

Differentially Private Stochastic Convex Optimization for Network Routing Applications

Matthew Tsao
Lyft Inc.
mtsao@lyft.com

Karthik Gopalakrishnan
Stanford University
gkarthik@stanford.edu

Kaidi Yang
National University of Singapore
kaidi.yang@nus.edu.sg

Marco Pavone
Stanford University
pavone@stanford.edu

October 27, 2022

Abstract

Network routing problems are common across many engineering applications. Computing optimal routing policies requires knowledge about network demand, i.e., the origin and destination (OD) of all requests in the network. However, privacy considerations make it challenging to share individual OD data that would be required for computing optimal policies. Privacy can be particularly challenging in standard network routing problems because sources and sinks can be easily identified from flow conservation constraints, making feasibility and privacy mutually exclusive.

In this paper, we present a differentially private algorithm for network routing problems. The main ingredient is a reformulation of network routing which moves all user data-dependent parameters out of the constraint set and into the objective function. We then present an algorithm for solving this formulation based on a differentially private variant of stochastic gradient descent. In this algorithm, differential privacy is achieved by injecting noise, and one may wonder if this noise injection compromises solution quality. We prove that our algorithm is both differentially private and asymptotically optimal as the size of the training set goes to infinity. We corroborate the theoretical results with numerical experiments on a road traffic network which show that our algorithm provides differentially private and near-optimal solutions in practice.

1 Introduction

Network routing problems appear in many important topics in engineering, including traffic routing in transportation systems, power routing in electrical grids, and packet routing in distributed computer systems. Network routing problems study settings where resources must be delivered to customers through a network with limited bandwidth. The goal is typically to route resources to their respective customers as efficiently as possible, or equivalently, with as little network congestion as possible.

One key challenge in network routing problems is that the requests (i.e., network demand) are often not known in advance. Indeed, it is difficult to know exactly how much power a neighborhood will need, or exactly how many visits a particular website will receive on any given day. Since information about the demand is often necessary to develop optimal or near-optimal network routing solutions, network routing algorithms often need a way of obtaining or estimating future demand. With the advent of big data and internet-of-things systems, crowd-sourcing has gained popularity as a demand forecasting approach for network routing systems. In crowd-sourcing, customers submit their request history to the network operator. The network operator uses this historical data to train a statistical or machine learning model to predict future demand from historical demand.

While crowd-sourcing provides a bountiful supply of data for training demand forecasting models, it can also introduce potential privacy concerns. Since crowd-sourcing uses individual-level customer data to train demand forecasting models, the model’s outputs may reveal sensitive user information, especially if it overfits to its training data [CTW⁺21]. Such privacy risks are problematic because they may deter users from sharing their data with network operators, hence reducing the supply of training data for demand forecasting models.

To address such concerns, the demand forecasting pipeline should be augmented with privacy-preserving mechanisms. Differential privacy [DMNS06] is a principled and popular method to occlude the influence a single user’s data on the result of a population study while also maintaining the study’s accuracy. This is done by carefully injecting noise into the desired computation so that data sets that differ by at most one data point will produce statistically indistinguishable results.

Providing differential privacy guarantees for the standard formulation of network routing is difficult because the constraints contain user data, meaning that in general feasibility and privacy become mutually exclusive. More specifically, in the standard network routing problem, the demand sources and sinks can be identified by checking for conservation of flow, and as a result, the presence of a user going to or from a very rare location can be detected from any feasible solution. Because differential privacy requires that the presence of any single user’s data be difficult to detect from the algorithm’s output, privacy and feasibility are at odds with one another in the standard formulation.

1.1 Statement of Contributions

In this paper we present a differentially private algorithm for network routing problems. The main ingredient is a reformulation of network routing which moves all user data dependent parameters out of the constraint set and into the objective function. We then present an algorithm for solving this formulation based on differentially private variants of stochastic gradient descent. In this algorithm, differential privacy is achieved by injecting noise, and one may wonder if this noise injection compromises solution quality. We prove that our algorithm is differentially private and under several reasonable regularity conditions, is also asymptotically optimal (as the size of the training set goes to infinity). We note that in exchange for becoming compatible with differentially private algorithms, this new formulation is more computationally expensive.

1.2 Related Work

Traffic assignment in transportation systems is one of the most well-known applications of network routing. Herein we focus our literature review on privacy research in transportation networks. Privacy works in transportation mainly focus on *location privacy*, where the objectives is to prevent untrusted and/or external entities from learning geographic locations or location sequences of an

individual [BS03]. Privacy-preserving approaches have been proposed for various location-based applications, e.g., trajectory publishing, mobile crowdsensing, traffic control, etc. These techniques are based on spatial cloaking [CML11] and differential privacy [Dwo08]. While not the setting of interest for this paper, there are many works that use Secure Multi-Party Computation (MPC) [GMW87] to achieve privacy in *distributed* mobility systems.

Spatial cloaking approaches aggregate users' exact locations into coarse information. These approaches are often based on k -anonymity [Swe02], where a mobility dataset is divided into equivalence classes based on data attributes (e.g., geological regions, time, etc.) so that each class contains at least k records [GFCA14, HC20]. These k -anonymity-based approaches can guarantee that every record in the dataset is indistinguishable from at least $k - 1$ other records. However, k -anonymity is generally considered to be a weak privacy guarantee, especially when k is small. Furthermore, very coarse data aggregation is required to address outliers or sparse data, and in these cases spatial cloaking-based approaches provide low data accuracy.

Differential privacy provides a principled privacy guarantee by producing randomized responses to queries, whereby two datasets that differ in only one entry produce statistically indistinguishable responses [DMNS06]. In other words, differential privacy ensures that an adversary with arbitrary background information (e.g., query responses, other entries) cannot infer individual entries with high confidence. Within transportation research, [WHL⁺18, YLL⁺19] share noisy versions of location data for mobile crowd-sensing applications. [GZFS15, GLTY18, AHFI⁺18, LYH20] use differential privacy to publish noisy versions of trajectory data. [DKBS15] and [HTP17] apply differential privacy to gradient descent algorithms for federated learning in mobility systems.

Many of the works mentioned in the previous paragraph establish differential privacy of their proposed algorithms by using composition properties of differential privacy. Composition theorems for differential privacy describe how well privacy is preserved when conducting several computations one after another. In [DKBS15] and [HTP17], composition theorems are applied as black boxes without considering the mathematical properties of the gradient descent algorithm. As a result, the privacy guarantees are overly conservative, meaning that large amounts of noise are added to the algorithm, leading to suboptimal behavior both in theory and in practice. Similarly, [GZFS15, GLTY18, AHFI⁺18, LYH20] use composition rules as a black box, and while privacy is achieved in this way, there are no accuracy guarantees for the algorithms presented in those works.

While blackbox applications of differential privacy can lead to impractical levels of noise injection, specialized applications of differential privacy were discovered that could provide privacy with much less noise. [WLK⁺17] show how a simple adjustment to stochastic gradient descent can give rise to an algorithm which is both differentially private, and under reasonable regularity assumptions, is also asymptotically optimal. [FMTT18] and [FKT20] refined this idea to develop stochastic gradient descent algorithms that are differentially private, computationally efficient, and have optimal convergence rates. These techniques cannot directly be used to solve the standard formulation of network routing because they study unconstrained optimization problems or optimization problems with public constraint sets (i.e., constraints that do not include any private data).

2 Model

In this section we define notations, network models, and assumptions that will allow us to formulate network routing as a data-driven convex optimization problem.

2.1 Preliminaries

Indicator Representation of Edge Sets: Let $G = (V, E)$ be a graph with vertices V and edges $E = \{e_1, \dots, e_m\}$. For any subset of edges $E' \subset E$, we define the indicator representation of E' as $\mathbb{1}_{[E']}$ as a boolean vector of length m in the following way: The i th entry of $\mathbb{1}_{[E']}$ is 1 if $e_i \in E'$, and is 0 otherwise.

Derivative Notation: For a scalar valued function $x \mapsto f(x)$, we use $\nabla f(x_0)$ to denote the gradient of the f with respect to x evaluated at the point $x = x_0$. For a vector valued function $x \mapsto g(x)$, and a variable z , we use $\mathcal{D}_z[g](x_0)$ to denote the derivative matrix of g with respect to z evaluated at the point $x = x_0$.

Projections: For a convex set $\mathcal{S} \subset \mathbb{R}^d$, we use $\Pi_{\mathcal{S}} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ to denote the projection operator for \mathcal{S} . For any $x \in \mathbb{R}^d$, $\Pi_{\mathcal{S}}(x) := \arg \min_{s \in \mathcal{S}} \|x - s\|_2$ to be the point in S that is closest to x .

2.2 Network, Demand, and Network Flow

In this section we will introduce a) a graph model of a network, b) a representation of network demand, c) the standard formulation for network routing and d) privacy requirements. The notation defined in this section is aggregated in table form in Section A for the reader's convenience.

Definition 1 (Network Representation). We use a directed graph $G = (V, E)$ to represent the network, where V and E represent the set of vertices and edges in the network, respectively. We will use $n := |V|$ and $m := |E|$ to denote the number of vertices and edges in the graph, respectively. For vertex pairs $(o, d) \in V \times V$ we will use $\mathcal{P}_G(o, d)$ to denote the set of simple paths from o to d in G .

Definition 2 (Operation Period). We use $\mathcal{T} := [t_{\text{start}}, t_{\text{end}}]$ to denote the operation period within which a network operator wants to optimize its routing decisions. We will also use T to denote the number of minutes in the operation period. For example, $t_{\text{start}} = 8 : 00\text{am}$, $t_{\text{end}} = 9 : 00\text{am}$ could represent a morning commute period where $T = 60$.

Definition 3 (Demand Representation). We study a stationary demand model where demand within the operation period \mathcal{T} is specified by a matrix $\Lambda \in \mathbb{R}_+^{n \times n}$. For each ordered pair of vertices $(o, d) \in V \times V$, $\Lambda(o, d)$ is the number of requests arriving per minute (i.e., the arrival rate) during \mathcal{T} that need routing from vertex o to vertex d .

Remark 1 (Estimating Λ from historical data). The arrival rates from historical demand are computed empirically, i.e., if Λ_t represents the demand for day t , then $\Lambda_t(o, d)$ is calculated by counting the number of (o, d) requests appearing on day t , and then dividing it by T to obtain requests per minute.

Definition 4 (Link Latency Functions). To model congestion effects, each link $e \in E$ has a latency function $f_e : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ which specifies the average travel time through the link as a function of the total flow on the link.

In this paper we study a setting where a network operator wants to route demand while minimizing the total travel time for the requests. With these definitions, the standard formulation of minimum travel time network routing is described in Definition 5.

Definition 5 (Standard Formulation of Network Flow). For a network $G = (V, E)$ with link latency functions $\{f_e\}_{e \in E}$ and demand Λ , the standard network flow problem is the following minimization program:

$$\underset{x}{\text{minimize}} \sum_{e \in E} y_e f_e(y_e) \quad (1a)$$

$$\text{subject to } x = \left\{ x^{(o,d)} \right\}_{(o,d) \in V \times V}, \quad (1b)$$

$$y_e = \sum_{(o,d) \in V \times V} x_e^{(o,d)} \text{ for all } e \in E, \quad (1c)$$

$$\sum_{v:(u,v) \in E} x_{(u,v)}^{(o,d)} - \sum_{w:(w,u) \in E} x_{(w,u)}^{(o,d)} = \Lambda(o,d) \mathbf{1}_{[u=o]} - \Lambda(o,d) \mathbf{1}_{[u=d]} \text{ for all } u \in V, (o,d) \in V \times V. \quad (1d)$$

In (1), the decision variable x is a collection of flows $\{x^{(o,d)}\}_{(o,d) \in V \times V}$, one for each non-zero entry of Λ , as represented by constraint (1b). (1d) are flow conservation constraints to ensure that $x^{(o,d)}$ sends $\Lambda(o,d)$ units of flow from o to d . Constraints (1c) ensure that $\{y_e\}_{e \in E}$ represents the total amount of flow on each edge. Finally, the objective function (1a) is to minimize the total travel time as a function of total network flow.

In the next subsection we will describe the rigorous privacy requirements that we will mandate while designing algorithms for network flow. We then describe in Section 2.4 why privacy and feasibility are mutually exclusive in the standard network flow formulation.

2.3 Privacy Requirements

We will use differential privacy to reason about the privacy-preserving properties of our algorithms. At a high level, changing one data point of the input to a differentially private algorithm should lead to a statistically indistinguishable change in the output. To make this concept concrete we will need to define data sets and adjacency.

Definition 6 (Data sets). Given a space of data points \mathcal{Z} , a data set L is any finite set of elements from \mathcal{Z} . In practice, each element of a data set is data collected from a single user, or data collected during a specific time period. We will use \mathcal{L} to denote the set of all data sets.

Definition 7 (Data Set Adjacency). Given a space of data sets \mathcal{L} , an adjacency relation Adj is a mapping $\text{Adj} : \mathcal{L} \times \mathcal{L} \rightarrow \{0, 1\}$. Two data sets L_1, L_2 are said to be adjacent if and only if $\text{Adj}(L_1, L_2) = 1$.

While the exact definition can vary across applications, adjacent data sets are sets that are very similar to one another. The most common definition of adjacency is the following: L, L' are adjacent if and only if L' can be obtained by adding, deleting, or modifying at most one data point from L , and vice versa. Thus comparing a function's output on adjacent data sets measures how sensitive the function is to changes in a single data point.

With these definitions in place, we are now ready to define differential privacy.

Definition 8 (Differential Privacy). For a given adjacency relation Adj , a function $M : \mathcal{L} \rightarrow \mathcal{X}$ is (ϵ, δ) -differentially private if for any $L_1, L_2 \in \mathcal{L}$ with $\text{Adj}(L_1, L_2) = 1$, the following inequality

$$\mathbb{P}[M(L_1) \in E] \leq e^\epsilon \mathbb{P}[M(L_2) \in E] + \delta$$

holds for every event $E \subset \mathcal{X}$.

The definition of differential privacy requires that changing one data point in the input to a differentially private algorithm cannot change the distribution of the algorithm's output by too much. Such a requirement ensures the following strong privacy guarantee: It is statistically impossible to reliably infer the properties of a single data point just by examining the algorithm's output, even if the observer knows all of the other data points in the input [DMNS06].

To proceed, we must specify what data set adjacency means in the context of network routing. For network routing problems, historical data is often a collection of routing requests through the network. To protect the privacy of those who submit requests, we want the routing policy we compute to not depend too much on any single request that appears in the historical data. This motivates the the following notion of data set adjacency that we will be using throughout this paper.

Definition 9 (Request Level Adjacency (RLA)). We define a function $\text{RLA} : \mathcal{L} \times \mathcal{L} \rightarrow \{0, 1\}$ which maps pairs of data sets to booleans. For two historical datasets of network demand $L := (\Lambda_1, \dots, \Lambda_N)$ and $L' := (\Lambda'_1, \dots, \Lambda'_N)$, we say L and L' are request-level-adjacent (RLA) if exactly one of the following statements is true:

1. L contains all of the requests in L' , and contains one extra request that is not present in L' .
2. L' contains all of the requests in L , and contains exactly one extra request that is not present in L .

Mathematically, L, L' are request-level-adjacent, i.e., $\text{RLA}(L, L') = 1$, if and only if they satisfy all of the following relations:

- There exists t so that $\Lambda_k = \Lambda'_k$ for all $k \neq t$.
- There exists two vertices o and d so that $\Lambda_t(o', d') = \Lambda'_t(o', d')$ for all $(o', d') \neq (o, d)$.
- $|\Lambda_t(o, d) - \Lambda'_t(o, d)| \leq \frac{1}{T}$.

Indeed, these relations dictate that one of the datasets contains an extra request from o to d which happened on the t th day. A difference of 1 request within a T minute operation period leads to a change of $\frac{1}{T}$ in the arrival rate. Aside from this difference, the datasets are otherwise identical.

2.4 Differential Privacy Challenges in Standard Network Flow

In the introduction we mentioned that privacy and feasibility can be mutually exclusive in the standard formulation of network flow described in (1). In this section we formally show that if Λ is constructed from a data set of trips as described in Remark 1, then trips to or from uncommon locations can be easily detected from any feasible solution to (1). As a result, announcing or releasing a feasible solution to (1) is not, in general, differentially private. Formally, we will prove the following theorem in this section.

Theorem 1 (Differential Privacy Impossibility for Standard Network Flow). *Let M be a function that takes as input a matrix Λ with non-negative entries and returns a feasible solution to the optimization problem (1) where Λ is used as a demand matrix. Then M cannot be (ϵ, δ) -differentially private for any $\delta < 1$.*

We further note that (ϵ, δ) -differential privacy only provides a meaningful privacy guarantee when $\epsilon < 1$ and $\delta < 1$ [Dwo08].

Proof of Theorem 1. Let $G = (V, E)$ be a network, and Λ be constructed from a historical data set of requests in G . Suppose there exist uncommon locations o and d for which Λ contains no trips to or from either o or d . Mathematically, this means that

$$\Lambda(o, u) = 0, \Lambda(u, o) = 0, \Lambda(d, u) = 0, \Lambda(u, d) = 0 \text{ for all } u \in V.$$

Such situations are not uncommon in transportation networks, if, for example, o and d are the homes of two different people who do not drive (perhaps they walk or bike to and from work).

If we now add a trip from o to d to the data set, and let Λ' be the resulting demand matrix, then we have

- i $\Lambda(o', d') = \Lambda'(o', d')$ for all $(o', d') \neq (o, d)$, and
- ii $\Lambda(o, d) = 0, \Lambda'(o, d) = \frac{1}{T}$.

Let $\text{Prob}_1, \text{Prob}_2$ be the optimization problem (1) using demand Λ, Λ' respectively. Because Λ, Λ' are request-level-adjacent, any differentially private algorithm must behave similarly when acting on Prob_1 and Prob_2 . However, this is impossible because the feasible sets of $\text{Prob}_1, \text{Prob}_2$ are disjoint. If we look at constraint (1d) with $u = d$ and (o, d) then any feasible solution to Prob_1 satisfies

$$\sum_{u:(u,d) \in E} x_{(u,d)}^{(o,d)} - \sum_{w:(d,w) \in E} x_{(u,d)}^{(o,d)} = \Lambda(o, d) \mathbf{1}_{[d=d]} = 0.$$

However, any feasible solution to Prob_2 satisfies

$$\sum_{u:(u,d) \in E} x_{(u,d)}^{(o,d)} - \sum_{w:(d,w) \in E} x_{(u,d)}^{(o,d)} = \Lambda(o, d) \mathbf{1}_{[d=d]} = \frac{1}{T}.$$

In other words, checking the net flow leaving node d will detect the presence of any trips going to or from d . We will now show that any algorithm which returns a feasible solution to (1) cannot be differentially private. To this end, define the event E to be the event that flow is conserved at node d . Then for any algorithm M that takes as input a demand matrix and returns a feasible solution to (1) with the specified demand matrix, we have $\mathbb{P}[M(\Lambda) \in E] = 1, \mathbb{P}[M(\Lambda') \in E] = 0$. Recalling Definition 8, M is (ϵ, δ) -differentially private only if $\mathbb{P}[M(\Lambda) \in E] \leq e^\epsilon \mathbb{P}[M(\Lambda') \in E] + \delta$. This equation can only be satisfied if $\delta \geq 1$. \square

We have two remarks regarding Theorem 1.

Remark 2. The same result holds if M returns the total flow y associated with a feasible solution (see (1c)), rather than returning the feasible solution itself. In other words, even total traffic measurements (without knowing the breakdown by (o, d) pairs) can already expose trips to or from uncommon locations.

Remark 3. The vulnerability of trips to and from uncommon locations is not a purely theoretical concern. A study on the New York City Taxi and Limousine data set was able to identify trips from residential areas to gentlemen’s club [Pan14]. Because the start locations were in residential areas, it was easy to re-identify the users who booked the taxi trips as the owners of the homes that the taxi trips began at. As a result, despite the data set being anonymized, users who had taken taxis to this gentlemen’s club were re-identified, and their reputations may have been negatively affected.

3 Routing Policy Formulation of Network Flow

To sidestep the impossibility result described in Theorem 1, we present an alternative formulation for the network flow problem in this section. The alternative formulation avoids the issues mentioned in Theorem 1 by moving all parameters related to user data from the constraints to the objective function, as described in (5). We note that (5) can only be solved if the demand Λ is known, which may not always be the case. For this reason, we present two variations of (5): (6) is the stochastic version of (5) where Λ is drawn from a distribution \mathcal{Q} , and (7) is the data driven approximation to (6) that one would solve if \mathcal{Q} is unknown.

Before formally defining the model, we provide a high level description of how this formulation works. In this formulation, a feasible solution $x = \{x^{(o,d)}\}_{(o,d) \in V \times V}$ specifies, for each $(o, d) \in V \times V$, a flow $x^{(o,d)}$ that routes 1 unit of flow from o to d . We note that a flow is specified for (o, d) even if there is no demand for this origin and destination in Λ , i.e., $\Lambda(o, d) = 0$. We refer to x as a *routing policy* due to its connection to randomized routing, which Remark 5 discusses in further detail. Given a feasible solution x , the objective function first calculates the total traffic on each edge by taking a linear combination of $\{x^{(o,d)}\}_{(o,d) \in V \times V}$ flows, where the coefficients of the linear combination are determined by the demand Λ , with high demand (o, d) pairs having larger coefficients. The total travel time can be computed from the total traffic in the same way as (1a). These ideas are formalized by the following definitions.

Definition 10 (Unit (o, d) flow). For a given origin-destination pair (o, d) , we say that a flow $x^{(o,d)} \in \mathbb{R}_+^m$ is a unit (o, d) flow if and only if it routes exactly 1 unit of flow from o to d . Formally, this condition is represented by the following constraints:

$$\sum_{v \in V: (v,u) \in E} x_{(v,u)}^{(o,d)} - \sum_{v \in V: (u,v) \in E} x_{(u,v)}^{(o,d)} = \begin{cases} -1 & \text{if } u = o \\ 1 & \text{if } u = d \\ 0 & \text{otherwise.} \end{cases} \quad \text{for all } u \in V \quad (2)$$

Indeed, (2) requires that the net flow entering o is -1 , the net flow entering d is 1, and that flow is conserved at all other vertices in the graph.

Definition 11 (Unit Network Flow). A unit network flow is a collection of flows $x = \{x^{(o,d)}\}_{(o,d) \in V \times V}$ so that $x^{(o,d)}$ is a unit (o, d) flow for each $(o, d) \in V \times V$. We use \mathcal{X} to denote the set of all unit network flows.

Remark 4. We can represent x as a concatenation of the vectors $\{x^{(o,d)}\}_{(o,d) \in V \times V}$. Since each unit (o, d) flow is a m dimensional vector, we have $x \in \mathbb{R}^{n^2 m}$.

We will refer to unit network flows as routing policies, due to their connection with randomized routing, described in the following remark.

Remark 5 (Routing demand using unit flow policies). A unit network flow represents a randomized routing policy. Due to the flow decomposition theorem [Tre11], every unit (o, d) flow can be written as a distribution of paths from o to d . Formally, for any unit (o, d) flow $x^{(o,d)}$, there exist paths $p_1^{(o,d)}, p_2^{(o,d)}, \dots, p_m^{(o,d)}$ from o to d and weights $w_1^{(o,d)}, \dots, w_m^{(o,d)}$ so that the following equations hold:

$$\begin{aligned} x^{(o,d)} &= \sum_{i=1}^m w_i^{(o,d)} \mathbb{1}_{[p_i^{(o,d)}]} & (3) \\ \sum_{i=1}^m w_i^{(o,d)} &= 1 \\ w_i^{(o,d)} &\geq 0 \text{ for all } 1 \leq i \leq m. \end{aligned}$$

Furthermore, $p_1^{(o,d)}, \dots, p_m^{(o,d)}$ and $w_1^{(o,d)}, \dots, w_m^{(o,d)}$ are efficiently computable. Defining the probability distribution $P_{x^{(o,d)}}$

$$P_{x^{(o,d)}}(\mathbb{1}_{[p_i^{(o,d)}]}) = w_i^{(o,d)} \text{ for all } 1 \leq i \leq m,$$

$x_e^{(o,d)}$ represents the probability that a random path chosen from $P_{x^{(o,d)}}$ contains e . x describes the expected behavior of the randomized routing policy that determines routes for (o, d) requests by drawing a path independently at random from $P_{x^{(o,d)}}$. In particular, when using this policy to serve demand Λ , by linearity of expectation, the expected number of requests from o to d whose assigned path contains e is exactly $\Lambda(o, d)x_e^{(o,d)}$. Furthermore, the average number of requests on each edge will be $\sum_{(o,d) \in V \times V} \Lambda(o, d)x_e^{(o,d)}$.

3.1 Minimum Total Travel Time Network Flow

In the minimum travel time network flow problem, the network operator wants to find a stationary routing policy for each (o, d) pair that will lead to small total travel times for the requests. Due to the equivalence between stationary (o, d) routing policies and unit (o, d) flows, the network operator can instead search over the space of unit (o, d) flows.

The total travel time of a flow y through G is given by $\sum_{e \in E} y_e f_e(y_e)$. The total flow resulting from a unit network flow x serving demand Λ is the sum of the flow contributions from each of the (o, d) pairs. With Remark 5 in mind, the total amount of flow on an edge $e \in E$ when serving Λ according to x is given by $\sum_{(o,d) \in V \times V} \Lambda(o, d)x_e^{(o,d)}$. We can thus define $F(x, \Lambda)$, the total travel time experienced by requests Λ when being routed by x , as follows:

$$F(x, \Lambda) := \sum_{e \in E} \left(\sum_{(o,d) \in V \times V} \Lambda(o, d)x_e^{(o,d)} \right) f_e \left(\sum_{(o,d) \in V \times V} \Lambda(o, d)x_e^{(o,d)} \right). \quad (4)$$

With these definitions in place, the unit network flow that minimizes the total travel time when serving the demand Λ can be found by solving the following optimization problem

$$\begin{aligned}
& \underset{x}{\text{minimize}} && F(x, \Lambda) && (5) \\
& \text{subject to } x = \left\{ x^{(o,d)} \right\}_{(o,d) \in V \times V}, \\
& && x^{(o,d)} \text{ is a unit } (o, d) \text{ flow for every } (o, d) \in V \times V,
\end{aligned}$$

3.2 Network Flow with Stochastic Demand

In practice, demand may vary from day to day, and such variations can be modeled by Λ being a random variable with distribution \mathcal{Q} . If \mathcal{Q} is known by the network operator, then rather than solving (5), the operator is interested in minimizing expected total travel time through the following stochastic optimization problem:

$$\begin{aligned}
& \underset{x}{\text{minimize}} && \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x, \Lambda)] && (6) \\
& \text{subject to } x = \left\{ x^{(o,d)} \right\}_{(o,d) \in V \times V}, \\
& && x^{(o,d)} \text{ is a unit } (o, d) \text{ flow for every } (o, d) \in V \times V,
\end{aligned}$$

We note that (6) is a generalization of (1) to the case when Λ is random.

In the more realistic case that \mathcal{Q} is not known, the optimization problem (6) can be approximated from historical data. We study a situation where the network operator has demand data $\Lambda_1, \Lambda_2, \dots, \Lambda_N \stackrel{\text{i.i.d.}}{\sim} \mathcal{Q}$ collected from operations of previous days. Using this data it can solve the following empirical approximation to (6):

$$\begin{aligned}
& \underset{x}{\text{minimize}} && \frac{1}{N} \sum_{k=1}^N F(x, \Lambda_k) && (7) \\
& \text{subject to } x = \left\{ x^{(o,d)} \right\}_{(o,d) \in V \times V}, \\
& && x^{(o,d)} \text{ is a unit } (o, d) \text{ flow for every } (o, d) \in V \times V,
\end{aligned}$$

The optimization problem in (7) uses historical data to estimate (6). In line with Assumption 1 we will assume that $\Lambda_t(o, d) \leq \lambda_{\max}$ for all values of t, o and d .

In Section 4 we show how (7) can be solved in a request-level differentially private way.

3.3 Assumptions on Travel Time functions

In this section we will introduce some assumptions that will help us establish our technical results. We will make the following assumptions on the network demand:

Assumption 1 (Bounded Demand). *We assume there exists a non-negative constant λ_{\max} so that $\Lambda(o, d) \leq \lambda_{\max}$ for every $(o, d) \in V \times V$. In practice, this constant can be estimated from historical data.*

The following are assumptions we make on the objective function F (see (4)). These assumptions are related to properties of the link latency functions $\{f_e\}_{e \in E}$. We present a typical model for link latency functions in Section 3.4 that satisfies all of the following assumptions.

Assumption 2 (Bounded Variance Gradients). *We assume there exists a non-negative constant K so that for every x , the variance of $\nabla F(x, \Lambda)$ is upper bounded by K^2 , i.e., $\mathbb{E}_{\Lambda \sim \mathcal{Q}} \left[\|\nabla F(x, \Lambda)\|_2^2 \right] \leq K^2$.*

Assumption 3 (Twice Differentiability). *We assume that $F(\cdot, \Lambda)$ is twice-differentiable for every Λ so that the hessian $H(x, \Lambda) := \frac{\partial^2}{\partial x^2} F(x, \Lambda)$ is defined on the entire domain of x .*

Assumption 4 (Strong Convexity). *We assume that there exists $\alpha > 0$ for which $F(\cdot, \Lambda)$ is α -strongly convex for every Λ . This means that for every Λ , and any unit network flows x, x' we have*

$$F(x', \Lambda) \geq F(x, \Lambda) + \nabla_x F(x, \Lambda)^\top (x' - x) + \frac{\alpha}{2} \|x' - x\|_2^2$$

Assumption 5 (Smoothness). *We assume that there exists $\beta > \alpha$ for which $F(\cdot, \Lambda)$ is β -smooth for every Λ . This means that for every Λ , and any unit network flows x, x' we have*

$$F(x', \Lambda) \leq F(x, \Lambda) + \nabla_x F(x, \Lambda)^\top (x' - x) + \frac{\beta}{2} \|x' - x\|_2^2.$$

Assumption 6 (Bounded second order partial derivative). *We assume that there exists $C > 0$ so that $\|\mathcal{D}_{\Lambda(o,d)} [\nabla_x F](x, \Lambda)\|_{op} \leq C$ for all x, Λ and $(o, d) \in V \times V$.*

Remark 6 (Satisfying Assumption 6). We can satisfy Assumption 6 for any positive C by re-scaling. In particular, letting $\theta := \max_{x', \Lambda'} \|\mathcal{D}_{\Lambda(o,d)} [\nabla_x F](x', \Lambda')\|_2$, then for any C we can define a re-scaled objective function

$$\tilde{F}(x, \Lambda) := \frac{C}{\theta} F(x, \Lambda).$$

By construction, \tilde{F} satisfies Assumption 6 with constant value C . Note, however, that the smoothness and strong convexity parameters for \tilde{F} will be rescaled accordingly to $\frac{\beta}{\theta}$ and $\frac{\alpha}{\theta}$ respectively.

3.4 Transportation Model satisfying assumptions from Section 3.3

In this section, we present a transportation network model that satisfies Assumptions 3,4,5 and 6. We study a network where the link latency functions are all affine where for each $e \in E$, there are non-negative constants q_e, c_e so that $f_e(y) = q_e y + c_e$. Let $Q \in \mathbb{R}^{m \times m}$ be defined so that $Q_{ee} = q_e$, and let $c \in \mathbb{R}^m$ be the concatenation of all of the zero order coefficients in the link latency functions.

As mentioned in Remark 4, we will represent x as a concatenation of $\{x^{(o,d)}\}_{(o,d) \in V \times V}$. As such, $x \in \mathbb{R}^{n^2 m}$. Let $\{(o_i, d_i)\}_{i=1}^{n^2}$ be the order in which the unit (o, d) flows are concatenated to produce x so that

$$x = \begin{bmatrix} x^{(o_1, d_1)} \\ x^{(o_2, d_2)} \\ \vdots \\ x^{(o_{n^2}, d_{n^2})} \end{bmatrix}.$$

The total flow on the links in the network when serving demand Λ according to x can then be written as:

$$y := \sum_{(o,d) \in V \times V} \Lambda(o,d)x^{(o,d)} = B_\Lambda x$$

where $B_\Lambda = \text{vec}(\Lambda) \otimes I_m$. Here $\text{vec}(\Lambda)$ is a n^2 dimensional row vector whose i th entry is $\Lambda(o_i, d_i)$, and \otimes represents the Kronecker product. Then when Λ is being routed according to x , the travel times on the links can be computed as

$$Qy + c = QB_\Lambda x + c,$$

which means that the total travel time can be written as

$$\begin{aligned} \widehat{F}(x, \Lambda) &:= \sum_{e \in E} y_e f_e(y_e) = y^\top (Qy + c) \\ &= x^\top B_\Lambda^\top (QB_\Lambda x + c) \\ &= x^\top B_\Lambda^\top QB_\Lambda x + c^\top B_\Lambda x. \end{aligned} \tag{8}$$

If we add a bit of ℓ_2 regularization, we obtain

$$\begin{aligned} F(x, \Lambda) &= \widehat{F}(x, \Lambda) + \frac{\alpha}{2} \|x\|_2^2 \\ &= x^\top B_\Lambda^\top QB_\Lambda x + c^\top B_\Lambda x + \frac{\alpha}{2} \|x\|_2^2 \\ &= x^\top \left(B_\Lambda^\top QB_\Lambda + \frac{\alpha}{2} I \right) x + c^\top B_\Lambda x. \end{aligned}$$

Recall that we use $H(\cdot, \Lambda)$ to denote the hessian of $F(\cdot, \Lambda)$ with respect to its first argument. We now make the following observations:

- The Hessian of $F(x, \Lambda)$ with respect to x is defined for all x and is equal to $2B_\Lambda^\top QB_\Lambda + \alpha I$. Hence Assumption 3 is satisfied.
- Since Q is a diagonal matrix with non-negative entries, $Q \succeq 0$. This implies that $B_\Lambda^\top QB_\Lambda \succeq 0$. Hence $H(x, \Lambda) \succeq \alpha I$. This implies that F is α -strongly convex, and hence Assumption 4 is satisfied.
- Note that $\|B_\Lambda^\top QB_\Lambda + (\alpha/2)I\|_{op} \leq \|Q\|_{op} \|B_\Lambda\|_{op}^2 + \alpha/2 \leq \|Q\|_{op} \|\text{vec}(\Lambda)\|_2^2 \|I_m\|_{op}^2 + \alpha/2$. Defining $\beta := n^2 (\max_e q_e) \lambda_{\max}^2 + \alpha/2$, we see that $\|H(x, \Lambda)\|_{op} \leq \beta$ for all $x \in \mathcal{X}$, meaning that F is β -smooth, and thus Assumption 5 is satisfied.
- By product rule,

$$\begin{aligned} \mathcal{D}_{\Lambda(o_i, d_i)} [\nabla F](x, \Lambda) &= \frac{\partial}{\partial \Lambda(o_i, d_i)} \left\{ 2 \left(B_\Lambda^\top QB_\Lambda + \alpha I \right) x + B_\Lambda^\top c \right\} \\ &= \frac{\partial}{\partial \Lambda(o_i, d_i)} \left\{ 2B_\Lambda^\top Qy + 2\alpha x + B_\Lambda^\top c \right\} \\ &= 2 \left((e_i \otimes I) Qy + B_\Lambda^\top Qx^{(o,d)} \right) + e_i \otimes c \end{aligned}$$

where e_i is the i th standard basis vector for \mathbb{R}^{n^2} . By triangle inequality we can then conclude that

$$\begin{aligned} \left\| \mathcal{D}_{\Lambda(o_i, d_i)} [\nabla F](x, \Lambda) \right\|_2 &\leq 2 \|Qy\|_2 + 2 \left\| B_\Lambda^\top Qx^{(o,d)} \right\|_2 + \|c\|_2 \\ &\leq 2 \|Q\|_{op} \|y\|_2 + 2 \|B_\Lambda\|_{op} \|Q\|_{op} \left\| x^{(o,d)} \right\|_2 + \|c\|_2 \\ &\stackrel{(a)}{\leq} 2 \|Q\|_{op} \|y\|_2 + 2n\lambda_{\max} \|Q\|_{op} \sqrt{m} + \|c\|_2 \\ &\stackrel{(b)}{\leq} 2\lambda_{\max} \|Q\|_{op} n^2 \sqrt{m} + 2\lambda_{\max} \|Q\|_{op} n\sqrt{m} + \|c\|_2. \end{aligned}$$

Here (a) is due to the fact that $x^{(o,d)}$ is a unit (o, d) flow, meaning $\|x^{(o,d)}\|_\infty \leq 1$ and thus $\|x^{(o,d)}\|_2 \leq \sqrt{m}$. Also, $B_\Lambda = \text{vec}(\Lambda) \otimes I_m$ implies that $\|B_\Lambda\|_{op} \leq \|\mathbb{1}_m^\top\|_2 \|I_m\|_{op} = n\lambda_{\max}$. (b) is due to $\|y\|_2 = \|B_\Lambda x^{(o,d)}\|_2 \leq \|B_\Lambda\|_{op} \|x^{(o,d)}\|_2 \leq n^2 \sqrt{m} \lambda_{\max}$.

Hence Assumption 6 is satisfied with $C = 2\lambda_{\max} \|Q\|_{op} \sqrt{mn}(n+1) + \|c\|_2$.

4 Differentially Private Network Routing Optimization

Given the setup from Section 3, our objective is to design a request-level differentially private algorithm that returns a near optimal solution to (6). Since the true distribution \mathcal{Q} of demand is unknown, we will design an algorithm for (7) and show that under the assumptions described in Section 3.2, the algorithm's solution is also near optimal for (6).

Computing a near-optimal solution to (7) while maintaining differential privacy may seem like a daunting task, but it turns out that a single modification to a well-known optimization algorithm gives an accurate and differentially private solution.

We present a Private Projected Stochastic Gradient Descent algorithm, which is described in Algorithm 1. As the name suggests, Algorithm 1 is a modified version of stochastic gradient descent. The algorithm conducts a single pass over the historical data, using each data point to perform a noisy gradient step (see line 6). The key difference between Algorithm 1 and standard stochastic gradient descent is in line 11, where instead of returning the final gradient descent iterate, Algorithm 1 returns a noisy version of the last iterate. Algorithm 1 has the following privacy and performance guarantees.

Theorem 2 (Privacy Guarantee for Algorithm 1). *Algorithm 1 is (ϵ, δ) -differentially private under request level adjacency defined in Definition 9.*

Theorem 3 (Performance Guarantee for Algorithm 1). *If $\Lambda_1, \dots, \Lambda_N \stackrel{i.i.d.}{\sim} \mathcal{Q}$, and x^* is a solution to (6), then the output x_{alg} of Algorithm 1 satisfies:*

$$\begin{aligned} \mathbb{E} [\|x_{alg} - x^*\|_2] &\leq \frac{1}{\sqrt{N}} \frac{KC \exp\left(\frac{\beta^2 \pi^2}{12\alpha^2}\right)}{\alpha} + \frac{n\sqrt{m}}{N} \left(\frac{\beta \exp\left(\frac{\beta^2 \pi^2}{12\alpha^2}\right)}{\alpha \min(1, 2\alpha)} + \frac{C}{\epsilon\alpha T} \sqrt{2 \ln\left(\frac{1.25}{\delta}\right)} \right) \\ &\leq O\left(\frac{1}{\sqrt{N}}\right) + O\left(\frac{1}{\epsilon N} \sqrt{\ln \frac{1}{\delta}}\right) \end{aligned}$$

In particular, x_{alg} is (ϵ, δ) -differentially private and converges to x^* as $N \rightarrow \infty$, meaning that privacy and asymptotic optimality are simultaneously achieved.

See Appendix B and Appendix C for proofs of Theorem 2 and Theorem 3 respectively.

Algorithm 1: Private Projected Stochastic Gradient Descent

- 1 **Input:** Historical demand data $\Lambda_1, \dots, \Lambda_N$, Desired privacy level (ϵ, δ) ;
 - 2 **Output:** Unit network flow $x \in \mathcal{X}$;
 - 3 Initialize $x_0 \in \mathcal{X}$ arbitrarily ;
 - 4 **for** $1 \leq k \leq N$ **do**
 - 5 $\eta_{k-1} \leftarrow \min\left(\frac{1}{\alpha k}, \frac{\min(1, 2\alpha)}{\beta}\right)$;
 - 6 $x_k \leftarrow \Pi_{\mathcal{X}}(x_{k-1} + \eta_{k-1} \nabla_x F(x_{k-1}, \Lambda_k))$;
 - 7 $s \leftarrow \frac{C}{T} \min\left(\frac{\min(1, 2\alpha)}{\beta}, \frac{1}{\alpha N}\right)$;
 - 8 $\sigma^2 \leftarrow 2 \frac{s^2}{\epsilon^2} \ln\left(\frac{1.25}{\delta}\right)$;
 - 9 $Z \sim \mathcal{N}(0, \sigma^2 I)$;
 - 10 $x_{\text{alg}} \leftarrow \Pi_{\mathcal{X}}(x_N + Z)$;
 - 11 **Return** x_{alg} ;
-

4.1 Discussion

Carefully adding noise to specific parts of existing algorithms is a principled approach for developing differentially private algorithms [DMNS06, FMTT18, FKT20]. The main challenge in such an approach is determining a) where and b) how much noise to add. Suppose the goal is to (approximately) compute a query function $f(L)$ on a data set L in a differentially private way. The latter question can be addressed by measuring the sensitivity of the desired query function.

Definition 12 (ℓ_2 sensitivity). Consider a function $f : \mathcal{L} \rightarrow \mathbb{R}^d$ which maps data sets to real vectors. For a given adjacency relation Adj , the ℓ_2 sensitivity of f , denoted $s_{\text{Adj}}(f)$, is the largest achievable difference in function value between adjacent data sets. Namely,

$$s_{\text{Adj}}(f) := \max_{\substack{L_1, L_2 \in \mathcal{L} \\ \text{Adj}(L_1, L_2)=1}} \|f(L_1) - f(L_2)\|_2.$$

Once the sensitivity of the query function is known, the required noise distribution can be determined using the Gaussian mechanism as described in Theorem 4.

Theorem 4 (From [DR14]). *Suppose $f : \mathcal{D} \rightarrow \mathbb{R}^p$ maps datasets to real vectors. If s_{Adj} is the ℓ_2 sensitivity of f , then $\hat{f}(D) := f(D) + Z$ where $Z \sim \mathcal{N}\left(0, 2 \frac{s_{\text{Adj}}^2}{\epsilon^2} \ln\left(\frac{1.25}{\delta}\right) I_p\right)$ is (ϵ, δ) -differentially private with respect to the adjacency relation Adj .*

Calculating the sensitivity of the simple query functions (e.g., counting, voting, selecting the maximum value) is relatively easy, making the Gaussian mechanism straightforward to apply. However, for more complicated functions, noise calibration becomes more involved.

Algorithm 1 is an application of the Gaussian mechanism. Moreover, Theorem 3 shows that the suboptimality of Algorithm 1 is $O\left(\frac{1}{\sqrt{N}}\right)$. The asymptotic optimality of Algorithm 1 comes from

the fact that the ℓ_2 sensitivity of the final gradient descent iterate is actually converging to zero as N approaches infinity. This fact enables us to add less noise as $N \rightarrow \infty$. Indeed, the Gaussian noise added to the final gradient descent iterate in Algorithm 1 has standard deviation which is $O\left(\frac{1}{N}\right)$.

In the remainder of this section, we sketch some of the mathematical ideas behind the perhaps non-intuitive result that the sensitivity of gradient descent converges to zero as the number of data points increases. For simplicity of exposition we will a) use scalar notation in place of vector notation for the sake of readability, and b) consider the simpler case of unconstrained gradient descent, which removes the need to perform projections. As a reminder, the full proof can be found in Appendix B.

Given a mobility data set, we use $x_t(L)$ to denote the t -th iterate of Algorithm 1 when using data set L . It is sufficient to show that the $\max_t \left| \frac{dx_N}{d\Lambda_t} \right|$ converges to zero for every t . This is because the sensitivity is obtained by integrating the derivative:

$$\|x_N(L_1) - x_N(L_2)\|_2 = \left\| \int_{L_1}^{L_2} \frac{dx_N}{dL} dL \right\|_2.$$

Because L_1, L_2 are adjacent, the distance between them is finite, and thus the above integral will converge to zero if its integrand converges to zero.

With this in mind, by chain rule, we can write

$$\frac{dx_N}{d\Lambda_t} = \frac{dx_t}{d\Lambda_t} \frac{dx_{t+1}}{dx_t} \frac{dx_{t+2}}{dx_{t+1}} \cdots \frac{dx_N}{dx_{N-1}}.$$

Next, by using properties of smooth and strongly convex functions, we show that $\left| \frac{dx_{t+1}}{dx_t} \right| \leq 1 - \theta$ for some positive constant θ . This result implies

$$\left| \frac{dx_N}{d\Lambda_t} \right| = \left| \frac{dx_t}{d\Lambda_t} \right| (1 - \theta)^{N-t}.$$

Next, recalling that $x_t = x_{t-1} - \eta_t \nabla F(x_{t-1}, \Lambda_t)$, note that x_{t-1} is computed from the first $t - 1$ gradient steps, which only depend on the first $t - 1$ data points $\Lambda_1, \dots, \Lambda_{t-1}$. Hence the $\frac{dx_{t-1}}{d\Lambda_t} = 0$. From this we see that

$$\begin{aligned} \frac{dx_t}{d\Lambda_t} &= \frac{d}{d\Lambda_t} (x_{t-1} - \eta_t \nabla F(x_{t-1}, \Lambda_t)) \\ &= -\eta_t \frac{d}{d\Lambda_t} \nabla F(x_{t-1}, \Lambda_t). \end{aligned}$$

By using Assumption 6 we have the following inequality:

$$\left| \frac{dx_t}{d\Lambda_t} \right| = \left| \eta_t \frac{d}{d\Lambda_t} \nabla F(x_{t-1}, \Lambda_t) \right| \leq \eta_t C.$$

Putting everything together, we have

$$\left| \frac{dx_N}{d\Lambda_t} \right| = C \eta_t (1 - \theta)^{N-t}.$$

The choice of stepsize in Algorithm 1 ensures three things: a) $\eta_t \leq 2$ for every t , b) η_t is non-increasing, and c) $\eta_t \rightarrow 0$ as $t \rightarrow \infty$. Given these facts, there are two cases to consider for $\left| \frac{dx_N}{d\Lambda_t} \right|$.

Case 1: $t \leq N/2$. In this case, we can upper bound $\eta_t \leq 2$ and $(1-\theta)^{N/2}$. Hence $\left| \frac{dx_N}{d\Lambda_t} \right| \leq 2C(1-\theta)^{N/2}$. Since $1-\theta < 1$, this converges to zero as $N \rightarrow \infty$.

Case 2: $t \geq N/2$. In this case, we simply upper bound $(1-\theta)^{N-t} \leq 1$ and $\eta_t \leq \eta_{N/2}$ which gives $\left| \frac{dx_N}{d\Lambda_t} \right| \leq C\eta_{N/2}$ which converges to zero as $N \rightarrow \infty$.

Finally, this shows that

$$\max_t \left| \frac{dx_N}{d\Lambda_t} \right| \leq \max(C\eta_{N/2}, 2C(1-\theta)^N),$$

and since both terms in the maximum are converging to zero, we have $\lim_{N \rightarrow \infty} \left| \frac{dx_N}{d\Lambda_t} \right| = 0$.

5 Experiments

Differentially private mechanisms add noise to provide a principled privacy guarantee for individual level user data. One immediate question is the degree to which the added noise degrades the quality of the obtained solution. In the previous section, we addressed this question theoretically in Theorem 3 by showing that Algorithm 1 is both differentially private and asymptotic optimal. In this section, we present empirical studies on privacy-performance trade-offs by comparing our algorithm's performance to that of a non-private network routing approach. To this end, we simulate a transportation network to evaluate the performance of our algorithm and the non-private optimal solution to the network flow problem (7).

We describe the dataset used for the experiments in Section 5.1. Next, the algorithms used in the experiments are described in Section 5.2. In Section 5.3, we evaluate the practical performance of our algorithm for different values of N . In particular, we study the effect of the number of samples, and quantify the loss in system performance we may experience due to the introduction of our privacy-preserving algorithm. Finally, in Section 5.4, we study the sensitivity of our algorithms to the simulation parameters. In particular, we study the convergence of the routing policy with increasing data for different demands, edge latency function, and the magnitude of the regularization loss introduced to convexify the edge latency function.

5.1 Data Set

We use data for the Sioux Fall network, which is available in the Transportation Network Test Problem (TNTP) dataset [Tra16]. This network has 24 nodes, 76 edges, and 528 OD pairs (see Figure 1 for an overview of the network topology). The distribution of mean hourly demand across different OD pairs is shown in Figure 1. The mean OD demand is 682 vehicles/hour. Furthermore, for more context on the scale of the network, the travel time on edges ranges from 2 to 10 minutes. Trip data is generated each hour, i.e., $T = 60$, from a Poisson distribution with a mean value that is given by the data. Our objective is to learn a routing policy for these trips that minimizes the total travel costs for all users. To model congestion, we use the link latency model described in Section 3.4. In particular, for each $e \in E$, we estimate the free flow latency as c_e directly from the data set. The sensitivity of the latency function to the traffic volume on the link, denoted by q_e is chosen such that the travel time on the link is doubled when the link flow equals the link capacity. In later experiments, we will change this factor to study the sensitivity of our algorithm.

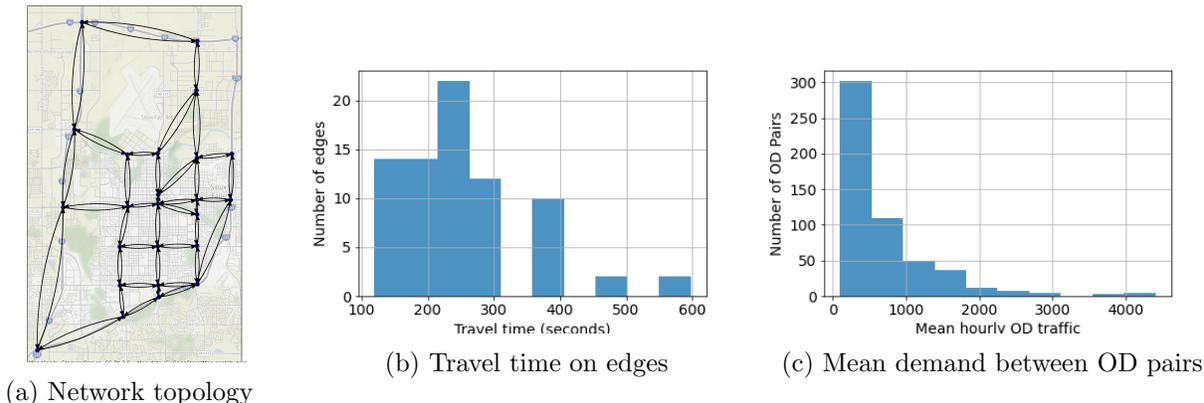


Figure 1: Description of the Sioux Falls dataset

5.2 Algorithms

Throughout Sections 5.3 and 5.4 we compare the performance of two different algorithms: a Baseline algorithm and Algorithm 1, which we describe below.

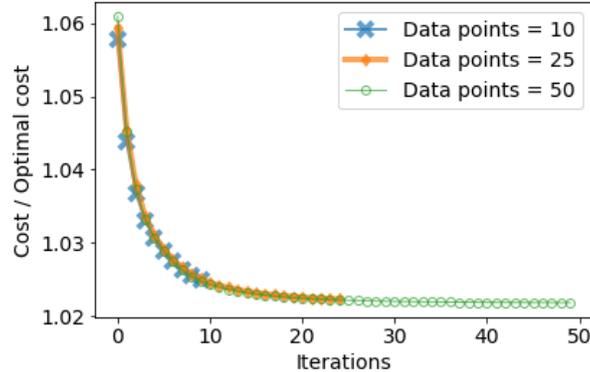
Baseline: This algorithm computes the minimum travel time solution to (7) for the Sioux Fall model described in Section 3.4. Recall that in the Section 3.4 model, (8) describes the travel time incurred when serving demand Λ with a routing policy x . Given a data set $\Lambda_1, \dots, \Lambda_N$, it computes the solution to the following minimization problem: $\min_{x \in \mathcal{X}} \frac{1}{N} \sum_{i=1}^N x^\top B_{\Lambda_i}^\top Q B_{\Lambda_i} x + c^\top B_{\Lambda_i} x$.

Algorithm 1: The travel time function described in (8) is not strongly convex because $B_\Lambda^\top Q B_\Lambda$ is rank deficient. In order to satisfy Assumption 4, we introduce an ℓ_2 regularization. Namely, given a data set $\Lambda_1, \dots, \Lambda_N$, we apply Algorithm 1 to the following minimization problem: $\min_{x \in \mathcal{X}} \alpha \|x\|_2^2 + \frac{1}{N} \sum_{i=1}^N x^\top B_{\Lambda_i}^\top Q B_{\Lambda_i} x + c^\top B_{\Lambda_i} x$.

For the Sioux Falls network, the parameters for implementing Algorithm 1 are set as follows. The smoothness parameter β is set to be the largest eigenvalue of $B_\Lambda^\top Q B_\Lambda$, which is equal to 2.08×10^7 . We set $C = \beta$, and the regularizer parameter $\alpha = 10^4$. It is easy to check that these values satisfy the assumptions described in Section 3.4. Note that several different values of α could have been used to convexify the latency function. However, our choice is governed by two factors. A small value of α will ensure that the regularized objective is a good estimate to the true objective, which is desirable. However, smaller values of α will lead to a larger condition number (which is β/α), resulting in slower convergence and necessitating more data for achieving a similar performance. Thus, the particular value of $\alpha = 10^4$ balances both these factors for our problem instance. In section 5.4, we present a sensitivity analysis with different values of α .

5.3 Performance of Algorithm 1

In this subsection, we study the convergence of the routing policy with each step of the gradient descent performed by our algorithm. Since the impact of a routing policy is directly reflected in the total travel time, we plot the travel cost induced by a the learned routing policy as a

Figure 2: Convergence dynamics for different values of N

function of the iterates. Recall that the number of iterations for a data set with N data points is N according to Algorithm 1. In our experiments we evaluate the cost of a routing policy x as $F(x, \frac{1}{N} \sum_{i=1}^N \Lambda_i)$ instead of using the sample average $\frac{1}{N} \sum_{i=1}^N F(x, \Lambda_i)$. This approximation is done solely for improving the run-time of our experiment (by up to 50X) and introduces less than $10^{-4}\%$ error in the evaluation of the system costs. Similarly, the optimal solution is also computed by minimizing the objective function $F(x, \frac{1}{N} \sum_{i=1}^N \Lambda_i)$ instead of the term described in Equation 7 to achieve a computational speed up of 30X. Evaluating the cost of this solution using the average-demand approximation results in errors less than $10^{-7}\%$ error. Thus, in these experiments, we define the optimal costs as the approximate costs obtained through this procedure. Further details justifying these approximation are presented in Appendix D

For our first set of experiments, we compare the objective values obtained by Algorithm 1 and the Baseline as a function of sample size N . In Figure 2, we plot the ratio of Algorithm 1’s cost to Baseline’s cost over the course of iterations for different values of N . For this set of experiments, we set the privacy parameters to $\epsilon = 0.1$ and $\delta = 0.1$. Note that Baseline’s cost is fixed for a given N and is computed offline to serve as a benchmark. For a given N , we only have N iterations since each data point is only used once in Algorithm 1 to maintain privacy. For all three experiments ($N = 10$, $N = 25$, $N = 50$), the cost decreases monotonically with additional iterations. It is therefore not surprising that the final costs for the $N = 50$ case is the lowest, as we expect the routing policy learned with 50 data points to be better than the routing policy learned from 10 data points. It is however very interesting that even with a random routing policy initialization, our algorithm finds solutions that are just around 2% away from the optimal policy. We suspect that this 2% gap is due to the fact that the two algorithms have slightly different objective functions. Indeed, Algorithm 1 has an ℓ_2 regularizer in the objective, but Baseline does not. Our results therefore show that although the convergence is guaranteed only in the limit $N \rightarrow \infty$, we can obtain practically useful solutions with a relatively small number of data points.

In the next set of experiments, we study the effect of different privacy parameters on total travel time. To this end, we compare the costs of the pre-noise and post-noise solutions x_N x_{alg} from Algorithm 1. We conduct this comparison for $\epsilon \in \{0.01, 0.1, 0.5\}$ and $\delta \in \{0.1, 0.5\}$. Table 1 presents the percentage increase in total travel time due to the addition of privacy noise. The results in indicate that the price of privacy, i.e., the increase in total travel time due to the introduction

ϵ	δ	Cost (% increase)
0.01	0.1	7.83×10^{-2}
0.01	0.5	3.97×10^{-3}
0.1	0.1	9.06×10^{-3}
0.1	0.5	5.96×10^{-3}
0.5	0.1	2.44×10^{-3}
0.5	0.5	2.05×10^{-3}

Table 1: Change in routing costs due to incorporating privacy.

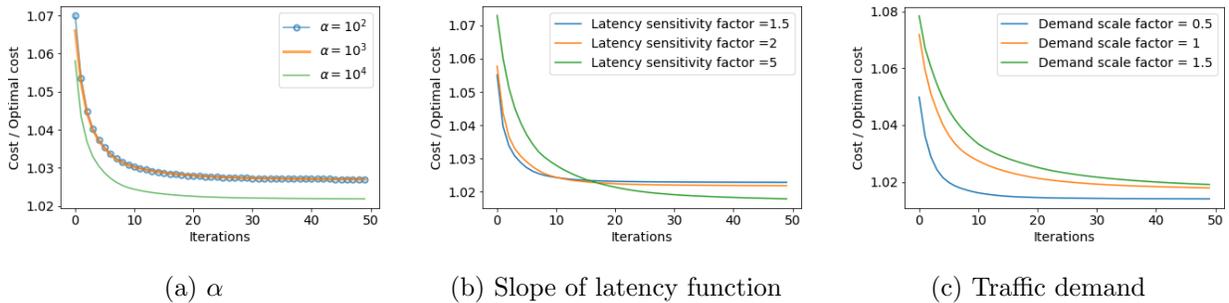


Figure 3: Sensitivity of the performance of Algorithm 1

of differential privacy noise is less than $7.8 \times 10^{-2}\%$ in the worst case. In fact, for more commonly used privacy parameters of $\epsilon = 0.1$ and $\delta = 0.1$, the cost of privacy is even smaller. One reason for this low cost of privacy is the high demand in the traffic network. From Figure 1c, we observe that every OD pair typically has a few hundred trips. With over 500 OD pairs, it is thus clear that there are tens of thousands of vehicles in the network contributing to trip information with every data point. Thus, with such a large number of vehicles, the noise required to protect the identity of one vehicle is not too high.

5.4 Sensitivity

We now study the sensitivity of our algorithm to input parameters. For these experiments, we fix the number of data points to $N = 50$, since prior results suggest that most of the cost benefits are obtained with 50 data points. First, we study the effect of the regularizer term by setting $\alpha \in \{10^2, 10^3, 10^4\}$ and plot the convergence of the normalized costs in Figure 3a. We see that for larger values of α , the costs decrease faster. This makes sense because larger α results in a lower condition number $\frac{\beta}{\alpha}$, which leads to faster convergence. We also note that the cost ratio does not go to 1 because Algorithm 1 is minimizing a regularized objective, while Baseline has no regularization.

In Figure 3b, we compare three scenarios with varying slope for the latency function q_e . Recall that our previous experiments set the slope q_e on each edge such that the travel time on the link doubles when the traffic is equal to the link capacity. This setting corresponds to a sensitivity factor of 2. We consider two more cases where where the travel time at capacity flow is 1.5 times and 5 times the free flow latency. Note that changing the latency sensitivity factor changes the matrix Q . Thus, for each of these experiments, we recompute the value of β and C and set it to

be equal to the largest eigenvalue of the appropriate $B_\Lambda^\top Q B_\Lambda$. The optimal cost also varies for all the three cases and is recomputed. The value of α is fixed at 10^4 . We observe that when the latency function is steeper, i.e., the sensitivity factor is higher, the algorithm takes more iterations to reduce the costs, but eventually ends up with the lowest costs. This is because a larger β leads to larger condition number $\frac{\beta}{\alpha}$, which makes convergence slow. However, a larger β means that the ℓ_2 regularizer is a smaller proportion of the total cost, meaning that the objectives of Algorithm 1 and Baseline become more similar, which is what we believe causes the cost ratio to improve as the latency function becomes steeper.

Finally, we present the sensitivity of our algorithm to varying traffic demand in Figure 3c. Elaborating further, in these experiments, we compare the nominal setting, where the demand is the mean demand with a scale factor of 1 to two cases. In the first case, we use a lower demand, where the mean traffic is 0.5 times the nominal traffic, and in the second case, the mean traffic is 1.5 times the nominal traffic. Again, as the demand changes, the matrix B changes, and we recompute β and C as before. We observe that for the same value of α , higher demand leads to better convergence and lower costs. This is because higher demand increases the travel time, making the ℓ_2 regularizer a smaller proportion of Algorithm 1’s objective. The objective functions becoming more similar leads to the cost ratio being closer to 1.

6 Conclusions

In this paper, we study the problem of learning network routing policies from sensitive user data. In particular, we consider the setting of a transportation network, where we want to learn and share a routing policy such that it does not reveal too much information about individual trips that may have contributed to learning this policy. Our paper presented a new approach to learn privacy-preserving routing policies by solving a reformulated network flow problem using a differentially private variant of the stochastic gradient descent algorithm. We prove that our algorithm is asymptotically optimal, meaning that the cost of the routing policy produced by our algorithm converges to the optimal non-private cost as the number of data points goes to infinity. Finally, our simulations on a Sioux Falls road network suggests that for realistic travel demands, we can learn differentially private routing policies that result in only a 2% suboptimality in terms of total travel time.

There are several interesting directions for future work. First, because differentially private algorithms are not allowed to be sensitive to single data points, they are naturally robust, and can be useful for tracking non-stationary demand distributions, as opposed to the stationary demand models we studied in this paper. In this paper, we studied request-level differential privacy where the goal is to occlude the influence of a single trip on the algorithm’s output. Another practical and important notion is user-level differential privacy where the goal is to occlude the influence of all trips belonging to the same person on the algorithm’s output. User-level privacy is harder to achieve, but is important in practice. Finally, generalizing the results to non-smooth objective functions would expand the domain of models that this technique can be applied to.

References

- [AHFI⁺18] Khalil Al-Hussaeni, Benjamin CM Fung, Farkhund Iqbal, Gaby G Dagher, and Eun G Park. Safepath: Differentially-private publishing of passenger trajectories in transportation systems. *Computer Networks*, 143:126–139, 2018.
- [BS03] Alastair R Beresford and Frank Stajano. Location privacy in pervasive computing. *IEEE Pervasive computing*, 2(1):46–55, 2003.
- [CML11] Chi-Yin Chow, Mohamed F Mokbel, and Xuan Liu. Spatial cloaking for anonymous location-based services in mobile peer-to-peer environments. *GeoInformatica*, 15(2):351–380, 2011.
- [CTW⁺21] Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom B. Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. Extracting training data from large language models. In *30th USENIX Security Symposium, USENIX Security 2021, August 11-13, 2021*, pages 2633–2650. USENIX Association, 2021.
- [DKBS15] Roy Dong, Walid Krichene, Alexandre M. Bayen, and S. Shankar Sastry. Differential privacy of populations in routing games. In *2015 54th IEEE Conference on Decision and Control (CDC)*, pages 2798–2803, 2015.
- [DMNS06] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam D. Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography Conference*, volume 3876, pages 265–284. Springer, 2006.
- [DR14] Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. *Found. Trends Theor. Comput. Sci.*, 9(3-4):211–407, 2014.
- [Dwo08] Cynthia Dwork. Differential privacy: A survey of results. In *International conference on theory and applications of models of computation*, pages 1–19. Springer, 2008.
- [FKT20] Vitaly Feldman, Tomer Koren, and Kunal Talwar. Private stochastic convex optimization: optimal rates in linear time. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing, STOC*, pages 439–449. ACM, 2020.
- [FMTT18] Vitaly Feldman, Ilya Mironov, Kunal Talwar, and Abhradeep Thakurta. Privacy amplification by iteration. In *59th IEEE Annual Symposium on Foundations of Computer Science, FOCS*, pages 521–532. IEEE Computer Society, 2018.
- [GFCA14] Moein Ghasemzadeh, Benjamin CM Fung, Rui Chen, and Anjali Awasthi. Anonymizing trajectory data for passenger flow analysis. *Transportation research part C: emerging technologies*, 39:63–79, 2014.
- [GLTY18] Mehmet Emre Gursoy, Ling Liu, Stacey Truex, and Lei Yu. Differentially private and utility preserving publication of trajectory data. *IEEE Transactions on Mobile Computing*, 18(10):2315–2329, 2018.

- [GMW87] Oded Goldreich, Silvio Micali, and Avi Wigderson. How to play any mental game or A completeness theorem for protocols with honest majority. In Alfred V. Aho, editor, *STOC 1987, New York, New York, USA*, pages 218–229. ACM, 1987.
- [GZFS15] Yanmin Gong, Chi Zhang, Yuguang Fang, and Jinyuan Sun. Protecting location privacy for task allocation in ad hoc mobile cloud computing. *IEEE Transactions on Emerging Topics in Computing*, 6(1):110–121, 2015.
- [HC20] Brian Yueshuai He and Joseph YJ Chow. Optimal privacy control for transport network data sharing. *Transportation Research Part C: Emerging Technologies*, 113, 2020.
- [HTP17] Shuo Han, Ufuk Topcu, and George J. Pappas. Differentially private distributed constrained optimization. *IEEE Transactions on Automatic Control*, 62(1):50–64, 2017.
- [LYH20] Yang Li, Dasen Yang, and Xianbiao Hu. A differential privacy-based privacy-preserving data publishing algorithm for transit smart card data. *Transportation Research Part C: Emerging Technologies*, 115:102634, 2020.
- [Pan14] Vijay Pandurangan. Riding with the stars: Passenger privacy in the nyc taxicab dataset. Available Online, 2014.
- [Swe02] Latanya Sweeney. k-anonymity: A model for protecting privacy. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 10(05):557–570, 2002.
- [Tra16] Transportation networks for research. github.com/bstabler/TransportationNetworks, 2016. Accessed January 20, 2021.
- [Tre11] Luca Trevisan. Lecture notes on optimization and algorithmic paradigms (lecture 11). Available online, 2011.
- [WHL⁺18] Zhibo Wang, Jiahui Hu, Ruizhao Lv, Jian Wei, Qian Wang, Dejun Yang, and Hairong Qi. Personalized privacy-preserving task allocation for mobile crowdsensing. *IEEE Transactions on Mobile Computing*, 18(6):1330–1341, 2018.
- [WLK⁺17] Xi Wu, Fengan Li, Arun Kumar, Kamalika Chaudhuri, Somesh Jha, and Jeffrey F. Naughton. Bolt-on differential privacy for scalable stochastic gradient descent-based analytics. In *Proceedings of the 2017 ACM International Conference on Management of Data, SIGMOD Conference 2017, Chicago, IL, USA, May 14-19, 2017*, pages 1307–1322. ACM, 2017.
- [YLL⁺19] Ke Yan, Guoming Lu, Guangchun Luo, Xu Zheng, Ling Tian, and Akshita Maradapu Vera Venkata Sai. Location privacy-aware task bidding and assignment for mobile crowd-sensing. *IEEE Access*, 7:131929–131943, 2019.

A Notation

The terms and concepts defined and used throughout the paper are all described here for the sake of convenience.

Notation	Definition
G	Graph representation of the network. (Definition 1)
V	The vertex set of the network G . (Definition 1)
E	The edge set of the network G . (Definition 1)
n	The number of vertices in G , i.e., $n = V $. (Definition 1)
m	The number of edges in G , i.e., $m = E $. (Definition 1)
\mathcal{T}	Operation period $\mathcal{T} = [t_{\text{start}}, t_{\text{end}}]$ which represents the time of day that the network operator is trying to optimize its decisions in. (Definition 2)
T	Number of minutes in the operation period. (Definition 2)
Λ	Demand matrix. $\Lambda(o, d)$ specifies the rate at which requests from $o \in V$ to $d \in V$ appear in the network. (Definition 3)
\mathcal{Q}	The distribution of Λ . (Defined in Section 3.2)
L	Dataset $L := (\Lambda_1, \dots, \Lambda_N)$ containing N days of historical request data for the operation period \mathcal{T} . (Defined in Section 3.2)
$x^{(o,d)}$	Unit (o, d) flow. A flow that routes exactly 1 unit of flow from o to d through the network. (Definition 10)
x	Unit network flow. Specifies a unit (o, d) flow for all pairs of vertices (o, d) in the network. (Definition 11)
\mathcal{X}	The set of all unit network flows. It is also the feasible set for optimization problems (5),(6) and (7). (Definition 11)
x_k	A unit network flow and the k th gradient descent iterate from Algorithm 1. $x_k(L)$ is used to show explicit dependence of x_k on the historical training data L . (Algorithm 1)
f_e	Link latency function for the edge $e \in E$. If x is the total flow on e , then $f(x)$ will be the average travel time through the edge. (Definition 4)
F	Total travel time function. $F(x, \Lambda)$ is the total travel time of requests in Λ experience when they are routed according to x . (Defined in (4))
(ϵ, δ)	Differential privacy parameters. As ϵ and δ get smaller, the privacy guarantees offered by differential privacy get stronger. (Definition 8)
K	Upper bound on the standard deviation of the gradient of F . (Assumption 2)
H	Hessian of F with respect to x . Concretely, $H(x, \Lambda)$ is the hessian of F with respect to x evaluated at (x, Λ) .
α	Strong convexity parameter. (Assumption 4)
β	Smoothness parameter. (Assumption 5)
C	Upper bound on $\ \mathcal{D}_\Lambda [\nabla F](x, \Lambda)\ _2$. (Assumption 6)

B Privacy Analysis (Proof of Theorem 2)

Recall the privacy guarantee for the Gaussian Mechanism from Theorem 4. To show that Algorithm 1 is (ϵ, δ) -differentially private, it suffices to show the following Lemma:

Lemma 1. *For every $1 \leq t \leq N$, any value of $L := (\Lambda_1, \dots, \Lambda_N)$, and any $(o, d) \in V \times V$ we have $\|\mathcal{D}_{\Lambda_t(o,d)}[x_N](L)\|_{op} \leq \min\left(\frac{C \min(1, 2\alpha)}{\beta}, \frac{C}{\alpha N}\right)$.*

Lemma 1 is sufficient to prove privacy, because it implies that the ℓ_2 sensitivity of x_N is at most $\frac{L}{\alpha N}$. To see why this is true, let $L_1 = (\Lambda_1, \dots, \Lambda'_t, \dots, \Lambda_N)$ and $L_2 = (\Lambda_1, \dots, \Lambda''_t, \dots, \Lambda_N)$ be request-level-adjacent data sets that differ only on the t th day. Furthermore, let (o', d') be the request for which Λ'_t and Λ''_t differ. Let $x_N(L_1)$ and $x_N(L_2)$ be the final iterates of gradient descent obtained by using the data sets L_1 and L_2 respectively. By the fundamental theorem of calculus,

$$\begin{aligned}
x_N(L_2) - x_N(L_1) &= \int_{L_1}^{L_2} \mathcal{D}_{\Lambda_t}[x_N](L) dL \\
&= \int_{\Lambda'_t}^{\Lambda''_t} \mathcal{D}_{\Lambda_t}[x_N](\Lambda_1, \dots, \Lambda_{t-1}, \Lambda, \Lambda_{t+1}, \dots, \Lambda_N) d\Lambda \\
\Rightarrow \left\| x_N^{(L')} - x_N^{(L)} \right\|_2 &= \left\| \int_{\Lambda'_t}^{\Lambda''_t} \mathcal{D}_{\Lambda_t}[x_N](\Lambda_1, \dots, \Lambda_{t-1}, \Lambda, \Lambda_{t+1}, \dots, \Lambda_N) d\Lambda \right\|_2 \\
&\leq \int_{\Lambda'_t}^{\Lambda''_t} \|\mathcal{D}_{\Lambda_t}[x_N](\Lambda_1, \dots, \Lambda_{t-1}, \Lambda, \Lambda_{t+1}, \dots, \Lambda_N) d\Lambda\|_2 \\
&= \int_{\Lambda'_t(o', d')}^{\Lambda''_t(o', d')} \|\mathcal{D}_{\Lambda_t(o', d')}[x_N](\Lambda_1, \dots, \Lambda_{t-1}, \Lambda, \Lambda_{t+1}, \dots, \Lambda_N) d\Lambda(o', d')\|_2 \\
&\leq \int_{\Lambda'_t(o', d')}^{\Lambda''_t(o', d')} \|\mathcal{D}_{\Lambda_t(o', d')}[x_N](\Lambda_1, \dots, \Lambda_{t-1}, \Lambda, \Lambda_{t+1}, \dots, \Lambda_N)\|_2 d\Lambda(o', d') \\
&\stackrel{(a)}{\leq} \int_{\Lambda'_t}^{\Lambda''_t} \min\left(\frac{C \min(1, 2\alpha)}{\beta}, \frac{C}{\alpha N}\right) d\Lambda(o', d') \\
&\leq \min\left(\frac{C \min(1, 2\alpha)}{\beta}, \frac{C}{\alpha N}\right) |\Lambda''_t(o', d') - \Lambda'_t(o', d')| \\
&\stackrel{(b)}{\leq} \min\left(\frac{C \min(1, 2\alpha)}{\beta T}, \frac{C}{\alpha T N}\right).
\end{aligned}$$

Here (a) due to Lemma 1. (b) is because $|\Lambda''_t(o, d) - \Lambda'_t(o, d)| \leq T^{-1}$ is a consequence of L_1, L_2 being request-level-adjacent.

Thus to establish (ϵ, δ) -differential privacy of Algorithm 1, all that remains is to prove Lemma 1.

Proof of Lemma 1. Define $L := (\Lambda_1, \dots, \Lambda_N)$. By Assumptions 3, 4 and 5, the functions $F(\cdot, \Lambda_1), \dots, F(\cdot, \Lambda_N)$ are all twice differentiable, α -strongly convex and β -smooth. For each $t \leq N$, and any $(o, d) \in V \times V$ by chain rule we have

$$\mathcal{D}_{\Lambda_t(o,d)}[x_N](L) = \left(\prod_{k=t+1}^{N-1} \mathcal{D}_{x_k}[x_{k+1}](L) \right) \mathcal{D}_{\Lambda_t(o,d)}[x_{t+1}](L).$$

Since $x_{k+1} = x_k - \eta_k \nabla_x F(x_k, \Lambda_k)$, differentiating both sides with respect to x_k gives

$$\mathcal{D}_{x_k}[x_{k+1}](L) = I - \eta_k H(x_k, \Lambda_k)$$

where $H(\cdot, \Lambda_k)$ is the Hessian of $F(\cdot, \Lambda_k)$. Since $F(\cdot, \Lambda_k)$ is α -strongly convex, we know $H(x_k, \Lambda_k) \succeq \alpha I$. By β -smoothness of $F(\cdot, \Lambda_k)$, we also know that $H(x_k, \Lambda_k) \preceq \beta I$. Therefore if $\eta_k \leq \frac{1}{\beta}$, we see that $\|\mathcal{D}_{x_k}[x_{k+1}](L)\|_{op} \leq 1 - \eta_k \alpha$. From this we can conclude that

$$\begin{aligned} \|\mathcal{D}_{\Lambda_t(o,d)}[x_N](L)\|_{op} &\leq \|\mathcal{D}_{\Lambda_t(o,d)}[x_{t+1}](L)\|_{op} \prod_{k=t+1}^{N-1} \|\mathcal{D}_{x_k}[x_{k+1}](L)\|_{op} \\ &\leq \|\mathcal{D}_{\Lambda_t(o,d)}[x_{t+1}](L)\|_{op} \prod_{k=t+1}^{N-1} (1 - \eta_k \alpha) \\ &\leq \|\mathcal{D}_{\Lambda_t(o,d)}[x_{t+1}](L)\|_{op} \prod_{k=t+1}^{N-1} \exp(-\eta_k \alpha) \\ &= \|\mathcal{D}_{\Lambda_t(o,d)}[x_{t+1}](L)\|_{op} \exp\left(-\sum_{k=t+1}^{N-1} \eta_k \alpha\right) \end{aligned} \quad (9)$$

Finally, note that x_k only depends on x_0 and $\Lambda_1, \dots, \Lambda_{k-1}$, and in particular it does not depend on Λ_k . Therefore differentiating both sides of $x_{k+1} = x_k - \eta_k \nabla_x F(x_k, \Lambda_k)$ with respect to $\Lambda_k(o, d)$ gives

$$\begin{aligned} \mathcal{D}_{\Lambda_t(o,d)}[x_{t+1}](L) &= -\eta_t \mathcal{D}_{\Lambda_t(o,d)}[\nabla_x F](x_t, \Lambda_t) \\ \implies \|\mathcal{D}_{\Lambda_t(o,d)}[x_{t+1}](L)\|_{op} &= \eta_t \|\mathcal{D}_{\Lambda_t(o,d)}[\nabla_x F](x_t, \Lambda_t)\|_{op} \stackrel{(a)}{\leq} C\eta_t, \end{aligned} \quad (10)$$

where (a) is due to Assumption 6. Combining inequalities (9) and (10) gives

$$\|\mathcal{D}_{\Lambda_t(o,d)}[x_N](L)\|_{op} \leq C\eta_t \exp\left(-\sum_{k=t+1}^{N-1} \eta_k \alpha\right). \quad (11)$$

Letting $t_0 := \max\left(\frac{\beta}{\alpha \min(1, 2\alpha)} - 1, 0\right)$ note that

$$\eta_t = \begin{cases} \frac{1}{\alpha(t+1)} & \text{if } t \geq t_0 \\ \frac{\min(1, 2\alpha)}{\beta} & \text{otherwise.} \end{cases}$$

If $t \geq t_0$, then $\eta_k = \frac{1}{\alpha(k+1)}$ for all $k \geq t$, and we see that

$$\begin{aligned} \|\mathcal{D}_{\Lambda_t(o,d)}[x_N](L)\|_{op} &\leq C\eta_t \exp\left(-\sum_{k=t+1}^{N-1} \eta_k \alpha\right) \\ &= \frac{C}{\alpha(t+1)} \exp\left(-\alpha \sum_{k=t+1}^{N-1} \frac{1}{\alpha(k+1)}\right) \end{aligned}$$

$$\begin{aligned}
&\leq \frac{C}{\alpha(t+1)} \exp\left(-\alpha \int_{t+1}^N \frac{1}{\alpha y} dy\right) \\
&= \frac{C}{\alpha(t+1)} \exp(-\ln N + \ln(t+1)) \\
&= \frac{C}{\alpha(t+1)} \frac{t+1}{N} = \frac{C}{\alpha N}
\end{aligned}$$

On the other hand, if $t \leq t_0$, then $\eta_t = \frac{\min(1, 2\alpha)}{\beta}$ and we see that

$$\begin{aligned}
\|\mathcal{D}_{\Lambda_t(o,d)}[x_N](L)\|_{op} &\leq C\eta_t \exp\left(-\sum_{k=t+1}^{N-1} \eta_k \alpha\right) \\
&= \frac{C \min(1, 2\alpha)}{\beta} \exp\left(-\alpha \sum_{k=t+1}^{N-1} \frac{1}{\alpha(k+1)}\right) \\
&\leq \frac{C \min(1, 2\alpha)}{\beta} \exp\left(-\alpha \sum_{k=\lfloor t_0 \rfloor + 1}^{N-1} \frac{1}{\alpha(k+1)}\right) \\
&\leq \frac{C \min(1, 2\alpha)}{\beta} \exp\left(-\alpha \int_{t_0+1}^N \frac{1}{\alpha y} dy\right) \\
&\leq \frac{C \min(1, 2\alpha)}{\beta} \frac{t_0 + 1}{N} \\
&= \frac{C \min(1, 2\alpha)}{\beta} \frac{\beta}{\alpha \min(1, 2\alpha) N} = \frac{C}{\alpha N}.
\end{aligned}$$

Thus in either case we have $\|\mathcal{D}_{\Lambda_t(o,d)}[x_N](L)\|_{op} \leq \frac{C}{\alpha N}$.

Finally, note that (11) implies $\|\mathcal{D}_{\Lambda_t(o,d)}[x_N](L)\|_{op} \leq C\eta_t \leq C \frac{\min(1, 2\alpha)}{\beta}$. Combining these two bounds gives the desired result:

$$\|\mathcal{D}_{\Lambda_t}[x_N](L)\|_{op} \leq \min\left(\frac{C \min(1, 2\alpha)}{\beta}, \frac{C}{\alpha N}\right).$$

□

C Performance Analysis (Proof of Theorem 3)

Let x^* be a solution to (6). Note that

$$\begin{aligned}
& \|x_{k+1} - x^*\|_2^2 \\
&= \left\| \prod_{\mathcal{X}} (x_k - \eta_k \nabla F(x_k, \Lambda_k)) - x^* \right\|_2^2 \\
&\stackrel{(a)}{=} \left\| \prod_{\mathcal{X}} (x_k - \eta_k \nabla F(x_k, \Lambda_k)) - \prod_{\mathcal{X}} (x^* - \eta_k \nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x^*, \Lambda)]) \right\|_2^2 \\
&\stackrel{(b)}{\leq} \|x_k - \eta_k \nabla F(x_k, \Lambda_k) - (x^* - \eta_k \nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x^*, \Lambda)])\|_2^2 \\
&= \|x_k - x^* - \eta_k (\nabla F(x_k, \Lambda_k) - \nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x^*, \Lambda)])\|_2^2 \\
&= \|x_k - x^*\|_2^2 + \eta_k^2 \|\nabla F(x_k, \Lambda_k) - \nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} F(x^*, \Lambda)\|_2^2 - 2\eta_k (\nabla F(x_k, \Lambda_k) - \nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} F(x^*, \Lambda))^\top (x_k - x^*)
\end{aligned}$$

Taking expectation of both sides conditioned on x_k , we see that

$$\begin{aligned}
& \mathbb{E} \left[\|x_{k+1} - x^*\|_2^2 \middle| x_k \right] \\
&\leq \mathbb{E} \left[\|x_k - x^*\|_2^2 + \underbrace{\eta_k^2 \|\nabla F(x_k, \Lambda_k) - \nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} F(x^*, \Lambda)\|_2^2}_{\text{Term 2}} - \underbrace{2\eta_k (\nabla F(x_k, \Lambda_k) - \nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} F(x^*, \Lambda))^\top (x_k - x^*)}_{\text{Term 2}} \middle| x_k \right]
\end{aligned}$$

To show that x_{k+1} is closer to x^* than x_k is, we will provide bounds for both Term 2 and Term 3.

C.1 Bounding Term 2

To upper bound Term 2, noting that $\mathbb{E}_{\Lambda_k}[\nabla f_k(x)] = \nabla f(x)$ for all $x \in \mathcal{X}$, we have

$$\begin{aligned}
& \eta_k^2 \mathbb{E} \left[\|\nabla F(x_k, \Lambda_k) - \nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x^*, \Lambda)]\|_2^2 \middle| x_k \right] \\
&= \eta_k^2 \mathbb{E} \left[\left\| \underbrace{\nabla F(x_k, \Lambda_k) - \nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x_k, \Lambda)]}_A + \underbrace{\nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x_k, \Lambda)] - \nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x^*, \Lambda)]}_B \right\|_2^2 \middle| x_k \right] \\
&= \eta_k^2 \mathbb{E} \left[\|A\|_2^2 + \|B\|_2^2 + 2A^\top B \middle| x_k \right] \tag{12}
\end{aligned}$$

By Assumption 6 and the dominated convergence theorem, we have $\nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x_k, \Lambda)] = \mathbb{E}_{\Lambda \sim \mathcal{Q}} [\nabla F(x_k, \Lambda)]$. Hence we see that

$$\mathbb{E} [A | x_k] = \mathbb{E}_{\Lambda_k \sim \mathcal{Q}} [F(x_k, \Lambda_k) | x_k] - \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x_k, \Lambda) | x_k] = 0.$$

From this observation we have the three following remarks:

1. Since A is a zero mean random vector, $\mathbb{E} \left[\|A\|_2^2 \middle| x_k \right]$ is the variance of $\nabla F(x_k, \Lambda)$ given x_k . By Assumption 2, $\mathbb{E}_{\Lambda \sim \mathcal{Q}} \left[\|\nabla F(x, \Lambda)\|_2^2 \right] \leq K^2$ for any x , which implies that $\mathbb{E} \left[\|A\|_2^2 \middle| x_k \right] \leq K^2$.
2. By Assumption 5, $F(\cdot, \Lambda)$ is β -smooth for every Λ . β -smoothness of $F(\cdot, \Lambda)$ implies that $\nabla F(\cdot, \Lambda)$ is β -lipschitz. Thus we have

$$\begin{aligned}
\mathbb{E} \left[\|B\|_2^2 \middle| x_k \right] &= \mathbb{E} \left[\|\mathbb{E}_{\Lambda \sim \mathcal{Q}} [\nabla F(x_k, \Lambda) - \nabla F(x^*, \Lambda)]\|_2^2 \middle| x_k \right] \\
&\stackrel{(a)}{\leq} \mathbb{E} \left[\mathbb{E}_{\Lambda \sim \mathcal{Q}} \left[\|\nabla F(x_k, \Lambda) - \nabla F(x^*, \Lambda)\|_2^2 \middle| x_k \right] \right] \\
&\stackrel{(b)}{\leq} \mathbb{E} \left[\mathbb{E}_{\Lambda \sim \mathcal{Q}} \left[\beta^2 \|x_k - x^*\|_2^2 \middle| x_k \right] \right] \\
&= \beta^2 \|x_k - x^*\|_2^2.
\end{aligned}$$

where (a) is due to Jensen's inequality and (b) is due to $\nabla F(\cdot, \Lambda)$ being β -Lipschitz.

3. Conditioned on x_k , A is zero mean and B is constant, meaning that $A^\top B$ is a zero mean random vector. Therefore $\mathbb{E}[A^\top B \mid x_k] = 0$.

Applying these three remarks to the inequality (12) we see that

$$\eta_k^2 \mathbb{E} \left[\|\nabla F(x_k, \Lambda_k) - \nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x^*, \Lambda)]\|_2^2 \middle| x_k \right] \leq \eta_k^2 \left(K^2 + \beta^2 \|x_k - x^*\|_2^2 \right). \quad (13)$$

C.2 Bounding Term 3

By linearity of expectation, Term 3 is equal to

$$\begin{aligned}
&-2\eta_k \mathbb{E} \left[(\nabla F(x_k, \Lambda_k) - \nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x^*, \Lambda)])^\top (x_k - x^*) \middle| x_k \right] \\
&= -2\eta_k \left(\mathbb{E} [\nabla f(x_k, \Lambda_k) \mid x_k] - \nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x^*, \Lambda)] \right)^\top (x_k - x^*) \\
&= -2\eta_k \left(\mathbb{E}_{\Lambda \sim \mathcal{Q}} [\nabla F(x_k, \Lambda) - \nabla F(x^*, \Lambda)] \right)^\top (x_k - x^*).
\end{aligned}$$

Define $f(x) := \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x, \Lambda)]$. We can re-write the above equation as

$$-2\eta_k \mathbb{E} \left[(\nabla F(x_k, \Lambda_k) - \nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x^*, \Lambda)])^\top (x_k - x^*) \middle| x_k \right] \leq -2\eta_k (\nabla f(x_k) - \nabla f(x^*))^\top (x_k - x^*)$$

Since $F(\cdot, \Lambda)$ is α -strongly convex for every Λ , we can conclude that f is also α -strongly convex. To upper bound Term 3, we use α -strong convexity of f to conclude that

$$\begin{aligned}
f(x_k) &\geq f(x^*) + \nabla f(x^*)^\top (x_k - x^*) + \frac{\alpha}{2} \|x_k - x^*\|_2^2 \\
&\text{and} \\
f(x^*) &\geq f(x_k) + \nabla f(x_k)^\top (x^* - x_k) + \frac{\alpha}{2} \|x_k - x^*\|_2^2.
\end{aligned}$$

Adding these inequalities together gives

$$f(x_k) + f(x^*) \geq f(x_k) + f(x^*) - (\nabla f(x_k) - \nabla f(x^*))^\top (x_k - x^*) + \alpha \|x_k - x^*\|_2^2$$

$$\implies -\alpha \|x_k - x^*\|_2^2 \geq -(\nabla f(x_k) - \nabla f(x^*))^\top (x_k - x^*).$$

This inequality implies the following bound on Term 3:

$$-2\eta_k \mathbb{E} \left[(\nabla F(x_k, \Lambda_k) - \nabla \mathbb{E}_{\Lambda \sim \mathcal{Q}} [F(x^*, \Lambda)])^\top (x_k - x^*) \middle| x_k \right] \leq -2\alpha\eta_k \|x_k - x^*\|_2^2. \quad (14)$$

C.3 Putting everything together

Applying the bounds (13) and (14) on Term 2 and Term 3 respectively, we see that

$$\begin{aligned} \mathbb{E} \left[\|x_{k+1} - x^*\|_2^2 \middle| x_k \right] &\leq \|x_k - x^*\|_2^2 + \eta_k^2 K^2 + \beta^2 \eta_k^2 \|x_k - x^*\|_2^2 - 2\alpha\eta_k \|x_k - x^*\|_2^2 \\ &= (1 - 2\alpha\eta_k + \beta^2 \eta_k^2) \|x_k - x^*\|_2^2 + \eta_k^2 K^2. \end{aligned} \quad (15)$$

By the tower property of expectation, we can write

$$\mathbb{E} \left[\|x_N - x^*\|_2^2 \right] = \mathbb{E} \left[\dots \mathbb{E} \left[\mathbb{E} \left[\|x_N - x^*\|_2^2 \middle| x_{N-1} \right] \middle| x_{N-2} \right] \dots \middle| x_0 \right]. \quad (16)$$

Combining the recursive relation from (15) with (16) we see that

$$\begin{aligned} \mathbb{E} \left[\|x_N - x^*\|_2^2 \right] &\leq \|x_0 - x^*\|_2^2 \left(\prod_{t=0}^{N-1} 1 - 2\alpha\eta_t + \beta^2 \eta_t^2 \right) + K^2 \sum_{t=0}^{N-1} \eta_t^2 \left(\prod_{k=t+1}^{N-1} 1 - 2\alpha\eta_k + \beta^2 \eta_k^2 \right) \\ &\leq \|x_0 - x^*\|_2^2 \exp \left(\sum_{t=0}^{N-1} -2\alpha\eta_t + \beta^2 \eta_t^2 \right) + K^2 \sum_{t=0}^{N-1} \eta_t^2 \exp \left(\sum_{k=t+1}^{N-1} -2\alpha\eta_k + \beta^2 \eta_k^2 \right). \end{aligned}$$

Since we chose $\eta_k := \min \left(\frac{1}{\alpha k}, \frac{\min(1, 2\alpha)}{\beta} \right)$, we have $\eta_k \leq \frac{1}{\alpha k}$. This means that $\sum_{k=1}^{\infty} \eta_k^2 \leq \sum_{k=1}^{\infty} \frac{1}{\alpha^2 k^2} = \frac{\pi^2}{6\alpha^2}$. Thus defining $C_{\alpha, \beta} := \exp \left(\frac{\beta^2 \pi^2}{6\alpha^2} \right)$, we have

$$\begin{aligned} \mathbb{E} \left[\|x_N - x^*\|_2^2 \right] &\leq C_{\alpha, \beta} \|x_0 - x^*\|_2^2 \exp \left(\sum_{t=0}^{N-1} -2\alpha\eta_t \right) + K^2 C_{\alpha, \beta} \sum_{t=0}^{N-1} \eta_t^2 \exp \left(\sum_{k=t+1}^{N-1} -2\alpha\eta_k \right) \\ &\stackrel{(a)}{\leq} C_{\alpha, \beta} \|x_0 - x^*\|_2^2 \exp \left(\sum_{t=0}^{N-1} -2\alpha\eta_t \right) + K^2 C_{\alpha, \beta} \sum_{t=0}^{N-1} \left(\frac{C}{\alpha N} \right)^2 \\ &= C_{\alpha, \beta} \|x_0 - x^*\|_2^2 \exp \left(\sum_{t=0}^{N-1} -2\alpha\eta_t \right) + \frac{K^2 C^2 C_{\alpha, \beta}}{\alpha^2 N} \end{aligned}$$

where (a) is because in Appendix B we showed that $\eta_t \exp \left(\sum_{k=t+1}^{N-1} -\alpha\eta_k \right) \leq \frac{C}{\alpha N}$ for all $0 \leq t \leq N$. Next, let $t_0 := \max \left(\frac{\beta}{\alpha \min(1, 2\alpha)} - 1, 0 \right)$ so that $\eta_t = \frac{1}{\alpha t}$ if $t \geq t_0$ and $\eta_t = \frac{\min(1, 2\alpha)}{\beta}$ otherwise. We then have

$$\mathbb{E} \left[\|x_N - x^*\|_2^2 \right] \leq C_{\alpha, \beta} \|x_0 - x^*\|_2^2 \exp \left(\sum_{t=0}^{N-1} -2\alpha\eta_t \right) + \frac{K^2 C^2 C_{\alpha, \beta}}{\alpha^2 N}$$

$$\begin{aligned}
&= C_{\alpha,\beta} \|x_0 - x^*\|_2^2 \exp\left(\sum_{t=0}^{t_0} -2\alpha\eta_t\right) \exp\left(\sum_{t=t_0+1}^{N-1} -2\alpha\eta_t\right) + \frac{K^2 C^2 C_{\alpha,\beta}}{\alpha^2 N} \\
&\leq C_{\alpha,\beta} \|x_0 - x^*\|_2^2 \exp\left(2\alpha \sum_{t=t_0+1}^{N-1} -\eta_t\right) + \frac{K^2 C^2 C_{\alpha,\beta}}{\alpha^2 N} \\
&\leq C_{\alpha,\beta} \|x_0 - x^*\|_2^2 \exp\left(-2\alpha \int_{t_0+1}^N \frac{1}{\alpha y} dy\right) + \frac{K^2 C^2 C_{\alpha,\beta}}{\alpha^2 N} \\
&= C_{\alpha,\beta} \|x_0 - x^*\|_2^2 \left(\frac{t_0+1}{N}\right)^2 + \frac{K^2 C^2 C_{\alpha,\beta}}{\alpha^2 N} \\
&= \frac{C_{\alpha,\beta} \|x_0 - x^*\|_2^2 \beta^2}{(\alpha \min(1, 2\alpha))^2 N^2} + \frac{K^2 C^2 C_{\alpha,\beta}}{\alpha^2 N}.
\end{aligned}$$

Finally, we have

$$\begin{aligned}
\mathbb{E} [\|x_{\text{alg}} - x^*\|_2] &= \mathbb{E} [\|\Pi_{\mathcal{X}}(x_N + Z) - x^*\|_2] \\
&= \mathbb{E} [\|\Pi_{\mathcal{X}}(x_N + Z) - \Pi_{\mathcal{X}}(x^*)\|_2] \\
&\leq \mathbb{E} [\|x_N + Z - x^*\|_2] \\
&\leq \mathbb{E} [\|x_N - x^*\|_2] + \mathbb{E} [\|Z\|_2] \\
&\stackrel{(a)}{\leq} \sqrt{\mathbb{E} [\|x_N - x^*\|_2^2]} + \sqrt{\mathbb{E} [\|Z\|_2^2]} \\
&\leq \sqrt{\frac{C_{\alpha,\beta} \|x_0 - x^*\|_2^2 \beta^2}{(\alpha \min(1, 2\alpha))^2 N^2} + \frac{K^2 C^2 C_{\alpha,\beta}}{\alpha^2 N}} + \sqrt{n^2 m \frac{2C^2}{\epsilon^2 \alpha^2 T^2 N^2} \ln\left(\frac{1.25}{\delta}\right)} \\
&\leq \frac{\sqrt{C_{\alpha,\beta}} \|x_0 - x^*\|_2 \beta}{(\alpha \min(1, 2\alpha)) N} + \frac{KC\sqrt{C_{\alpha,\beta}}}{\alpha\sqrt{N}} + \frac{Cn\sqrt{m}}{\epsilon\alpha TN} \sqrt{2 \ln\left(\frac{1.25}{\delta}\right)} \\
&\stackrel{(b)}{\leq} \frac{\sqrt{C_{\alpha,\beta}} n^2 m \beta}{(\alpha \min(1, 2\alpha)) N} + \frac{KC\sqrt{C_{\alpha,\beta}}}{\alpha\sqrt{N}} + \frac{Cn\sqrt{m}}{\epsilon\alpha TN} \sqrt{2 \ln\left(\frac{1.25}{\delta}\right)}
\end{aligned}$$

where (a) is due to Jensen's inequality and (b) is due to the fact that $\|x' - x''\|_2 \leq n\sqrt{m}$ for any pair of unit network flows x', x'' .

	Error #1	Error #2	Error #3
N	$\frac{\langle F_\alpha(x_\alpha^{opt}, \Lambda_i) \rangle - F_\alpha(x_\alpha^{opt}, \langle \Lambda_i \rangle)}{\langle F_\alpha(x_\alpha^{opt}, \Lambda_i) \rangle}$	$\frac{\langle F_\alpha(x_\alpha, \Lambda_i) \rangle - \langle F_\alpha(x_\alpha^{opt}, \Lambda_i) \rangle}{\langle F_\alpha(x_\alpha^{opt}, \Lambda_i) \rangle}$	$\frac{\langle F_{\alpha=0}(x_\alpha, \Lambda_i) \rangle - \langle F_{\alpha=0}(x_{\alpha=0}, \Lambda_i) \rangle}{\langle F_{\alpha=0}(x_{\alpha=0}, \Lambda_i) \rangle}$
10	8.97×10^{-7}	$< 10^{-10}$	8.17×10^{-3}
25	1.02×10^{-6}	$< 10^{-10}$	8.16×10^{-3}
50	9.84×10^{-7}	$< 10^{-10}$	8.17×10^{-3}

Table 2: Approximation errors.

D Computational approximations

In our experiments, we make several approximations for computational tractability. In this section, we provide empirical evidence that these approximations are reasonable and do not introduce significant errors. For ease of discussion, we define the following notations.

- x_α^{opt} : Solution obtained by solving the optimization problem (7) with regularizer α
- x_α : Solution obtained by solving the optimization problem with the average demand and a regularizer α
- $F_\alpha(x, \Lambda_i)$: Evaluating the routing policy x on demand Λ_i with a regularizer α
- $\langle F_\alpha(x, \Lambda_i) \rangle$: Evaluating the average cost of the routing policy x on the set of demands $\{\Lambda_1 \dots, \Lambda_N\}$ with a regularizer α . More precisely, $\langle F_\alpha(x, \Lambda_i) \rangle = \frac{1}{N} \sum_{i=1}^N F_\alpha(x, \Lambda_i)$
- $F_\alpha(x, \langle \Lambda_i \rangle)$: Evaluating the cost of the routing policy x on the average demands $\langle \Lambda_i \rangle = \frac{1}{N} \sum_{i=1}^N \Lambda_i$ with a regularizer α

Table 2 presents errors from three different approximations. Our first approximation is to compute the system costs for a given policy by using the average demand instead of averaging the costs over every observed demand. The first column (denoted as Error #) lists the fractional error introduced by this approximation for different values of N when using the optimal routing policy. We note that the error is less than 10^{-6} , and we observe a 30-50X improvement in computational time when evaluating the costs using the average flow. This justifies the use of the average demand for estimating costs. Our second approximation is in solving an easier optimization problem to compute the optimal routing policy. In this case, the exact approach would be to solve the optimization problem described in Equation 7. However, the size of this problem grows rapidly with the number of data points N . Our approximation involves solving the easier optimization problem of maximizing $F_\alpha(x, \langle \Lambda_i \rangle)$ to obtain the routing policy x_α instead of solving the original optimization problem to obtain x_α^{opt} . The second column of the table (titled Error #2) presents the error introduced due to this approximation on the travel costs. The small errors indicate that this assumption is reasonable, and helps us obtain upto a 30X speedup in solving the optimization problem. Finally, we show through numerical evaluations that the addition of the $\alpha = 10^3$ regularizer does not change the travel costs significantly (third column, denoted as Error #3), and results in less than a 0.01% error in the total travel cost.