 q^2).$$

1502 For every query, with probability at least $1 - \delta_{\text{apx}}/q$, the return values are correct.

1503 Furthermore, total time cost over all queries is at most

1504

$$1505 \tilde{O}((\epsilon_{\text{apx},x}^{-2} \zeta_x^2 + \epsilon_{\text{apx},s}^{-2} \zeta_s^2) q^2 \tau^2 m).$$

- 1506
- 1507 • UPDATE($\delta_{\tilde{x}} \in \mathbb{R}^{n_{\text{tot}}}, \delta_h \in \mathbb{R}^{n_{\text{tot}}}, \delta_{\tilde{h}} \in \mathbb{R}^{n_{\text{tot}} \times m}, \delta_{\epsilon_x}, \delta_{\epsilon_s}, \delta_{H_{w,\tilde{x}}^{1/2} \hat{x}}, \delta_{H_{w,\tilde{x}}^{-1/2} \hat{s}}, \delta_{c_s} \in \mathbb{R}^{n_{\text{tot}}}$):
1508 Update sketches of $H_{w,\tilde{x}^{(\ell)}}^{1/2} x^{(\ell+1)}$ and $H_{w,\tilde{x}^{(\ell)}}^{-1/2} s^{(\ell+1)}$ by calling BATCHSKETCH.UPDATE.
1509

1512 This effectively moves $\bar{x}^{(\ell)}$ to $\bar{x}^{(\ell+1)}$ while keeping $(x^{(\ell+1)}, s^{(\ell+1)})$ unchanged. Then ad-
 1513 vance timestamp ℓ .
 1514 Each update costs
 1515
$$\tilde{O}(\tau^2(\|\delta_{\bar{x}}\|_0 + \|\delta_h\|_0 + \|\delta_{\tilde{h}}\|_0 + \|\delta_{\epsilon_x}\|_0 + \|\delta_{\epsilon_s}\|_0) + \|\delta_{H_{w,\bar{x}}^{1/2}\hat{x}}\|_0 + \|\delta_{H_{w,\bar{x}}^{-1/2}\hat{s}}\|_0 + \|\delta_{c_s}\|_0)$$

 1516
 1517
 1518 time.
 1519

1520 *Proof.* The proof is essentially the same as proof of (Gu & Song, 2022, Theorem 4.18). For the
 1521 running time claims, we plug in Theorem D.8 when necessary. \square
 1522

1523 D.3.3 BATCHSKETCH

1524 In this section we present the data structure BATCHSKETCH. It maintains a sketch of $H_{\bar{x}}^{1/2}x$ and
 1525 $H_{\bar{x}}^{-1/2}s$. It is a variation of BATCHSKETCH in Gu & Song (2022).
 1526

1527 We recall the following definition from Gu & Song (2022).
 1528

1529 **Definition D.7** (Partition tree). A partition tree (\mathcal{S}, χ) of \mathbb{R}^n is a constant degree rooted tree
 1530 $\mathcal{S} = (V, E)$ and a labeling of the vertices $\chi : V \rightarrow 2^{[n]}$, such that

- 1531 • $\chi(\text{root}) = [n]$;
- 1532
- 1533 • if v is a leaf of \mathcal{S} , then $|\chi(v)| = 1$;
- 1534
- 1535 • for any non-leaf node $v \in V$, the set $\{\chi(c) : c \text{ is a child of } v\}$ is a partition of $\chi(v)$.

1537 **Algorithm 6** The BATCHSKETCH data structure used by Algorithm 4 and 5.

1538 1: **data structure** BATCHSKETCH ▷ Theorem D.8
 1539 2: **members**
 1540 3: $\Phi \in \mathbb{R}^{r \times n_{\text{tot}}}$ ▷ All sketches need to share the same sketching matrix
 1541 4: \mathcal{S}, χ partition tree
 1542 5: $\ell \in \mathbb{N}$ ▷ Current timestamp
 1543 6: BALANCEDSKETCH sketch $\mathcal{W}^\top h$, sketch $\mathcal{W}^\top \tilde{h}$, sketch $\mathcal{W}^\top \epsilon_x$, sketch $\mathcal{W}^\top \epsilon_s$ ▷ Algorithm 10
 1544 7: VECTORSKETCH sketch $H_{w,\bar{x}}^{1/2}\hat{x}$, sketch $H_{w,\bar{x}}^{-1/2}\hat{s}$, sketch c_s ▷ Algorithm 9
 1545 8: $\beta_x, \beta_s, \beta_{c_s} \in \mathbb{R}, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$
 1546 9: (history $[t]_{t \geq 0}$) ▷ Snapshot of data at timestamp t . See Remark D.9.
 1547 10: **end members**
 1548 11: **procedure** INITIALIZE($\bar{x} \in \mathbb{R}^{n_{\text{tot}}}, h \in \mathbb{R}^{n_{\text{tot}}}, \tilde{h} \in \mathbb{R}^{n_{\text{tot}} \times m}, \epsilon_x, \epsilon_s, H_{w,\bar{x}}^{1/2}\hat{x}, H_{w,\bar{x}}^{-1/2}\hat{s}, c_s \in$
 1549 $\mathbb{R}^{n_{\text{tot}}}, \beta_x, \beta_s, \beta_{c_s} \in \mathbb{R}, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m, \delta_{\text{apx}} \in \mathbb{R}$)
 1550 12: Construct partition tree (\mathcal{S}, χ) as in Definition D.11
 1551 13: $r \leftarrow \Theta(\log^3(n_{\text{tot}}) \log(1/\delta_{\text{apx}}))$
 1552 14: Initialize $\Phi \in \mathbb{R}^{r \times n_{\text{tot}}}$ with iid $\mathcal{N}(0, \frac{1}{r})$
 1553 15: $\beta_x \leftarrow \beta_x, \beta_s \leftarrow \beta_s, \beta_{c_s} \leftarrow \beta_{c_s}, \tilde{\beta}_x \leftarrow \tilde{\beta}_x, \tilde{\beta}_s \leftarrow \tilde{\beta}_s$
 1554 16: sketch $\mathcal{W}^\top h$.INITIALIZE($\mathcal{S}, \chi, \Phi, \bar{x}, h$) ▷ Algorithm 10
 1555 17: sketch $\mathcal{W}^\top \tilde{h}$.INITIALIZE($\mathcal{S}, \chi, \Phi, \bar{x}, \tilde{h}$) ▷ Algorithm 10
 1556 18: sketch $\mathcal{W}^\top \epsilon_x$.INITIALIZE($\mathcal{S}, \chi, \Phi, \bar{x}, \epsilon_x$) ▷ Algorithm 10
 1557 19: sketch $\mathcal{W}^\top \epsilon_s$.INITIALIZE($\mathcal{S}, \chi, \Phi, \bar{x}, \epsilon_s$) ▷ Algorithm 10
 1558 20: sketch $H_{w,\bar{x}}^{1/2}\hat{x}$.INITIALIZE($\mathcal{S}, \chi, \Phi, H_{w,\bar{x}}^{1/2}\hat{x}$) ▷ Algorithm 9
 1559 21: sketch $H_{w,\bar{x}}^{-1/2}\hat{s}$.INITIALIZE($\mathcal{S}, \chi, \Phi, H_{w,\bar{x}}^{-1/2}\hat{s}$) ▷ Algorithm 9
 1560 22: sketch c_s .INITIALIZE($\mathcal{S}, \chi, \Phi, c_s$) ▷ Algorithm 9
 1561 23: $\ell \leftarrow 0$. Make snapshot history $[\ell]$ ▷ Remark D.9
 1562 24: **end procedure**
 1563 25: **end data structure**

1564 **Theorem D.8.** Data structure BATCHSKETCH (Algorithm 6, 8) supports the following operations:
 1565

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

Algorithm 7 BATCHSKETCH Algorithm 6 continued.

1: **data structure** BATCHSKETCH ▷ Theorem D.8
2: **procedure** MOVE($\beta_x, \beta_s, \beta_{c_s} \in \mathbb{R}, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$)
3: $\beta_x \leftarrow \beta_x, \beta_s \leftarrow \beta_s, \beta_{c_s} \leftarrow \beta_{c_s}, \tilde{\beta}_x \leftarrow \tilde{\beta}_x, \tilde{\beta}_s \leftarrow \tilde{\beta}_s$ ▷ Do not update ℓ yet
4: **end procedure**
5: **procedure** UPDATE($\delta_{\bar{x}} \in \mathbb{R}^{n_{\text{tot}}}, \delta_h \in \mathbb{R}^{n_{\text{tot}}}, \delta_{\tilde{h}} \in \mathbb{R}^{n_{\text{tot}} \times m}, \delta_{\epsilon_x}, \delta_{\epsilon_s}, \delta_{H_{w, \bar{x}}^{1/2} \hat{x}}, \delta_{H_{w, \bar{x}}^{-1/2} \hat{s}}, \delta_{c_s} \in \mathbb{R}^{n_{\text{tot}}}$)
6: sketch $\mathcal{W}^\top h$.UPDATE($\delta_{\bar{x}}, \delta_h$) ▷ Algorithm 11
7: sketch $\mathcal{W}^\top \tilde{h}$.UPDATE($\delta_{\bar{x}}, \delta_{\tilde{h}}$) ▷ Algorithm 11
8: sketch $\mathcal{W}^\top \epsilon_x$.UPDATE($\delta_{\bar{x}}, \delta_{\epsilon_x}$) ▷ Algorithm 11
9: sketch $\mathcal{W}^\top \epsilon_s$.UPDATE($\delta_{\bar{x}}, \delta_{\epsilon_s}$) ▷ Algorithm 11
10: sketch $H_{w, \bar{x}}^{1/2} \hat{x}$.UPDATE($\delta_{H_{w, \bar{x}}^{1/2} \hat{x}}$) ▷ Algorithm 9
11: sketch $H_{w, \bar{x}}^{-1/2} \hat{s}$.UPDATE($\delta_{H_{w, \bar{x}}^{-1/2} \hat{s}}$) ▷ Algorithm 9
12: sketch c_s .UPDATE(δ_{c_s}) ▷ Algorithm 9
13: $\ell \leftarrow \ell + 1$
14: Make snapshot history[ℓ] ▷ Remark D.9
15: **end procedure**
16: **end data structure**

• INITIALIZE($\bar{x} \in \mathbb{R}^{n_{\text{tot}}}, h \in \mathbb{R}^{n_{\text{tot}}}, \tilde{h} \in \mathbb{R}^{n_{\text{tot}} \times m}, \epsilon_x, \epsilon_s, H_{w, \bar{x}}^{1/2} \hat{x}, H_{w, \bar{x}}^{-1/2} \hat{s}, c_s \in \mathbb{R}^{n_{\text{tot}}}, \beta_x, \beta_s, \beta_{c_s} \in \mathbb{R}, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m, \delta_{\text{apx}} \in \mathbb{R}$): *Initialize the data structure in $\tilde{O}(n\tau^{\omega-1} + n\tau m)$ time.*

• MOVE($\beta_x, \beta_s, \beta_{c_s} \in \mathbb{R}, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$): *Update values of $\beta_x, \beta_s, \beta_{c_s}, \tilde{\beta}_x, \tilde{\beta}_s$ in $O(m)$ time. This effectively moves $(x^{(\ell)}, s^{(\ell)})$ to $(x^{(\ell+1)}, s^{(\ell+1)})$ while keeping $\bar{x}^{(\ell)}$ unchanged.*

• UPDATE($\delta_{\bar{x}} \in \mathbb{R}^{n_{\text{tot}}}, \delta_h \in \mathbb{R}^{n_{\text{tot}}}, \delta_{\tilde{h}} \in \mathbb{R}^{n_{\text{tot}} \times m}, \delta_{\epsilon_x}, \delta_{\epsilon_s}, \delta_{H_{w, \bar{x}}^{1/2} \hat{x}}, \delta_{H_{w, \bar{x}}^{-1/2} \hat{s}}, \delta_{c_s} \in \mathbb{R}^{n_{\text{tot}}}$): *Update sketches of $H_{w, \bar{x}^{(\ell)}}^{1/2} x^{(\ell+1)}$ and $H_{w, \bar{x}^{(\ell)}}^{-1/2} s^{(\ell+1)}$. This effectively moves $\bar{x}^{(\ell)}$ to $\bar{x}^{(\ell+1)}$ while keeping $(x^{(\ell+1)}, s^{(\ell+1)})$ unchanged. Then advance timestamp ℓ .*

Each update costs

$$\tilde{O}(\tau^2(\|\delta_{\bar{x}}\|_0 + \|\delta_h\|_0 + \|\delta_{\tilde{h}}\|_0 + \|\delta_{\epsilon_x}\|_0 + \|\delta_{\epsilon_s}\|_0) + \|\delta_{H_{w, \bar{x}}^{1/2} \hat{x}}\|_0 + \|\delta_{H_{w, \bar{x}}^{-1/2} \hat{s}}\|_0 + \|\delta_{c_s}\|_0)$$

time.

• QUERY($x^{(\ell') \in \mathbb{N}, \epsilon \in \mathbb{R}$): *Given timestamp ℓ' , return a set $S \subseteq [n]$ where*

$$S \supseteq \{i \in [n] : \|H_{w, \bar{x}^{(\ell')}}^{1/2} x_i^{(\ell')} - H_{w, \bar{x}^{(\ell')}}^{1/2} x_i^{(\ell+1)}\|_2 \geq \epsilon\},$$

and

$$|S| = O(\epsilon^{-2}(\ell - \ell' + 1) \sum_{\ell' \leq t \leq \ell} \|H_{w, \bar{x}^{(t)}}^{1/2} x^{(t)} - H_{w, \bar{x}^{(t)}}^{1/2} x^{(t+1)}\|_2^2 + \sum_{\ell' \leq t \leq \ell-1} \|\bar{x}^{(t)} - \bar{x}^{(t+1)}\|_{2,0})$$

where ℓ is the current timestamp.

For every query, with probability at least $1 - \delta$, the return values are correct, and costs at most

$$\tilde{O}(\tau^2 \cdot (\epsilon^{-2}(\ell - \ell' + 1) \sum_{\ell' \leq t \leq \ell} \|H_{\bar{x}^{(t)}}^{1/2} x^{(t)} - H_{\bar{x}^{(t)}}^{1/2} x^{(t+1)}\|_2^2 + \sum_{\ell' \leq t \leq \ell-1} \|\bar{x}^{(t)} - \bar{x}^{(t+1)}\|_{2,0}))$$

running time.

• QUERY($s^{(\ell') \in \mathbb{N}, \epsilon \in \mathbb{R}$): *Given timestamp ℓ' , return a set $S \subseteq [n]$ where*

$$S \supseteq \{i \in [n] : \|H_{w, \bar{x}^{(\ell')}}^{-1/2} s_i^{(\ell')} - H_{w, \bar{x}^{(\ell')}}^{-1/2} s_i^{(\ell+1)}\|_2 \geq \epsilon\}$$

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

Algorithm 8 BATCHSKETCH Algorithm 6, 7 continued.

```

1: data structure BATCHSKETCH ▷ Theorem D.8
2: private:
3: procedure QUERYxSKETCH( $v \in S$ ) ▷ Return the value of  $\Phi_{\chi(v)}(H_{w,\bar{x}}^{1/2}x)_{\chi(v)}$ 
4:   return sketch $H_{w,\bar{x}}^{1/2}\hat{x}$ .QUERY( $v$ ) + sketch $\mathcal{W}^\top h$ .QUERY( $v$ ) ·  $\beta_x$  − sketch $\mathcal{W}^\top \tilde{h}$ .QUERY( $v$ ) ·
5:    $\tilde{\beta}_x$  + sketch $\mathcal{W}^\top \epsilon_x$ .QUERY( $v$ ) ▷ Algorithm 9, 10
6: end procedure
7: procedure QUERYsSKETCH( $v \in S$ ) ▷ Return the value of  $\Phi_{\chi(v)}(H_{w,\bar{x}}^{-1/2}s)_{\chi(v)}$ 
8:   return sketch $H_{w,\bar{x}}^{-1/2}\hat{s}$ .QUERY( $v$ ) + sketch $c_s$ .QUERY( $v$ ) ·  $\beta_{c_s}$  − sketch $\mathcal{W}^\top h$ .QUERY( $v$ ) ·
9:    $\beta_s$  + sketch $\mathcal{W}^\top \tilde{h}$ .QUERY( $v$ ) ·  $\tilde{\beta}_s$  − sketch $\mathcal{W}^\top \epsilon_s$ .QUERY( $v$ ) ▷ Algorithm 9, 10
10: end procedure
11: public:
12: procedure QUERYx( $\ell' \in \mathbb{N}, \epsilon \in \mathbb{R}$ )
13:    $L_0 = \{\text{root}(S)\}$ 
14:    $S \leftarrow \emptyset$ 
15:   for  $d = 0 \rightarrow \infty$  do
16:     if  $L_d = \emptyset$  then
17:       return  $S$ 
18:     end if
19:      $L_{d+1} \leftarrow \emptyset$ 
20:     for  $v \in L_d$  do
21:       if  $v$  is a leaf node then
22:          $S \leftarrow S \cup \{v\}$ 
23:       else
24:         for  $u$  child of  $v$  do
25:           if  $\|\text{QUERY}_x\text{SKETCH}(u) - \text{history}[\ell']\text{.QUERY}_x\text{SKETCH}(u)\|_2 > 0.9\epsilon$  then
26:              $L_{d+1} \leftarrow L_{d+1} \cup \{u\}$ 
27:           end if
28:         end for
29:       end if
30:     end for
31:   end procedure
32: procedure QUERYs( $\ell' \in \mathbb{N}, \epsilon \in \mathbb{R}$ )
33:   Same as QUERYx, except for replacing QUERYxSKETCH in Line 23 with QUERYsSKETCH.
34: end procedure
35: end structure

```

and

$$|S| = O(\epsilon^{-2}(\ell - \ell' + 1) \sum_{\ell' \leq t \leq \ell} \|H_{w,\bar{x}^{(t)}}^{-1/2}s^{(t)} - H_{w,\bar{x}^{(t)}}^{-1/2}s^{(t+1)}\|_2^2 + \sum_{\ell' \leq t \leq \ell-1} \|\bar{x}^{(t)} - \bar{x}^{(t+1)}\|_{2,0})$$

where ℓ is the current timestamp.

For every query, with probability at least $1 - \delta$, the return values are correct, and costs at most

$$\tilde{O}(\tau^2 \cdot (\epsilon^{-2}(\ell - \ell' + 1) \sum_{\ell' \leq t \leq \ell} \|H_{\bar{x}^{(t)}}^{1/2}s^{(t)} - H_{\bar{x}^{(t)}}^{1/2}s^{(t+1)}\|_2^2 + \sum_{\ell' \leq t \leq \ell-1} \|\bar{x}^{(t)} - \bar{x}^{(t+1)}\|_{2,0}))$$

running time.

Proof. The proof is essentially the same as proof of (Gu & Song, 2022, Theorem 4.21). For the running time claims, we plug in Lemma D.10 and D.12 when necessary. \square

Remark D.9 (Snapshot). As in previous works, we use persistent data structures (e.g., Driscoll et al. (1989)) to keep a snapshot of the data structure after every update. This allows us to support query to

1674 historical data. This incurs an $O(\log n_{\text{tot}}) = \tilde{O}(1)$ multiplicative factor in all running times, which
 1675 we ignore in our analysis.
 1676

1677 D.3.4 VECTORSKETCH

1678 VECTORSKETCH is a data structure used to maintain sketches of sparsely-changing vectors. It is a
 1679 direct application of segment trees. For completeness, we include code (Algorithm 9) from (Gu &
 1680 Song, 2022, Algorithm 10).
 1681

1682 **Algorithm 9** (Gu & Song, 2022, Algorithm 10). Used in Algorithm 6, 7, 8.

```

1684 1: data structure VECTORSKETCH ▷ Lemma D.10
1685 2: private: members
1686 3:    $\Phi \in \mathbb{R}^{r \times n_{\text{tot}}}$ 
1687 4:   Partition tree  $(\mathcal{S}, \chi)$ 
1688 5:    $x \in \mathbb{R}^{n_{\text{tot}}}$ 
1689 6:   Segment tree  $\mathcal{T}$  on  $[n]$  with values in  $\mathbb{R}^r$ 
1690 7: end members
1691 8: procedure INITIALIZE( $\mathcal{S}, \chi$  : partition tree,  $\Phi \in \mathbb{R}^{r \times n_{\text{tot}}}, x \in \mathbb{R}^{n_{\text{tot}}}$ )
1692 9:    $(\mathcal{S}, \chi) \leftarrow (\mathcal{S}, \chi), \Phi \leftarrow \Phi$ 
1693 10:   $x \leftarrow x$ 
1694 11:  Order leaves of  $\mathcal{S}$  (variable blocks) such that every node  $\chi(v)$  corresponds to a contiguous
1695 12:  interval  $\subseteq [n]$ .
1696 13:  Build a segment tree  $\mathcal{T}$  on  $[n]$  such that each segment tree interval  $I \subseteq [n]$  maintains
1697 14:   $\Phi_I x_I \in \mathbb{R}^r$ .
1698 15: end procedure
1699 16: procedure UPDATE( $\delta_x \in \mathbb{R}^{n_{\text{tot}}}$ )
1700 17:   for all  $i \in [n_{\text{tot}}]$  such that  $\delta_{x,i} \neq 0$  do
1701 18:     Let  $j \in [n]$  be such that  $i$  is in  $j$ -th block
1702 19:     Update  $\mathcal{T}$  at  $j$ -th coordinate  $\Phi_j x_j \leftarrow \Phi_j x_j + \Phi_i \cdot \delta_{x,i}$ .
1703 20:      $x_i \leftarrow x_i + \delta_{x,i}$ 
1704 21:   end for
1705 22: end procedure
1706 23: procedure QUERY( $v \in V(\mathcal{S})$ )
1707 24:   Find interval  $I$  corresponding to  $\chi(v)$ 
1708 25:   return range sum of  $\mathcal{T}$  on interval  $I$ 
1709 26: end procedure
1710 27: end data structure

```

1710 **Lemma D.10** ((Gu & Song, 2022, Lemma 4.23)). *Given a partition tree (\mathcal{S}, χ) of \mathbb{R}^n , and a*
 1711 *JL sketching matrix $\Phi \in \mathbb{R}^{r \times n_{\text{tot}}}$, the data structure VECTORSKETCH (Algorithm 9) maintains*
 1712 *$\Phi_{\chi(v)} x_{\chi(v)}$ for all nodes v in the partition tree implicitly through the following functions:*

- 1713 • INITIALIZE(\mathcal{S}, χ, Φ): *Initializes the data structure in $O(r n_{\text{tot}})$ time.*
- 1714 • UPDATE($\delta_x \in \mathbb{R}^{n_{\text{tot}}}$): *Maintains the data structure for $x \leftarrow x + \delta_x$ in $O(r \|\delta_x\|_0 \log n)$*
 1715 *time.*
- 1716 • QUERY($v \in V(\mathcal{S})$): *Outputs $\Phi_{\chi(v)} x_{\chi(v)}$ in $O(r \log n)$ time.*

1719 D.3.5 BALANCEDSKETCH

1720 In this section, we present data structure BALANCEDSKETCH. It is a data structure for maintaining a
 1721 sketch of a vector of form $\mathcal{W}^\top h$, where $\mathcal{W} = L_{\bar{x}}^{-1} H_{w, \bar{x}}^{1/2}$ and $h \in \mathbb{R}^{n_{\text{tot}}}$ is a sparsely-changing vector.
 1722 This is a variation of BLOCKBALANCEDSKETCH in Gu & Song (2022).
 1723

1724 We use the following construction of a partition tree.
 1725

1726 **Definition D.11** (Construction of Partition Tree). *We fix an ordering π of $[n]$ using the heavy-light*
 1727 *decomposition (Lemma A.10). Let \mathcal{S} be a complete binary tree with leaf set $[n]$ and ordering π . Let χ*
map a node to the set of leaves in its subtree. Then (\mathcal{S}, χ) is a valid partition tree.

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

Algorithm 10 The BALANCEDSKETCH data structure is used in Algorithm 6, 7, 8.

1: **data structure** BALANCEDSKETCH ▷ Lemma D.12
2: **private: members**
3: $\Phi \in \mathbb{R}^{r \times n_{\text{tot}}}$
4: Partition tree (\mathcal{S}, χ) with balanced binary tree \mathcal{B}
5: $t \in \mathbb{N}$
6: $h \in \mathbb{R}^{n_{\text{tot}}}, \bar{x} \in \mathbb{R}^{n_{\text{tot}}}, H_{w, \bar{x}} \in \mathbb{R}^{n_{\text{tot}} \times n_{\text{tot}}}$
7: $\{L[t] \in \mathbb{R}^{n_{\text{tot}} \times n_{\text{tot}}}\}_{t \geq 0}$
8: $\{J_v \in \mathbb{R}^{r \times n_{\text{tot}}}\}_{v \in \mathcal{S}}$
9: $\{Z_v \in \mathbb{R}^{r \times n_{\text{tot}}}\}_{v \in \mathcal{B}}$
10: $\{y_v^\nabla \in \mathbb{R}^r\}_{v \in \mathcal{B}}$
11: $\{t_v \in \mathbb{N}\}_{v \in \mathcal{B}}$
12: **end members**
13: **procedure** INITIALIZE(\mathcal{S}, χ : partition tree, $\Phi \in \mathbb{R}^{r \times n_{\text{tot}}}, \bar{x} \in \mathbb{R}^{n_{\text{tot}}}, h \in \mathbb{R}^{n_{\text{tot}} \times k}$)
14: $(\mathcal{S}, \chi) \leftarrow (\mathcal{S}, \chi), \Phi \leftarrow \Phi$
15: $t \leftarrow 0, h \leftarrow h$
16: $H_{w, \bar{x}} \leftarrow \nabla^2 \phi(\bar{x}), B_{\bar{x}} \leftarrow Q + \bar{t} H_{w, \bar{x}}$
17: Compute lower Cholesky factor $L_{\bar{x}}[t]$ of $B_{\bar{x}}$
18: **for all** $v \in \mathcal{S}$ **do**
19: $J_v \leftarrow \Phi_{\chi(v)} H_{w, \bar{x}}^{1/2}$
20: **end for**
21: **for all** $v \in \mathcal{B}$ **do**
22: $Z_v \leftarrow J_v L_{\bar{x}}[t]^{-\top}$
23: $y_v^\nabla \leftarrow Z_v (I - I_{\Lambda(v)}) h$
24: $t_v \leftarrow t$
25: **end for**
26: **end procedure**
27: **procedure** QUERY($v \in \mathcal{S}$)
28: **if** $v \in \mathcal{S} \setminus \mathcal{B}$ **then**
29: **return** $J_v \cdot L_{\bar{x}}[t]^{-\top} h$
30: **end if**
31: $\Delta_{L_{\bar{x}}} \leftarrow (L_{\bar{x}}[t] - L_{\bar{x}}[t_v]) \cdot I_{\Lambda(v)}$
32: $\delta_{Z_v} \leftarrow -(L_{\bar{x}}[t]^{-1} \cdot \Delta_{L_{\bar{x}}} \cdot Z_v^\top)^\top$
33: $Z_v \leftarrow Z_v + \delta_{Z_v}$
34: $\delta_{y_v^\nabla} \leftarrow \delta_{Z_v} \cdot (I - I_{\Lambda(v)}) h$
35: $y_v^\nabla \leftarrow y_v^\nabla + \delta_{y_v^\nabla}$
36: $t_v \leftarrow t$
37: $y_v^\Delta \leftarrow Z_v \cdot I_{\Lambda(v)} \cdot h$
38: **return** $y_v^\Delta + y_v^\nabla$
39: **end procedure**
40: **end data structure**

Lemma D.12. Given an elimination tree \mathcal{T} with height η , a JL matrix $\Phi \in \mathbb{R}^{r \times n_{\text{tot}}}$, and a partition tree (\mathcal{S}, χ) constructed as in Definition D.11 with height $\tilde{O}(1)$, the data structure BALANCEDSKETCH (Algorithm 10, 11, 12), maintains $\Phi_{\chi(v)} (\mathcal{W}^\top h)_{\chi(v)}$ for each $v \in V(\mathcal{S})$ through the following operations

- INITIALIZE((\mathcal{S}, χ) : partition tree, $\Phi \in \mathbb{R}^{n_{\text{tot}}}, \bar{x} \in \mathbb{R}^{n_{\text{tot}}}, h \in \mathbb{R}^{n_{\text{tot}} \times k}$): Initializes the data structure in $\tilde{O}(r(n\tau^{\omega-1} + n\tau k))$ time.
- UPDATE($\delta_{\bar{x}} \in \mathbb{R}^{n_{\text{tot}}}, \delta_h \in \mathbb{R}^{n_{\text{tot}} \times k}$): Updates all sketches in \mathcal{S} implicitly to reflect (\mathcal{W}, h) updating to $(\mathcal{W}^{\text{new}}, h^{\text{new}})$ in $\tilde{O}(r\tau^2 k)$ time.
- QUERY($v \in \mathcal{S}$): Outputs $\Phi_{\chi(v)} (\mathcal{W}^\top h)_{\chi(v)}$ in $\tilde{O}(r\tau^2 k)$ time.

Proof. The proof is almost same as the proof of (Gu & Song, 2022, Lemma 4.24). (In fact, our \mathcal{W} is simpler than the one used in Gu & Song (2022).)

Algorithm 11 BALANCEDSKETCH Algorithm 10 continued. This is used in Algorithm 6, 7, 8.

```

1782 1: data structure BALANCEDSKETCH
1783 2: procedure UPDATE( $\delta_{\bar{x}} \in \mathbb{R}^{n_{\text{tot}}}, \delta_h \in \mathbb{R}^{n_{\text{tot}} \times k}$ )
1784 3:   for  $i \in [n]$  where  $\delta_{\bar{x},i} \neq 0$  do
1785 4:     UPDATE $\bar{x}(\delta_{\bar{x},i})$ 
1786 5:   end for
1787 6:   for all  $\delta_{h,i} \neq 0$  do
1788 7:      $v \leftarrow \Lambda^\circ(i)$ 
1789 8:     for all  $u \in \mathcal{P}^{\mathcal{B}}(v)$  do
1790 9:        $y_u^\nabla \leftarrow y_v^\nabla + Z_u \cdot I_{\{i\}} \cdot \delta_h$ 
1791 10:    end for
1792 11:  end for
1793 12:   $h \leftarrow h + \delta_h$ 
1794 13: end procedure
1795 14: procedure UPDATE $\bar{x}(\delta_{\bar{x},i} \in \mathbb{R}^{n_i})$ 
1796 15:   $t \leftarrow t + 1$ 
1797 16:   $\bar{x}_i \leftarrow \bar{x}_i + \delta_{\bar{x},i}$ 
1798 17:   $\Delta_{H_{w,\bar{x}}(i,i)} \leftarrow \nabla^2 \phi_i(\bar{x}_i) - H_{w,\bar{x}}(i,i)$ 
1799 18:  Compute  $\Delta_{L_{\bar{x}}}$  such that  $L_{\bar{x}}[t] \leftarrow L_{\bar{x}}[t-1] + \Delta_{L_{\bar{x}}}$  is the lower Cholesky factor of  $A(H_{w,\bar{x}} +$ 
1800  $\Delta_{H_{w,\bar{x}}})^{-1} A^\top$ 
1801 19:   $S \leftarrow \mathcal{P}^{\mathcal{B}}(\Lambda^\circ(\text{low}^\top(i)))$ 
1802 20:  UPDATE $L(S, \Delta_{L_{\bar{x}}})$ 
1803 21:  UPDATE $H(i, \Delta_{H_{w,\bar{x}}(i,i)})$ 
1804 22: end procedure
1805 23: end data structure

```

For INITIALIZE running time, we note that computing Z_v for all $v \in \mathcal{B}$ takes $\tilde{O}(rn\tau^{\omega-1})$ time by (Gu & Song, 2022, Lemma 8.3). Because Z_v is supported on the path from v to the root in \mathcal{T} , we know that $\text{nnz}(Z) = O(rn\tau)$. Therefore computing y_v^∇ for all $v \in \mathcal{B}$ takes $\tilde{O}(rn\tau k)$ time.

Remaining claims follow from combining proof of (Gu & Song, 2022, Lemma 4.24) and (Gu & Song, 2022, Lemma 8.3). \square

D.4 ANALYSIS OF CENTRALPATHMAINTENANCE

Lemma D.13 (Correctness of CENTRALPATHMAINTENANCE). *Algorithm 1 implicitly maintains the primal-dual solution pair (x, s) via representation Eq. (8)(9). It also explicitly maintains $(\bar{x}, \bar{s}) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$ such that $\|\bar{x}_i - x_i\|_{\bar{x}_i} \leq \bar{\epsilon}$ and $\|\bar{s}_i - s_i\|_{\bar{x}_i}^* \leq t\bar{\epsilon}w_i$ for all $i \in [n]$ with probability at least 0.9.*

Proof. We correctly maintain the implicit representation because of correctness of exact.UPDATE (Theorem D.3).

We show that $\|\bar{x}_i - x_i\|_{\bar{x}_i} \leq \bar{\epsilon}$ and $\|\bar{s}_i - s_i\|_{\bar{x}_i}^* \leq t\bar{\epsilon}w_i$ for all $i \in [n]$ (c.f. Algorithm 20, Line 16). approx maintains an ℓ_∞ approximation of $H_{w,\bar{x}}^{1/2}x$. For $\ell \leq q$, we have

$$\|H_{w,\bar{x}}^{1/2}x^{(\ell+1)} - H_{w,\bar{x}}^{1/2}x^{(\ell)}\|_2 = \|\delta_x\|_{w,\bar{x}} \leq \frac{9}{8}\alpha \leq \zeta_x$$

where the first step from definition of $\|\cdot\|_{w,\bar{x}}$, the second step follows from Lemma F.11, the third step follows from definition of ζ_x .

By Theorem D.6, with probability at least $1 - \delta_{\text{apx}}$, approx correctly maintains \bar{x} such that $\|H_{w,\bar{x}}^{1/2}\bar{x} - H_{w,\bar{x}}^{1/2}x\|_\infty \leq \epsilon_{\text{apx},x} \leq \bar{\epsilon}$. Then

$$\|\bar{x}_i - x_i\|_{\bar{x}_i} \leq w_i^{-1/2} \|H_{w,\bar{x}}^{1/2}\bar{x} - H_{w,\bar{x}}^{1/2}x\|_\infty \leq w_i^{-1/2}\bar{\epsilon} \leq \bar{\epsilon}.$$

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

Algorithm 12 BALANCEDSKETCH Algorithm 10, 11 continued. This is used in Algorithm 6, 7, 8.

1: **data structure** BALANCEDSKETCH ▷ Lemma D.12
2: **private:**
3: **procedure** UPDATE $L(S \subset \mathcal{B}, \Delta_{L_{\bar{x}}} \in \mathbb{R}^{n_{\text{tot}} \times n_{\text{tot}}})$
4: **for all** $v \in S$ **do**
5: $\delta_{Z_v} \leftarrow -(L_{\bar{x}}[t-1]^{-1}(L_{\bar{x}}[t-1] - L_{\bar{x}}[t_v]) \cdot I_{\Lambda(v)} \cdot Z_v^\top)^\top$
6: $\delta'_{Z_v} \leftarrow -(L_{\bar{x}}[t]^{-1} \cdot \Delta_{L_{\bar{x}}} \cdot (Z_v + \delta_{Z_v})^\top)^\top$
7: $Z_v \leftarrow Z_v + \delta_{Z_v} + \delta'_{Z_v}$
8: $\delta_{y_v^\nabla} \leftarrow (\delta_{Z_v} + \delta'_{Z_v})(I - I_{\Lambda(v)})h$
9: $y_v^\nabla \leftarrow y_v^\nabla + \delta_{y_v^\nabla}$
10: $t_v \leftarrow t$
11: **end for**
12: **end procedure**
13: **private:**
14: **procedure** UPDATE $H(i \in [n], \Delta_{H_{w,\bar{x}}(i,i)} \in \mathbb{R}^{n_i \times n_i})$
15: Find u such that $\chi(u) = \{i\}$
16: $\Delta_{H_{w,\bar{x}}(i,i)}^{1/2} \leftarrow (H_{w,\bar{x}}(i,i) + \Delta_{H_{w,\bar{x}}(i,i)})^{1/2} - H_{w,\bar{x}}^{1/2}(i,i)$
17: $\delta_{J_u} \leftarrow \Phi_i \cdot \Delta_{H_{w,\bar{x}}(i,i)}^{1/2}$
18: **for all** $v \in \mathcal{P}^S(u)$ **do**
19: $J_v \leftarrow J_v + \delta_{J_u}$
20: **if** $v \in \mathcal{B}$ **then**
21: $\delta_{Z_v} \leftarrow \delta_{J_v} \cdot L_{\bar{x}}[t_v]^{-\top}$
22: $Z_v \leftarrow Z_v + \delta_{Z_v}$
23: $\delta_{y_v^\nabla} \leftarrow \delta_{Z_v} \cdot (I - I_{\Lambda(v)}) \cdot h$
24: $y_v^\nabla \leftarrow y_v^\nabla + \delta_{y_v^\nabla}$
25: **end if**
26: **end for**
27: $H_{w,\bar{x}} \leftarrow H_{w,\bar{x}} + \Delta_{H_{w,\bar{x}}(i,i)}$
28: **end procedure**
29: **end data structure**

Note that the last step is loose by a factor of $w_i^{1/2}$. When w_i s are large, we could improve running time by using a tighter choice of $\epsilon_{\text{apx},x}$, as did in Gu & Song (2022). Here we use a loose bound for simplicity of presentation. Same remark applies to s .

The proof for s is similar. We have

$$\|H_{w,\bar{x}}^{-1/2} \delta_s\|_2 = \|\delta_s\|_{w,\bar{x}}^* \leq \frac{17}{8} \alpha \cdot t \leq \zeta_s$$

and

$$\|\bar{s}_i - s_i\|_{\bar{x}_i}^* \leq w_i^{1/2} \|H_{w,\bar{x}}^{-1/2} \bar{s} - H_{w,\bar{x}}^{-1/2} s\|_\infty \leq w_i^{1/2} \epsilon_{\text{apx},s} \leq \bar{\epsilon} \cdot \bar{t} \cdot w_i. \quad \square$$

Lemma D.14. We bound the running time of CENTRALPATHMAINTENANCE as following.

- CENTRALPATHMAINTENANCE.INITIALIZE takes $\tilde{O}(n\tau^{\omega-1} + n\tau m + nm^{\omega-1})$ time.
- If CENTRALPATHMAINTENANCE.MULTIPLYANDMOVE is called N times, then it has total running time

$$\tilde{O}((Nn^{-1/2} + \log(t_{\max}/t_{\min})) \cdot n(\tau^2 m + \tau m^2)^{1/2} (\tau^{\omega-1} + \tau m + m^{\omega-1})^{1/2}).$$

- CENTRALPATHMAINTENANCE.OUTPUT takes $\tilde{O}(n\tau m)$ time.

Proof. INITIALIZE part: By Theorem D.3 and D.6.

OUTPUT part: By Theorem D.3.

MULTIPLYANDMOVE part: Between two restarts, the total size of $|L_x|$ returned by approx.QUERY is bounded by $\tilde{O}(q^2 \zeta_x^2 / \epsilon_{\text{apx},x}^2)$ by Theorem D.6. By plugging in $\zeta_x = 2\alpha$, $\epsilon_{\text{apx},x} = \bar{\epsilon}$, we have $\sum_{\ell \in [q]} |L_x^{(\ell)}| = \tilde{O}(q^2)$. Similarly, for s we have $\sum_{\ell \in [q]} |L_s^{(\ell)}| = \tilde{O}(q^2)$.

Update time: By Theorem D.3 and D.6, in a sequence of q updates, total cost for update is $\tilde{O}(q^2(\tau^2 m + \tau m^2))$. So the amortized update cost per iteration is $\tilde{O}(q(\tau^2 m + \tau m^2))$. The total update cost is

$$\text{number of iterations} \cdot \text{time per iteration} = \tilde{O}(Nq(\tau^2 m + \tau m^2)).$$

Init/restart time: We restart the data structure whenever $k > q$ or $|\bar{t} - t| > \bar{t}\epsilon_t$, so there are $O(N/q + \log(t_{\max}/t_{\min})\epsilon_t^{-1})$ restarts in total. By Theorem D.3 and D.6, time cost per restart is $\tilde{O}(n(\tau^{\omega-1} + \tau m + m^{\omega-1}))$. So the total initialization time is

$$\text{number of restarts} \cdot \text{time per restart} = \tilde{O}((N/q + \log(t_{\max}/t_{\min})\epsilon_t^{-1}) \cdot n(\tau^{\omega-1} + \tau m + m^{\omega-1})).$$

Combine everything: Overall running time is

$$\tilde{O}(Nq(\tau^2 m + \tau m^2) + (N/q + \log(t_{\max}/t_{\min})\epsilon_t^{-1}) \cdot n(\tau^{\omega-1} + \tau m + m^{\omega-1})).$$

Taking $\epsilon_t = \frac{1}{2}\bar{\epsilon}$, the optimal choice for q is

$$q = n^{1/2}(\tau^2 m + \tau m^2)^{-1/2}(\tau^{\omega-1} + \tau m + m^{\omega-1})^{1/2},$$

achieving overall running time

$$\tilde{O}((Nn^{-1/2} + \log(t_{\max}/t_{\min})) \cdot n(\tau^2 m + \tau m^2)^{1/2}(\tau^{\omega-1} + \tau m + m^{\omega-1})^{1/2}). \quad \square$$

Proof of Theorem D.2. Combining Lemma D.13 and D.14. \square

D.5 PROOF OF MAIN STATEMENT

Proof of Theorem D.1. Use CENTRALPATHMAINTENANCE (Algorithm 1) as the maintenance data structure in Algorithm 20. Combining Theorem D.2 and Theorem F.1 finishes the proof. \square

E ALGORITHM FOR LOW-RANK QP

In this section we present a nearly-linear time algorithm for solving low-rank QP with small number of linear constraints. We briefly describe the outline of this section.

- In Section E.1, we present the main statement of Section E.
- In Section E.2, we present the main data structure CENTRALPATHMAINTENANCE.
- In Section E.3, we present several data structures used in CENTRALPATHMAINTENANCE, including EXACTDS (Section E.3.1), APPROXDS (Section E.3.2), BATCHSKETCH (Section E.3.3).
- In Section E.4, we prove correctness and running time of CENTRALPATHMAINTENANCE data structure.
- In Section E.5, we prove the main result (Theorem E.1).

E.1 MAIN STATEMENT

We consider programs of the form (16), i.e.,

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} x^\top Q x + c^\top x$$

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

$$\begin{aligned} \text{s.t. } Ax &= b \\ x_i &\in \mathcal{K}_i \quad \forall i \in [n] \end{aligned}$$

where $Q \in \mathcal{S}^{n_{\text{tot}}}$, $c \in \mathbb{R}^{n_{\text{tot}}}$, $A \in \mathbb{R}^{m \times n_{\text{tot}}}$, $b \in \mathbb{R}^m$, $\mathcal{K}_i \subset \mathbb{R}^{n_i}$ is a convex set. For simplicity, we assume that $n_i = O(1)$ for all $i \in [n]$.

Theorem E.1. Consider the convex program (16). Let $\phi_i : \mathcal{K}_i \rightarrow \mathbb{R}$ be a ν_i -self-concordant barrier for all $i \in [n]$. Suppose the program satisfies the following properties:

- Inner radius r : There exists $z \in \mathbb{R}^{n_{\text{tot}}}$ such that $Az = b$ and $B(z, r) \in \mathcal{K}$.
- Outer radius R : $\mathcal{K} \subseteq B(0, R)$ where $0 \in \mathbb{R}^{n_{\text{tot}}}$.
- Lipschitz constant L : $\|Q\|_{2 \rightarrow 2} \leq L$, $\|c\|_2 \leq L$.
- Low rank: We are given a factorization $Q = UV^\top$ where $U, V \in \mathbb{R}^{n_{\text{tot}} \times k}$.

Let $(w_i)_{i \in [n]} \in \mathbb{R}_{\geq 1}^n$ and $\kappa = \sum_{i \in [n]} w_i \nu_i$. Given any $0 < \epsilon \leq \frac{1}{2}$, we can find an approximate solution $x \in \mathcal{K}$ satisfying

$$\begin{aligned} \frac{1}{2} x^\top Q x + c^\top x &\leq \min_{Ax=b, x \in \mathcal{K}} \left(\frac{1}{2} x^\top Q x + c^\top x \right) + \epsilon LR(R+1), \\ \|Ax - b\|_1 &\leq 3\epsilon(R\|A\|_1 + \|b\|_1), \end{aligned}$$

in expected time

$$\tilde{O}((\sqrt{\kappa}n^{-1/2} + \log(R/(r\epsilon))) \cdot n(k+m)^{(\omega+1)/2}).$$

When $\max_{i \in [n]} \nu_i = \tilde{O}(1)$, $w_i = 1$, the running time simplifies to

$$\tilde{O}(n(k+m)^{(\omega+1)/2} \log(R/(r\epsilon))).$$

E.2 ALGORITHM STRUCTURE AND CENTRAL PATH MAINTENANCE

Similar to the low-treewidth case, our algorithm is based on the robust IPM. Details of the robust IPM will be given in Section F. During the algorithm, we maintain a primal-dual solution pair $(x, s) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$ on the robust central path. In addition, we maintain a sparsely-changing approximation $(\bar{x}, \bar{s}) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$ to (x, s) . In each iteration, we implicitly perform update

$$\begin{aligned} x &\leftarrow x + \bar{t} B_{w, \bar{x}, \bar{t}}^{-1/2} (I - P_{w, \bar{x}, \bar{t}}) B_{w, \bar{x}, \bar{t}}^{-1/2} \delta_\mu \\ s &\leftarrow s + \bar{t} \delta_\mu - \bar{t}^2 H_{w, \bar{x}} B_{w, \bar{x}, \bar{t}}^{-1/2} (I - P_{w, \bar{x}, \bar{t}}) B_{w, \bar{x}, \bar{t}}^{-1/2} \delta_\mu \end{aligned}$$

where

$$H_{w, \bar{x}} = \nabla^2 \phi_w(\bar{x}) \quad (\text{see Eq. (24)})$$

$$B_{w, \bar{x}, \bar{t}} = Q + \bar{t} H_{w, \bar{x}} \quad (\text{see Eq. (25)})$$

$$P_{w, \bar{x}, \bar{t}} = B_{w, \bar{x}, \bar{t}}^{-1/2} A^\top (AB_{w, \bar{x}, \bar{t}}^{-1} A^\top)^{-1} AB_{w, \bar{x}, \bar{t}}^{-1/2} \quad (\text{see Eq. (26)})$$

and explicitly maintain (\bar{x}, \bar{s}) such that they remain close to (x, s) in ℓ_∞ -distance.

This task is handled by the CENTRALPATHMAINTENANCE data structure, which is our main data structure. The robust IPM algorithm (Algorithm 19, 20) directly calls it in every iteration.

The CENTRALPATHMAINTENANCE data structure (Algorithm 13) has two main sub data structures, EXACTDS (Algorithm 14, 15) and APPROXDS (Algorithm 16). EXACTDS is used to maintain (x, s) , and APPROXDS is used to maintain (\bar{x}, \bar{s}) .

Theorem E.2. Data structure CENTRALPATHMAINTENANCE (Algorithm 13) implicitly maintains the central path primal-dual solution pair $(x, s) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$ and explicitly maintains its approximation $(\bar{x}, \bar{s}) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$ using the following functions:

- INITIALIZE($x \in \mathbb{R}^{n_{\text{tot}}}$, $s \in \mathbb{R}^{n_{\text{tot}}}$, $t_0 \in \mathbb{R}_{>0}$, $\epsilon \in (0, 1)$): Initializes the data structure with initial primal-dual solution pair $(x, s) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$, initial central path timestamp $t_0 \in \mathbb{R}_{>0}$ in $\tilde{O}(n(k^{\omega-1} + m^{\omega-1}))$ time.

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

Algorithm 13 Main algorithm for low-rank QP.

```

1: data structure CENTRALPATHMAINTENANCE ▷ Theorem E.2
2: private : members
3:   EXACTDS exact ▷ Algorithm 14, 15
4:   APPROXDS approx ▷ Algorithm 16
5:    $\ell \in \mathbb{N}$ 
6: end members
7: procedure INITIALIZE( $x, s \in \mathbb{R}^{n_{\text{tot}}}, t \in \mathbb{R}_+, \bar{\epsilon} \in (0, 1)$ ) ▷ Algorithm 14
8:   exact.INITIALIZE( $x, s, x, s, t$ )
9:    $\ell \leftarrow 0$ 
10:   $w \leftarrow \nu_{\max}, N \leftarrow \sqrt{\kappa} \log n \log \frac{n\kappa R}{\bar{\epsilon}r}$ 
11:   $q \leftarrow n^{1/2}(k^2 + m^2)^{-1/2}(d^{\omega-1} + m^{\omega-1})^{1/2}$ 
12:   $\epsilon_{\text{apx},x} \leftarrow \bar{\epsilon}, \zeta_x \leftarrow 2\alpha, \delta_{\text{apx}} \leftarrow \frac{1}{N}$ 
13:   $\epsilon_{\text{apx},s} \leftarrow \bar{\epsilon} \cdot \bar{t}, \zeta_s \leftarrow 3\alpha\bar{t}$ 
14:
      approx.INITIALIZE( $x, s, h, \tilde{h}, H_{w,\bar{x}}^{1/2}\hat{x}, H_{w,\bar{x}}^{-1/2}\hat{s}, \beta_x, \beta_s, \hat{\beta}_x, \hat{\beta}_s, \tilde{\beta}_x, \tilde{\beta}_s, q, \&\text{exact},$ 
       $\epsilon_{\text{apx},x}, \epsilon_{\text{apx},s}, \delta_{\text{apx}}$ )
15:   ▷ Algorithm 16. Parameters from  $x$  to  $\tilde{\beta}_s$  come from exact.  $\&\text{exact}$  is pointer to exact
16: end procedure
17: procedure MULTIPLYANDMOVE( $t \in \mathbb{R}_+$ )
18:    $\ell \leftarrow \ell + 1$ 
19:   if  $|\bar{t} - t| > \bar{t} \cdot \epsilon_t$  or  $\ell > q$  then
20:      $x, s \leftarrow \text{exact.OUTPUT}()$  ▷ Algorithm 15
21:     INITIALIZE( $x, s, t, \bar{\epsilon}$ )
22:   end if
23:    $\beta_x, \beta_s, \hat{\beta}_x, \hat{\beta}_s, \tilde{\beta}_x, \tilde{\beta}_s \leftarrow \text{exact.MOVE}()$  ▷ Algorithm 14
24:    $\delta_{\bar{x}}, \delta_{\bar{s}} \leftarrow \text{approx.MOVEANDQUERY}(\beta_x, \beta_s, \hat{\beta}_x, \hat{\beta}_s, \tilde{\beta}_x, \tilde{\beta}_s)$  ▷ Algorithm 16
25:    $\delta_h, \delta_{\tilde{h}}, \delta_{\tilde{h}}, \delta_{H_{w,\bar{x}}^{1/2}\hat{x}}, \delta_{H_{w,\bar{x}}^{-1/2}\hat{s}} \leftarrow \text{exact.UPDATE}(\delta_{\bar{x}}, \delta_{\bar{s}})$  ▷ Algorithm 15
26:   approx.UPDATE( $\delta_{\bar{x}}, \delta_h, \delta_{\tilde{h}}, \delta_{\tilde{h}}, \delta_{H_{w,\bar{x}}^{1/2}\hat{x}}, \delta_{H_{w,\bar{x}}^{-1/2}\hat{s}}$ ) ▷ Algorithm 16
27: end procedure
28: procedure OUTPUT()
29:   return exact.OUTPUT() ▷ Algorithm 15
30: end procedure
31: end data structure

```

- MULTIPLYANDMOVE($t \in \mathbb{R}_{>0}$): *It implicitly maintains*

$$\begin{aligned}
x &\leftarrow x + \bar{t} B_{w,\bar{x},\bar{t}}^{-1/2} (I - P_{w,\bar{x},\bar{t}}) B_{w,\bar{x},\bar{t}}^{-1/2} \delta_{\mu}(\bar{x}, \bar{s}, \bar{t}) \\
s &\leftarrow s + \bar{t} \delta_{\mu} - \bar{t}^2 H_{w,\bar{x}} B_{w,\bar{x},\bar{t}}^{-1/2} (I - P_{w,\bar{x},\bar{t}}) B_{w,\bar{x},\bar{t}}^{-1/2} \delta_{\mu}(\bar{x}, \bar{s}, \bar{t})
\end{aligned}$$

where $H_{w,\bar{x}}, B_{w,\bar{x},\bar{t}}, P_{w,\bar{x},\bar{t}}$ are defined in Eq. (24)(25)(26) respectively, and \bar{t} is some timestamp satisfying $|\bar{t} - t| \leq \epsilon_t \cdot \bar{t}$.

It also explicitly maintains $(\bar{x}, \bar{s}) \in \mathbb{R}^{n_{\text{tot}} \times n_{\text{tot}}}$ such that $\|\bar{x}_i - x_i\|_{\bar{x}_i} \leq \bar{\epsilon}$ and $\|\bar{s}_i - s_i\|_{\bar{x}_i}^* \leq \bar{t} \bar{\epsilon} w_i$ for all $i \in [n]$ with probability at least 0.9.

Assuming the function is called at most N times and t decreases from t_{\max} to t_{\min} , the total running time is

$$\tilde{O}((Nn^{-1/2} + \log(t_{\max}/t_{\min})) \cdot n(k^{(\omega+1)/2} + m^{(\omega+1)/2})).$$

- OUTPUT: Computes $(x, s) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$ exactly and outputs them in $\tilde{O}(n(k+m))$ time.

2052 E.3 DATA STRUCTURES USED IN CENTRALPATHMAINTENANCE

2053 In this section we present several data structures used in CENTRALPATHMAINTENANCE, including:

- 2054 • EXACTDS (Section E.3.1): This data structure maintains an implicit representation of the
- 2055 primal-dual solution pair (x, s) . This is directly used by CENTRALPATHMAINTENANCE.
- 2056 • APPROXDS (Section E.3.2): This data structure explicitly maintains an approximation (\bar{x}, \bar{s})
- 2057 of (x, s) . This data structure is directly used by CENTRALPATHMAINTENANCE.
- 2058 • BATCHSKETCH (Section E.3.3): This data structure maintains a sketch of (x, s) . This data
- 2059 structure is used by APPROXDS.

2060 E.3.1 EXACTDS

2061 In this section we present the data structure EXACTDS. It maintains an implicit representation of the

2062 primal-dual solution pair (x, s) by maintaining several sparsely-changing vectors (see Eq. (12)(13)).

2063 **Theorem E.3.** *Data structure EXACTDS (Algorithm 14, 15) implicitly maintains the primal-dual*

2064 *pair $(x, s) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$, computable via the expression*

$$2065 x = \hat{x} + H_{w,\bar{x}}^{-1/2} h \beta_x + H_{w,\bar{x}}^{-1/2} \hat{h} \hat{\beta}_x + H_{w,\bar{x}}^{-1/2} \tilde{h} \tilde{\beta}_x, \quad (12)$$

$$2066 s = \hat{s} + H_{w,\bar{x}}^{1/2} h \beta_s + H_{w,\bar{x}}^{1/2} \hat{h} \hat{\beta}_s + H_{w,\bar{x}}^{1/2} \tilde{h} \tilde{\beta}_s, \quad (13)$$

2067 where $\hat{x}, \hat{s} \in \mathbb{R}^{n_{\text{tot}}}$, $h = H_{w,\bar{x}}^{-1/2} \bar{\delta}_\mu \in \mathbb{R}^{n_{\text{tot}}}$, $\hat{h} = H_{w,\bar{x}}^{-1/2} U^\top \in \mathbb{R}^{n_{\text{tot}} \times k}$, $\tilde{h} = H_{w,\bar{x}}^{-1/2} A^\top \in \mathbb{R}^{n_{\text{tot}} \times m}$,

2068 $\beta_x, \beta_s \in \mathbb{R}$, $\hat{\beta}_x, \hat{\beta}_s \in \mathbb{R}^k$, $\tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$.

2069 The data structure supports the following functions:

- 2070 • INITIALIZE($x, s, \bar{x}, \bar{s} \in \mathbb{R}^{n_{\text{tot}}}, \bar{t} \in \mathbb{R}_{>0}$): *Initializes the data structure in $\tilde{O}(n(k^\omega + m^\omega))$*
- 2071 *time, with initial value of the primal-dual pair (x, s) , its initial approximation (\bar{x}, \bar{s}) , and*
- 2072 *initial approximate timestamp \bar{t} .*

- 2073 • MOVE(): *Performs robust central path step*

$$2074 x \leftarrow x + \bar{t} B_{\bar{x}}^{-1} \delta_\mu - \bar{t} B_{\bar{x}}^{-1} A^\top (A B_{\bar{x}}^{-1} A^\top)^{-1} A B_{\bar{x}}^{-1} \delta_\mu, \quad (14)$$

$$2075 s \leftarrow s + \bar{t} \delta_\mu - \bar{t}^2 B_{\bar{x}}^{-1} \delta_\mu + \bar{t}^2 B_{\bar{x}}^{-1} A^\top (A B_{\bar{x}}^{-1} A^\top)^{-1} A B_{\bar{x}}^{-1} \delta_\mu \quad (15)$$

2076 *in $O(k^\omega + m^\omega)$ time by updating its implicit representation.*

- 2077 • UPDATE($\delta_{\bar{x}}, \delta_{\bar{s}} \in \mathbb{R}^{n_{\text{tot}}}$): *Updates the approximation pair (\bar{x}, \bar{s}) to $(\bar{x}^{\text{new}} = \bar{x} + \delta_{\bar{x}} \in$*
- 2078 *$\mathbb{R}^{n_{\text{tot}}}, \bar{s}^{\text{new}} = \bar{s} + \delta_{\bar{s}} \in \mathbb{R}^{n_{\text{tot}}})$ in $\tilde{O}((k^2 + m^2)(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0))$ time, and output the*
- 2079 *changes in variables $h, \hat{h}, \tilde{h}, H_{w,\bar{x}}^{1/2} \hat{x}, H_{w,\bar{x}}^{-1/2} \hat{s}$.*

2080 *Furthermore, $h, H_{w,\bar{x}}^{1/2} \hat{x}, H_{w,\bar{x}}^{-1/2} \hat{s}$ changes in $O(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0)$ coordinates, \hat{h} changes in*

2081 *$O(k(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0))$ coordinates, \tilde{h} changes in $O(m(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0))$ coordinates.*

- 2082 • OUTPUT(): *Output x and s in $\tilde{O}(n(k + m))$ time.*

- 2083 • QUERY $x(i \in [n])$: *Output x_i in $\tilde{O}(k + m)$ time. This function is used by APPROXDS.*

- 2084 • QUERY $s(i \in [n])$: *Output s_i in $\tilde{O}(k + m)$ time. This function is used by APPROXDS.*

2085 *Proof of Theorem E.3.* By combining Lemma E.4 and E.5. □

2086 **Lemma E.4.** EXACTDS correctly maintains an implicit representation of (x, s) , i.e., invariant

$$2087 x = \hat{x} + H_{w,\bar{x}}^{-1/2} h \beta_x + H_{w,\bar{x}}^{-1/2} \hat{h} \hat{\beta}_x + H_{w,\bar{x}}^{-1/2} \tilde{h} \tilde{\beta}_x,$$

$$2088 s = \hat{s} + H_{w,\bar{x}}^{1/2} h \beta_s + H_{w,\bar{x}}^{1/2} \hat{h} \hat{\beta}_s + H_{w,\bar{x}}^{1/2} \tilde{h} \tilde{\beta}_s,$$

$$2089 h = H_{w,\bar{x}}^{-1/2} \bar{\delta}_\mu \in \mathbb{R}^{n_{\text{tot}}}, \hat{h} = H_{w,\bar{x}}^{-1/2} U^\top \in \mathbb{R}^{n_{\text{tot}} \times d}, \tilde{h} = H_{w,\bar{x}}^{-1/2} A^\top \in \mathbb{R}^{n_{\text{tot}} \times m},$$

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159

Algorithm 14 This is used in Algorithm 13.

▷ Theorem E.3

```

1: data structure EXACTDS
2: members
3:    $\bar{x}, \bar{s} \in \mathbb{R}^{n_{\text{tot}}}, \bar{t} \in \mathbb{R}_+, H_{w,\bar{x}} \in \mathbb{R}^{n_{\text{tot}} \times n_{\text{tot}}}$ 
4:    $\hat{x}, \hat{s}, \in \mathbb{R}^{n_{\text{tot}}}, \hat{h} \in \mathbb{R}^{n_{\text{tot}} \times k}, \tilde{h} \in \mathbb{R}^{n_{\text{tot}} \times m}, \beta_x, \beta_s \in \mathbb{R}, \hat{\beta}_x, \hat{\beta}_s \in \mathbb{R}^d, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$ 
5:    $u_1, u_2 \in \mathbb{R}^{k \times m}, u_3 \in \mathbb{R}^{m \times m}, u_4 \in \mathbb{R}^m, u_5 \in \mathbb{R}^d, u_6 \in \mathbb{R}^{k \times k}$ 
6:    $\bar{\alpha} \in \mathbb{R}, \bar{\delta}_\mu \in \mathbb{R}^n$ 
7:    $K \in \mathbb{N}$ 
8: end members
9: procedure INITIALIZE( $x, s, \bar{x}, \bar{s} \in \mathbb{R}^{n_{\text{tot}}}, \bar{t} \in \mathbb{R}_+$ )
10:   $\bar{x} \leftarrow x, \bar{s} \leftarrow s, \bar{t} \leftarrow \bar{t}$ 
11:   $\hat{x} \leftarrow x, \hat{s} \leftarrow s, \beta_x \leftarrow 0, \beta_s \leftarrow 0, \hat{\beta}_x \leftarrow 0, \hat{\beta}_s \leftarrow 0, \tilde{\beta}_x \leftarrow 0, \tilde{\beta}_s \leftarrow 0$ 
12:   $H_{w,\bar{x}} \leftarrow \nabla^2 \phi_w(\bar{x})$ 
13:  INITIALIZE $h(\bar{x}, \bar{s}, H_{w,\bar{x}})$ 
14: end procedure
15: procedure INITIALIZE $h(\bar{x}, \bar{s} \in \mathbb{R}^{n_{\text{tot}}}, H_{w,\bar{x}} \in \mathbb{R}^{n_{\text{tot}} \times n_{\text{tot}}})$ 
16:  for  $i \in [n]$  do
17:     $(\bar{\delta}_\mu)_i \leftarrow -\frac{\alpha \sinh(\frac{\lambda}{w_i} \gamma_i(\bar{x}, \bar{s}, \bar{t}))}{\gamma_i(\bar{x}, \bar{s}, \bar{t})} \cdot \mu_i(\bar{x}, \bar{s}, \bar{t})$ 
18:     $\bar{\alpha} \leftarrow \bar{\alpha} + w_i^{-1} \cosh^2(\frac{\lambda}{w_i} \gamma_i(\bar{x}, \bar{s}, \bar{t}))$ 
19:  end for
20:   $h \leftarrow H_{w,\bar{x}}^{-1/2} \bar{\delta}_\mu, \hat{h} \leftarrow H_{w,\bar{x}}^{-1/2} U^\top, \tilde{h} \leftarrow H_{w,\bar{x}}^{-1/2} A^\top$ 
21:   $u_1 \leftarrow U H_{w,\bar{x}}^{-1} A^\top, u_2 \leftarrow V H_{w,\bar{x}}^{-1} A^\top, u_3 \leftarrow A H_{w,\bar{x}}^{-1} A^\top$ 
22:   $u_4 \leftarrow A H_{w,\bar{x}}^{-1} \bar{\delta}_\mu, u_5 \leftarrow V H_{w,\bar{x}}^{-1} \bar{\delta}_\mu, u_6 \leftarrow V H_{w,\bar{x}}^{-1} U^\top$ 
23: end procedure
24: procedure MOVE()
25:   $v_0 \leftarrow I + \bar{t}^{-1} u_6 \in \mathbb{R}^{k \times k}$ 
26:   $v_1 \leftarrow \bar{t}^{-1} u_3 - \bar{t}^{-2} u_1^\top v_0^{-1} u_2 \in \mathbb{R}^{m \times m}$ 
27:   $v_2 \leftarrow \bar{t}^{-1} u_4 - \bar{t}^{-2} u_1^\top v_0^{-1} u_5 \in \mathbb{R}^m$ 
28:   $\beta_x \leftarrow \beta_x + (\bar{\alpha})^{-1/2}$ 
29:   $\tilde{\beta}_x \leftarrow \tilde{\beta}_x - (\bar{\alpha})^{-1/2} \cdot \bar{t}^{-1} v_0^{-1} u_5 + (\bar{\alpha})^{-1/2} \cdot \bar{t}^{-1} v_0^{-1} u_2 v_1^{-1} v_2$ 
30:   $\hat{\beta}_x \leftarrow \hat{\beta}_x - (\bar{\alpha})^{-1/2} \cdot v_1^{-1} v_2$ 
31:   $\beta_s \leftarrow \beta_s$ 
32:   $\tilde{\beta}_s \leftarrow \tilde{\beta}_s + (\bar{\alpha})^{-1/2} \cdot v_0^{-1} u_5 - (\bar{\alpha})^{-1/2} \cdot v_0^{-1} u_2 v_1^{-1} v_2$ 
33:   $\hat{\beta}_s \leftarrow \hat{\beta}_s + (\bar{\alpha})^{-1/2} \cdot \bar{t} v_1^{-1} v_2$ 
34:  return  $\beta_x, \beta_s, \hat{\beta}_x, \hat{\beta}_s, \tilde{\beta}_x, \tilde{\beta}_s$ 
35: end procedure
36: end data structure

```

$$u_1 = U H_{w,\bar{x}}^{-1} A^\top \in \mathbb{R}^{d \times m}, u_2 = V H_{w,\bar{x}}^{-1} A^\top \in \mathbb{R}^{d \times m}, u_3 = A H_{w,\bar{x}}^{-1} A^\top \in \mathbb{R}^{m \times m},$$

$$u_4 = A H_{w,\bar{x}}^{-1} \bar{\delta}_\mu \in \mathbb{R}^m, u_5 = V H_{w,\bar{x}}^{-1} \bar{\delta}_\mu \in \mathbb{R}^d, u_6 = V H_{w,\bar{x}}^{-1} U^\top \in \mathbb{R}^{d \times d},$$

$$\bar{\alpha} = \sum_{i \in [n]} w_i^{-1} \cosh^2\left(\frac{\lambda}{w_i} \gamma_i(\bar{x}, \bar{s}, \bar{t})\right),$$

$$\bar{\delta}_\mu = \bar{\alpha}^{1/2} \delta_\mu(\bar{x}, \bar{s}, \bar{t})$$

always holds after every external call, and return values of the queries are correct.

Proof. INITIALIZE: By checking the definitions we see that all invariants are satisfied after INITIALIZE.

MOVE: By the invariants, we have

$$v_0 = I + \bar{t}^{-1} V H_{w,\bar{x}}^{-1} U^\top,$$

2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213

Algorithm 15 Algorithm 14 continued.

▷ Theorem E.3

```

1: data structure EXACTDS
2: procedure OUTPUT()
3:   return  $\hat{x} + H_{w,\bar{x}}^{-1/2} h \beta_x + H_{w,\bar{x}}^{-1/2} \hat{h} \hat{\beta}_x + H_{w,\bar{x}}^{-1/2} \tilde{h} \tilde{\beta}_x, \hat{s} + H_{w,\bar{x}}^{1/2} h \beta_s + H_{w,\bar{x}}^{1/2} \hat{h} \hat{\beta}_s + H_{w,\bar{x}}^{1/2} \tilde{h} \tilde{\beta}_s$ 
4: end procedure
5: procedure QUERY  $x(i \in [n])$ 
6:   return  $\hat{x}_i + H_{w,\bar{x}}^{-1/2} h_{i,*} \beta_x + H_{w,\bar{x}}^{-1/2} \hat{h}_{i,*} \hat{\beta}_x + H_{w,\bar{x}}^{-1/2} \tilde{h}_{i,*} \tilde{\beta}_x$ 
7: end procedure
8: procedure QUERY  $s(i \in [n])$ 
9:   return  $\hat{s}_i + H_{w,\bar{x}}^{1/2} h_{i,*} \beta_s + H_{w,\bar{x}}^{1/2} \hat{h}_{i,*} \hat{\beta}_s + H_{w,\bar{x}}^{1/2} \tilde{h}_{i,*} \tilde{\beta}_s$ 
10: end procedure
11: procedure UPDATE( $\delta_{\bar{x}}, \delta_{\bar{s}} \in \mathbb{R}^{n_{\text{tot}}}$ )
12:    $\Delta_{H_{w,\bar{x}}} \leftarrow \nabla^2 \phi_w(\bar{x} + \delta_{\bar{x}}) - H_{w,\bar{x}}$  ▷  $\Delta_{H_{w,\bar{x}}}$  is non-zero only for diagonal blocks  $(i, i)$  for which  $\delta_{\bar{x},i} \neq 0$ 
13:    $S \leftarrow \{i \in [n] \mid \delta_{\bar{x},i} \neq 0 \text{ or } \delta_{\bar{s},i} \neq 0\}$ 
14:    $\delta_{\bar{\delta}_\mu} \leftarrow 0$ 
15:   for  $i \in S$  do
16:     Let  $\gamma_i = \gamma_i(\bar{x}, \bar{s}, \bar{t}), \gamma_i^{\text{new}} = \gamma_i(\bar{x} + \delta_{\bar{x}}, \bar{s} + \delta_{\bar{s}}, \bar{t}), \mu_i^{\text{new}} = \mu_i(\bar{x} + \delta_{\bar{x}}, \bar{s} + \delta_{\bar{s}}, \bar{t})$ 
17:      $\bar{\alpha} \leftarrow \bar{\alpha} - w_i^{-1} \cosh^2(\frac{\lambda}{w_i} \gamma_i) + w_i^{-1} \cosh^2(\frac{\lambda}{w_i} \gamma_i^{\text{new}})$ 
18:      $\delta_{\bar{\delta}_{\mu,i}} \leftarrow -\alpha \sinh(\frac{\lambda}{w_i} \gamma_i^{\text{new}}) \cdot \frac{1}{\gamma_i^{\text{new}}} \cdot \mu_i^{\text{new}} - \delta_{\bar{\delta}_{\mu,i}}$ 
19:   end for
20:    $\delta_h \leftarrow \Delta_{H_{w,\bar{x}}^{-1/2}} (\bar{\delta}_\mu + \delta_{\bar{\delta}_\mu}) + H_{w,\bar{x}}^{-1/2} \delta_{\bar{\delta}_\mu}$ 
21:    $\delta_{\hat{h}} \leftarrow \Delta_{H_{w,\bar{x}}^{-1/2}} U^\top$ 
22:    $\delta_{\tilde{h}} \leftarrow \Delta_{H_{w,\bar{x}}^{-1/2}} A^\top$ 
23:    $\delta_{\hat{x}} \leftarrow -(\delta_h \beta_x + \delta_{\hat{h}} \hat{\beta}_x + \delta_{\tilde{h}} \tilde{\beta}_x)$ 
24:    $\delta_{\hat{s}} \leftarrow -(\delta_h \beta_s + \delta_{\hat{h}} \hat{\beta}_s + \delta_{\tilde{h}} \tilde{\beta}_s)$ 
25:    $h \leftarrow h + \delta_h, \hat{h} \leftarrow \hat{h} + \delta_{\hat{h}}, \tilde{h} \leftarrow \tilde{h} + \delta_{\tilde{h}}, \hat{x} \leftarrow \hat{x} + \delta_{\hat{x}}, \hat{s} \leftarrow \hat{s} + \delta_{\hat{s}}$ 
26:    $u_1 \leftarrow u_1 + U \Delta_{H_{w,\bar{x}}^{-1}} A^\top$ 
27:    $u_2 \leftarrow u_2 + V \Delta_{H_{w,\bar{x}}^{-1}} A^\top$ 
28:    $u_3 \leftarrow u_3 + A \Delta_{H_{w,\bar{x}}^{-1}} A^\top$ 
29:    $u_4 \leftarrow u_4 + A (\Delta_{H_{w,\bar{x}}^{-1}} (\bar{\delta}_\mu + \delta_{\bar{\delta}_\mu}) + H_{w,\bar{x}}^{-1} \delta_{\bar{\delta}_\mu})$ 
30:    $u_5 \leftarrow u_5 + V (\Delta_{H_{w,\bar{x}}^{-1}} (\bar{\delta}_\mu + \delta_{\bar{\delta}_\mu}) + H_{w,\bar{x}}^{-1} \delta_{\bar{\delta}_\mu})$ 
31:    $u_6 \leftarrow u_6 + V \Delta_{H_{w,\bar{x}}^{-1}} U^\top$ 
32:    $\bar{x} \leftarrow \bar{x} + \delta_{\bar{x}}, \bar{s} \leftarrow \bar{s} + \delta_{\bar{s}}$ 
33:    $H_{w,\bar{x}} \leftarrow H_{w,\bar{x}} + \Delta_{H_{w,\bar{x}}}$ 
34:   return  $\delta_h, \delta_{\hat{h}}, \delta_{\tilde{h}}, \delta_{H_{w,\bar{x}}^{1/2} \hat{x}}, \delta_{H_{w,\bar{x}}^{-1/2} \hat{s}}$ 
35: end procedure
36: end data structure

```

$$\begin{aligned}
v_1 &= \bar{t}^{-1} A H_{w,\bar{x}}^{-1} A^\top - \bar{t}^{-1} A H_{w,\bar{x}}^{-1} U^\top (I + \bar{t}^{-1} V H_{w,\bar{x}}^{-1} U^\top)^{-1} V H_{w,\bar{x}} A^\top \\
&= A B_{\bar{x}}^{-1} A^\top \\
v_2 &= \bar{t}^{-1} A H_{w,\bar{x}}^{-1} \bar{\delta}_\mu - \bar{t}^{-1} A H_{w,\bar{x}}^{-1} U^\top (I + \bar{t}^{-1} V H_{w,\bar{x}}^{-1} U^\top)^{-1} V H_{w,\bar{x}} \bar{\delta}_\mu \\
&= A B_{\bar{x}}^{-1} \bar{\delta}_\mu.
\end{aligned}$$

By implicit representation (12),

$$\begin{aligned}
\delta_x &= H_{w,\bar{x}}^{-1/2} h \delta_{\beta_x} + H_{w,\bar{x}}^{-1/2} \hat{h} \delta_{\hat{\beta}_x} + H_{w,\bar{x}}^{-1/2} \tilde{h} \delta_{\tilde{\beta}_x} \\
&= H_{w,\bar{x}}^{-1} \bar{\delta}_\mu \cdot (\bar{\alpha})^{-1/2}
\end{aligned}$$

$$\begin{aligned}
& + H_{w,\bar{x}}^{-1} U^\top \cdot (\bar{\alpha})^{-1/2} \bar{t}^{-1} v_0^{-1} (-u_5 + u_2 v_1^{-1} v_2) \\
& - H_{w,\bar{x}}^{-1} A^\top \cdot (\bar{\alpha})^{-1/2} v_1^{-1} v_2 \\
& = H_{w,\bar{x}}^{-1} \delta_\mu \\
& + H_{w,\bar{x}}^{-1} U^\top \bar{t}^{-1} (I + \bar{t}^{-1} V H_{w,\bar{x}}^{-1} U^\top)^{-1} (-V H_{w,\bar{x}}^{-1} \delta_\mu + V H_{w,\bar{x}}^{-1} A^\top (A B_{\bar{x}}^{-1} A^\top)^{-1} A B_{\bar{x}}^{-1} \delta_\mu) \\
& - H_{w,\bar{x}}^{-1} A^\top (A B_{\bar{x}}^{-1} A^\top)^{-1} A B_{\bar{x}}^{-1} \delta_\mu \\
& = \bar{t} \cdot (\bar{t}^{-1} H_{w,\bar{x}}^{-1} - \bar{t}^{-2} H_{w,\bar{x}}^{-1} U^\top (I + \bar{t}^{-1} V H_{w,\bar{x}}^{-1} U^\top)^{-1} V H_{w,\bar{x}}^{-1}) \delta_\mu \\
& - \bar{t} (\bar{t}^{-1} H_{w,\bar{x}}^{-1} - \bar{t}^{-2} \bar{t}^{-2} H_{w,\bar{x}}^{-1} U^\top (I + \bar{t}^{-1} V H_{w,\bar{x}}^{-1} U^\top)^{-1} V H_{w,\bar{x}}^{-1}) A^\top (A B_{\bar{x}}^{-1} A^\top)^{-1} A B_{\bar{x}}^{-1} \delta_\mu \\
& = \bar{t} B_{\bar{x}}^{-1} \delta_\mu - \bar{t} B_{\bar{x}}^{-1} A^\top (A B_{\bar{x}}^{-1} A^\top)^{-1} A B_{\bar{x}}^{-1} \delta_\mu.
\end{aligned}$$

Comparing with the robust central path step (14), we see that x is updated correctly.

For s , from implicit representation 13 we have

$$\begin{aligned}
\delta_s & = H_{w,\bar{x}}^{1/2} h \delta_{\beta_x} + H_{w,\bar{x}}^{1/2} \hat{h} \delta_{\hat{\beta}_x} + H_{w,\bar{x}}^{1/2} \tilde{h} \delta_{\tilde{\beta}_x} \\
& = -U^\top \cdot (\bar{\alpha})^{-1/2} \cdot v_0^{-1} (-u_5 + u_2 v_1^{-1} v_2) + A^\top \cdot (\bar{\alpha})^{-1/2} \cdot \bar{t} v_1^{-1} v_2 \\
& = \bar{t} \delta_\mu - \bar{t}^2 B_{\bar{x}}^{-1} \delta_\mu + \bar{t}^2 B_{\bar{x}}^{-1} A^\top (A B_{\bar{x}}^{-1} A^\top)^{-1} A B_{\bar{x}}^{-1} \delta_\mu.
\end{aligned}$$

Comparing with robust central path step (15), we see that s is updated correctly.

UPDATE: We would like to prove that UPDATE correctly updates the values of $\hat{x}, \hat{s}, h, \hat{h}, \tilde{h}, u_1, u_2, u_3, u_4, u_5, u_6, \bar{\alpha}, \bar{\delta}_\mu$, while preserving the values of (x, s) . In fact, by checking the definitions, it is easy to see that $h, \hat{h}, \tilde{h}, u_1, u_2, u_3, u_4, u_5, u_6, \bar{\alpha}, \bar{\delta}_\mu$ are updated correctly. Furthermore

$$\begin{aligned}
\delta_x & = \delta_{\hat{x}} + \delta_h \beta_x + \delta_{\hat{h}} \hat{\beta}_x + \delta_{\tilde{h}} \tilde{\beta}_x = 0, \\
\delta_s & = \delta_{\hat{s}} + \delta_h \beta_s + \delta_{\hat{h}} \hat{\beta}_s + \delta_{\tilde{h}} \tilde{\beta}_s = 0.
\end{aligned}$$

So values of (x, s) are preserved. \square

Lemma E.5. *We bound the running time of EXACTDS as following.*

- (i) EXACTDS.INITIALIZE (Algorithm 14) runs in $\tilde{O}(n(k^{\omega-1} + m^{\omega-1}))$ time.
- (ii) EXACTDS.MOVE (Algorithm 14) runs in $\tilde{O}(k^\omega + m^\omega)$ time.
- (iii) EXACTDS.OUTPUT (Algorithm 15) runs in $\tilde{O}(n(k + m))$ time and correctly outputs (x, s) .
- (iv) EXACTDS.QUERY x and EXACTDS.QUERY s (Algorithm 15) runs in $\tilde{O}(k + m)$ time and returns the correct answer.
- (v) EXACTDS.UPDATE (Algorithm 15) runs in $\tilde{O}((k^2 + m^2)(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0))$ time. Furthermore, $\|\delta_h\|_0, \|\delta_{\hat{x}}\|_0, \|\delta_{\hat{s}}\|_0 = O(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0)$, $\text{nnz}(\hat{h}) = O(d(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0))$, $\text{nnz}(\tilde{h}) = O(m(\|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0))$.

Proof. (i) EXACTDS.INITIALIZE: Computing u_1 and u_2 takes $\mathcal{T}_{\text{mat}}(k, n, m) = \tilde{O}(n(k^{\omega-1} + m^{\omega-1}))$ time. Computing u_3 takes $\mathcal{T}_{\text{mat}}(m, n, m) = \tilde{O}(nm^{\omega-1})$ time. Computing u_4 takes $O(nm)$ time. Computing u_5 takes $O(nk)$ time. Computing u_6 takes $\mathcal{T}_{\text{mat}}(k, n, k) = \tilde{O}(nk^{\omega-1})$ time. All other computations are cheaper.

(ii) EXACTDS.MOVE: Computing v_0^{-1} takes $\tilde{O}(k^\omega)$ time. Computing v_1^{-1} takes $\tilde{O}(m^\omega)$ time. All other computations are cheaper.

(iii) EXACTDS.OUTPUT: Takes $\tilde{O}(n(k + m))$ time.

(iv) EXACTDS.QUERY x and EXACTDS.QUERY s : Takes $\tilde{O}(k + m)$ time.

2268 (v) EXACTDS.UPDATE: For simplicity, write $t = \|\delta_{\bar{x}}\|_0 + \|\delta_{\bar{s}}\|_0$. Computing δ_h takes $\tilde{O}(t)$
2269 time. Computing $\delta_{\hat{h}}$ takes $\tilde{O}(tk)$ time. Computing $\delta_{\tilde{h}}$ takes $\tilde{O}(tm)$ time. Computing $\delta_{\bar{x}}$ and
2270 $\delta_{\bar{s}}$ takes $\tilde{O}(t(k+m))$ time. The sparsity statements follow directly. Computing u_1 and u_2
2271 takes $\tilde{O}(tkm)$ time. Computing u_3 takes $\tilde{O}(tm^2)$ time. Computing u_4 takes $\tilde{O}(tm)$ time.
2272 Computing u_5 takes $\tilde{O}(tk)$ time. Computing u_6 takes $\tilde{O}(tk^2)$ time. \square
2273
2274

2275 E.3.2 APPROXDS

2276 In this section we present the data structure APPROXDS. Given BATCHSKETCH, a data structure
2277 maintaining a sketch of the primal-dual pair $(x, s) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$, APPROXDS maintains a
2278 sparsely-changing ℓ_∞ -approximation of (x, s) .
2279

2280 **Algorithm 16** This is used in Algorithm 13.

2281 1: **data structure** APPROXDS ▷ Theorem E.6
2282 2: **private : members**
2283 3: $\epsilon_{\text{apx},x}, \epsilon_{\text{apx},s} \in \mathbb{R}$
2284 4: $\ell \in \mathbb{N}$
2285 5: BATCHSKETCH bs ▷ This maintains a sketch of $H_{w,\bar{x}}^{1/2}x$ and $H_{w,\bar{x}}^{-1/2}s$. See Algorithm 17 and 18.
2286 6: EXACTDS* exact ▷ This is a pointer to the EXACTDS (Algorithm 14, 15) we maintain in parallel to
2287 APPROXDS.
2288 7: $\tilde{x}, \tilde{s} \in \mathbb{R}^{n_{\text{tot}}}$ ▷ (\tilde{x}, \tilde{s}) is a sparsely-changing approximation of (x, s) . They have the same value as
2289 (\bar{x}, \bar{s}) , but for these local variables we use (\tilde{x}, \tilde{s}) to avoid confusion.
2290 8: **end members**
2291 9: **procedure** INITIALIZE($x, s \in \mathbb{R}^{n_{\text{tot}}}, h \in \mathbb{R}^{n_{\text{tot}}}, \hat{h} \in \mathbb{R}^{n_{\text{tot}} \times k}, \tilde{h} \in \mathbb{R}^{n_{\text{tot}} \times m}, H_{w,\bar{x}}^{1/2}\hat{x}, H_{w,\bar{x}}^{-1/2}\hat{s} \in$
2292 $\mathbb{R}^{n_{\text{tot}}}, \beta_x, \beta_s \in \mathbb{R}, \hat{\beta}_x, \hat{\beta}_s \in \mathbb{R}^d, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m, q \in \mathbb{N}, \text{EXACTDS* exact}, \epsilon_{\text{apx},x}, \epsilon_{\text{apx},s}, \delta_{\text{apx}} \in \mathbb{R}$)
2293 10: $\ell \leftarrow 0, q \leftarrow q$
2294 11: $\epsilon_{\text{apx},x} \leftarrow \epsilon_{\text{apx},x}, \epsilon_{\text{apx},s} \leftarrow \epsilon_{\text{apx},s}$
2295 12: bs.INITIALIZE($x, h, \hat{h}, \tilde{h}, H_{w,\bar{x}}^{1/2}\hat{x}, H_{w,\bar{x}}^{-1/2}\hat{s}, \beta_x, \beta_s, \hat{\beta}_x, \hat{\beta}_s, \tilde{\beta}_x, \tilde{\beta}_s, \delta_{\text{apx}}/q$) ▷ Algorithm 17
2296 13: $\tilde{x} \leftarrow x, \tilde{s} \leftarrow s$
2297 14: exact \leftarrow exact
2298 15: **end procedure**
2299 16: **procedure** UPDATE($\delta_{\bar{x}} \in \mathbb{R}^{n_{\text{tot}}}, \delta_h \in \mathbb{R}^{n_{\text{tot}}}, \delta_{\hat{h}} \in \mathbb{R}^{n_{\text{tot}} \times k}, \delta_{\tilde{h}} \in \mathbb{R}^{n_{\text{tot}} \times m}, \delta_{H_{w,\bar{x}}^{1/2}\hat{x}}, \delta_{H_{w,\bar{x}}^{-1/2}\hat{s}} \in \mathbb{R}^{n_{\text{tot}}}$)
2300 17: bs.UPDATE($\delta_{\bar{x}}, \delta_h, \delta_{\hat{h}}, \delta_{\tilde{h}}, \delta_{H_{w,\bar{x}}^{1/2}\hat{x}}, \delta_{H_{w,\bar{x}}^{-1/2}\hat{s}}$) ▷ Algorithm 17
2301 18: $\ell \leftarrow \ell + 1$
2302 19: **end procedure**
2303 20: **procedure** MOVEANDQUERY($\beta_x, \beta_s \in \mathbb{R}, \hat{\beta}_x, \hat{\beta}_s \in \mathbb{R}^d, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$)
2304 21: bs.MOVE($\beta_x, \beta_s, \hat{\beta}_x, \hat{\beta}_s, \tilde{\beta}_x, \tilde{\beta}_s$) ▷ Algorithm 17. Do not update ℓ yet
2305 22: $\delta_{\bar{x}} \leftarrow \text{QUERY}x(\epsilon_{\text{apx},x}/(2 \log q + 1))$ ▷ Algorithm 16
2306 23: $\delta_{\bar{s}} \leftarrow \text{QUERY}s(\epsilon_{\text{apx},s}/(2 \log q + 1))$ ▷ Algorithm 16
2307 24: $\tilde{x} \leftarrow \tilde{x} + \delta_{\bar{x}}, \tilde{s} \leftarrow \tilde{s} + \delta_{\bar{s}}$
2308 25: **return** $(\delta_{\bar{x}}, \delta_{\bar{s}})$
2309 26: **end procedure**
2310 27: **procedure** QUERY($x \in \mathbb{R}$)
2311 28: Same as Algorithm 5, QUERY x .
2312 29: **end procedure**
2313 30: **procedure** QUERY($s \in \mathbb{R}$)
2314 31: Same as Algorithm 5, QUERY s .
2315 32: **end procedure**
2316 33: **end data structure**

2317 **Theorem E.6.** Given parameters $\epsilon_{\text{apx},x}, \epsilon_{\text{apx},s} \in (0, 1), \delta_{\text{apx}} \in (0, 1), \zeta_x, \zeta_s \in \mathbb{R}$ such that

$$2318 \quad \left\| H_{w,\bar{x}^{(\ell)}}^{1/2}x^{(\ell)} - H_{w,\bar{x}^{(\ell+1)}}^{1/2}x^{(\ell+1)} \right\|_2 \leq \zeta_x, \quad \left\| H_{w,\bar{x}^{(\ell)}}^{-1/2}s^{(\ell)} - H_{w,\bar{x}^{(\ell+1)}}^{-1/2}s^{(\ell+1)} \right\|_2 \leq \zeta_s$$

2319 for all $\ell \in \{0, \dots, q-1\}$, data structure APPROXDS (Algorithm 16) supports the following
2320 operations:

- 2321 • INITIALIZE($x, s \in \mathbb{R}^{n_{\text{tot}}}, h \in \mathbb{R}^{n_{\text{tot}}}, \hat{h} \in \mathbb{R}^{n_{\text{tot}} \times k}, \tilde{h} \in$
 $\mathbb{R}^{n_{\text{tot}} \times m}, H_{w,\bar{x}}^{1/2}\hat{x}, H_{w,\bar{x}}^{-1/2}\hat{s} \in \mathbb{R}^{n_{\text{tot}}}, \beta_x, \beta_s \in \mathbb{R}, \hat{\beta}_x, \hat{\beta}_s \in \mathbb{R}^d, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m, q \in$

2322 \mathbb{N} , EXACTDS* exact, $\epsilon_{\text{apx},x}$, $\epsilon_{\text{apx},s}$, $\delta_{\text{apx}} \in \mathbb{R}$: Initialize the data structure in
 2323 $\tilde{O}(n(k+m))$ time.

2325 • MOVEANDQUERY($\beta_x, \beta_s \in \mathbb{R}, \hat{\beta}_x, \hat{\beta}_s \in \mathbb{R}^d, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$): Update values of
 2326 $\beta_x, \beta_s, \hat{\beta}_x, \hat{\beta}_s, \tilde{\beta}_x, \tilde{\beta}_s$ by calling BATCHSKETCH.MOVE. This effectively moves $(x^{(\ell)}, s^{(\ell)})$
 2327 to $(x^{(\ell+1)}, s^{(\ell+1)})$ while keeping $\bar{x}^{(\ell)}$ unchanged.

2328 Then return two sets $L_x^{(\ell)}, L_s^{(\ell)} \subset [n]$ where

2329
$$L_x^{(\ell)} \supseteq \{i \in [n] : \|H_{w, \bar{x}^{(\ell)}}^{1/2} x_i^{(\ell)} - H_{w, \bar{x}^{(\ell)}}^{1/2} x_i^{(\ell+1)}\|_2 \geq \epsilon_{\text{apx},x}\},$$

2330
$$L_s^{(\ell)} \supseteq \{i \in [n] : \|H_{w, \bar{x}^{(\ell)}}^{-1/2} s_i^{(\ell)} - H_{w, \bar{x}^{(\ell)}}^{-1/2} s_i^{(\ell+1)}\|_2 \geq \epsilon_{\text{apx},s}\},$$

2331 satisfying

2332
$$\sum_{0 \leq \ell \leq q-1} |L_x^{(\ell)}| = \tilde{O}(\epsilon_{\text{apx},x}^{-2} \zeta_x^2 q^2),$$

2333
$$\sum_{0 \leq \ell \leq q-1} |L_s^{(\ell)}| = \tilde{O}(\epsilon_{\text{apx},s}^{-2} \zeta_s^2 q^2).$$

2334 For every query, with probability at least $1 - \delta_{\text{apx}}/q$, the return values are correct.

2335 Furthermore, total time cost over all queries is at most

2336
$$\tilde{O}((\epsilon_{\text{apx},x}^{-2} \zeta_x^2 + \epsilon_{\text{apx},s}^{-2} \zeta_s^2) q^2 (k+m)).$$

2337 • UPDATE($\delta_{\bar{x}} \in \mathbb{R}^{n_{\text{tot}}}, \delta_h \in \mathbb{R}^{n_{\text{tot}}}, \delta_{\hat{h}} \in \mathbb{R}^{n_{\text{tot}} \times d}, \delta_{\tilde{h}} \in \mathbb{R}^{n_{\text{tot}} \times m}, \delta_{H_{w, \bar{x}}^{1/2} \hat{x}}, \delta_{H_{w, \bar{x}}^{-1/2} \hat{s}} \in \mathbb{R}^{n_{\text{tot}}}$):

2338 Update sketches of $H_{w, \bar{x}^{(\ell)}}^{1/2} x^{(\ell+1)}$ and $H_{w, \bar{x}^{(\ell)}}^{-1/2} s^{(\ell+1)}$ by calling BATCHSKETCH.UPDATE.
 2339 This effectively moves $\bar{x}^{(\ell)}$ to $\bar{x}^{(\ell+1)}$ while keeping $(x^{(\ell+1)}, s^{(\ell+1)})$ unchanged. Then ad-
 2340 vance timestamp ℓ .

2341 Each update costs

2342
$$\tilde{O}(\|\delta_h\|_0 + \text{nnz}(\delta_{\hat{h}}) + \text{nnz}(\delta_{\tilde{h}}) + \|H_{w, \bar{x}}^{1/2} \hat{x}\|_0 + \|H_{w, \bar{x}}^{-1/2} \hat{s}\|_0)$$

2343 time.

2344 *Proof.* The proof is essentially the same as proof of (Gu & Song, 2022, Theorem 4.18). For the
 2345 running time claims, we plug in Theorem E.7 when necessary. \square

2346 E.3.3 BATCHSKETCH

2347 In this section we present the data structure BATCHSKETCH. It maintains a sketch of $H_{\bar{x}}^{1/2} x$ and
 2348 $H_{\bar{x}}^{-1/2} s$. It is a variation of BATCHSKETCH in Gu & Song (2022).

2349 **Theorem E.7.** Data structure BATCHSKETCH (Algorithm 17, 18) supports the following operations:

2350 • INITIALIZE($\bar{x} \in \mathbb{R}^{n_{\text{tot}}}, h \in \mathbb{R}^{n_{\text{tot}}}, \hat{h} \in \mathbb{R}^{n_{\text{tot}} \times k}, \tilde{h} \in \mathbb{R}^{n_{\text{tot}} \times m}, H_{w, \bar{x}}^{1/2} \hat{x}, H_{w, \bar{x}}^{-1/2} \hat{s} \in$
 2351 $\mathbb{R}^{n_{\text{tot}}}, \beta_x, \beta_s \in \mathbb{R}, \hat{\beta}_x, \hat{\beta}_s \in \mathbb{R}^k, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m, \delta_{\text{apx}} \in \mathbb{R}$): Initialize the data structure
 2352 in $\tilde{O}(n(k+m))$ time.

2353 • MOVE($\beta_x, \beta_s \in \mathbb{R}, \hat{\beta}_x, \hat{\beta}_s \in \mathbb{R}^k, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$): Update values of $\beta_x, \beta_s, \hat{\beta}_x, \hat{\beta}_s, \tilde{\beta}_x, \tilde{\beta}_s$ in
 2354 $O(k+m)$ time. This effectively moves $(x^{(\ell)}, s^{(\ell)})$ to $(x^{(\ell+1)}, s^{(\ell+1)})$ while keeping $\bar{x}^{(\ell)}$
 2355 unchanged.

2356 • UPDATE($\delta_{\bar{x}} \in \mathbb{R}^{n_{\text{tot}}}, \delta_h \in \mathbb{R}^{n_{\text{tot}}}, \delta_{\hat{h}} \in \mathbb{R}^{n_{\text{tot}} \times k}, \delta_{\tilde{h}} \in \mathbb{R}^{n_{\text{tot}} \times m}, \delta_{H_{w, \bar{x}}^{1/2} \hat{x}}, \delta_{H_{w, \bar{x}}^{-1/2} \hat{s}} \in \mathbb{R}^{n_{\text{tot}}}$):
 2357 Update sketches of $H_{w, \bar{x}^{(\ell)}}^{1/2} x^{(\ell+1)}$ and $H_{w, \bar{x}^{(\ell)}}^{-1/2} s^{(\ell+1)}$. This effectively moves $\bar{x}^{(\ell)}$ to $\bar{x}^{(\ell+1)}$
 2358 while keeping $(x^{(\ell+1)}, s^{(\ell+1)})$ unchanged. Then advance timestamp ℓ .

2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429

Algorithm 17 This is used by Algorithm 16.

1: **data structure** BATCHSKETCH ▷ Theorem E.7
2: **members**
3: $\Phi \in \mathbb{R}^{r \times n_{\text{tot}}}$ ▷ All sketches need to share the same sketching matrix
4: \mathcal{S}, χ partition tree
5: $\ell \in \mathbb{N}$ ▷ Current timestamp
6: **VECTORSKETCH** $\text{sketch}H_{w,\bar{x}}^{1/2}\hat{x}, \text{sketch}H_{w,\bar{x}}^{-1/2}\hat{s}, \text{sketch}h, \text{sketch}\hat{h}, \text{sketch}\tilde{h}$ ▷ Algorithm 9
7: $\beta_x, \beta_s \in \mathbb{R}, \hat{\beta}_x, \hat{\beta}_s \in \mathbb{R}^d, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$
8: $(\text{history}[t])_{t \geq 0}$ ▷ Snapshot of data at timestamp t . See Remark D.9.
9: **end members**
10: **procedure** INITIALIZE($\bar{x} \in \mathbb{R}^{n_{\text{tot}}}, h \in \mathbb{R}^{n_{\text{tot}}}, \hat{h} \in \mathbb{R}^{n_{\text{tot}} \times k}, \tilde{h} \in \mathbb{R}^{n_{\text{tot}} \times m}, H_{w,\bar{x}}^{1/2}\hat{x}, H_{w,\bar{x}}^{-1/2}\hat{s} \in \mathbb{R}^{n_{\text{tot}}}, \beta_x, \beta_s \in \mathbb{R}, \hat{\beta}_x, \hat{\beta}_s \in \mathbb{R}^d, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m, \delta_{\text{apx}} \in \mathbb{R}$)
11: Construct any valid partition tree (\mathcal{S}, χ)
12: $r \leftarrow \Theta(\log^3(n_{\text{tot}}) \log(1/\delta_{\text{apx}}))$
13: Initialize $\Phi \in \mathbb{R}^{r \times n_{\text{tot}}}$ with iid $\mathcal{N}(0, \frac{1}{r})$
14: $\beta_x \leftarrow \beta_x, \beta_s \leftarrow \beta_s, \hat{\beta}_x \leftarrow \hat{\beta}_x, \hat{\beta}_s \leftarrow \hat{\beta}_s, \tilde{\beta}_x \leftarrow \tilde{\beta}_x, \tilde{\beta}_s \leftarrow \tilde{\beta}_s$
15: $\text{sketch}H_{w,\bar{x}}^{1/2}\hat{x}$.INITIALIZE($\mathcal{S}, \chi, \Phi, H_{w,\bar{x}}^{1/2}\hat{x}$) ▷ Algorithm 9
16: $\text{sketch}H_{w,\bar{x}}^{-1/2}\hat{s}$.INITIALIZE($\mathcal{S}, \chi, \Phi, H_{w,\bar{x}}^{-1/2}\hat{s}$) ▷ Algorithm 9
17: $\text{sketch}h$.INITIALIZE($\mathcal{S}, \chi, \Phi, h$) ▷ Algorithm 9
18: $\text{sketch}\hat{h}$.INITIALIZE($\mathcal{S}, \chi, \Phi, \hat{h}$) ▷ Algorithm 9. Here we construct one sketch for $\hat{h}_{*,i}$ for every $i \in [k]$.
19: $\text{sketch}\tilde{h}$.INITIALIZE($\mathcal{S}, \chi, \Phi, \tilde{h}$) ▷ Algorithm 9. Here we construct one sketch for $\tilde{h}_{*,i}$ for every $i \in [m]$.
20: $\ell \leftarrow 0$
21: Make snapshot history $[\ell]$ ▷ Remark D.9
22: **end procedure**
23: **procedure** MOVE($\beta_x, \beta_s \in \mathbb{R}, \hat{\beta}_x, \hat{\beta}_s \in \mathbb{R}^k, \tilde{\beta}_x, \tilde{\beta}_s \in \mathbb{R}^m$)
24: $\beta_x \leftarrow \beta_x, \beta_s \leftarrow \beta_s, \hat{\beta}_x \leftarrow \hat{\beta}_x, \hat{\beta}_s \leftarrow \hat{\beta}_s, \tilde{\beta}_x \leftarrow \tilde{\beta}_x, \tilde{\beta}_s \leftarrow \tilde{\beta}_s$ ▷ Do not update ℓ yet
25: **end procedure**
26: **procedure** UPDATE($\delta_{\bar{x}} \in \mathbb{R}^{n_{\text{tot}}}, \delta_h \in \mathbb{R}^{n_{\text{tot}}}, \delta_{\hat{h}} \in \mathbb{R}^{n_{\text{tot}} \times k}, \delta_{\tilde{h}} \in \mathbb{R}^{n_{\text{tot}} \times m}, \delta_{H_{w,\bar{x}}^{1/2}\hat{x}}, \delta_{H_{w,\bar{x}}^{-1/2}\hat{s}} \in \mathbb{R}^{n_{\text{tot}}}$)
27: $\text{sketch}H_{w,\bar{x}}^{1/2}\hat{x}$.UPDATE($\delta_{H_{w,\bar{x}}^{1/2}\hat{x}}$) ▷ Algorithm 9
28: $\text{sketch}H_{w,\bar{x}}^{-1/2}\hat{s}$.UPDATE($\delta_{H_{w,\bar{x}}^{-1/2}\hat{s}}$) ▷ Algorithm 9
29: $\text{sketch}h$.UPDATE(δ_h) ▷ Algorithm 9
30: $\text{sketch}\hat{h}$.UPDATE($\delta_{\hat{h}}$) ▷ Algorithm 9
31: $\text{sketch}\tilde{h}$.UPDATE($\delta_{\tilde{h}}$) ▷ Algorithm 9
32: $\ell \leftarrow \ell + 1$
33: Make snapshot history $[\ell]$ ▷ Remark D.9
34: **end procedure**
35: **end data structure**

Each update costs

$$\tilde{O}(\|\delta_h\|_0 + \text{nnz}(\delta_{\hat{h}}) + \text{nnz}(\delta_{\tilde{h}}) + \|H_{w,\bar{x}}^{1/2}\hat{x}\|_0 + \|H_{w,\bar{x}}^{-1/2}\hat{s}\|_0).$$

• QUERY($x(\ell' \in \mathbb{N}, \epsilon \in \mathbb{R})$): Given timestamp ℓ' , return a set $S \subseteq [n]$ where

$$S \supseteq \{i \in [n] : \|H_{w,\bar{x}(\ell')}^{1/2}x_i^{(\ell')} - H_{w,\bar{x}(\ell)}^{1/2}x_i^{(\ell+1)}\|_2 \geq \epsilon\},$$

and

$$|S| = O(\epsilon^{-2}(\ell - \ell' + 1) \sum_{\ell' \leq t \leq \ell} \|H_{w,\bar{x}^{(t)}}^{1/2}x^{(t)} - H_{w,\bar{x}^{(t)}}^{1/2}x^{(t+1)}\|_2^2 + \sum_{\ell' \leq t \leq \ell-1} \|\bar{x}^{(t)} - \bar{x}^{(t+1)}\|_{2,0})$$

2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451

Algorithm 18 BATCHSKETCH Algorithm 17 continued. This is used by Algorithm 16.

1: **data structure** BATCHSKETCH ▷ Theorem E.7
2: **private:**
3: **procedure** QUERY_xSKETCH($v \in \mathcal{S}$) ▷ Return the value of $\Phi_{\chi(v)}(H_{w,\bar{x}}^{1/2}x)_{\chi(v)}$
4: **return** sketch $H_{w,\bar{x}}^{1/2}\hat{x}$.QUERY(v) + sketch h .QUERY(v) · β_x + sketch \hat{h} .QUERY(v) · $\hat{\beta}_x$ +
sketch \tilde{h} .QUERY(v) · $\tilde{\beta}_x$ ▷ Algorithm 9
5: **end procedure**
6: **procedure** QUERY_sSKETCH($v \in \mathcal{S}$) ▷ Return the value of $\Phi_{\chi(v)}(H_{w,\bar{x}}^{-1/2}s)_{\chi(v)}$
7: **return** sketch $H_{w,\bar{x}}^{-1/2}\hat{s}$.QUERY(v) + sketch h .QUERY(v) · β_s + sketch \hat{h} .QUERY(v) · $\hat{\beta}_s$ +
sketch \tilde{h} .QUERY(v) · $\tilde{\beta}_s$ ▷ Algorithm 9
8: **end procedure**
9: **public:**
10: **procedure** QUERY_x($\ell' \in \mathbb{N}, \epsilon \in \mathbb{R}$)
11: Same as Algorithm 7, QUERY_x, using QUERY_xSKETCH defined here instead of the one in
Algorithm 7.
12: **end procedure**
13: **procedure** QUERY_s($\ell' \in \mathbb{N}, \epsilon \in \mathbb{R}$)
14: Same as Algorithm 7, QUERY_s, using QUERY_sSKETCH defined here instead of the one in
Algorithm 7.
15: **end procedure**
16: **end structure**

2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476

where ℓ is the current timestamp.

For every query, with probability at least $1 - \delta$, the return values are correct, and costs at most

$$\tilde{O}((k+m) \cdot (\epsilon^{-2}(\ell - \ell' + 1) \sum_{\ell' \leq t \leq \ell} \|H_{\bar{x}(t)}^{1/2}x^{(t)} - H_{\bar{x}(t)}^{1/2}x^{(t+1)}\|_2^2 + \sum_{\ell' \leq t \leq \ell-1} \|\bar{x}^{(t)} - \bar{x}^{(t+1)}\|_{2,0}))$$

running time.

- QUERY_s($\ell' \in \mathbb{N}, \epsilon \in \mathbb{R}$): Given timestamp ℓ' , return a set $S \subseteq [n]$ where

$$S \supseteq \{i \in [n] : \|H_{w,\bar{x}(\ell')}^{-1/2}s_i^{(\ell')} - H_{w,\bar{x}(\ell)}^{-1/2}s_i^{(\ell+1)}\|_2 \geq \epsilon\}$$

and

$$|S| = O(\epsilon^{-2}(\ell - \ell' + 1) \sum_{\ell' \leq t \leq \ell} \|H_{w,\bar{x}(t)}^{-1/2}s^{(t)} - H_{w,\bar{x}(t)}^{-1/2}s^{(t+1)}\|_2^2 + \sum_{\ell' \leq t \leq \ell-1} \|\bar{x}^{(t)} - \bar{x}^{(t+1)}\|_{2,0})$$

where ℓ is the current timestamp.

For every query, with probability at least $1 - \delta$, the return values are correct, and costs at most

$$\tilde{O}((k+m) \cdot (\epsilon^{-2}(\ell - \ell' + 1) \sum_{\ell' \leq t \leq \ell} \|H_{\bar{x}(t)}^{1/2}s^{(t)} - H_{\bar{x}(t)}^{1/2}s^{(t+1)}\|_2^2 + \sum_{\ell' \leq t \leq \ell-1} \|\bar{x}^{(t)} - \bar{x}^{(t+1)}\|_{2,0}))$$

running time.

2477
2478

Proof. The proof is essentially the same as proof of (Gu & Song, 2022, Theorem 4.21). □

2479
2480

E.4 ANALYSIS OF CENTRALPATHMAINTENANCE

2481
2482
2483

Lemma E.8 (Correctness of CENTRALPATHMAINTENANCE). *Algorithm 13 implicitly maintains the primal-dual solution pair (x, s) via representation Eq. (12)(13). It also explicitly maintains $(\bar{x}, \bar{s}) \in \mathbb{R}^{n_{\text{tot}}} \times \mathbb{R}^{n_{\text{tot}}}$ such that $\|\bar{x}_i - x_i\|_{\bar{x}_i} \leq \bar{\epsilon}$ and $\|\bar{s}_i - s_i\|_{\bar{x}_i}^* \leq t\bar{\epsilon}w_i$ for all $i \in [n]$ with probability at least 0.9.*

2484 *Proof.* Same as proof of Lemma D.13. □
 2485

2486 **Lemma E.9.** *We bound the running time of CENTRALPATHMAINTENANCE as following.*
 2487

- 2488 • CENTRALPATHMAINTENANCE.INITIALIZE takes $\tilde{O}(n(k^{\omega-1} + m^{\omega-1}))$ time.
- 2489
- 2490 • If CENTRALPATHMAINTENANCE.MULTIPLYANDMOVE is called N times, then it has
- 2491 total running time

$$2492 \quad \tilde{O}((Nn^{-1/2} + \log(t_{\max}/t_{\min})) \cdot n(k+m)^{(\omega+1)/2}).$$

- 2493
- 2494
- 2495 • CENTRALPATHMAINTENANCE.OUTPUT takes $\tilde{O}(n(k+m))$ time.
- 2496

2497

2498 *Proof.* INITIALIZE part: By Theorem E.3 and E.6.

2499 OUTPUT part: By Theorem E.3.

2500 MULTIPLYANDMOVE part: Between two restarts, the total size of $|L_x|$ returned by approx.QUERY
 2501 is bounded by $\tilde{O}(q^2 \zeta_x^2 / \epsilon_{\text{apx},x}^2)$ by Theorem E.6. By plugging in $\zeta_x = 2\alpha$, $\epsilon_{\text{apx},x} = \bar{\epsilon}$, we have
 2502 $\sum_{\ell \in [q]} |L_x^{(\ell)}| = \tilde{O}(q^2)$. Similarly, for s we have $\sum_{\ell \in [q]} |L_s^{(\ell)}| = \tilde{O}(q^2)$.

2503 **Update time:** By Theorem E.3 and E.6, in a sequence of q updates, total cost for update is $\tilde{O}(q^2(k^2 +$
 2504 $m^2))$. So the amortized update cost per iteration is $\tilde{O}(q(k^2 + m^2))$. The total update cost is

$$2505 \quad \text{number of iterations} \cdot \text{time per iteration} = \tilde{O}(Nq(k^2 + m^2)).$$

2506

2507 **Init/restart time:** We restart the data structure whenever $K > q$ or $|\bar{t} - t| > \bar{t}\epsilon_t$, so there are
 2508 $O(N/q + \log(t_{\max}/t_{\min})\epsilon_t^{-1})$ restarts in total. By Theorem E.3 and E.6, time cost per restart is
 2509 $\tilde{O}(n(k^{\omega-1} + m^{\omega-1}))$. So the total initialization time is

$$2510 \quad \text{number of restarts} \cdot \text{time per restart} = \tilde{O}((N/q + \log(t_{\max}/t_{\min})\epsilon_t^{-1}) \cdot n(k^{\omega-1} + m^{\omega-1})).$$

2511 **Combine everything:** Overall running time is

$$2512 \quad \tilde{O}(Nq(k^2 + m^2) + (N/q + \log(t_{\max}/t_{\min})\epsilon_t^{-1}) \cdot n(k^{\omega-1} + m^{\omega-1})).$$

2513 Taking $\epsilon_t = \frac{1}{2}\bar{\epsilon}$, the optimal choice for q is

$$2514 \quad q = n^{1/2}(k^2 + m^2)^{-1/2}(k^{\omega-1} + m^{\omega-1})^{1/2},$$

2515 achieving overall running time

$$2516 \quad \tilde{O}((Nn^{-1/2} + \log(t_{\max}/t_{\min})) \cdot n(k^2 + m^2)^{1/2}(k^{\omega-1} + m^{\omega-1})^{1/2})$$

$$2517 \quad = \tilde{O}((Nn^{-1/2} + \log(t_{\max}/t_{\min})) \cdot n(k+m)^{(\omega+1)/2}). \quad \square$$

2518 *Proof of Theorem E.2.* Combining Lemma E.8 and E.9. □
 2519

2520 E.5 PROOF OF MAIN STATEMENT

2521

2522 *Proof of Theorem E.1.* Use CENTRALPATHMAINTENANCE (Algorithm 13) as the maintenance data
 2523 structure in Algorithm 20. Combining Theorem E.2 and Theorem F.1 finishes the proof. □
 2524

2538 F ROBUST IPM ANALYSIS

2539
2540 In this section we present a robust IPM algorithm for quadratic programming. The algorithm is a
2541 modification of previous robust IPM algorithms for linear programming Lee et al. (2019); Lee &
2542 Vempala (2021).

2543 Convention: Variables are in n blocks of dimension n_i ($i \in [n]$). Total dimension is $n_{\text{tot}} = \sum_{i \in [n]} n_i$.
2544 We write $x = (x_1, \dots, x_n) \in \mathbb{R}^{n_{\text{tot}}}$ where $x_i \in \mathbb{R}^{n_i}$. We consider programs of the following form:
2545

$$2546 \min_{x \in \mathbb{R}^n} \frac{1}{2} x^\top Q x + c^\top x \quad (16)$$

$$2547 \text{s.t. } Ax = b$$

$$2548 x_i \in \mathcal{K}_i \quad \forall i \in [n]$$

2551 where $Q \in \mathcal{S}^{n_{\text{tot}}}$, $c \in \mathbb{R}^{n_{\text{tot}}}$, $A \in \mathbb{R}^{m \times n_{\text{tot}}}$, $b \in \mathbb{R}^m$, $\mathcal{K}_i \subset \mathbb{R}^{n_i}$ is a convex set. Let $\mathcal{K} = \prod_{i \in [n]} \mathcal{K}_i$.

2552 **Theorem F.1.** Consider the convex program (16). Let $\phi_i : \mathcal{K}_i \rightarrow \mathbb{R}$ be a ν_i -self-concordant barrier
2553 for all $i \in [n]$. Suppose the program satisfies the following properties:
2554

- 2555 • Inner radius r : There exists $z \in \mathbb{R}^{n_{\text{tot}}}$ such that $Az = b$ and $B(z, r) \in \mathcal{K}$.
- 2556 • Outer radius R : $\mathcal{K} \subseteq B(0, R)$ where $0 \in \mathbb{R}^{n_{\text{tot}}}$.
- 2557 • Lipschitz constant L : $\|Q\|_{2 \rightarrow 2} \leq L$, $\|c\|_2 \leq L$.

2560 Let $(w_i)_{i \in [n]} \in \mathbb{R}_{\geq 1}^n$ and $\kappa = \sum_{i \in [n]} w_i \nu_i$. For any $0 < \epsilon \leq \frac{1}{2}$, Algorithm 19 outputs an approximate
2561 solution x in $O(\sqrt{\kappa} \log n \log \frac{n\kappa R}{\epsilon r})$ steps, satisfying
2562

$$2563 \frac{1}{2} x^\top Q x + c^\top x \leq \min_{Ax=b, x \in \mathcal{K}} \left(\frac{1}{2} x^\top Q x + c^\top x \right) + \epsilon LR(R+1),$$

$$2564 \|Ax - b\|_1 \leq 3\epsilon(R\|A\|_1 + \|b\|_1),$$

$$2565 x \in \mathcal{K}.$$

2569 Algorithm 19 Our main algorithm

- 2570 1: **procedure** ROBUSTQPIPM($Q \in \mathcal{S}^{n_{\text{tot}}}$, $c \in \mathbb{R}^{n_{\text{tot}}}$, $A \in \mathbb{R}^{m \times n_{\text{tot}}}$, $b \in \mathbb{R}^m$, $(\phi_i : \mathcal{K}_i \rightarrow$
2571 $\mathbb{R})_{i \in [n]}$, $w \in \mathbb{R}^n$)
 - 2572 2: */* Initial point reduction */*
 - 2573 3: $\rho \leftarrow LR(R+1)$, $x^{(0)} \leftarrow \arg \min_x \sum_{i \in [n]} w_i \phi_i(x_i)$, $s^{(0)} \leftarrow \epsilon \rho (c + Qx^{(0)})$
 - 2574 4: $\bar{x} \leftarrow \begin{bmatrix} x^{(0)} \\ 1 \end{bmatrix}$, $\bar{s} \leftarrow \begin{bmatrix} s^{(0)} \\ 1 \end{bmatrix}$, $\bar{Q} \leftarrow \begin{bmatrix} \epsilon \rho Q & 0 \\ 0 & 0 \end{bmatrix}$, $\bar{A} \leftarrow [A \mid b - Ax^{(0)}]$
 - 2575 5: $\bar{w} \leftarrow \begin{bmatrix} w \\ 1 \end{bmatrix}$, $\bar{\phi}_i = \phi_i \forall i \in [n]$, $\bar{\phi}_{n+1}(x) := -\log x - \log(2-x)$
 - 2576 6: $(x, s) \leftarrow \text{CENTERING}(\bar{Q}, \bar{A}, (\bar{\phi}_i)_{i \in [n+1]}, \bar{w}, \bar{x}, \bar{s}, t_{\text{start}} = 1, t_{\text{end}} = \frac{\epsilon^2}{4\kappa})$
 - 2577 7: **return** $(x_{1:n}, s_{1:n})$
 - 2578 8: **end procedure**
-

2583 F.1 PRELIMINARIES

2584
2585 Previous works on linear programming (e.g. Lee et al. (2019), Lee & Vempala (2021)) use the
2586 following path:

$$2587 s/t + \nabla \phi_w(x) = \mu,$$

$$2588 Ax = b,$$

$$2589 A^\top y + s = c$$

2590 where $\phi_w(x) := \sum_{i=1}^n w_i \phi_i(x_i)$.

2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645

Algorithm 20 Subroutine used by Algorithm 19

1: **procedure** CENTERING($Q \in \mathcal{S}^{n_{\text{tot}}}, A \in \mathbb{R}^{m \times n_{\text{tot}}}, (\phi_i : \mathcal{K}_i \rightarrow \mathbb{R})_{i \in [n]}, w \in \mathbb{R}^n, x \in \mathbb{R}^{n_{\text{tot}}}, s \in \mathbb{R}^{n_{\text{tot}}}, t_{\text{start}} \in \mathbb{R}_{>0}, t_{\text{end}} \in \mathbb{R}_{>0}$)

2: /* Parameters */

3: $\lambda = 64 \log(256n \sum_{i \in [n]} w_i), \bar{\epsilon} = \frac{1}{1440} \lambda, \alpha = \frac{\bar{\epsilon}}{2}$

4: $\epsilon_t = \frac{\bar{\epsilon}}{4} (\min_{i \in [n]} \frac{w_i}{w_i + \nu_i}), h = \frac{\alpha}{64\sqrt{\kappa}}$

5: /* Definitions */

6: $\phi_w(x) := \sum_{i \in [n]} w_i \phi_i(x_i)$

7: $\mu_i(x, s, t) := s/t + w_i \nabla \phi_i(x_i), \forall i \in [n]$ ▷ Eq. (17)

8: $\gamma_i(x, s, t) \leftarrow \|\mu_i^t(x, s)\|_{x_i}^*, \forall i \in [n]$ ▷ Eq. (18)

9: $c_i(x, s, t) := \frac{\sinh(\frac{\lambda}{w_i} \gamma_i(x, s, t))}{\gamma_i(x, s, t) \sqrt{\sum_{j \in [n]} w_j^{-1} \cosh^2(\frac{\lambda}{w_j} \gamma_j(x, s, t))}}, \forall i \in [n]$ ▷ Eq. (22)

10: $H_{w,x} := \nabla^2 \phi_w(x)$ ▷ Eq. (24)

11: $B_{w,x,t} := Q + tH_{w,x}$ ▷ Eq. (25)

12: $P_{w,x,t} := B_{w,x,t}^{-1/2} A^\top (AB_{w,x,t}^{-1} A^\top)^{-1} AB_{w,x,t}^{-1/2}$ ▷ Eq. (26)

13: /* Main loop */

14: $\bar{t} \leftarrow t \leftarrow t_{\text{start}}, \bar{x} \leftarrow x, \bar{s} \leftarrow s$

15: **while** $t > t_{\text{end}}$ **do**

16: Maintain $\bar{x}, \bar{s}, \bar{t}$ such that $\|\bar{x}_i - x_i\|_{\bar{x}_i} \leq \bar{\epsilon}, \|\bar{s}_i - s_i\|_{\bar{x}_i}^* \leq t\bar{\epsilon}w_i$ and $|\bar{t} - t| \leq \epsilon_t \bar{t}$

17: $\delta_{\mu,i} \leftarrow -\alpha \cdot c_i(\bar{x}, \bar{s}, \bar{t}) \cdot \mu_i(\bar{x}, \bar{s}, \bar{t}), \forall i \in [n]$ ▷ Eq. (21)

18: Pick δ_x and δ_s such that $A\delta_x = 0, \delta_s - Q\delta_x \in \text{Range}(A^\top)$ and

$$\|\delta_x - \bar{t} B_{w,\bar{x},\bar{t}}^{-1/2} (I - P_{w,\bar{x},\bar{t}}) B_{w,\bar{x},\bar{t}}^{-1/2} \delta_\mu\|_{w,\bar{x}} \leq \bar{\epsilon} \alpha,$$

$$\|\bar{t}^{-1} \delta_s - (\delta_\mu - \bar{t} H_{w,\bar{x}} B_{w,\bar{x},\bar{t}}^{-1/2} (I - P_{w,\bar{x},\bar{t}}) B_{w,\bar{x},\bar{t}}^{-1/2} \delta_\mu)\|_{w,\bar{x}}^* \leq \bar{\epsilon} \alpha.$$

19: $t \leftarrow \max\{(1-h)t, t_{\text{end}}\}, x \leftarrow x + \delta_x, s \leftarrow s + \delta_s$

20: **end while**

21: **return** (x, s)

22: **end procedure**

For quadratic programming, we modify the above central path as following:

$$\begin{aligned} s/t + \nabla \phi_w(x) &= \mu, \\ Ax &= b, \\ -Qx + A^\top y + s &= c. \end{aligned}$$

We make the following definitions.

Definition F.2. For each $i \in [n]$, we define the i -th coordinate error

$$\mu_i(x, s, t) := \frac{s_i}{t} + w_i \nabla \phi_i(x_i) \tag{17}$$

We define μ_i 's norm as

$$\gamma_i(x, s, t) := \|\mu_i(x, s, t)\|_{x_i}^*. \tag{18}$$

We define the soft-max function by

$$\Psi_\lambda(r) := \sum_{i=1}^m \cosh(\lambda \frac{r_i}{w_i}) \tag{19}$$

for some $\lambda > 0$ and the potential function is the soft-max of the norm of the error of each coordinate

$$\Phi(x, s, t) = \Psi_\lambda(\gamma(x, s, t)) \tag{20}$$

2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699

We choose the step direction δ_μ as

$$\delta_{\mu,i} := -\alpha \cdot c_i(x, s, t) \cdot \mu_i(x, s, t) \quad (21)$$

where

$$c_i(x, s, t) := \frac{\sinh(\frac{\lambda}{w_i} \gamma_i(x, s, t))}{\gamma_i(x, s, t) \sqrt{\sum_{j \in [n]} w_j^{-1} \cosh^2(\frac{\lambda}{w_j} \gamma_j(x, s, t))}} \quad (22)$$

We define induced norms as following. Note that we include the weight vector w in the subscript to avoid confusion.

Definition F.3. For each block \mathcal{K}_i , we define

$$\begin{aligned} \|v\|_{x_i} &:= \|v\|_{\nabla^2 \phi_i(x_i)}, \\ \|v\|_{x_i}^* &:= \|v\|_{(\nabla^2 \phi_i(x_i))^{-1}} \end{aligned}$$

for $v \in \mathbb{R}^{n_i}$.

For the whole domain $\mathcal{K} = \prod_{i=1}^n \mathcal{K}_i$, we define

$$\begin{aligned} \|v\|_{w,x} &:= \|v\|_{\nabla^2 \phi_w(x)} = \left(\sum_{i=1}^n w_i \|v_i\|_{x_i}^2 \right)^{1/2}, \\ \|v\|_{w,x}^* &:= \|v\|_{(\nabla^2 \phi_w(x))^{-1}} = \left(\sum_{i=1}^n w_i^{-1} (\|v_i\|_{x_i}^*)^2 \right)^{1/2} \end{aligned}$$

for $v \in \mathbb{R}^{n_{\text{tot}}}$.

The Hessian matrices of the barrier functions appear a lot in the computation.

Definition F.4. We define matrices $H_{x,i} \in \mathbb{R}^{n_i \times n_i}$ and $H_{w,x} \in \mathbb{R}^{n_{\text{tot}} \times n_{\text{tot}}}$ as

$$H_{x,i} := \nabla^2 \phi_i(x_i), \quad (23)$$

$$H_{w,x} := \nabla^2 \phi_w(x). \quad (24)$$

From the definition, we see that

$$H_{w,x,(i,i)} = w_i H_{x,i}.$$

The following equations are immediate from definition.

Claim F.5. Let $H_{w,x} \in \mathbb{R}^{n_{\text{tot}} \times n_{\text{tot}}}$ be defined as Definition F.4. For $v \in \mathbb{R}^{n_{\text{tot}}}$, we have

$$\begin{aligned} \|v\|_{w,x} &= \|H_{w,x}^{1/2} v\|_2, \\ \|v\|_{w,x}^* &= \|H_{w,x}^{-1/2} v\|_2. \end{aligned}$$

Claim F.6. For each $i \in [n]$, let $H_{x,i}$ be defined as Definition F.4. For $v \in \mathbb{R}^{n_i}$, $i \in [n]$, we have

$$\begin{aligned} \|v\|_{x_i} &= \|H_{x,i}^{1/2} v\|_2, \\ \|v\|_{x_i}^* &= \|H_{x,i}^{-1/2} v\|_2. \end{aligned}$$

We define matrices B and P used in the algorithm.

Definition F.7. Let A, Q denote two fixed matrices. Let $H_{w,x} \in \mathbb{R}^{n_{\text{tot}} \times n_{\text{tot}}}$ be defined as Definition F.4. We define matrix $B_{w,x,t} \in \mathbb{R}^{n_{\text{tot}} \times n_{\text{tot}}}$ as

$$B_{w,x,t} := Q + t \cdot H_{w,x} \quad (25)$$

We define projection matrix $P_{w,x,t} \in \mathbb{R}^{n_{\text{tot}} \times n_{\text{tot}}}$ as

$$P_{w,x,t} \leftarrow B_{w,x,t}^{-1/2} A^\top (A B_{w,x,t}^{-1} A^\top)^{-1} A B_{w,x,t}^{-1/2}. \quad (26)$$

2700 F.2 DERIVING THE CENTRAL PATH STEP
2701

2702 In this section we explain how to derive the central path step.

2703 We follow the central path
2704

$$\begin{aligned} 2705 \quad & s/t + \nabla\phi_w(x) = \mu \\ 2706 \quad & Ax = b \\ 2707 \quad & -Qx + A^\top y + s = c \end{aligned}$$

2708
2709
2710 We perform gradient descent on μ with step δ_μ . Then Newton step gives

$$2711 \quad \frac{1}{t}\delta_s + \nabla^2\phi_w(x)\delta_x = \delta_\mu \quad (27)$$

$$2712 \quad A\delta_x = 0 \quad (28)$$

$$2713 \quad -Q\delta_x + A^\top\delta_y + \delta_s = 0 \quad (29)$$

2714 where δ_x (resp. δ_y, δ_s) is the step taken by x (resp. y, s).

2715 For simplicity, we define $H \in \mathbb{R}^{n_{\text{tot}} \times n_{\text{tot}}}$ to represent $\nabla^2\phi_w(x)$.⁸

2716 From Eq. (27) we get

$$2717 \quad \delta_s = t\delta_\mu - tH\delta_x. \quad (30)$$

2718 Plug the above equation into Eq. (29) we get

$$2719 \quad -Q\delta_x + A^\top\delta_y + t\delta_\mu - tH\delta_x = 0. \quad (31)$$

2720 Let $B = Q + tH$, multiply by AB^{-1} we get

$$2721 \quad -A\delta_x + AB^{-1}A^\top\delta_y + tAB^{-1}\delta_\mu = 0.$$

2722 Using Eq. (28) we get

$$2723 \quad AB^{-1}A^\top\delta_y + tAB^{-1}\delta_\mu = 0.$$

2724 Solve for δ_y (assuming that $AB^{-1}A$ is invertible), we get

$$2725 \quad \delta_y = -t(AB^{-1}A^\top)^{-1}AB^{-1}\delta_\mu.$$

2726 Plug into Eq. (31) we get

$$2727 \quad -B\delta_x - tA^\top(AB^{-1}A^\top)^{-1}AB^{-1}\delta_\mu + t\delta_\mu = 0.$$

2728 Solve for δ_x we get

$$\begin{aligned} 2729 \quad \delta_x &= tB^{-1}\delta_\mu - tB^{-1}A^\top(AB^{-1}A^\top)^{-1}AB^{-1}\delta_\mu \\ 2730 \quad &= tB^{-1/2}(I - P)B^{-1/2}\delta_\mu \end{aligned}$$

2731 where $P = B^{-1/2}A^\top(AB^{-1}A^\top)^{-1}AB^{-1/2}$ is the projection matrix. Solve for δ_s in Eq. (30) we get

$$2732 \quad \delta_s = t\delta_\mu - t^2HB^{-1/2}(I - P)B^{-1/2}\delta_\mu.$$

2733 In summary, we have

$$\begin{aligned} 2734 \quad \delta_x &= tB^{-1/2}(I - P)B^{-1/2}\delta_\mu, \\ 2735 \quad \delta_y &= -t(AB^{-1}A^\top)^{-1}AB^{-1}\delta_\mu, \\ 2736 \quad \delta_s &= t\delta_\mu - t^2HB^{-1/2}(I - P)B^{-1/2}\delta_\mu, \\ 2737 \quad P &= B^{-1/2}A^\top(AB^{-1}A^\top)^{-1}AB^{-1/2}. \end{aligned}$$

2738 These equations will guide the design of our actual algorithm.

2739 ⁸In this section, and in this section only, we omit the subscript in H, B, P for simplicity.

F.3 BOUNDING MOVEMENT OF POTENTIAL FUNCTION

The goal of this section is to bound the movement of potential function during the robust IPM algorithm.

In robust IPM, we do not need to follow the ideal central path exactly over the entire algorithm. Instead, we only use an approximate version. For convenience of analysis we state two assumptions (see Algorithm 20, Line 18).

Assumption F.8. *We make the following assumptions on $\delta_x \in \mathbb{R}^{n_{\text{tot}}}$ and $\delta_s \in \mathbb{R}^{n_{\text{tot}}}$.*

$$\begin{aligned} \|\delta_x - \bar{t} B_{w, \bar{x}, \bar{t}}^{-1/2} (I - P_{w, \bar{x}, \bar{t}}) B_{w, \bar{x}, \bar{t}}^{-1/2} \delta_\mu\|_{w, \bar{x}} &\leq \bar{\epsilon} \alpha, \\ \|\bar{t}^{-1} \delta_s - (\delta_\mu - \bar{t} H_{w, \bar{x}} B_{w, \bar{x}, \bar{t}}^{-1/2} (I - P_{w, \bar{x}, \bar{t}}) B_{w, \bar{x}, \bar{t}}^{-1/2} \delta_\mu)\|_{w, \bar{x}}^* &\leq \bar{\epsilon} \alpha. \end{aligned}$$

The following lemma bounds the movement of potential function Ψ assuming bound on δ_γ .

Lemma F.9 (Ye, 2020, Lemma A.5). *For any $r \in \mathbb{R}^{n_{\text{tot}}}$, and $w \in \mathbb{R}_{\geq 1}^{n_{\text{tot}}}$. Let α and λ denote the parameters that are satisfying $0 \leq \alpha \leq \frac{1}{8\lambda}$.*

Let $\epsilon_r \in \mathbb{R}^{n_{\text{tot}}}$ denote a vector satisfying

$$\left(\sum_{i=1}^n w_i^{-1} \epsilon_{r,i}^2 \right)^{1/2} \leq \alpha/8.$$

Suppose that vector $\bar{r} \in \mathbb{R}^{n_{\text{tot}}}$ is satisfying the following property

$$|r_i - \bar{r}_i| \leq \frac{w_i}{8\lambda}, \quad \forall i \in [n]$$

We define vector $\delta_r \in \mathbb{R}^{n_{\text{tot}}}$ as follows:

$$\delta_{r,i} := \frac{-\alpha \cdot \sinh\left(\frac{\lambda}{w_i} \bar{r}_i\right)}{\sqrt{\sum_{j=1}^n w_j^{-1} \cosh^2\left(\frac{\lambda}{w_j} \bar{r}_j\right)}} + \epsilon_{r,i}.$$

Then, we have that

$$\Psi_\lambda(r + \delta_r) \leq \Psi_\lambda(r) - \frac{\alpha\lambda}{2} \left(\sum_{i=1}^n w_i^{-1} \cosh^2\left(\lambda \frac{r_i}{w_i}\right) \right)^{1/2} + \alpha\lambda \left(\sum_{i=1}^n w_i^{-1} \right)^{1/2}$$

The following lemma bounds the norm of δ_μ .

Lemma F.10 (Bounding norm of δ_μ).

$$\|\delta_\mu(\bar{x}, \bar{s}, \bar{t})\|_{w, \bar{x}}^* \leq \alpha.$$

Proof.

$$\begin{aligned} (\|\delta_\mu(\bar{x}, \bar{s}, \bar{t})\|_{w, \bar{x}}^*)^2 &= \sum_{i=1}^n w_i^{-1} (\|\delta_{\mu,i}(\bar{x}, \bar{s}, \bar{t})\|_{\bar{x}_i}^*)^2 \\ &= \alpha^2 \sum_{i \in [n]} w_i^{-1} c_i^2(\bar{x}, \bar{s}, \bar{t}) \cdot \|\mu_i(\bar{x}, \bar{s}, \bar{t})\|_{\bar{x}_i}^2 \\ &= \alpha^2 \sum_{i \in [n]} w_i^{-1} c_i^2(\bar{x}, \bar{s}, \bar{t}) \cdot \|H_{\bar{x}, i}^{-1/2} \mu_i(\bar{x}, \bar{s}, \bar{t})\|_2^2 \\ &= \alpha^2 \sum_{i \in [n]} w_i^{-1} c_i^2(\bar{x}, \bar{s}, \bar{t}) \cdot \gamma_i^2(\bar{x}, \bar{s}, \bar{t}) \\ &= \alpha^2 \sum_{i \in [n]} \frac{w_i^{-1} \sinh^2\left(\frac{\lambda}{w_i} \gamma_i(\bar{x}, \bar{s}, \bar{t})\right)}{\gamma_i^2(\bar{x}, \bar{s}, \bar{t}) \cdot \sum_{j \in [n]} w_j^{-1} \cosh^2\left(\frac{\lambda}{w_j} \gamma_j(\bar{x}, \bar{s}, \bar{t})\right)} \cdot \gamma_i^2(\bar{x}, \bar{s}, \bar{t}) \end{aligned}$$

2808

2809

2810

2811

2812

2813

2814

2815

2816

2817

2818

2819

2820

2821

2822

2823

2824

2825

2826

2827

2828

2829

2830

2831

2832

2833

2834

2835

2836

2837

2838

2839

2840

2841

2842

2843

2844

2845

2846

2847

2848

2849

2850

2851

2852

2853

2854

2855

2856

2857

2858

2859

2860

2861

$$\begin{aligned}
&= \alpha^2 \frac{\sum_{j \in [n]} w_j^{-1} \sinh^2(\frac{\lambda}{w_j} \gamma_j(\bar{x}, \bar{s}, \bar{t}))}{\sum_{j \in [n]} w_j^{-1} \cosh^2(\frac{\lambda}{w_j} \gamma_j(\bar{x}, \bar{s}, \bar{t}))} \\
&\leq \alpha^2.
\end{aligned}$$

where the first step follows from Definition F.3, the second step follows from $\delta_{\mu,i}(\bar{x}, \bar{s}, \bar{t}) = -\alpha \cdot c_i(\bar{x}, \bar{s}, \bar{t}) \cdot \mu_i(\bar{x}, \bar{s}, \bar{t})$, the third step follows from norm of \bar{x}_i (see Definition F.3), the fourth step follows from $\gamma_i(\bar{x}, \bar{s}, \bar{t}) = \|H_{\bar{x},i}^{-1/2} \mu_i(\bar{x}, \bar{s}, \bar{t})\|_2$ (see Eq. (18)), the fifth step follows from $c_i(\bar{x}, \bar{s}, \bar{t})^2 = \frac{\sinh^2(\frac{\lambda}{w_i} \gamma_i(\bar{x}, \bar{s}, \bar{t}))}{\gamma_i^2(\bar{x}, \bar{s}, \bar{t}) \sum_{j \in [n]} w_j^{-1} \cosh^2(\frac{\lambda}{w_j} \gamma_j(\bar{x}, \bar{s}, \bar{t}))}$ (see Eq. (22)), the sixth step follows from canceling the term $\gamma_i^2(\bar{x}, \bar{s}, \bar{t})$, and the last step follows from $\cosh^2(x) \geq \sinh^2(x)$ for all x . \square

The following lemma bounds the norm of δ_x and δ_s .

Lemma F.11. *For each $i \in [n]$, we define $\alpha_i := \|\delta_{x,i}\|_{\bar{x}_i}$. Then, we have*

$$\|\delta_x\|_{w,\bar{x}} = \left(\sum_{i \in [n]} w_i \alpha_i^2 \right)^{1/2} \leq \frac{9}{8} \alpha. \quad (32)$$

In particular, we have $\alpha_i \leq \frac{9}{8} \alpha$. Similarly, for δ_s , we have

$$\|\delta_s\|_{w,\bar{x}}^* = \sqrt{\sum_{i \in [n]} w_i^{-1} (\|\delta_{s,i}\|_{\bar{x}_i}^*)^2} \leq \frac{17}{8} \alpha \cdot t. \quad (33)$$

Proof. For δ_x , we have

$$\begin{aligned}
\|\delta_x\|_{w,\bar{x}} &\leq \|\bar{t} H_{w,\bar{x}}^{1/2} B_{w,\bar{x},\bar{t}}^{-1/2} (I - P_{w,\bar{x},\bar{t}}) B_{w,\bar{x},\bar{t}}^{-1/2} \delta_\mu\|_2 + \bar{\epsilon} \alpha \\
&\leq \|\bar{t}^{1/2} (I - P_{w,\bar{x},\bar{t}}) B_{w,\bar{x},\bar{t}}^{-1/2} \delta_\mu\|_2 + \bar{\epsilon} \alpha \\
&\leq \|\bar{t}^{1/2} B_{w,\bar{x},\bar{t}}^{-1/2} \delta_\mu\|_2 + \bar{\epsilon} \alpha \\
&\leq \|H_{w,\bar{x}}^{-1/2} \delta_\mu\|_2 + \bar{\epsilon} \alpha \\
&\leq \alpha + \bar{\epsilon} \alpha \\
&\leq \frac{9}{8} \alpha.
\end{aligned}$$

First step follows from Assumption F.8. Second step is because $\bar{t} H_{w,\bar{x}} \preceq B_{w,\bar{x},\bar{t}}$. Third step is because $P_{w,\bar{x},\bar{t}}$ is a projection matrix. Fourth step is because $\bar{t} H_{w,\bar{x}} \preceq B_{w,\bar{x},\bar{t}}$. Fifth step is by Lemma F.10. Sixth step is because $\bar{\epsilon} \leq \frac{1}{8}$.

For δ_s , we have

$$\begin{aligned}
\|\delta_s\|_{w,\bar{x}}^* &\leq \|\bar{t} \delta_\mu\|_{w,\bar{x}}^* + \|\bar{t}^2 H_{w,\bar{x}} B_{w,\bar{x},\bar{t}}^{-1/2} (I - P_{w,\bar{x},\bar{t}}) B_{w,\bar{x},\bar{t}}^{-1/2} \delta_\mu\|_{w,\bar{x}}^* + \bar{\epsilon} \alpha \bar{t} \\
&\leq \alpha \bar{t} + \alpha \bar{t} + \bar{\epsilon} \alpha \bar{t} \\
&\leq \frac{17}{8} \alpha \cdot t.
\end{aligned}$$

First step is by triangle inequality and the assumption that

$$\delta_s \approx \bar{t} \delta_\mu - \bar{t}^2 H_{w,\bar{x}} B_{w,\bar{x},\bar{t}}^{-1/2} (I - P_{w,\bar{x},\bar{t}}) B_{w,\bar{x},\bar{t}}^{-1/2} \delta_\mu.$$

Second step is by same analysis as the analysis for δ_x . Third step is by $\bar{t} \leq \frac{33}{32} t$ and $\bar{\epsilon} \leq \frac{1}{32}$. \square

The following lemma shows that μ^{new} is close to $\mu + \delta_\mu$ under an approximate step.

Lemma F.12 (Variation of (Ye, 2020, Lemma A.9)). *For each $i \in [n]$, we define*

$$\beta_i := \|\epsilon_{\mu,i}\|_{x_i}^*$$

2862 For each $i \in [n]$, let

$$2863 \mu_i(x^{\text{new}}, s^{\text{new}}, t) = \mu_i(x, s, t) + \delta_{\mu,i} + \epsilon_{\mu,i}.$$

2864 Then, we have

$$2865 \left(\sum_{i=1}^n w_i^{-1} \beta_i^2 \right)^{1/2} \leq 15\bar{\epsilon}\alpha.$$

2866 *Proof.* The proof is similar as (Ye, 2020, Lemma A.9), except for changing the definitions of ϵ_1 and ϵ_2 :

$$2867 \epsilon_1 := H_{w,\bar{x}}^{1/2} \delta_x - \bar{t} \cdot H_{w,\bar{x}}^{1/2} B_{w,\bar{x},\bar{t}}^{-1/2} (I - P_{w,\bar{x},\bar{t}}) B_{w,\bar{x},\bar{t}}^{-1/2} \delta_\mu,$$

$$2868 \epsilon_2 := \bar{t}^{-1} H_{w,\bar{x}}^{-1/2} \delta_s - H_{w,\bar{x}}^{-1/2} (\delta_\mu - \bar{t} H_{w,\bar{x}} B_{w,\bar{x},\bar{t}}^{-1/2} (I - P_{w,\bar{x},\bar{t}}) B_{w,\bar{x},\bar{t}}^{-1/2} \delta_\mu).$$

2869 One key step in the proof of Ye (2020) is the following property:

$$2870 \delta_{\mu,i} = \bar{t}^{-1} \cdot \delta_{s,i} + H_{w,\bar{x}} \delta_{x,i} - H_{w,\bar{x}}^{1/2} (\epsilon_1 + \epsilon_2).$$

2871 Under our new definition of ϵ_1 and ϵ_2 , the above property still holds. Remaining parts of the proof are similar and we omit the details here. \square

2872 The following lemma shows that error $\mu(\bar{x}, \bar{s}, \bar{t})$ on the robust central path is close to error $\mu(x, s, t)$ on the ideal central path. Furthermore, norms of errors $\gamma_i(x, s, t)$ and $\gamma_i(\bar{x}, \bar{s}, \bar{t})$ are also close to each other.

2873 **Lemma F.13** ((Ye, 2020, Lemma A.10)). Assume that $\gamma_i(x, s, t) \leq w_i$ for all i . For all $i \in [n]$, we have

$$2874 \|\mu_i(x, s, t) - \mu_i(\bar{x}, \bar{s}, \bar{t})\|_{x_i}^* \leq 3\bar{\epsilon}w_i.$$

2875 Furthermore, we have that

$$2876 |\gamma_i(x, s, t) - \gamma_i(\bar{x}, \bar{s}, \bar{t})| \leq 5\bar{\epsilon}w_i.$$

2877 *Proof.* Same as proof of (Ye, 2020, Lemma A.10). \square

2878 The following lemma bounds the change of γ under one robust IPM step.

2879 **Lemma F.14** ((Ye, 2020, Lemma A.12)). Assume $\Phi(x, s, t) \leq \cosh(\lambda)$. For all $i \in [n]$, we define

$$2880 \epsilon_{r,i} := \gamma_i(x^{\text{new}}, s^{\text{new}}) - \gamma_i(x, s, t) + \alpha \cdot c_i(\bar{x}, \bar{s}, \bar{t}) \cdot \gamma_i(\bar{x}, \bar{s}, \bar{t}).$$

2881 Then, we have

$$2882 \left(\sum_{i=1}^n w_i^{-1} \epsilon_{r,i}^2 \right)^{1/2} \leq 90 \cdot \bar{\epsilon} \cdot \lambda \alpha + 4 \cdot \max_{i \in [n]} (w_i^{-1} \gamma_i(x, s, t)) \cdot \alpha.$$

2883 *Proof.* The proof is similar to the proof of (Ye, 2020, Lemma A.12). By replacing corresponding references in Ye (2020) by our versions (Lemma F.11, F.12, F.13) we get proof of this lemma. \square

2884 Finally, the following theorem bounds the movement of potential function Φ under one robust IPM step.

2885 **Theorem F.15** (Variation of (Ye, 2020, Theorem A.15)). Assume $\Phi(x, s, t) \leq \cosh(\lambda/64)$. Then for any $0 \leq h \leq \frac{\alpha}{64 \sqrt{\sum_{i \in [n]} w_i \nu_i}}$, we have

$$2886 \Phi(x^{\text{new}}, s^{\text{new}}, t^{\text{new}}) \leq \left(1 - \frac{\alpha \lambda}{\sqrt{\sum_{i \in [n]} w_i}} \right) \cdot \Phi(x, s, t) + \alpha \lambda \sqrt{\sum_{i \in [n]} w_i^{-1}}.$$

2887 In particular, for any $\cosh(\lambda/128) \leq \Phi(x, s, t) \leq \cosh(\lambda/64)$, we have that

$$2888 \Phi(x^{\text{new}}, s^{\text{new}}, t^{\text{new}}) \leq \Phi(x, s, t).$$

2889 *Proof.* Similar to the proof of (Ye, 2020, Theorem A.15), but replacing lemmas with the corresponding QP versions. \square

2916 F.4 INITIAL POINT REDUCTION
2917

2918 In this section, we propose an initial point reduction scheme for quadratic programming. Our scheme
2919 is closer to Lee et al. (2019) rather than Ye (2020); Lee & Vempala (2021). The reason is that Lee
2920 & Vempala (2021)'s initial point reduction requires an efficient algorithm for finding the optimal
2921 solution to an unconstrained program, which may be difficult in quadratic programming.

2922 **Lemma F.16** ((Nesterov, 1998, Theorem 4.1.7 and Lemma 4.2.4)). *Let ϕ be a ν -self-concordant*
2923 *barrier. Then for any $x, y \in \text{dom}(\phi)$, we have*

$$2924 \langle \nabla \phi(x), y - x \rangle \leq \nu,$$

$$2925$$

$$2926 \langle \nabla \phi(y) - \nabla \phi(x), y - x \rangle \geq \frac{\|y - x\|_x^2}{1 + \|y - x\|_x}.$$

2928 Let $x^* = \arg \min_x \phi(x)$. For any $x \in \mathbb{R}^n$ such that $\|x - x^*\|_{x^*} \leq 1$, we have that $x \in \text{dom}(\phi)$.

2929 **Lemma F.17** (QP version of (Lee et al., 2019, Lemma D.2)). *Work under the setting of Theorem F.1.*
2930 *Let $x^{(0)} = \arg \min_x \sum_{i \in [n]} w_i \phi_i(x_i)$. Let $\rho = \frac{1}{LR(R+1)}$. For any $0 < \epsilon \leq \frac{1}{2}$, the modified program*

$$2932 \min_{\substack{\bar{A}\bar{x}=\bar{b}, \bar{x} \in \mathcal{K} \times \mathbb{R}_{\geq 0}}} \left(\frac{1}{2} \bar{x}^\top \bar{Q} \bar{x} + \bar{c}^\top \bar{x} \right)$$

2935 with

$$2936 \bar{Q} = \begin{bmatrix} \epsilon \rho Q & 0 \\ 0 & 0 \end{bmatrix}, \quad \bar{A} = [A \mid b - Ax^{(0)}], \quad \bar{b} = b, \quad \bar{c} = \begin{bmatrix} \epsilon \rho c \\ 1 \end{bmatrix}$$

2939 satisfies the following:

- 2941 • $\bar{x} = \begin{bmatrix} x^{(0)} \\ 1 \end{bmatrix}$, $\bar{y} = 0 \in \mathbb{R}^m$ and $\bar{s} = \begin{bmatrix} \epsilon \rho(c + Qx^{(0)}) \\ 1 \end{bmatrix}$ are feasible primal dual vectors with
- 2942 $\|\bar{s} + \nabla \bar{\phi}_w(\bar{x})\|_{\bar{x}}^* \leq \epsilon$ where $\bar{\phi}_w(\bar{x}) = \sum_{i=1}^n w_i \phi_i(\bar{x}_i) - \log(\bar{x}_{n+1})$.
- 2944 • For any $\bar{x} \in \mathcal{K} \times \mathbb{R}_{\geq 0}$ satisfying $\bar{A}\bar{x} = \bar{b}$ and

$$2946 \frac{1}{2} \bar{x}^\top \bar{Q} \bar{x} + \bar{c}^\top \bar{x} \leq \min_{\substack{\bar{A}\bar{x}=\bar{b}, \bar{x} \in \mathcal{K} \times \mathbb{R}_{\geq 0}}} \left(\frac{1}{2} \bar{x}^\top \bar{Q} \bar{x} + \bar{c}^\top \bar{x} \right) + \epsilon^2, \quad (34)$$

2949 the vector $\bar{x}_{1:n}$ ($\bar{x}_{1:n}$ is the first n coordinates of \bar{x}) is an approximate solution to the
2950 original convex program in the following sense:

$$2951 \frac{1}{2} \bar{x}_{1:n}^\top Q \bar{x}_{1:n} + c^\top \bar{x}_{1:n} \leq \min_{Ax=b, x \in \mathcal{K}} \left(\frac{1}{2} x^\top Q x + c^\top x \right) + \epsilon \rho^{-1},$$

$$2952$$

$$2953 \|A\bar{x}_{1:n} - b\|_1 \leq 3\epsilon \cdot (R\|A\|_1 + \|b\|_1),$$

$$2954 \bar{x}_{1:n} \in \mathcal{K}.$$

2956 *Proof.* **First bullet point:** Direct computation shows that $(\bar{x}, \bar{y}, \bar{s})$ is feasible.

2958 Let us compute $\|\bar{s} + \nabla \bar{\phi}_w(\bar{x})\|_{\bar{x}}^*$. We have

$$2959 \|\bar{s} + \nabla \bar{\phi}_w(\bar{x})\|_{\bar{x}}^* = \|\epsilon \rho(c + Qx^{(0)})\|_{\nabla^2 \phi_w(x^{(0)})^{-1}}$$

2962 Lemma F.16 says that for all $x \in \mathbb{R}^n$ with $\|x - x^{(0)}\|_{w, x^{(0)}} \leq 1$, we have $x \in \mathcal{K}$, because
2963 $x^{(0)} = \arg \min_x \phi_w(x)$. Therefore for any v such that $v^\top \nabla^2 \phi_w(x^{(0)}) v \leq 1$, we have $x^{(0)} \pm v \in \mathcal{K}$
2964 and hence $\|x^{(0)} \pm v\|_2 \leq R$. This implies $\|v\|_2 \leq R$ for any $v^\top \nabla^2 \phi_w(x^{(0)}) v \leq 1$. Hence
2965 $(\nabla^2 \phi_w(x^{(0)}))^{-1} \preceq R^2 \cdot I$. So we have

$$2966 \|\bar{s} + \nabla \bar{\phi}_w(\bar{x})\|_{\bar{x}}^* = \|\epsilon \rho(c + Qx^{(0)})\|_{\nabla^2 \phi_w(x^{(0)})^{-1}}$$

$$2967 \leq \epsilon \rho R \|c + Qx^{(0)}\|_2$$

$$2968 \leq \epsilon \rho R (\|c\|_2 + \|Q\|_{2 \rightarrow 2} \|x^{(0)}\|_2)$$

2969

2970

2971

2972

2973

2974

Second bullet point: We define

2975

2976

2977

2978

2979

2980

2981

2982

2983

2984

2985

2986

2987

2988

2989

2990

2991

2992

2993

2994

2995

2996

2997

2998

2999

3000

3001

3002

3003

3004

3005

3006

3007

3008

3009

3010

3011

3012

3013

3014

3015

3016

3017

3018

3019

3020

3021

3022

3023

$$\begin{aligned} &\leq \epsilon \rho R(L + LR) \\ &\leq \epsilon. \end{aligned}$$

$$\text{OPT} := \min_{Ax=b, x \in \mathcal{K}} \left(\frac{1}{2} x^\top Q x + c^\top x \right), \quad (35)$$

$$\overline{\text{OPT}} := \min_{A\bar{x}=b, \bar{x} \in \mathcal{K} \times \mathbb{R}_{\geq 0}} \left(\frac{1}{2} \bar{x}^\top Q \bar{x} + \bar{c}^\top \bar{x} \right). \quad (36)$$

For any feasible x in the original problem (35), $\bar{x} = \begin{bmatrix} x \\ 0 \end{bmatrix}$ is feasible in the modified problem (36).

Therefore we have

$$\overline{\text{OPT}} \leq \epsilon \rho \left(\frac{1}{2} x^\top Q x + c^\top x \right) = \epsilon \rho \cdot \text{OPT}.$$

Given a feasible \bar{x} satisfying (34), we write $\bar{x} = \begin{bmatrix} \bar{x}_{1:n} \\ \tau \end{bmatrix}$ for some $\tau \geq 0$. Then we have

$$\epsilon \rho \left(\frac{1}{2} \bar{x}_{1:n}^\top Q \bar{x}_{1:n} + c^\top \bar{x}_{1:n} \right) + \tau \leq \overline{\text{OPT}} + \epsilon^2 \leq \epsilon \rho \cdot \text{OPT} + \epsilon^2.$$

Therefore

$$\frac{1}{2} \bar{x}_{1:n}^\top Q \bar{x}_{1:n} + c^\top \bar{x}_{1:n} \leq \text{OPT} + \epsilon \rho^{-1}.$$

We have

$$\tau \leq -\epsilon \rho \left(\frac{1}{2} \bar{x}_{1:n}^\top Q \bar{x}_{1:n} + c^\top \bar{x}_{1:n} \right) + \epsilon \rho \cdot \text{OPT} + \epsilon^2 \leq 3\epsilon$$

because $|\frac{1}{2} x^\top Q x + c^\top x| \leq LR(R+1)$ for all $x \in \mathcal{K}$.

Note that \bar{x} satisfies $A\bar{x}_{1:n} + (b - Ax^{(0)})\tau = b$. So

$$\|A\bar{x}_{1:n} - b\|_1 \leq \|b - Ax^{(0)}\|_1 \cdot \tau.$$

This finishes the proof. \square

The following lemma is a generalization of (Lee et al., 2019, Lemma D.3) to quadratic program, and with weight vector w .

Lemma F.18 (QP version of (Lee et al., 2019, Lemma D.3)). *Work under the setting of Theorem F.1. Suppose we have $\frac{s_i}{t} + w_i \nabla \phi_i(x_i) = \mu_i$ for all $i \in [n]$, $-Qx + A^\top y + s = c$ and $Ax = b$. If $\|\mu_i\|_{x_i}^* \leq w_i$ for all $i \in [n]$, then we have*

$$\frac{1}{2} x^\top Q x + c^\top x \leq \frac{1}{2} x^*{}^\top Q x^* + c^\top x^* + 4t\kappa$$

where $x^* = \arg \min_{Ax=b, x \in \mathcal{K}} (\frac{1}{2} x^\top Q x + c^\top x)$.

Proof. Let $x_\alpha = (1-\alpha)x + \alpha x^*$ for some α to be chosen. By Lemma F.16, we have $\langle \nabla \phi_w(x_\alpha), x^* - x_\alpha \rangle \leq \kappa$. (Note that ϕ_w is a κ -self-concordant barrier for \mathcal{K} .) Therefore we have

$$\begin{aligned} \frac{\kappa \alpha}{1-\alpha} &\geq \langle \nabla \phi_w(x_\alpha), x_\alpha - x \rangle \\ &= \langle \nabla \phi_w(x_\alpha) - \nabla \phi_w(x), x_\alpha - x \rangle + \langle \mu - \frac{s}{t}, x_\alpha - x \rangle \\ &\geq \sum_{i \in [n]} w_i \frac{\|x_{\alpha,i} - x_i\|_{x_i}^2}{1 + \|x_{\alpha,i} - x_i\|_{x_i}} + \langle \mu, x_\alpha - x \rangle - \frac{1}{t} \langle c - A^\top y + Qx, x_\alpha - x \rangle \end{aligned}$$

$$\geq \sum_{i \in [n]} w_i \frac{\alpha^2 \|x_i^* - x_i\|_{x_i}^2}{1 + \alpha \|x_i^* - x_i\|_{x_i}} - \alpha \sum_{i \in [n]} \|\mu_i\|_{x_i}^* \|x_i^* - x_i\|_{x_i} - \frac{\alpha}{t} \langle c + Qx, x^* - x \rangle.$$

First step is because $\langle \nabla \phi_w(x_\alpha), x^* - x_\alpha \rangle \leq \nu$. Second step is because $\mu = \frac{s}{t} + \nabla \phi_w(x)$. Third step is by Lemma F.16 and $c = -Qx + A^\top y + s$. Fourth step is by Cauchy-Schwarz and $Ax_\alpha = Ax$.

So we get

$$\begin{aligned} & \frac{1}{t} (x^\top Qx + c^\top x) \\ & \leq \frac{1}{t} (x^\top Qx^* + c^\top x^*) + \frac{\kappa}{1 - \alpha} + \sum_{i \in [n]} \|\mu_i\|_{x_i}^* \|x_i^* - x_i\|_{x_i} - \sum_{i \in [n]} w_i \frac{\alpha \|x_i^* - x_i\|_{x_i}^2}{1 + \alpha \|x_i^* - x_i\|_{x_i}} \\ & \leq \frac{1}{t} \left(\frac{1}{2} x^\top Qx + \frac{1}{2} x^{*\top} Qx^* + c^\top x^* \right) + \frac{\kappa}{1 - \alpha} + \sum_{i \in [n]} w_i \|x_i^* - x_i\|_{x_i} - \sum_{i \in [n]} w_i \frac{\alpha \|x_i^* - x_i\|_{x_i}^2}{1 + \alpha \|x_i^* - x_i\|_{x_i}} \\ & = \frac{1}{t} \left(\frac{1}{2} x^\top Qx + \frac{1}{2} x^{*\top} Qx^* + c^\top x^* \right) + \frac{\kappa}{1 - \alpha} + \sum_{i \in [n]} w_i \frac{\|x_i^* - x_i\|_{x_i}}{1 + \alpha \|x_i^* - x_i\|_{x_i}} \\ & \leq \frac{1}{t} \left(\frac{1}{2} x^\top Qx + \frac{1}{2} x^{*\top} Qx^* + c^\top x^* \right) + \frac{\kappa}{1 - \alpha} + \sum_{i \in [n]} \frac{w_i}{\alpha} \\ & \leq \frac{1}{t} \left(\frac{1}{2} x^\top Qx + \frac{1}{2} x^{*\top} Qx^* + c^\top x^* \right) + \frac{\kappa}{\alpha(1 - \alpha)}. \end{aligned}$$

First step is by rearranging terms in the previous inequality. Second step is by AM-GM inequality and $\|\mu_i\|_{x_i}^* \leq w_i$. Third step is by merging the last two terms. Fourth step is by bounding the last term. Fifth step is by $\sum_{i \in [n]} w_i \leq \sum_{i \in [n]} w_i \nu_i = \kappa$.

Finally,

$$\begin{aligned} \frac{1}{2} x^\top Qx + c^\top x & \leq \frac{1}{2} x^{*\top} Qx^* + c^\top x^* + \frac{\kappa t}{\alpha(1 - \alpha)} \\ & \leq \frac{1}{2} x^{*\top} Qx^* + c^\top x^* + 4\kappa t. \end{aligned}$$

First step is by rearranging terms in the previous inequality. Second step is by taking $\alpha = \frac{1}{2}$. This finishes the proof. \square

F.5 PROOF OF THEOREM F.1

In this section we combine everything and prove Theorem F.1.

Proof of Theorem F.1. Lemma F.17 shows that the initial x and s satisfies

$$\|\mu\|_{w,x}^* \leq \epsilon.$$

This implies $w_i^{-1} \|\mu_i\|_{x_i}^* \leq \epsilon$ because $w_i \geq 1$.

Because $\epsilon \leq \frac{1}{\lambda}$, we have

$$\Phi(x, s, t) = \sum_{i \in [n]} \cosh(\lambda w_i^{-1} \|\mu_i\|_{x_i}^*) \leq n \cosh(1) \leq \cosh(\lambda/64)$$

for the initial x and s , by the choice of λ .

Using Theorem F.15, we see that

$$\Phi(x, s, t) \leq \cosh(\lambda/64)$$

during the entire algorithm.

So at the end of the algorithm, we have $w_i^{-1} \|\mu_i\|_{x_i}^* \leq \frac{1}{64}$ for all $i \in [n]$. In particular, $\|\mu_i\|_{x_i}^* \leq w_i$ for all $i \in [n]$.

3078
3079
3080
3081
3082
3083
3084
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131

Therefore, applying Lemma F.18 we get

$$\begin{aligned} \frac{1}{2}x^\top Qx + c^\top x &\leq \frac{1}{2}x^{*\top} Qx^* + c^\top x^* + 4t\kappa \\ &\leq \frac{1}{2}x^{*\top} Qx^* + c^\top x^* + \epsilon^2 \end{aligned}$$

where we used the stop condition for t at the end.

So Lemma F.17 shows how to get an approximate solution for the original quadratic program with error $\epsilon LR(R + 1)$.

The number of iterations is because we decrease t by a factor of $1 - h$ every iteration, and the choice $h = \frac{\alpha}{64\sqrt{\kappa}}$. \square

G GAUSSIAN KERNEL SVM: ALMOST-LINEAR TIME ALGORITHM AND HARDNESS

In this section, we provide both algorithm and hardness for Gaussian kernel SVM problem. For the algorithm, we utilize a result due to Aggarwal & Alman (2022) in conjunction with our low-rank QP solver to obtain an $O(n^{1+o(1)} \log(1/\epsilon))$ time algorithm. For the hardness, we build upon the framework outlined in Backurs et al. (2017) and improve their results in terms of dependence on dimension d .

We start by proving a simple lemma that shows that if $K = UV^\top$ for low-rank U, V , then the quadratic objective $K \circ (yy^\top)$ also admits such a factorization via a simple scaling.

Lemma G.1. *Let $U, V \in \mathbb{R}^{n \times k}$ and $y \in \mathbb{R}^n$. Then, there exists a pair of matrices $\tilde{U}, \tilde{V} \in \mathbb{R}^{n \times k}$ such that*

$$\tilde{U}\tilde{V}^\top = (UV^\top) \circ (yy^\top)$$

moreover, \tilde{U}, \tilde{V} can be computed in time $O(nk)$.

Proof. The proof relies on the following identity for Hadamard product: for any matrix A and conforming vectors x, y (all real), one has

$$A \circ (yx^\top) = D_y A D_x$$

where $D_y, D_x \in \mathbb{R}^{n \times n}$ are diagonal matrices that put y, x on their diagonals. Thus, we can simply compute \tilde{U}, \tilde{V} as follows:

$$\begin{aligned} \tilde{U} &= D_y U, \\ \tilde{V} &= D_y V, \end{aligned}$$

consequently,

$$\begin{aligned} \tilde{U}\tilde{V}^\top &= D_y UV^\top D_y \\ &= (yy^\top) \circ (UV^\top) \\ &= (UV^\top) \circ (yy^\top), \end{aligned}$$

as desired. Moreover, the diagonal scaling of U, V can be indeed performed in $O(nk)$ time, as advertised. \square

Throughout this section, we will let B denote the squared radius of the dataset.

G.1 ALMOST-LINEAR TIME ALGORITHM FOR GAUSSIAN KERNEL SVM

We state a result due to Aggarwal & Alman (2022), in which they present an optimal-degree polynomial approximation to the function e^{-x} and consequentially, this produces an efficient approximate scheme to the Batch Gaussian Kernel Density Estimation problem.

We start by introducing a notion that captures the minimum degree polynomial that well-approximates e^{-x} :

3132 **Definition G.2.** Let $f : [0, B] \rightarrow \mathbb{R}$, we let $q_{B;\epsilon}(f) \in \mathbb{N}$ denote the minimum degree of a non-constant
3133 polynomial $p(x)$ such that
3134

$$\sup_{x \in [0, B]} |p(x) - f(x)| \leq \epsilon$$

3135 Utilizing the Chebyshev polynomial machinery together with the orthogonal polynomial families, Ag-
3136 garwal & Alman (2022) provides the following characterization on $q_{B;\epsilon}(f)$:
3137

3138 **Theorem G.3** (Theorem 1.2 of Aggarwal & Alman (2022)). Let $B \geq 1$ and $\epsilon \in (0, 1)$. Then
3139

$$q_{B;\epsilon}(e^{-x}) = \Theta(\max\{\sqrt{B \log(1/\epsilon)}, \frac{\log(1/\epsilon)}{\log(B^{-1} \log(1/\epsilon))}\})$$

3140 **Theorem G.4** (Corollary 1.7 of Aggarwal & Alman (2022)). Let $x_1, \dots, x_n \in \mathbb{R}^d$ be a dataset with
3141 squared radius B and $\epsilon \in (0, 1)$. Let $q = q_{B;\epsilon}(e^{-x})$. Let $K \in \mathbb{R}^{n \times n}$ be the Gaussian kernel matrix
3142 formed by x_1, \dots, x_n . Finally, let $k = \binom{2d+2q}{2q}$. Then, there exists a deterministic algorithm that
3143 computes a pair of matrices $U, V \in \mathbb{R}^{n \times k}$ such that for any vector $v \in \mathbb{R}^n$,
3144

$$\|Kv - UV^\top v\|_\infty \leq \epsilon \|v\|_1.$$

3145 Moreover, matrices U, V can be computed in time $O(nkd)$.
3146

3147 Even though ℓ_∞ error in terms of ℓ_1 norm of vector v seems quite weak, it can be conveniently
3148 translated into more standard guarantees, e.g., spectral norm error. The following lemma provides
3149 a conversion of errors that come in handy later when integrating the kernel approximation to our
3150 low-rank QP solver.
3151

3152 **Lemma G.5.** Let $K \in \mathbb{R}^{n \times n}$ be a PSD kernel matrix and $\epsilon \in (0, 1)$ be a parameter. Let $\tilde{K} \in \mathbb{R}^{n \times n}$
3153 be an approximation to K with the guarantee that for any $v \in \mathbb{R}^n$,
3154

$$\|Kv - \tilde{K}v\|_\infty \leq \epsilon \|v\|_1,$$

3155 then
3156

$$|v^\top Kv - v^\top \tilde{K}v| \leq \epsilon \|v\|_1^2 \leq \epsilon n \|v\|_2^2.$$

3157 *Proof.* The proof is a simple application of Hölder's inequality:
3158

$$\begin{aligned} |v^\top (Kv - \tilde{K}v)| &= |\langle v, Kv - \tilde{K}v \rangle| \\ &\leq \|v\|_1 \|Kv - \tilde{K}v\|_\infty \\ &\leq \epsilon \|v\|_1^2 \\ &\leq \epsilon n \|v\|_2^2, \end{aligned}$$

3159 where the second step is by Hölder's inequality, and the last step is by Cauchy-Schwarz. This
3160 completes the proof. \square
3161

3162 We can now combine the Gaussian kernel low-rank decomposition with our low-rank QP solver
3163 to provide an almost-linear time algorithm for Gaussian kernel SVM. We restate the kernel SVM
3164 formulation here.
3165

3166 **Definition G.6** (Restatement of Definition 1.3). Given a data matrix $X \in \mathbb{R}^{n \times d}$ and labels $y \in \mathbb{R}^n$.
3167 Let $Q \in \mathbb{R}^{n \times n}$ denote a matrix where $Q_{i,j} = K(x_i, x_j) \cdot y_i y_j$ for $i, j \in [n]$. The hard-margin kernel
3168 SVM problem with bias asks to solve the following program.
3169

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^n} \quad & \mathbf{1}_n^\top \alpha - \frac{1}{2} \alpha^\top Q \alpha \\ \text{s.t.} \quad & \alpha^\top y = 0 \\ & \alpha \geq 0. \end{aligned}$$

3170 **Theorem G.7.** Let Gaussian kernel SVM training problem be defined as above with kernel function
3171 $K(x_i, x_j) = \exp(-\|x_i - x_j\|_2^2)$. Suppose the dataset has squared radius $B \geq 1$, and let $\epsilon \in (0, 1)$
3172 be the precision parameter. Suppose the program satisfies the following:
3173

3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212
3213
3214
3215
3216
3217
3218
3219
3220
3221
3222
3223
3224
3225
3226
3227
3228
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239

- There exists a point $z \in \mathbb{R}^n$ such that there is an Euclidean ball with radius r centered at z that is contained in the constraint set.
- The constraint set is enclosed by an Euclidean ball of radius R , centered at the origin.

Then, there exists a randomized algorithm that outputs an approximate solution $\hat{\alpha} \in \mathbb{R}^n$ such that $\hat{\alpha} \geq 0$, moreover,

$$\begin{aligned} \mathbf{1}_n^\top \hat{\alpha} - \frac{1}{2} \hat{\alpha}^\top Q \hat{\alpha} &\geq \text{OPT} - \epsilon, \\ \|\hat{\alpha}^\top y\|_1 &\leq 3\epsilon, \end{aligned}$$

where OPT denote the optimal cost of the objective function. Let $q = q_{B; \Theta(\epsilon/nR^2)}(e^{-x})$ and $k = \binom{2d+2q}{2q}$. Then, the vector $\hat{\alpha}$ can be computed in expected time

$$\tilde{O}(nk^{(\omega+1)/2} \log(nR/(\epsilon r))).$$

Proof. Throughout the proof, we set $\epsilon_1 = O(\epsilon/(nR^2))$. We will craft an algorithm that first computes an approximate Gaussian kernel together with a proper low-rank factorization, then use this proxy kernel matrix to solve the quadratic program. We will use K to denote the exact Gaussian kernel matrix, Q to denote the exact quadratic matrix.

Approximate the Gaussian kernel matrix with finer granularity. We start by invoking Theorem G.4 using data matrix X with accuracy parameter ϵ_1 . We let $\tilde{K} = UV^\top$ to denote this approximate kernel matrix, and we let $\tilde{Q} = D_y UV^\top D_y$ to denote the approximate quadratic matrix. Owing to Lemma G.5, we know that for any vector $x \in \mathbb{R}^n$,

$$\begin{aligned} |x^\top (Q - \tilde{Q})x| &= |(D_y x)^\top (K - \tilde{K})(D_y x)| \\ &\leq \epsilon_1 n \|D_y x\|_2^2 \\ &= \epsilon_1 n \|x\|_2^2, \end{aligned}$$

where we use the fact that $y \in \{\pm 1\}^n$. This also implies that

$$\|Q - \tilde{Q}\| \leq \epsilon_1 n \tag{37}$$

this simple bound will come in handy later.

Solving the approximate program to high precision. Given \tilde{Q} , we solve the following program:

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^n} \quad &\mathbf{1}_n^\top \alpha - \frac{1}{2} \alpha^\top \tilde{Q} \alpha \\ \text{s.t.} \quad &\alpha^\top y = 0 \\ &\alpha \geq 0 \end{aligned}$$

by invoking Theorem E.1. To do so, we need a bound on the Lipschitz constant of the program, i.e., the spectral norm of \tilde{Q} and ℓ_2 norm of $\mathbf{1}$. The latter is clearly \sqrt{n} , we shall show the first term is at most $(1 + \epsilon_1) \cdot n$.

Note that

$$\begin{aligned} \|Q\| &= \|D_y K D_y\| \\ &\leq \text{tr}[D_y K D_y] \\ &= \text{tr}[K] \\ &\leq n, \end{aligned}$$

where we use K is PSD. Combining with Eq. (37) and triangle inequality, we have

$$\begin{aligned} \|\tilde{Q}\| &\leq \|Q\| + \|Q - \tilde{Q}\| \\ &\leq (1 + \epsilon_1) \cdot n. \end{aligned}$$

3240 With these Lipschitz constants, we examine the error guarantee provided by Theorem E.1: it produces
 3241 a vector $\hat{\alpha} \in \mathbb{R}^n$ such that
 3242

$$3243 \mathbf{1}_n^\top \hat{\alpha} - \frac{1}{2} \hat{\alpha}^\top \tilde{Q} \hat{\alpha} \geq \max_{\alpha^\top y=0, \alpha \geq 0} (\mathbf{1}_n^\top \alpha - \frac{1}{2} \alpha^\top \tilde{Q} \alpha) - O(\epsilon_1 n R^2),$$

$$3244 \|\hat{\alpha}^\top y\|_1 \leq O(\epsilon_1 n R),$$

3246 we mainly focus on the first error bound, as we need to understand the quality of \hat{x} when plug into
 3247 the program with Q .
 3248

3249 We will follow a chain of triangle inequalities, so we first bound

$$3250 |\hat{\alpha}^\top (\tilde{Q} - Q) \hat{\alpha}| \leq \epsilon n \|\hat{\alpha}\|_2^2$$

$$3251 \leq \epsilon n R^2.$$

3253 Next, let

$$3254 \alpha' := \arg \max_{\alpha^\top y=0, \alpha \geq 0} \mathbf{1}_n^\top \alpha - \frac{1}{2} \alpha^\top \tilde{Q} \alpha,$$

$$3255 \alpha^* := \arg \max_{\alpha^\top y=0, \alpha \geq 0} \mathbf{1}_n^\top \alpha - \frac{1}{2} \alpha^\top Q \alpha,$$

3258 then we have the following

$$3260 \mathbf{1}_n^\top \alpha' - \frac{1}{2} \alpha'^\top \tilde{Q} \alpha' \geq \mathbf{1}_n^\top \alpha^* - \frac{1}{2} (\alpha^*)^\top \tilde{Q} \alpha^*$$

$$3261 \geq \mathbf{1}_n^\top \alpha^* - \frac{1}{2} (\alpha^*)^\top Q \alpha^* - O(\epsilon_1 n R^2)$$

$$3262 = \text{OPT} - O(\epsilon_1 n R^2),$$

3263 where the second step is by applying Lemma G.5 to α^* . Now we are ready to bound the final error:

$$3264 \mathbf{1}_n^\top \hat{\alpha} - \frac{1}{2} \hat{\alpha}^\top Q \hat{\alpha} \geq \mathbf{1}_n^\top \hat{\alpha} - \frac{1}{2} \hat{\alpha}^\top \tilde{Q} \hat{\alpha} - O(\epsilon_1 n R^2)$$

$$3265 \geq \mathbf{1}_n^\top \alpha' - \frac{1}{2} \alpha'^\top \tilde{Q} \alpha' - O(\epsilon_1 n R^2)$$

$$3266 \geq \text{OPT} - O(\epsilon_1 n R^2).$$

3267 The final error guarantee follows by the choice of ϵ_1 , and we indeed design an algorithm that outputs
 3268 a vector $\hat{\alpha}$ with

$$3269 \mathbf{1}_n^\top \hat{\alpha} - \frac{1}{2} \hat{\alpha}^\top Q \hat{\alpha} \geq \text{OPT} - \epsilon,$$

$$3270 \|\hat{\alpha}^\top y\|_1 \leq \epsilon.$$

3271 **Runtime analysis.** It remains to analyze the runtime of our proposed algorithm. We first compute an
 3272 approximate kernel \tilde{K} with parameter ϵ_1 , owing to Theorem G.4, we have

$$3273 q_{B, \epsilon_1}(e^{-x}) = \Theta(\max\{\sqrt{B \log(nR/\epsilon)}, \frac{\log(nR/\epsilon)}{\log(B^{-1} \log(nR/\epsilon))}\})$$

3274 then by setting $k = \binom{2d+2q}{2q}$, the matrix \tilde{K} can be computed in time $O(nkd)$. Given this rank- k
 3275 factorization, the program can then be solved with precision ϵ_1 in time

$$3276 \tilde{O}(nk^{(\omega+1)/2} \log(nR/(\epsilon r))),$$

3277 as desired. □

3278 **Remark G.8.** To understand the value range of k , let us consider the set of parameters:

$$3279 d = O(\log n), \epsilon = 1/\text{poly } n, R = \text{poly } n, B = \Theta(1),$$

3294 under this setting, $O(\log(nR/\epsilon)) = O(\log n)$ and the degree q is

3295
3296
$$q = \Theta(\sqrt{\log n})$$

3297 the rank k is then

3298
3299
$$k = \binom{2d + 2q}{2q}$$

3300
3301
$$\leq \Theta((\log n)^{\frac{1}{2}} \sqrt{\log n})$$

3302
$$= \Theta(2^{\Theta(\log \log n \sqrt{\log n})})$$

3303
3304
$$= n^{o(1)},$$

3305 consequentially, our algorithm runs in almost-linear time in n :

3306
$$\tilde{O}(n^{1+o(1)} \log n).$$

3308 It is worth noting to achieve the almost-linear runtime, the data radius B can be further relaxed. In
3309 fact, as long as

3310
$$B = o\left(\frac{\log n}{\log \log n}\right),$$

3311 we can ensure that $k = n^{o(1)}$ and subsequently the almost-linear runtime.

3312
3313
3314 The runtime we obtain can be viewed as a consequence of the “phase transition” phenomenon
3315 observed in Aggarwal & Alman (2022), in which they prove that if $B = \omega(\log n)$, then quadratic
3316 time in n is essentially needed to approximate the Gaussian kernel assuming **SETH**.

3318 G.2 HARDNESS OF GAUSSIAN KERNEL SVM WITH LARGE RADIUS

3319
3320 In this section, we show that for $d = O(\log n)$, any algorithm that solves the program associated to
3321 hard-margin Gaussian kernel SVM would require $\Omega(n^{2-o(1)})$ time for $B = \omega(\log n)$. This justifies
3322 the choice of B in Remark G.8. To prove the hardness result, we need to introduce the approximate
3323 Hamming nearest neighbor problem.

3324 **Definition G.9.** For $\delta > 0$ and $n, d \in \mathbb{N}$, let $\{a_1, \dots, a_n\}, \{b_1, \dots, b_n\} \subseteq \{0, 1\}^d$ be two sets
3325 of vectors, and let $t \in \{0, 1, \dots, d\}$ be a threshold. The $(1 + \delta)$ -Approximate Hamming Nearest
3326 Neighbor problem asks to distinguish the following two cases:

- 3327
3328
 - If there exists some a_i and b_j such that $\|a_i - b_j\|_1 \leq t$, output “Yes”;
 - If for any $i, j \in [n]$, we have $\|a_i - b_j\|_1 > (1 + \delta) \cdot t$, output “No”.

3329
3330 Note that the algorithm can output any value if the datasets fall in neither of these two cases. We will
3331 utilize the following hardness result due to Rubinfeld.

3332
3333 **Theorem G.10 (Rubinfeld (2018)).** Assuming **SETH**, for every $q > 0$, there exists $\delta > 0$ and $C > 0$
3334 such that $(1 + \delta)$ -Approximate Hamming Nearest Neighbor in dimension $d = C \log n$ requires time
3335 $\Omega(n^{2-q})$.

3336
3337 A final ingredient is a rewriting of the dual SVM into its primal form, without resorting to optimize
3338 over an infinite-dimensional hyperplane.

3339 **Lemma G.11.** Consider the dual hard-margin kernel SVM defined as

3340
$$\max_{\alpha \in \mathbb{R}^n} \mathbf{1}^\top \alpha - \frac{1}{2} \sum_{i,j \in [n] \times [n]} \alpha_i \alpha_j y_i y_j K(w_i, w_j),$$

3341
3342 s.t. $\alpha^\top y = 0,$
3343 $\alpha \geq 0.$

3344
3345 The primal program can be written as

3346
$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{2} \sum_{i,j \in [n] \times [n]} \alpha_i \alpha_j y_i y_j K(w_i, w_j),$$

3347

3348

$$\begin{aligned} \text{s.t. } & y_i f(w_i) \geq 1, \\ & \alpha \geq 0, \end{aligned}$$

3349

3350

3351

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is defined as

3352

3353

3354

3355

$$f(w) = \sum_{j=1}^n \alpha_j y_j \mathbf{K}(w_j, w) - b.$$

3356

Moreover, the primal and dual program has no duality gap and the optimal solution α to both programs are the same.

3357

3358

3359

Proof. Recall that the primal hard-margin SVM is the following program:

3360

3361

3362

3363

$$\begin{aligned} \min_v & \frac{1}{2} \|v\|_2^2, \\ \text{s.t. } & y_i (v^\top \phi(w_i) - b) \geq 1, \end{aligned}$$

3364

3365

3366

3367

where $b \in \mathbb{R}$ is the bias term and $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^K$ is the feature mapping corresponding to the kernel in the sense that $\mathbf{K}(w_i, w_j) = \phi(w_i)^\top \phi(w_j)$. The optimal weight $v = \sum_{i=1}^n \alpha_i y_i \phi(w_i)$ where $\alpha \in \mathbb{R}^n$ is the optimal solution to the dual program. Consequently,

3368

3369

3370

3371

3372

3373

3374

3375

3376

$$\begin{aligned} \|v\|_2^2 &= \left(\sum_{i=1}^n \alpha_i y_i \phi(w_i) \right)^2 \\ &= \sum_{i,j \in [n] \times [n]} \alpha_i \alpha_j y_i y_j \phi(w_i)^\top \phi(w_j) \\ &= \sum_{i,j \in [n] \times [n]} \alpha_i \alpha_j y_i y_j \mathbf{K}(w_i, w_j) \\ &= \alpha^\top Q \alpha, \end{aligned}$$

3377

where the matrix Q is the usual

3378

3379

$$Q = (y y^\top) \circ K,$$

3380

3381

the constraint can be rewritten as

3382

3383

3384

3385

3386

3387

3388

3389

3390

3391

3392

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is defined as

3393

3394

3395

3396

3397

$$f(w) = \sum_{j=1}^n \alpha_j y_j \mathbf{K}(w_j, w) - b.$$

Thus, we can alternatively write the primal as

3398

3399

3400

3401

$$\begin{aligned} \min_{\alpha \in \mathbb{R}^n} & \frac{1}{2} \alpha^\top Q \alpha, \\ \text{s.t. } & y_i f(w_i) \geq 1. \end{aligned}$$

For the strong duality and optimal solution, see, e.g., [Muller et al. \(2001\)](#). □

We will now prove the almost-quadratic lower bound for Gaussian kernel SVM. Our proof strategy is similar to that of [Backurs et al. \(2017\)](#) with different set of parameters. It is also worth noting that the [Backurs et al. \(2017\)](#) construction

- Requires the dimension $d = \Theta(\log^3 n)$;
- Requires the squared dataset radius $B = \Theta(\log^4 n)$.

We will improve both of these results.

Theorem G.12. *Assuming **SETH**, for every $q > 0$, there exists a hard-margin Gaussian kernel SVM without the bias term as defined in [Definition 1.3](#) with $d = \Theta(\log n)$ and error $\epsilon = \exp(-\Theta(\log^2 n))$ for inputs whose squared radius is at most $B = \Theta(\log^2 n)$ requiring time $\Omega(n^{2-q})$ to solve.*

Proof. Let $l = \sqrt{2(c'\delta)^{-1} \log n}$. We will provide a reduction from $(1 + \delta)$ -Approximate Hamming Nearest Neighbor to Gaussian kernel SVM. Let $A := \{a_1, \dots, a_n\}, B := \{b_1, \dots, b_n\} \subseteq \{0, 1\}^d$ be the datasets, we assign label 1 to all vectors a_i and label -1 to all vectors b_j , moreover, we scale both A and B by l , this results in two datasets with points in $\{0, l\}^d$. The squared radius of this dataset is then

$$\begin{aligned} B &= \max\{\max_{i,j} \|la_i - la_j\|_2^2, \max_{i,j} \|lb_i - lb_j\|_2^2, \max_{i,j} \|la_i - lb_j\|_2^2\} \\ &\leq l^2 d \\ &= \Theta(\delta^{-1} \log^2 n). \end{aligned}$$

To simplify the notation, we will implicitly assume A and B are scaled by l without explicitly writing out la_i, lb_j . Now consider the following three programs:

- Classifying A :

$$\begin{aligned} \min_{\alpha \in \mathbb{R}_{\geq 0}^n} \quad & \frac{1}{2} \sum_{i,j \in [n] \times [n]} \alpha_i \alpha_j \mathbf{K}(a_i, a_j), \\ \text{s.t.} \quad & \sum_{j=1}^n \alpha_j \mathbf{K}(a_i, a_j) \geq 1, \quad \forall i \in [n] \end{aligned} \quad (38)$$

- Classifying B :

$$\begin{aligned} \min_{\beta \in \mathbb{R}_{\geq 0}^n} \quad & \frac{1}{2} \sum_{i,j \in [n] \times [n]} \beta_i \beta_j \mathbf{K}(b_i, b_j), \\ \text{s.t.} \quad & -\sum_{j=1}^n \beta_j \mathbf{K}(b_i, b_j) \leq -1, \quad \forall i \in [n] \end{aligned} \quad (39)$$

- Classifying both A and B :

$$\begin{aligned} \min_{\alpha, \beta \in \mathbb{R}_{\geq 0}^n} \quad & \frac{1}{2} \sum_{i,j \in [n] \times [n]} \alpha_i \alpha_j \mathbf{K}(a_i, a_j) + \frac{1}{2} \sum_{i,j \in [n] \times [n]} \beta_i \beta_j \mathbf{K}(b_i, b_j) - \sum_{i,j \in [n] \times [n]} \alpha_i \beta_j \mathbf{K}(a_i, b_j), \\ \text{s.t.} \quad & \sum_{j=1}^n \alpha_j \mathbf{K}(a_i, a_j) - \sum_{j=1}^n \beta_j \mathbf{K}(a_i, b_j) \geq 1, \quad \forall i \in [n], \\ & \sum_{j=1}^n \alpha_j \mathbf{K}(b_i, a_j) - \sum_{j=1}^n \beta_j \mathbf{K}(b_i, b_j) \leq -1, \quad \forall i \in [n] \end{aligned} \quad (40)$$

We will prove that the optimal solution α_i^* 's and β_i^* 's are both lower and upper bounded. Use $\text{Val}(A)$, $\text{Val}(B)$ and $\text{Val}(A, B)$ to denote the value of program (38), (39) and (40) respectively, then note that

$$\text{Val}(A) \leq \frac{n^2}{2}$$

3456 by plugging in $\alpha = \mathbf{1}$ and setting all vectors to be the same. On the other hand,
 3457

$$\begin{aligned}
 3458 \quad \text{Val}(A) &\geq \frac{1}{2} \sum_{i=1}^n (\alpha_i^*)^2 \mathsf{K}(a_i, a_i) \\
 3459 \\
 3460 \\
 3461 \quad &= \frac{1}{2} \sum_{i=1}^n (\alpha_i^*)^2. \\
 3462 \\
 3463
 \end{aligned}$$

3464 Combining these two, we can conclude that for any α_i^* , it must be $\alpha_i^* \leq n$. For the lower bound,
 3465 consider the inequality constraint for the i -th point:

$$\begin{aligned}
 3466 \quad \alpha_i^* + \sum_{j \neq i} \alpha_j^* \mathsf{K}(a_i, a_j) &\geq 1, \\
 3467 \\
 3468
 \end{aligned}$$

3469 to estimate $\mathsf{K}(a_i, a_j)$, note that $\|a_i - a_j\|_2^2 = \|a_i - a_j\|_1 \geq 1$ for $j \neq i$,⁹ and

$$\begin{aligned}
 3470 \quad \mathsf{K}(a_i, a_j) &= \exp(-l^2 \|a_i - a_j\|_2^2) \\
 3471 \quad &= \exp(-2(c'\delta)^{-1} \log n \|a_i - a_j\|_1) \\
 3472 \quad &\leq \exp(-2(c'\delta)^{-1} \log n) \\
 3473 \quad &\leq n^{-10}/100, \\
 3474 \\
 3475
 \end{aligned}$$

3476 combining with $\alpha_j^* \leq n$, we have

$$\begin{aligned}
 3477 \quad \alpha_i^* &\geq 1 - \sum_{j \neq i} \alpha_j^* \mathsf{K}(a_i, a_j) \\
 3478 \\
 3479 \quad &\geq 1 - n \cdot n \cdot n^{-10}/100 \\
 3480 \quad &\geq 1/2. \\
 3481
 \end{aligned}$$

3482 This lower bound on α_i^* is helpful when we attempt to lower bound $\text{Val}(A, B)$ with $\text{Val}(A) + \text{Val}(B)$.
 3483 Following the outline of [Backurs et al. \(2017\)](#), we consider the three dual programs:

- 3484
 3485 • Dual of classifying A :

$$\begin{aligned}
 3486 \\
 3487 \quad \max_{\alpha \in \mathbb{R}_{\geq 0}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j \mathsf{K}(a_i, a_j) \quad (41) \\
 3488 \\
 3489
 \end{aligned}$$

- 3490 • Dual of classifying B :

$$\begin{aligned}
 3491 \\
 3492 \quad \max_{\beta \in \mathbb{R}_{\geq 0}^n} \sum_{i=1}^n \beta_i - \frac{1}{2} \sum_{i,j} \beta_i \beta_j \mathsf{K}(b_i, b_j) \quad (42) \\
 3493 \\
 3494
 \end{aligned}$$

- 3495 • Dual of classifying A and B :

$$\begin{aligned}
 3496 \\
 3497 \quad \max_{\alpha, \beta \in \mathbb{R}_{\geq 0}^n} \sum_{i=1}^n \alpha_i + \sum_{i=1}^n \beta_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j \mathsf{K}(a_i, a_j) - \frac{1}{2} \sum_{i,j} \beta_i \beta_j \mathsf{K}(b_i, b_j) + \sum_{i,j} \alpha_i \beta_j \mathsf{K}(a_i, b_j) \\
 3498 \\
 3499 \quad (43) \\
 3500
 \end{aligned}$$

3501 as the SVM program exhibits strong duality, we know that the optimal value of the primal equals
 3502 to the dual, so we can alternatively bound $\text{Val}(A, B)$ using the dual program. Plug in α^*, β^* to
 3503 program (43), we have

$$\begin{aligned}
 3504 \quad \text{Val}(A, B) &\geq \sum_{i=1}^n \alpha_i^* + \sum_{i=1}^n \beta_i^* - \frac{1}{2} \sum_{i,j} \alpha_i^* \alpha_j^* \mathsf{K}(a_i, a_j) - \frac{1}{2} \sum_{i,j} \beta_i^* \beta_j^* \mathsf{K}(b_i, b_j) + \sum_{i,j} \alpha_i^* \beta_j^* \mathsf{K}(a_i, b_j) \\
 3505 \\
 3506 \quad &= \text{Val}(A) + \text{Val}(B) + \sum_{i,j} \alpha_i^* \beta_j^* \mathsf{K}(a_i, b_j), \\
 3507 \\
 3508 \\
 3509
 \end{aligned}$$

⁹We without loss of generality that during preprocess, we have remove duplicates in A and B .

3510 to bound the third term, we consider the pair (a_{i_0}, b_{j_0}) such that $\|a_{i_0} - b_{j_0}\|_1 \leq t - 1$, and note that

$$3511$$

$$3512 \quad \sum_{i,j} \alpha_i^* \beta_j^* \mathsf{K}(a_i, b_j) \geq \alpha_{i_0}^* \beta_{j_0}^* \mathsf{K}(a_{i_0}, b_{j_0})$$

$$3513$$

$$3514 \quad \geq \frac{1}{4} \exp(-2(c'\delta)^{-1} \log n \cdot (t - 1)).$$

$$3515$$

3516 To wrap up, we have

$$3517 \quad \text{Val}(A, B) \geq \text{Val}(A) + \text{Val}(B) + \frac{1}{4} \exp(-2(c'\delta)^{-1} \log n \cdot (t - 1))$$

$$3518$$

$$3519$$

3520 We now prove the “No” case, where for any a_i, b_j , $\|a_i - b_j\|_1 \geq t$. We have

$$3521 \quad \mathsf{K}(a_i, b_j) = \exp(-l^2 \|a_i - b_j\|_2^2)$$

$$3522 \quad \leq \exp(-2(c'\delta)^{-1} \log n \cdot t),$$

$$3523$$

3524 we let $m := \exp(-2(c'\delta)^{-1} \log n \cdot t)$, set $\alpha' = \alpha^* + 10n^2m$ and $\beta' = \beta^* + 10n^2m$, we let V to

$$3525 \quad \text{denote the value when evaluating program (40) with } \alpha', \beta'. \text{ We will essentially show that}$$

$$3526 \quad \text{Val}(A, B) \leq V$$

$$3527$$

3528 and

$$3529 \quad V \leq \text{Val}(A) + \text{Val}(B) + 400n^6m$$

$$3530$$

3531 chaining these two gives us a certificate for the “No” case. To prove the first assertion, we show that

$$3532 \quad \alpha', \beta' \text{ are feasible solution to program (40) since}$$

$$3533 \quad \sum_{j=1}^n \alpha'_j \mathsf{K}(a_i, a_j) = \sum_{j=1}^n (\alpha_j^* \mathsf{K}(a_i, a_j) + 10n^2m \mathsf{K}(a_i, a_j))$$

$$3534$$

$$3535 \quad = \alpha_i^* + 10n^2m + \sum_{j \neq i} (\alpha_j^* + 10n^2m) \mathsf{K}(a_i, a_j)$$

$$3536$$

$$3537 \quad \geq \alpha_i^* + \sum_{j \neq i} \alpha_j^* \mathsf{K}(a_i, a_j) + 10n^2m$$

$$3538$$

$$3539 \quad = 10n^2m + \sum_{j=1}^n \alpha_j^* \mathsf{K}(a_i, a_j)$$

$$3540$$

$$3541 \quad \geq 10n^2m + 1$$

$$3542$$

$$3543$$

$$3544$$

3545 where we use α_i^* satisfy the inequality constraint of program (38). We compute an upper bound on

$$3546 \quad \sum_{j=1}^n \beta'_j \mathsf{K}(a_i, b_j):$$

$$3547$$

$$3548 \quad \sum_{j=1}^n \beta'_j \mathsf{K}(a_i, b_j) \leq \sum_{j=1}^n 2nm$$

$$3549$$

$$3550 \quad = 2n^2m,$$

$$3551$$

3552 where we use the fact that $m = \exp(-2(c'\delta)^{-1} \log n \cdot t) \leq n^{-10}/10$ therefore $\beta^* + 10n^2m \leq 2n$.

3553 Thus, it must be the case that

$$3554 \quad \sum_{j=1}^n \alpha'_j \mathsf{K}(a_i, a_j) - \sum_{j=1}^n \beta'_j \mathsf{K}(a_i, b_j) \geq 8n^2m + 1$$

$$3555$$

$$3556 \quad \geq 1,$$

$$3557$$

$$3558$$

3559 as desired. The other linear constraint follows by a symmetric argument. This indeed shows that

3560 α', β' are feasible solutions to program (40) and $\text{Val}(A, B) \leq V$.

3561 To prove an upper bound on V , we note that

$$3562 \quad V = \frac{1}{2} \sum_{i,j} \alpha'_i \alpha'_j \mathsf{K}(a_i, a_j) + \frac{1}{2} \sum_{i,j} \beta'_i \beta'_j \mathsf{K}(b_i, b_j) - \sum_{i,j} \alpha'_i \beta'_j \mathsf{K}(a_i, b_j)$$

$$3563$$

3564
3565
3566
3567
3568
3569
3570
3571
3572
3573
3574
3575
3576
3577
3578
3579
3580
3581
3582
3583
3584
3585
3586
3587
3588
3589
3590
3591
3592
3593
3594
3595
3596
3597
3598
3599
3600
3601
3602
3603
3604
3605
3606
3607
3608
3609
3610
3611
3612
3613
3614
3615
3616
3617

$$\leq \frac{1}{2} \sum_{i,j} \alpha'_i \alpha'_j \mathsf{K}(a_i, a_j) + \frac{1}{2} \sum_{i,j} \beta'_i \beta'_j \mathsf{K}(b_i, b_j),$$

we bound the first quantity, as the second follows similarly:

$$\begin{aligned} \frac{1}{2} \sum_{i,j} \alpha'_i \alpha'_j \mathsf{K}(a_i, a_j) &= \frac{1}{2} \sum_{i,j} (\alpha_i^* \alpha_j^* + 10n^2 m (\alpha_i^* + \alpha_j^*) + 100n^4 m^2) \mathsf{K}(a_i, a_j) \\ &\leq \text{Val}(A) + \sum_{i,j} 10n^3 m \mathsf{K}(a_i, a_j) + \sum_{i,j} 100n^4 m^2 \mathsf{K}(a_i, a_j) \\ &\leq \text{Val}(A) + 10n^5 m + 100n^6 m^2 \\ &\leq \text{Val}(A) + 200n^6 m, \end{aligned}$$

we can thus conclude

$$V \leq \text{Val}(A) + \text{Val}(B) + 400n^6 m.$$

Chaining these two, we obtain the following threshold for the “No” case:

$$\text{Val}(A, B) \leq \text{Val}(A) + \text{Val}(B) + 400n^6 m.$$

Finally, we observe that

$$400n^6 \exp(-2(c'\delta)^{-1} \log n \cdot t) \ll \frac{1}{4} \exp(-2(c'\delta)^{-1} \log n \cdot (t - 1)),$$

we can therefore distinguish these two cases.

Note that when one considers solving the program with additive error, we need to make sure that the error is smaller than the distinguishing threshold, i.e.,

$$\begin{aligned} \epsilon &\leq \frac{1}{4} \exp(-2(c'\delta)^{-1} \log n \cdot (t - 1)) \\ &\leq \frac{1}{4} \exp(-2(c'\delta)^{-1} d \log n) \\ &= \exp(-\Theta(\log^2 n)), \end{aligned}$$

where we use $t \leq d$ and $d = \Theta(\log n)$. This concludes the proof. \square

Remark G.13. *Our proof can be interpreted as using a stronger complexity theoretical tool in place of the one used by [Backurs et al. \(2017\)](#), to obtain a better dependence on dimension d and B . We also note that the construction due to [Backurs et al. \(2017\)](#) has the relation that $B = \Theta(d \log n)$, this is because in order to lower bound $\text{Val}(A, B)$, one has to lower bound the optimal values of α_i^* 's and β_j^* 's. To do so, one needs to further scale up a_i 's and b_j 's so that within datasets A and B , the radius is at least $\Theta(\log n)$. This is in contrast to the Batch Gaussian KDE studied in [Aggarwal & Alman \(2022\)](#), where they show the almost-quadratic lower bound can be achieved for both $d, B = \Theta(\log n)$.*

Similar to [Backurs et al. \(2017\)](#), we obtain hardness results for hard-margin kernel SVM with bias, and soft-margin kernel SVM with bias.

Corollary G.14. *Assuming SETH , for every $q > 0$, there exists a hard-margin Gaussian kernel SVM with the bias term with $d = \Theta(\log n)$ and error $\epsilon = \exp(-\Theta(\log^2 n))$ for inputs whose squared radius is at most $B = \Theta(\log^6 n)$ requiring time $\Omega(n^{2-q})$ to solve.*

Proof. The proof is similar to [Backurs et al. \(2017\)](#). Given a hard instance of [Theorem G.12](#), except we append $\Theta(\log n)$ entries with magnitude $\Theta(\log^2 n)$ instead of distributing the values across $\Theta(\log^3 n)$ entries. Rest of the proof follows exactly the same as [Backurs et al. \(2017\)](#). \square

Corollary G.15. *Assuming SETH , for every $q > 0$, there exists a soft-margin Gaussian kernel SVM with the bias term with $d = \Theta(\log n)$ and error $\epsilon = \exp(-\Theta(\log^2 n))$ for inputs whose squared radius is at most $B = \Theta(\log^6 n)$ requiring time $\Omega(n^{2-q})$ to solve.*

3618 **Remark G.16.** Compared to the construction of *Backurs et al. (2017)* in which they distribute a
3619 total mass of $\Theta(\log^3 n)$ across $\Theta(\log^3 n)$ entries so that they ensure after the reduction, the vectors
3620 take values in $\{-1, 0, 1\}$, we instead distribute the mass across $\Theta(\log n)$ entries so that each entry
3621 has magnitude $\Theta(\log^2 n)$. To make the reduction work, the total mass of $\Theta(\log^3 n)$ is needed, and
3622 for *Backurs et al. (2017)*, it is fine to append an extra $\Theta(\log^3 n)$ entries as their hardness result for
3623 hard-margin SVM without bias does require $d = \Theta(\log^3 n)$. For us, we need to restrict $d = \Theta(\log n)$
3624 at the price of each entry has a larger magnitude of $\Theta(\log^2 n)$. This blows up the squared radius
3625 from $\log^2 n$ to $\log^6 n$. We note that the *Backurs et al. (2017)* construction has squared radius $\log^4 n$.
3626

3627
3628
3629
3630
3631
3632
3633
3634
3635
3636
3637
3638
3639
3640
3641
3642
3643
3644
3645
3646
3647
3648
3649
3650
3651
3652
3653
3654
3655
3656
3657
3658
3659
3660
3661
3662
3663
3664
3665
3666
3667
3668
3669
3670
3671