

# TIPS: Mining Top-K Locations to Minimize User-Inconvenience for Trajectory-Aware Services

Shubhadip Mitra, Priya Saraf, Arnab Bhattacharya

Dept. of Computer Science and Engineering, Indian Institute of Technology, Kanpur, India.

{smitr,priyas,arnabb}@cse.iitk.ac.in



**Abstract**—*Facility location* problems aim to identify the best locations to set up new services. Majority of the existing works typically assume that the users are *static*. However, there exists a wide array of services such as fuel stations, ATMs, food joints, etc., that are widely accessed by mobile users besides the static ones. Such *trajectory-aware services* should, therefore, factor in the trajectories of its users rather than simply their static locations. In this work, we introduce the problem of optimal placement of facility locations for such trajectory-aware services that *minimize the user inconvenience*. The inconvenience of a user is the *extra* distance traveled by her from her regular path to avail a service. We call this the TIPS problem (*Trajectory-aware Inconvenience-minimizing Placement of Services*) and consider two variants of it. The goal of the first variant, MAX-TIPS, is to minimize the *maximum* inconvenience faced by any user, while that of the second, AVG-TIPS, is to minimize the *average* inconvenience over all the users. We show that both these problems are NP-hard, and propose multiple efficient heuristics to solve them. Empirical evaluation on real urban-scale road networks validate the efficiency and effectiveness of the proposed heuristics.

**Index Terms**—Trajectory-aware service, User inconvenience, TIPS.

## 1 INTRODUCTION

*Facility location* problems identify the best locations to set up new facilities (or services) for its users [1], [2]. This has been also studied as *optimal location* problems [3], [4], [5], [6], [7]. Majority of the existing works, however, assume that the users of the service are *static* [8], [9], [10]. However, services such as fuel stations, automobile service stations, ATMs, food joints, convenience stores, etc., are widely accessed by mobile users, besides the static users [11]. For example, it is common for many users to make their daily purchases while returning from their workplaces. This practice is increasingly becoming common because of rapid expansion of cities due to growing urban population and consequent longer work commute trips. The placement of such services should therefore take into account the *mobility patterns* or the *trajectories* of its users, rather than simply their static home and office locations. The necessity for such *mobility-aware location* selection has also been highlighted in recent studies [11], [12], [13]. Moreover, our experimental findings in Sec. 6.4 suggest that there is 10-40 % cost-savings if we factor in the user-trajectories instead of their static locations.

A *user trajectory* is a sequence of spatial points that lie on the path of a user while travelling. It is important to note that trajectories strictly generalize the static user scenario as static users can always be modeled as trajectories with a single location. In general, however, trajectories capture user location patterns more effectively and realistically.

In this work, we extend two key optimal location problems, namely the *MinMax Location Query* [6], [7], [10], [14] and the *Min-Dist Location Query* [8], [9], [12], [15] (also referred to as the *MinSum Location Query* [16]). Given a set of customers (or users)  $\mathcal{C}$ , a set of existing facilities  $\mathcal{F}$ , and a set of candidate locations that can host a new facility  $\mathcal{S}$ , the goal of the MinMax Location Query (respectively, Min-Dist Location Query) is to identify a facility location in  $\mathcal{S}$  that minimizes the maximum (respectively, average) distance of any user to its nearest facility.

Majority of the existing works assume that the users are *static* and ignore their mobile behavior. Further, most of these works also restrict themselves to reporting a *single* facility location, which is polynomially solvable. Motivated by these two limitations, in this work, we introduce two novel optimal location problems, *MAX-TIPS* and *AVG-TIPS*, that factor the user trajectories and report any desired number of facility locations.

Formally, given a set of trajectories  $\mathcal{T}$ , a set of existing facilities  $\mathcal{F}$ , a set of candidate sites  $\mathcal{S}$ , an integer  $k$ , and a user-fraction  $\gamma \in [0, 1]$ , the MAX-TIPS problem seeks to report a set  $\mathcal{Q} \subseteq \mathcal{S}$  of  $k$  locations that minimizes the *maximum inconvenience* over any  $\gamma$  fraction of the trajectories  $\mathcal{T}$ , while the AVG-TIPS problem aims to identify the  $k$  locations that minimize the *average inconvenience* faced by any user. The *inconvenience* of a user on a trajectory is defined as the *extra distance* traveled by her with respect to her normal trajectory to avail the service. The proposed problems are NP-hard.

While for critical services such as ambulance or fire stations, it is desirable to minimize the *maximum inconvenience*, for other services such as ATMs, fuel stations or convenience stores, it is desirable to minimize the *average inconvenience*. Both these TIPS problems have applications in various resource planning scenarios [15], [16], [17]. Some of the direct applications are:

*Placement of drop-boxes for crowd-sourced taxi shipment service:* Chen et al. [18] presented a novel scheme for city-wide shipment of items using the regular passenger-carrying taxis in a crowd-

sourced manner. The taxis participating in this service are required to collect and drop the shipments at a nearby drop-box whenever they are idle, i.e., not carrying any passenger. The drop-boxes must be located in a manner such that the maximum inconvenience faced by majority of the taxis is minimal.

*Locating multiple ATMs of a given bank:* Suppose a bank plans to set up  $k$  ATMs in a given city. Since mobile users often access nearby ATMs, they must be placed such that the average user-inconvenience is minimized.

*Placement of food trucks of a given chain:* Extending the example given in [17], consider a restaurant chain that wants to place its  $k$  food trucks on the road network, so as to serve the mobile customers (who have shared their trips) such that the average inconvenience faced by any user is minimal.

We illustrate the TIPS problems through an example shown in Fig. 1. There are 6 trajectories  $T_1, \dots, T_6$  (shown with blue dashed lines with arrows indicating the directions of the respective trip), one existing facility at  $s_0$  (marked in green) and 4 candidate sites  $s_1, \dots, s_4$  (shown in red) to host a new facility. The road segments are marked in black with corresponding distances (assumed to be the same on both ways). The nodes  $v_1$  and  $v_2$  (shown in black) are general points on the road network that do not host a facility. If the user on trajectory  $T_1$  (that passes through  $v_1$  and  $s_1$ ) wishes to access the facility at  $s_0$ , she needs to detour from  $v_1$ , visit  $s_0$ , and join her regular path at  $s_1$ . As a result, her inconvenience (i.e., the extra distance traveled) is  $1 + 2 - 2 = 1$  unit. Now if another facility comes up at  $s_1$ , her inconvenience reduces to 0, as there is no detour. For trajectories  $T_2, T_3$ , the inconvenience w.r.t.  $s_0$  is  $7 + 2 + 2 + 7 = 18$  units (as they need to take a round trip via  $s_2, s_1, s_0, s_1, s_2$ , in that order). Similarly, for trajectories  $T_4, T_5$  and  $T_6$ , it is 30, 32 and 20 respectively. Note that the inconvenience is measured w.r.t. the nearest facility from any point on the trajectory.

Assume that the service-provider wants to set up 2 new facilities besides the existing one at  $s_0$ , with the objective to minimize the *maximum* inconvenience faced by any user. If the new facilities are hosted at  $\{s_1, s_2\}$ , the inconvenience of trajectories  $T_1, T_2, T_3$  are 0 units each. For trajectories  $T_4$  and  $T_5$  the nearest facility is  $s_2$  and, therefore, their inconveniences are 12 units each. Similarly, for trajectory  $T_6$ , the nearest facility is  $s_1$  and, thus, its inconvenience is 16 units. Therefore, the maximum inconvenience among all the trajectories due to the selection  $\{s_0, s_1, s_2\}$  is 16 units. The maximum inconvenience for all such selections of 2 new sites (along with  $s_0$ ) are listed in Fig. 1 under the column  $\gamma = 1$ . (We will shortly explain the meaning of  $\gamma$ .) The selection  $\{s_0, s_3, s_4\}$  offers the optimal maximum inconvenience of 12 units. Importantly, although most number of trajectories pass through  $s_2$ , it is *not* part of the optimal solution.

Next, suppose the objective is to minimize the *average* (or equivalently, the *total*) inconvenience over all the user trajectories. The last column in the table in Fig. 1 lists the total inconvenience for all possible selections. The selection  $\{s_0, s_2, s_3\}$  offers the optimal total inconvenience of 21 units. Thus, the optimal average inconvenience is  $21/6 = 3.5$  units. (The optimal maximum inconvenience was 12 units.)

The maximum inconvenience problem suffers from the issue of outlier trajectories where a trajectory is very different from all the other ones and, therefore, accommodating for it becomes harder. Thus, instead of considering all the trajectories, the service provider may choose a fraction of the trajectories, over which the maximum inconvenience will be minimized. We call this fraction

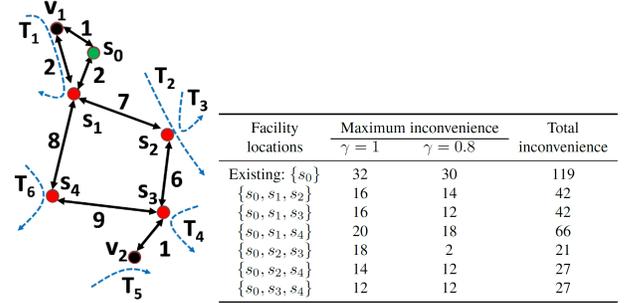


Fig. 1: Illustration of the need for minimizing user-inconvenience for trajectory-aware services.

the *user-fraction*,  $\gamma$ .

Referring to Fig. 1, when  $\gamma = 0.8$ , the goal is to minimize the maximum inconvenience over at least  $0.8 \times 6 = 4.8$  or 5 trajectories. The values of the maximum inconvenience for all the selections for  $\gamma = 0.8$  are listed in the table. Note that the optimal selection for  $\gamma = 0.8$  is  $\{s_0, s_2, s_3\}$  which is different from the optimal selection for  $\gamma = 1$ . Note that the optimal inconvenience for  $\gamma = 0.8$  falls to 2 units as compared to 12 units for  $\gamma = 1$ .

A naive approach to solve either of the TIPS problems involves enumerating all  $k$ -sized subsets of  $\mathcal{S}$ , computing the maximum or the average inconvenience (depending on the objective) for each subset, and returning the subset with the minimal objective value. This requires an exponential time and space complexity, which is infeasible for almost all datasets. Thus, the major challenges in solving TIPS are as follows:

- 1) **Quality:** Since both the problems are NP-hard (proved later), optimal algorithms are impractical. Thus, we need efficient heuristics that offer high quality solutions.
  - 2) **Scalability:** Any basic approach to solve the above problems would typically need to compute and store pairwise distances between the sets of candidate sites and trajectories. However, for any city-scale datasets, this time and storage requirement is prohibitively large (of the order of 100s of GB). Thus, it is necessary to design solutions that are practical and scalable.
- In this paper, we propose efficient heuristics for both MAX-TIPS and AVG-TIPS that overcome the above challenges. To summarize, our major contributions are:
- 1) To the best of our knowledge, this is the *first* work, that factors in user-mobility to identify the best  $k$  facility locations that minimize the maximum and the average user-inconvenience, in presence or absence of *existing facilities*. In particular, we introduce two *facility location problems* over user trajectories, namely, MAX-TIPS and AVG-TIPS (Sec. 3).
  - 2) We show that both these problems are *NP-hard* (Th. 1), and propose one exact algorithm and two polynomial-time efficient *heuristics* to solve each of them (Sec. 4 and Sec. 5). The exact algorithms are based on integer linear programming (Sec.4.1 and 5.1). For MAX-TIPS, while the first heuristic is an index-free greedy algorithm (Sec. 4.2), the second one (Sec. 4.3) uses an index structure based on multi-resolution clustering of the road network. For AVG-TIPS, the first heuristic is a hybrid of two standard clustering algorithms that are based on local search techniques (Sec. 5.2), while the second one (Sec. 5.3) uses a simple greedy approach.
  - 3) Empirical evaluation on urban scale datasets show that our heuristics are *effective* in terms of quality, and *efficient* in terms of space and running time (Sec. 6).

## 2 RELATED WORK

The related work is divided into the following main classes.

**Location Selection (LS) Queries:** Location Selection (LS) queries identify best locations to set up new facilities. They are broadly of two types: Optimal Location (OL) queries [3] and Facility Location (FL) Queries [2]. An OL query typically has three inputs, a set of existing facilities, a set of users, and a set of candidate sites, and the aim is to find a candidate site to host a new facility that optimizes an objective function based on the distances among the facilities and the users. On the other hand, an FL query considers a set of users, and a set of candidate sites, and seeks to identify  $k$  ( $k \geq 1$ ) facility locations that optimize certain objective function which is usually based on the distances among the users and the facilities. While OL queries are polynomially solvable, FL queries are NP-hard. In Table 1, we summarize the key related works in the area of location selection based on the following attributes.

**(1) Type of Users:** Majority of the early works assumed that the users are fixed to a single location such as home, and did not consider their mobile behavior. However, many recent studies including ours factor in user-mobility.

**(2) Type of Output:** While many works report a single location to set up a new facility, others report top- $k$  facility locations that collectively optimize the desired objective. While the former is polynomially solvable, the latter is NP-hard. Both the TIPS problems return  $k$  locations, and are NP-hard as well.

**(3) Type of Objective:** Based on the objective function, the related works can be classified into two categories: *distance minimizing* and *influence maximizing*. In distance minimizing queries, the goal is to minimize the maximum or the average distance of any user to its nearest facility. While the former is referred to as the MinMax Location Query [6], [7], [10], [14], the latter is known as the Min-Dist Location Query or MinSum Location Query [8], [9], [12], [15], [16]. MAX-TIPS problem is a generalization of the MinMax Location query, while AVG-TIPS is a generalization of the Min-Dist Location query. These generalizations are in terms of reporting  $k$  locations instead of a single location, and considering mobile users on a road network instead of static ones. The aim of influence maximizing queries is to find a candidate site that has maximal *influence* over its users. The influence of a site  $s$  is the number of users for which  $s$  is the nearest facility. These problems are usually modeled as *Reverse Nearest Neighbor* queries. The key difference between these works and ours is that they assume that the new facility is *competing* with the existing ones; in our model, both the new and the existing facilities (if any) belong to a given service provider and, hence, complement each other. This is especially true for public services such as ATMs, post offices, hospitals, gas stations, parking spots, etc.

**(4) Type of Underlying Space:** Many earlier models assume that the underlying space is Euclidean. Since user movements are typically restricted by a road network, and network distances can significantly vary from corresponding Euclidean distances, recent works base their studies on road networks. Our TIPS formulation is also based on the road network. The massive distance computations, and absence of geometric properties make the latter problems more challenging.

Referring to Table 1, we note that this is the *first* work that factors in *user-mobility* to mine the *top-k* facility locations that minimizes the *inconvenience* (defined in terms of the distances between the users and facilities) caused to the users traveling on

a *road network*. Next, we discuss the important LS works that are closely related to our proposed TIPS query.

• **Key related works factoring in user-mobility:** In [13], the authors studied the PRIME-LS problem that aims to find an optimal location which can influence the most number of moving objects. Two algorithms were proposed based on two pruning techniques and optimization strategies to filter out unpromising candidate sites. However, since the goal is to maximize the influence rather than to minimize the distance, the attributes of this problem are quite different to ours (as explained earlier). Thus, these algorithms are inapplicable to solve the TIPS problems. The studies in [11], [22] consider the FL problem over user trajectories moving on a road network. They assume that a user is attracted to a facility if its trajectory lies within a specified distance threshold. While both these works aim to maximize the user coverage, [22] reports a *single* optimal road segment, and [11] reports the  $k$  best facility locations.

MinMax Location Query and Min-Dist Location Query have been studied in [4], [6], [7], [8], [9], [10], [12], [14], [15], [16], [17], [20]. All other works except [12], [17] assume the users to be static. However, in contrast to our work, both these works find a *single* optimal location for the min-dist location problem over user trajectories. Moreover, since [17] study the problem over Euclidean space, their techniques are not applicable to our model which is based on a road network. The study of [12] is however based on a road network. Based on reference location transformation, they propose two groups of algorithms. While the first group uses spatial locality based index structures, the other group does not use any index structure but computes from scratch. In our empirical study, we consider this algorithm as a baseline.

Majority of the works in the FL literature also assume that the users are static [1], [2]. Works that consider human mobility include [23], [24], [25], [26], [27], [28], [29]. These works, however, assume a flow model to characterize mobility instead of using real trajectories. The proposed models are mostly theoretical and are not scalable for real city-scale road networks. In particular, all these approaches require extensive distance computations which leads to large memory overhead and are, hence, infeasible [11]. A fairly comprehensive literature survey is available in [30].

Hodgson [23] posed the first FL problem that minimizes the average user inconvenience as a generalization of  $k$ -medians problem. In [26], it was shown that the problem does not admit constant factor approximation unless  $P = NP$ . However, no approximation algorithm was proposed. In contrast, we propose two heuristics for this problem.

**Clustering Problems:** The following clustering problems are related to the current work.

**(1)  $k$ -Center Problem:** Given a set  $\mathcal{S}$  of  $n$  points, the  $k$ -center problem is to determine a set  $\mathcal{Q} \subseteq \mathcal{S}$  of size  $k$ , referred to as *centers*, such that the maximum distance of any point in  $\mathcal{S}$  to its nearest center is minimized. This problem was introduced and proved to be NP-hard in [31]. They also proposed a greedy heuristic that offers a factor of 2 approximation. Our proposed MAX-TIPS problem is a generalization of the  $k$ -center problem as trajectories generalize static users. We have extended the greedy algorithm in [31] to design a heuristic (with bounded quality guarantees) to solve MAX-TIPS (Sec. 4).

**(2)  $k$ -Medoids Problem:** Given a set  $\mathcal{S}$  of  $n$  points, the  $k$ -medoids problem is to determine a set  $\mathcal{Q} \subseteq \mathcal{S}$  of size  $k$ , referred to as *medoids*, such that the sum of distances of each point in  $\mathcal{S}$  to its nearest medoid is minimized [32]. This is also referred to as the

Property	[3]	[4]	[5]	[6]	[7]	[8]	[9]	[10]	[11]	[12]	[13]	[14]	[15]	[16]	[17]	[19]	[20]	[21]	[22]
Users	Stat	Stat	Stat	Stat	Stat	Stat	Stat	Stat	Mob	Mob	Mob	Stat	Stat	Stat	Mob	Stat	Stat	Stat	Mob
Output	Sngl	Sngl	Sngl	Sngl	Sngl+Top-k	Sngl	Sngl	Top-k	Top-k	Sngl	Sngl	Sngl+Top-k	Sngl	Sngl	Sngl	Sngl	Sngl	Top-k	Sngl
Objective function	Infl	Dist	Infl	Dist	Dist	Dist	Dist	Dist	Infl	Dist	Infl	Dist	Dist	Infl	Dist	Infl	Dist	Dist	Infl
Underlying space	Eucl	Eucl	Road	Road	Road	Eucl	Eucl	Road	Road	Road	Eucl	Road	Eucl	Road	Eucl	Eucl	Road+Eucl	Eucl	Road

TABLE 1: Summary of Related Work. (Stat: static, Mob: mobile; Sngl: single, Top-k: top-k; Dist: distance, Infl: influence; Eucl: Euclidean space, Road: road network.)

*k*-median problem. The first constant factor algorithm for the *k*-median problem in general metric space, with an approximation ratio of  $6\frac{2}{3}$  was proposed in [33]. Later, [34] improved this factor to 6. Subsequently, [35] designed a 4-approximation algorithm that runs in  $O(n^3)$  time. Korupolu et al. [36] proposed a local search based approximation scheme by allowing a constant factor blow-up in *k*. Arya et al. [37] improved the approximation bound for the metric *k*-medians problem to  $3 + 2/p$  where *p* is the number of medians swapped simultaneously. Three popular local search based techniques for the *k*-medoids problem are PAM [32], CLARA [32] and CLARANS [38]. These schemes are detailed in Sec. 5.2. Most of these works, however, fail to scale on large datasets. Moreover, all these algorithms require the distance matrix between each pair of points. This quadratic space and time requirement renders them infeasible for real datasets. Our proposed AVG-TIPS problem generalizes the *k*-medoids problem as explained in Sec. 3. The HCC heuristic to solve AVG-TIPS (Sec. 5.2) builds on the ideas of CLARA and CLARANS. More importantly, to address the challenge of high space overhead, we design sampling techniques to reduce the sizes of the sets.

### 3 THE TIPS PROBLEM

Consider a road network  $G = \{V, E\}$  over a geographical area where  $V = \{v_1, \dots, v_N\}$  denotes the set of road intersections, and  $E$  denotes the road segments between two adjacent road intersections. To model the direction of the underlying traffic that passes over a road segment, we assume that the edges are directed. Assume a set of candidate sites  $\mathcal{S} = \{s_1, \dots, s_n\}$  where a certain service or facility can be set up. The set  $\mathcal{S}$  can be in addition to the existing facility locations  $\mathcal{F}$ . Without loss of generality, we can augment the vertices  $V$  to include all the sites. Thus,  $\mathcal{S} \subseteq V$ . Further, we also assume that the set of existing facilities  $\mathcal{F} \subseteq V$ .

The set of trajectories is denoted by  $\mathcal{T} = \{T_1, \dots, T_m\}$  where each trajectory  $T_j = \{v_{j_1}, \dots, v_{j_l}\}$ ,  $v_{j_i} \in V$ , is a sequence of locations that the user passes through. Usually any service or facility is used by a mix of static users and mobile users. As the above definition of trajectory allows both single and multiple locations of a single user, it captures both static and mobile users simultaneously. The trajectories are usually recorded as GPS traces and may contain arbitrary spatial points on the road network. For our purpose, each trajectory is map-matched [39] to form a sequence of road intersections through which it passes.

We also assume that each trajectory belongs to a separate user. The framework can be easily generalized to multiple trajectories belonging to a single user. The union of road intersections that the user passes through in any of her trajectories will be treated as the nodes in her trajectory. In effect, this will minimize her inconvenience from any of her trajectories.

Suppose  $d(v_i, v_j)$  denotes the shortest road network distance along a directed path from node  $v_i$  to  $v_j$ , and  $d_r(v_i, v_j)$  denotes the shortest distance of a round-trip starting at node  $v_i$ , visiting  $v_j$ , and returning to  $v_i$ , i.e.,  $d_r(v_i, v_j) = d(v_i, v_j) + d(v_j, v_i)$ . In general,  $d(v_i, v_j) \neq d(v_j, v_i)$ , but  $d_r(v_i, v_j) = d_r(v_j, v_i)$ .

The extra distance traveled by any user on trajectory  $T_j$  to avail a service at site  $s_i \in V$ , denoted by  $d_r(T_j, s_i)$ , is defined as follows:  $d_r(T_j, s_i) = \min_{v_k, v_l \in T_j} \{d(v_k, s_i) + d(s_i, v_l) - d(v_k, v_l)\}$ , i.e., it deviates from its trajectory at  $v_k$ , reach site  $s_i$  and then return to  $v_l \in T_j$  such that the deviation is minimum.

The round-trip distance between two trajectories  $T_i$  and  $T_j$  is defined as the minimum pairwise distance among its sites:  $d_r(T_i, T_j) = \min_{v_i \in T_i, v_j \in T_j} \{d_r(v_i, v_j)\}$ . Henceforth, *distance* implies *round-trip distance* unless mentioned otherwise.

It is inconvenient for a user to avail a service if the nearest service location is far off from her trajectory. We define this inconvenience as follows. Given a set of service locations  $\mathcal{Q} \subseteq \mathcal{S}$ , the *inconvenience* of a user on trajectory  $T_j$ , denoted by  $I_j$ , is the extra distance travelled to avail a service at the nearest service location in  $\mathcal{Q}$ . Formally,  $I_j = \min_{s_i \in \mathcal{Q}} \{d_r(T_j, s_i)\}$ .

Using the above setting, we introduce two novel problems, namely MAX-TIPS and AVG-TIPS, described as follows.

The MAX-TIPS problem aims to report a set of *k* service locations that minimizes the maximum inconvenience over a given *user-fraction* of the set of trajectories.

**Problem 1 (MAX-TIPS).** *Given a set of trajectories  $\mathcal{T}$ , a set of existing facilities  $\mathcal{F}$ , a set of candidate sites  $\mathcal{S}$  that can host the services, a positive integer *k*, and a user-fraction  $\gamma$  ( $0 < \gamma \leq 1$ ), the MAX-TIPS problem seeks to report a set  $\mathcal{Q} \subseteq \mathcal{S}$ ,  $|\mathcal{Q}| = k$ , that minimizes the maximum inconvenience over any set  $\mathcal{T}' \subseteq \mathcal{T}$  such that  $|\mathcal{T}'| \geq \gamma \times |\mathcal{T}|$ , i.e., it minimizes  $MI(\mathcal{Q}) = \max_{T_j \in \mathcal{T}'} \{I_j\}$ , where  $I_j = \min_{s_i \in \mathcal{Q} \cup \mathcal{F}} \{d_r(T_j, s_i)\}$ .*

Intuitively, as the user-fraction increases, the optimal value of  $MI(\mathcal{Q})$  increases because of the need to serve more number of users with the same number of *k* facilities. When  $\gamma = 1$ , the goal is to serve *all* the trajectories such that the maximum inconvenience faced by any trajectory is minimized.

The AVG-TIPS problem seeks to identify *k* service locations that minimizes the expected or average inconvenience across all the trajectories. Since the number of trajectories is fixed, minimizing the average inconvenience is equivalent to minimizing the total inconvenience over all the trajectories.

**Problem 2 (AVG-TIPS).** *Given a set of trajectories  $\mathcal{T}$ , a set of existing facilities  $\mathcal{F}$ , a set of candidate sites  $\mathcal{S}$  that can host the services, and a positive integer *k*, the AVG-TIPS problem seeks to report a set  $\mathcal{Q} \subseteq \mathcal{S}$ ,  $|\mathcal{Q}| = k$ , that minimizes the total inconvenience over  $\mathcal{T}$ , i.e., it minimizes  $TI(\mathcal{Q}) = \sum_{T_j \in \mathcal{T}} I_j$ , where  $I_j = \min_{s_i \in \mathcal{Q} \cup \mathcal{F}} \{d_r(T_j, s_i)\}$ .*

We next show that both the TIPS problems are NP-hard.

**Theorem 1 (NP-hardness of TIPS).** *MAX-TIPS and AVG-TIPS are NP-hard problems.*

*Proof.* Since the *k*-center problem is NP-hard [31] and it reduces to the MAX-TIPS problem with each trajectory being a single user location, the set of existing facilities  $\mathcal{F} = \emptyset$  and  $\gamma = 1$ , MAX-TIPS is also NP-hard.

Since the *k*-medoids problem is NP-hard [32] and it reduces to the AVG-TIPS problem with each trajectory being a single user location and  $\mathcal{F} = \emptyset$ , AVG-TIPS is also NP-hard.  $\square$

## 4 ALGORITHMS FOR MAX-TIPS

In this section, we present an optimal algorithm and two heuristics to solve the MAX-TIPS problem.

### 4.1 Optimal Algorithm

We present an *optimal* solution to the MAX-TIPS problem in the form of an *integer linear program* (ILP). For the ease of representation of the ILP, we assume that the set of candidate sites  $\mathcal{S} = \{s_i | 1 \leq i \leq n\}$  is augmented with the set of existing facilities  $\mathcal{F}$ .

$$\text{minimize } Z \quad \text{such that} \quad (1)$$

$$\forall 1 \leq i \leq n, Z \geq z_i \quad (2)$$

$$\forall 1 \leq i \leq n, \forall 1 \leq j \leq m, z_i \geq d_r(T_j, s_i) \times y_{ij} \quad (3)$$

$$\sum_{i=1}^n x_i \leq k + |\mathcal{F}|, \quad (4)$$

$$\forall 1 \leq j \leq m, \sum_{i=1}^n y_{ij} \leq 1 \quad (5)$$

$$\sum_{j=1}^m \sum_{i=1}^n y_{ij} \geq \gamma \times |\mathcal{T}| \quad (6)$$

$$\forall 1 \leq i \leq n, \forall 1 \leq j \leq m, y_{ij} \leq x_i \quad (7)$$

$$\forall 1 \leq i \leq n, x_i \in \{0, 1\} \quad (8)$$

$$\forall s_i \in \mathcal{F}, x_i = 1 \quad (9)$$

$$\forall 1 \leq i \leq n, \forall 1 \leq j \leq m, y_{ij} \in \{0, 1\} \quad (10)$$

The Boolean variable  $x_i = 1$  if and only if the site  $s_i$  is selected, or it is an existing facility. The Boolean variable  $y_{ij} = 1$  if and only if the site  $s_i$  is a serving facility (either existing or new), and the trajectory  $T_j$  is served by  $s_i$ . The variable  $z_i$  captures the maximum inconvenience offered to any trajectory served by the site  $s_i$  (Ineq. (3)). The constraint in Ineq. (2), along with the objective function in Eq. (1) ensure that the maximum inconvenience is minimized over the set of selected sites. The constraint in Ineq. (4) ensures that at most  $k$  new sites are selected in the answer set. Since  $Z$  monotonically decreases with  $\sum_{i=1}^n x_i$ ,  $Z$  will attain its optimal value only when  $\sum_{i=1}^n x_i = k + |\mathcal{F}|$ . The constraint in Ineq. (5) guarantees that each trajectory is served by at most one service location. The constraint in Ineq. (6) ensures that at least  $\gamma \cdot |\mathcal{T}|$  trajectories are served. The constraint in Ineq. (7) guarantees that if  $x_i = 0$ , then  $\forall j, y_{ij} = 0$ , i.e., no trajectory is served by the site  $s_i$  as it is not a serving facility location.

The optimal algorithm is impractical except for very small datasets (demonstrated in Sec. 6.2). Therefore, we next present a couple of polynomial-time heuristics to solve MAX-TIPS.

### 4.2 MIF Algorithm

MAX-TIPS problem being a generalization of the  $k$ -center problem (Sec. 2), the most natural approach to solve MAX-TIPS is to extend the greedy heuristic for the  $k$ -center problem [31]. We refer to this adaptation as *MIF* (*Most-Inconvenient-First*). It iteratively selects a site and the corresponding most inconvenient trajectory in each of the  $k$  iterations. The details are as follows.

Initially, we set  $\mathcal{Q} = \mathcal{F}$ , i.e., the existing set of facilities. The algorithm maintains a map, called the *Nearest Facility* map, denoted by  $\mathcal{NF}$ . This map keeps the trajectories in  $\mathcal{T}$  in a sorted order based on their distance to the nearest facility in  $\mathcal{Q}$ . Let  $\mathcal{RT}$  be an empty set of representative trajectories. The algorithm

runs in iterations. At the beginning of each iteration, a trajectory  $T_i \in \mathcal{NF}$  is chosen whose rank is  $\lfloor \gamma \times |\mathcal{T}| \rfloor$  in the sorted ordering. The reason behind this choice is that  $T_i$  faces the maximum inconvenience among the first  $\lfloor \gamma \times |\mathcal{T}| \rfloor$  trajectories in  $\mathcal{NF}$  in the sorted ordering.  $T_i$  is then added to the set  $\mathcal{RT}$ . In case of no existing facilities,  $\mathcal{NF}$  is initially empty. In such a scenario, any random trajectory is added to  $\mathcal{RT}$  in the first iteration. Next, we choose a candidate site  $s_i \in \mathcal{S} \setminus \mathcal{Q}$  that is nearest to  $T_i$ , and add it to the set  $\mathcal{Q}$ . If there are multiple such candidate sites, then the tie is broken arbitrarily. The above process is repeated until  $k$  new facility locations are selected.

Let us evaluate this algorithm on the example in Fig. 1 with  $k = 2$  and  $\gamma = 1$ . Initially,  $\mathcal{Q}$  is set to  $\{s_0\}$ . Since  $T_5$  is the farthest trajectory from  $s_0$ , it is added to  $\mathcal{RT}$ . Next,  $s_3$  is chosen and added to  $\mathcal{Q}$ , because it is the nearest site to  $T_5$ . Subsequently, in the next iteration,  $T_6$  is added to  $\mathcal{RT}$ . As a result,  $s_4$  is added to the answer set. Finally, the algorithm concludes with the selection  $\mathcal{Q} = \{s_0, s_3, s_4\}$ , which is also the optimal answer.

Now, consider the same example with  $k = 2$  and  $\gamma = 0.8$ . Once again,  $\mathcal{Q}$  is initialized with  $s_0$ . In iteration 1,  $T_4$  and  $s_3$  are chosen, and in the next iteration,  $T_3$  and  $s_2$  are chosen. The final selection is thus,  $\mathcal{Q} = \{s_0, s_2, s_3\}$ , which is again the optimal solution.

We observe that for any set of sites  $\mathcal{Q}$ , and two sets of trajectories,  $\mathcal{T}', \mathcal{T}$  such that  $\mathcal{T}' \subseteq \mathcal{T}$ , the maximum inconvenience faced by any trajectory in  $\mathcal{T}'$  due to the set  $\mathcal{Q}$  is at most the maximum inconvenience faced by any trajectory in  $\mathcal{T}$  due to the set  $\mathcal{Q}$ . Therefore, any approximation bound that holds for the MAX-TIPS problem with user-fraction  $\gamma = 1$  will also hold for  $\gamma < 1$ . Hence, we next discuss the approximation results only for  $\gamma = 1$ . Further, for ease of analysis, we assume that all the nodes in the road network are candidate sites, i.e.,  $\mathcal{S} = V$ .

**Theorem 2.** *Let  $d$  and  $d^*$  be the maximum inconvenience offered by the answer sets returned by the MIF algorithm, and the optimal algorithm for MAX-TIPS, respectively. Then,  $d \leq 2d^*$  for  $k = 1$ , and  $d \leq 2d^* + L$  for  $k \geq 2$  where  $L$  is the length of the longest trajectory.*

The proof is given in Appendix A.1.

**Theorem 3.** *The time and space complexities of MIF algorithm are  $O(k.l.n \log n + k.m.l^2 + k.m \log m)$  and  $O(l(n + m))$  respectively, where  $n = |V|$  is the total number of nodes in the road network,  $l$  is the maximum number of nodes in any trajectory,  $m$  is the total number of trajectories in  $\mathcal{T}$ , and  $k$  is the total number of iterations.*

The proof is stated in Appendix A.2.

Although MIF offers bounded quality guarantees, it is quite slow. This is because it does not use any pre-computed distances. The next scheme, however, leverages on pre-computed distances, and offers significantly faster response times.

### 4.3 Algorithm using NetClus

We observe that MIF takes significant computation time for calculating node-to-node distances, and node-to-trajectory distances. This is because we cannot afford to pre-compute and store all pairs node-to-trajectory distances, which is overwhelmingly large. However, an indexing scheme can be used that pre-computes and stores only a *small* set of node-to trajectory distances.

In this section, we propose a heuristic that uses the NetClus indexing framework [11] that was originally designed to solve the

following *TOPS* problem [11]:

Given a set of trajectories  $\mathcal{T}$ , a set of candidate sites  $\mathcal{S}$  that can host the services, the *TOPS* problem with query parameters  $(k, \tau)$  seeks to report the best  $k$  sites,  $\mathcal{Q} \subseteq \mathcal{S}$ ,  $|\mathcal{Q}| = k$ , that cover maximum number of trajectories. It is assumed that a site  $s_i$  covers a trajectory  $T_j$ , if and only if  $d_r(T_j, s_i) \leq \tau$ , where  $\tau$  is referred to as the coverage threshold.

NetClus performs multi-resolution clustering of the nodes in the road network,  $V$ . NetClus maintains  $t$  instances of index structures  $\mathcal{I}_0, \dots, \mathcal{I}_{t-1}$  of varying cluster radii. A particular index instance is useful for a particular range of query coverage thresholds. From one instance to the next, the radius is increased by a factor of  $(1 + \epsilon)$  for some  $\epsilon > 0$ . Assume that the normal range of query coverage threshold  $\tau$  is  $[\tau_{min}, \tau_{max}]$ . Then the total number of index instances is  $t = \lceil \log_{(1+\epsilon)}(\tau_{max}/\tau_{min}) \rceil + 1$ . For each index instance, NetClus maps the trajectories to the sequence of clusters that they pass through.

Intuitively, as the coverage threshold  $\tau$  increases, the number of trajectories covered by any set of candidate sites  $\mathcal{Q}$  also increases. This was validated empirically in [11].

Exploiting the general monotonic behavior of the trajectory coverage with respect to the coverage threshold  $\tau$ , we propose the following heuristic to answer the MAX-TIPS problem. Our goal is to identify the smallest value of  $\tau$  such that there exists a set  $\mathcal{Q} \subseteq \mathcal{S}$  of size  $k$  that covers at least  $\gamma \times |\mathcal{T}|$  number of trajectories in  $\mathcal{T}$ . To guess this desired value of  $\tau$ , we perform a binary search over the range of  $\tau$ , i.e.,  $[\tau_{min}, \tau_{max}]$ .

The algorithm proceeds in iterations. In each iteration, it computes the value of the coverage threshold,  $\tau = \frac{\tau_{min}^c + \tau_{max}^c}{2}$  where  $\tau_{min}^c$  and  $\tau_{max}^c$  denote the current ranges. Initially,  $\tau_{min}^c = \tau_{min}$  and  $\tau_{max}^c = \tau_{max}$ . Next, the *TOPS* query with parameters  $(k, \tau)$  is computed. While doing so, the existing facilities  $\mathcal{F}$  must be taken into consideration [11]. The trajectories that lie within the coverage threshold  $\tau$  of any existing facility, are deemed to be covered. If the trajectory coverage value, i.e., the number of trajectories covered by the set  $\mathcal{Q} \cup \mathcal{F}$ , is lower than  $\gamma \times |\mathcal{T}|$ , then  $\tau_{min}^c$  is set to  $\tau$ , else  $\tau_{max}^c$  is set to  $\tau$ . Consequently, in the next iteration, *TOPS* query is computed with the revised value of  $\tau$ . Since this process can continue forever, it is stopped when the difference between  $\tau_{max}^c$  and  $\tau_{min}^c$  falls below a desired precision. Suppose the final iteration executed the *TOPS* query with parameters  $(k, \tau')$  and returned the set  $\mathcal{Q}$ . Then the answer to the MAX-TIPS problem is also  $\mathcal{Q}$  with maximum inconvenience as  $\tau'$ . We call this algorithm simply NetClus.

As the monotonicity of the trajectory coverage w.r.t. the coverage threshold  $\tau$  is not guaranteed theoretically, the quality of this heuristic cannot be bounded. Empirically, however, it performs the best in terms of both running time and quality (Sec. 6).

**Theorem 4.** *The time and space complexities of NetClus algorithm are  $O(\log_2(\tau_{max}/\tau_{min}) \cdot t_{TOPS})$  and  $O(t \cdot (n + m \cdot l))$  respectively, where  $O(t_{TOPS})$  is the time required to answer a *TOPS* query by NetClus and  $\tau_{min}$  and  $\tau_{max}$  are the ranges of  $\tau$  values indexed by NetClus. Further,  $t = 1 + \lceil \log_{1+\epsilon}(\tau_{max}/\tau_{min}) \rceil$  is the number of index instances,  $\epsilon > 0$  is the index resolution parameter,  $m = |\mathcal{T}|$ ,  $n = |V|$  and  $l$  is the maximum number of nodes in any trajectory.*

The proof is given in Appendix A.3.

While the NetClus approach is index-based, MIF is non-index based. In case of NetClus, the distance of the trajectories and sites to their respective cluster centres is pre-computed, thereby making

it efficient. For MIF, all the necessary distance computations are performed online and, hence, it is slower.

## 5 ALGORITHMS FOR AVG-TIPS

This section presents an optimal algorithm and two heuristics for the AVG-TIPS problem.

### 5.1 Optimal Algorithm

The following *integer linear program* solves the AVG-TIPS problem. As in Sec. 4.1, we assume that the set of candidate sites  $\mathcal{S} = \{s_i | 1 \leq i \leq n\}$  is augmented with the set of existing facilities  $\mathcal{F}$ .

$$\text{minimize } Z = \sum_{j=1}^m \sum_{i=1}^n [d_r(T_j, s_i) \times y_{ij}] \quad \text{such that} \quad (11)$$

$$\sum_{i=1}^n x_i \leq k + |\mathcal{F}| \quad (12)$$

$$\forall 1 \leq j \leq m, \sum_{i=1}^n y_{ij} = 1 \quad (13)$$

$$\forall 1 \leq i \leq n, \forall 1 \leq j \leq m, y_{ij} \leq x_i \quad (14)$$

$$\forall 1 \leq i \leq n, x_i \in \{0, 1\} \quad (15)$$

$$\forall s_i \in \mathcal{F}, x_i = 1 \quad (16)$$

$$\forall 1 \leq i \leq n, \forall 1 \leq j \leq m, y_{ij} \in \{0, 1\} \quad (17)$$

The objective function in Eq. (11) ensures that the total (equivalently, mean) inconvenience is minimized over the set of selected sites. The constraint in Ineq. (13) guarantees that each trajectory is served by exactly one service location. The semantics of rest of the constraints are same as those stated in Sec. 4.1.

This optimal algorithm is impractical except for extremely small datasets (shown in Sec. 6.3). Therefore, we next present the following heuristics for the AVG-TIPS problem.

### 5.2 HCC Algorithm

Recall that the AVG-TIPS problem generalizes the  $k$ -medoids problem (Sec. 2). Our first heuristic, HCC, therefore, builds on the three popular approaches for the  $k$ -medoids problem, PAM [32], CLARA [32], and CLARANS [38]. We first describe these approaches, and then discuss the proposed HCC algorithm.

**PAM:** The basic idea of PAM is as follows. Given a set of  $n$  objects, it starts by choosing  $k$  random objects, called as *medoids*. Each non-medoid object is assigned to its nearest medoid. The cost of a particular clustering is the sum of the distances of each non-medoid object to its nearest medoid. The PAM algorithm proceeds in iterations. In each iteration, it swaps one of the existing medoids with a non-medoid object such that the cost of the resulting clustering decreases. To realize the swap with minimal cost, it computes the cost of each possible swap which are as many as  $k(n - k)$ . This step is computationally very expensive. PAM terminates when it reaches a local minima, i.e., there is no possible swap resulting in a solution with a lower cost.

**CLARA:** For large datasets, the PAM algorithm is infeasible owing to its high running time. To address this limitation, the CLARA algorithm [32] that relies on sampling was proposed. The idea is to create random samples of size much smaller than  $n$ , and execute the PAM algorithm on each of these samples, and return the clustering that offers the minimal cost over all the samples.

**CLARANS:** To improve the quality of clustering of CLARA, another algorithm called CLARANS was proposed [38]. The basic

idea of CLARANS is to avoid scanning all possible swaps in each iteration of PAM. Essentially, a small fraction of the total  $k(n-k)$  possible swaps are scanned and the swap that offers the minimal clustering cost is executed. To increase the robustness, this scheme is repeated a few number of times, and finally the clustering with the minimal cost is reported.

Inspired by the above approaches, we propose a new algorithm, *HCC (Hybrid-CLARA-CLARANS)*, for the AVG-TIPS problem that combines the ideas of sampling (from CLARA) and scanning a small number of swaps (from CLARANS). The basic idea is to consider a few samples of sufficiently small size, and for each sample, to scan a sufficiently small number of swaps. The details are described next.

### 5.2.1 Details

First, we describe the basic algorithm, and later, discuss how to make it scalable. Given a set  $\mathcal{T}$  of  $m$  trajectories, and a set  $\mathcal{S}$  of  $n$  candidate sites, initially, we compute the distances between each pair of trajectory and site. Then we execute the following steps.

Initialize  $\mathcal{Q}$  to  $\mathcal{F}$ , i.e., the set of existing facilities. Choose a random set of  $k$  sites in  $\mathcal{S}$ , referred to as *medoids*, and add it to  $\mathcal{Q}$ . The total inconvenience of the set  $\mathcal{Q}$  is  $TI(\mathcal{Q}) = \sum_{T_j \in \mathcal{T}} I_j$ , where  $I_j = \min_{s_i \in \mathcal{Q}} d_r(T_j, s_i)$ . To efficiently compute the value of  $TI(\mathcal{Q})$ , we use the  $\mathcal{NF}$  map (discussed in Sec. 4.2) which tracks the trajectories based on the distance to their nearest facility in  $\mathcal{Q}$ .

Subsequently, the algorithm proceeds in iterations. In each iteration, it scans a fraction of the total possible swaps between a medoid and a non-medoid site in  $\mathcal{S} \setminus \mathcal{Q}$ , as discussed above, and executes the swap that results in the lowest total inconvenience  $TI(\mathcal{Q})$ . Out of the total number of possible swaps,  $k(n-k)$ , HCC scans only a fraction, referred to as the *swap-fraction*, and denoted by  $sf$ . This process terminates when there is no possible swap that leads to a solution with a lower total inconvenience, or after a pre-specified number of maximum iterations  $\eta$ , whatever is encountered earlier. While performing the swaps, we ensure that we do not swap any existing facility in  $\mathcal{Q}$ . The algorithm makes  $t$  ( $t \geq 1$ ) trials of the above steps to account for the randomly chosen initial set of medoids. Finally, it returns the set  $\mathcal{Q}$  with the lowest total inconvenience achieved over the  $t$  trials.

Let us see the working of this algorithm for the AVG-TIPS problem on the example shown in Fig. 1 with  $k = 2$ . We will consider a single trial with swap-fraction  $sf = 1$ . Initially,  $\mathcal{Q}$  is set to  $\{s_0\}$ . Next, let us start with the following set of initial medoids  $\{s_1, s_2\}$  that is added to  $\mathcal{Q}$ . We note that  $TI(\mathcal{Q}) = 42$ . Since  $sf = 1$ , we examine all possible swaps resulting in the following sets:  $\{s_0, s_1, s_3\}$ ,  $\{s_0, s_1, s_4\}$ ,  $\{s_0, s_2, s_3\}$  and  $\{s_0, s_2, s_4\}$ . Since  $TI(\{s_0, s_2, s_3\}) = 21$  is minimal, it executes this swap. Since there is no further possible swap that results in lower total inconvenience, the algorithm terminates, reporting the set  $\{s_2, s_3\}$  as the answer. Incidentally, this happens to be the optimal solution.

This heuristic requires the distance values between each pair of trajectory  $T_j \in \mathcal{T}$  and site  $s_i \in \mathcal{S}$ . Computing and storing these distances for large datasets may be infeasible when  $|\mathcal{S}|$  or  $|\mathcal{T}|$  is large. Thus, under such circumstances, we propose to sample the set of candidate sites and trajectories, using the following schemes.

### 5.2.2 Site Sampling

The sampling technique is based on a simple clustering idea that clusters the set of nodes  $V$  in the road network and samples at most a single candidate site from each cluster. The details are as

follows. A random node  $v_p \in V$  (that is not yet clustered) is chosen as the cluster-center of a cluster  $Clus(v_p)$  that consists of all nodes  $v_q \in V$  that are not yet clustered and are within some distance threshold  $R$  from  $v_p$ . This process is repeated until every node in  $V$  is clustered. Finally, a sample  $\mathcal{S}' \subseteq \mathcal{S}$  is created by selecting at most a single site from each cluster that is closest to the cluster-center, which we refer to as its cluster-representative. If a cluster has no candidate site, then it has no cluster-representative either. As the chosen sites belong to different clusters, they are not expected to be close to each other. Thus, the sampled sites are typically nicely distributed over the road network.

### 5.2.3 Trajectory Sampling

Consider any trajectory  $T_j \in \mathcal{T}$ . For each node  $v_p \in T_j$ , let  $v'_p$  be the cluster-center of the cluster that contains  $v_p$ . Each trajectory  $T_j$  is mapped to a set of cluster-centers  $T'_j = \{v'_p\}$ . This transforms the trajectories into a coarser representation as nodes that are close to each other in the trajectory are likely to fall in the same cluster. This transformed set of trajectories  $\{T'_j\}$  is denoted by  $\mathcal{T}'$ .

Following this, the set of trajectories  $\mathcal{T}'$  is clustered based on their Jaccard similarity measure, as proposed in [40]. The high level overview of this method is as follows. Suppose, we are required to create a trajectory sample of size  $s$ . Initially, each trajectory  $T'_j \in \mathcal{T}'$  is a cluster by itself with  $T'_j$  being the cluster-representative. For any two trajectories,  $T'_p$  and  $T'_q$ , their Jaccard similarity is given by  $J(T'_p, T'_q) = |T'_p \cap T'_q| / |T'_p \cup T'_q|$ .

The clustering follows an iterative algorithm where in each iteration it fuses a pair of clusters with cluster-representatives  $T'_p$  and  $T'_q$  that have the maximum Jaccard similarity. After fusing, either of the two trajectories,  $T'_p$  or  $T'_q$  is deemed as the cluster-representative of the new cluster. The algorithm continues fusing a pair of clusters in each iteration in this manner, until there are exactly  $s$  clusters. The cluster-representatives of these  $s$  clusters are mapped back to their original representation as sequence of nodes, and reported as the desired trajectory sample.

**Theorem 5.** *The time and space complexities of HCC are  $O(t \cdot \eta \cdot sf \cdot k^2(n-k) \cdot m')$  and  $O(m' \cdot (n' + l))$  respectively, where  $n'$  and  $m'$  are the number of sites and trajectories produced after sampling of the sites and the trajectories, respectively. Here,  $k$  is the number of medoids,  $sf$  is the swap-fraction,  $t$  is the number of trials,  $\eta$  is the maximum number of iterations per trial, and  $l$  is the maximum number of nodes in any trajectory.*

The proof is available in Appendix A.4.

Unfortunately, HCC does not have any approximation bound, since the original algorithms, PAM, CLARA and CLARANS do not have any quality guarantees either.

## 5.3 GREAT Algorithm

We next propose a greedy heuristic called GREAT (GREedy Avg-Tips) that offers bounded quality guarantees.

It is an iterative algorithm that works on the principle of maximizing the marginal gain in each successive iteration. It starts with the set  $\mathcal{Q} \leftarrow \mathcal{F}$ . In each iteration  $\theta = \{1, \dots, k\}$ , it selects a site  $s_\theta \in \mathcal{S} \setminus \mathcal{Q}$  such that the total inconvenience of the resulting set,  $TI(\mathcal{Q} \cup \{s_\theta\})$ , is minimized. The site  $s_\theta$  is added to the set  $\mathcal{Q}$ . The algorithm terminates after  $k$  iterations.

Similar to HCC, the GREAT algorithm also assumes that the distances between each pair of trajectory  $T_j \in \mathcal{T}$  and site  $s_i \in \mathcal{S}$  are pre-computed and available. We also use the  $\mathcal{NF}$  map (as in

the case of HCC) to track the distance between each trajectory and its nearest facility. Whenever a new site is chosen to be added into  $\mathcal{Q}$ , we check whether it would be the nearest facility for each trajectory  $T_j$ , and update it accordingly.

Let us see the working of this algorithm for the AVG-TIPS problem on the example shown in Fig. 1 with  $k = 2$ . Before the algorithm begins, we set  $\mathcal{Q} \leftarrow \{s_0\}$ . In the first iteration,  $s_3$  is selected as the set  $\{s_0, s_3\}$  offers the least total inconvenience of 43 units. In the next iteration,  $s_2$  is selected, resulting in optimal total inconvenience of 21 units.

### 5.3.1 Analysis of GREAT

We, next, state an important property of AVG-TIPS problem, which in turn, helps us to bound the quality of GREAT.

A function  $f$  defined on any subset of a set  $\mathcal{S}$  is *sub-modular* if for any pair of subsets  $\mathcal{Q}, \mathcal{R} \subseteq \mathcal{S}$ ,  $f(\mathcal{Q}) + f(\mathcal{R}) \geq f(\mathcal{Q} \cup \mathcal{R}) + f(\mathcal{Q} \cap \mathcal{R})$  [41]. A function  $f$  is *super-modular* if its negative ( $-f$ ) is sub-modular. The following result shows that the function  $TI(\mathcal{Q})$  is super-modular.

**Theorem 6.** *For any set of candidate sites  $\mathcal{Q} \subseteq \mathcal{S}$ , the total inconvenience  $TI(\mathcal{Q})$  is a non-increasing super-modular function.*

*Proof.* Consider any pair of sets  $\mathcal{Q}, \mathcal{R} \subseteq \mathcal{S}$  such that  $\mathcal{Q} \subseteq \mathcal{R}$ .

First, we show that  $TI(\mathcal{Q})$  is a non-increasing function. Since  $I_j(\mathcal{Q}) = \min_{s_i \in \mathcal{Q}} \{d_r(T_j, s_i)\}$  is a minimum function over the set  $\mathcal{Q}$ , it follows that  $I_j(\mathcal{R}) \leq I_j(\mathcal{Q})$ . Thus,  $TI(\mathcal{R}) = \sum_{j=1}^m I_j(\mathcal{R}) \leq \sum_{j=1}^m I_j(\mathcal{Q}) = TI(\mathcal{Q})$ . Hence,  $TI(\mathcal{Q})$  is a non-increasing function.

To show that  $TI(\mathcal{Q})$  is super-modular, it is sufficient to show that for any site  $s \in \mathcal{S} \setminus \mathcal{R}$ , the following holds [41]:

$$TI(\mathcal{Q} \cup \{s\}) - TI(\mathcal{Q}) \leq TI(\mathcal{R} \cup \{s\}) - TI(\mathcal{R}) \quad (18)$$

Since  $TI(\mathcal{Q}) = \sum_{j=1}^m I_j$ , it is, therefore, enough to prove that for any trajectory  $T_j \in \mathcal{T}$ ,

$$I_j(\mathcal{Q} \cup \{s\}) - I_j(\mathcal{Q}) \leq I_j(\mathcal{R} \cup \{s\}) - I_j(\mathcal{R}) \quad (19)$$

Suppose the site  $s^* \in \mathcal{R} \cup \{s\}$  is the nearest facility to trajectory  $T_j$  in the set  $\mathcal{R} \cup \{s\}$ . There can be two cases:

(a)  $s^* = s$ :  $I_j(\mathcal{Q} \cup \{s\}) = I_j(\{s\}) = I_j(\mathcal{R} \cup \{s\})$ . Further, since  $I_j(\mathcal{Q}) \geq I_j(\mathcal{R})$  (using the non-increasing property), Ineq. (19) follows.

(b)  $s^* \neq s, s^* \in \mathcal{R}$ : Here,  $I_j(\mathcal{R} \cup \{s\}) - I_j(\mathcal{R}) = 0$ . Using the non-increasing property of  $I_j(\mathcal{Q})$ ,  $I_j(\mathcal{Q} \cup \{s\}) - I_j(\mathcal{Q}) \leq 0$ . Thus, Ineq. (19) follows.  $\square$

The next result bounds the quality of GREAT.

**Theorem 7.** *Let  $OPT \subseteq \mathcal{S}, |OPT| = k$  denotes an optimal solution to the AVG-TIPS problem. Let  $\mathcal{Q} \subseteq \mathcal{S}, |\mathcal{Q}| = k$  be the solution reported by the GREAT algorithm. Then,*

$$TI(\mathcal{Q}) = TI(OPT) \quad \text{for } k = 1$$

$$TI(\mathcal{Q}) \leq (1 - 1/e)TI(OPT) + TI(\mathcal{F})/e \quad \text{for } k \geq 2$$

where  $TI(\mathcal{F})$  refers to the initial total inconvenience offered by the existing facilities  $\mathcal{F}$ . We assume  $TI(\mathcal{F})$  to be a non-negative constant.

The proof is stated in Appendix A.5.

The next result bounds the complexity of GREAT.

**Theorem 8.** *The space and time complexities of GREAT are  $O(m \cdot (n + l))$  and  $O(k \cdot m \cdot n)$  respectively, where  $m = |\mathcal{T}|$ ,  $n = |\mathcal{S}|$  and  $l$  is the maximum length of any trajectory.*

Dataset	Type	# Trajectories	# Sites
Beijing-Small (BS)	Real	8,083	30
Beijing-Medium (BM)	Real	8,083	47,125
Beijing-Large (BL)	Real	123,179	269,686
Beijing-Medium-Sampled (BMS)	Real	1,278	1,883
Beijing-Large-Sampled (BLS)	Real	9,701	15,775
Bengaluru	Synthetic	9,950	61,563
New York	Synthetic	9,950	355,930
Atlanta	Synthetic	9,950	389,680

TABLE 2: Summary of datasets.

The proof is available in Appendix A.6.

The key drawback of the above scheme is its space requirement of  $O(m \cdot n)$  which is prohibitively large for city-scale datasets. To alleviate this problem, we work with sampled set of trajectories and candidate sites, as described in Sec. 5.2.2 and Sec. 5.2.3.

## 6 EXPERIMENTAL EVALUATION

In this section, we perform extensive experiments to assess the quality, scalability and practicality of the different heuristics. Since both the TIPS problems, MAX-TIPS and AVG-TIPS, are introduced in this work, there is *no* existing *practical* baseline technique to compare against. Nevertheless, we do compare with the nearest available baseline techniques for both *static* and *mobile* users' scenarios. Further, we compare with the optimal algorithms for both MAX-TIPS and AVG-TIPS, on small datasets. The experiments were conducted using Java (version 1.7.0) code on an Intel(R) Core i7-4770 CPU @3.40GHz machine with 32GB RAM running Ubuntu 14.04.2 LTS OS.

### 6.1 Datasets

We conducted experiments on both real and synthetic datasets, whose details are shown in Table 2. For simplicity, we assume that the set of candidate sites is same as the set of nodes in the road network, unless otherwise stated.

**Real datasets:** We use GPS traces from Beijing consisting of user trajectories generated by tracking taxis for a week [42], [43]. To generate trajectories as sequences of road intersections, the raw GPS-traces were map-matched [39] to the Beijing road network extracted from OpenStreetMap (<http://www.openstreetmap.org/>). This dataset, referred to as Beijing-Large (BL) here, is the largest and most widely used publicly available trajectory datasets.

For AVG-TIPS, all the algorithms require distances between each pair of site and trajectory. For a large dataset such as BL, computing and storing such pairwise distances is infeasible. Therefore, for thorough evaluation, we use a medium sized dataset, Beijing-Medium (BM).

To assess the quality of sampling techniques, representative datasets, referred to as Beijing-Medium-Sampled (BMS) and Beijing-Large-Sampled (BLS) respectively, are derived from the BM and BL datasets. These datasets are generated using the sampling schemes described in Sec. 5.2.2 and Sec. 5.2.3.

Since both the TIPS problems are NP-hard, the optimal algorithm requires exponential time and, therefore, can be run only on a small dataset. Hence, we evaluate all the algorithms against the optimal on Beijing-Small (BS) dataset which consists of the same set of trajectories as in BM, but has only 30 candidate sites, which are sampled randomly. To increase the robustness of the results,

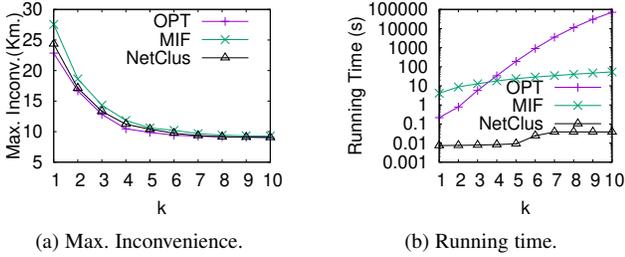


Fig. 2: MAX-TIPS: Comparison with optimal at  $\gamma = 100\%$ .

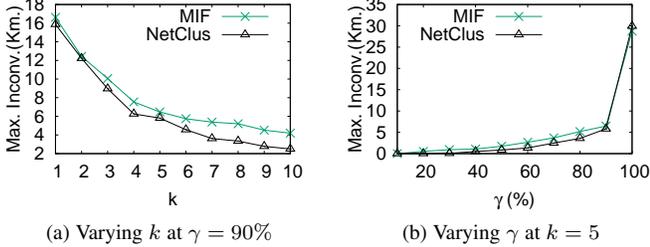


Fig. 3: MAX-TIPS: Quality results.

such sampling was performed 5 times, and the experiments were conducted 10 times for each sample.

**Synthetic datasets:** To study the impact of city geographies, we generated three synthetic datasets that emulate trajectory patterns followed in New York, Atlanta and Bengaluru. We use an online traffic generator tool, MNTG (<http://mntg.cs.umn.edu/tg/index.php>) to generate the traffic traces, and map-match them to generate trajectories in the desired format.

## 6.2 MAX-TIPS Results

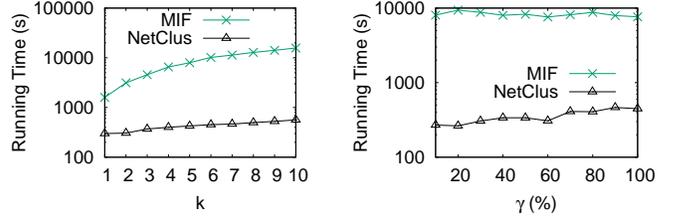
We evaluate the performance of three different algorithms to solve MAX-TIPS, Opt, MIF, and NetClus, on the two basic parameters: (i) desired number of service locations  $k$ , varied in the range  $[1, 10]$ , and (ii) user-fraction  $\gamma$ , varied in the range 10% to 100%. The default values of  $k$  and  $\gamma$  are 5 and 90% respectively. The metrics evaluated are (i) maximum inconvenience,  $MI$ , and (ii) running time. The dataset used is BL, unless otherwise stated.

**Comparison with Optimal:** The optimal algorithm used, is the integer linear program (ILP), given in Sec. 4. Since the optimal algorithm requires exponential running times, we run it only on the BS dataset. Fig. 2 shows that  $MI$  values offered by MIF and NetClus are very close to that of Opt although the running times are much better. Opt requires several hours to complete even for this small dataset and, therefore, is not practical at all. Consequently, we do not experiment with Opt any further.

**Quality Results:** Fig. 3a shows the  $MI$  values of MIF and NetClus on the BL dataset. NetClus offers the least  $MI$  value, beating MIF by more than 20% on an average.

Fig. 3b shows that as user-fraction increases from 90% to 100%, there is a sharp rise in  $MI$  value for both MIF and NetClus. This is because there are generally trajectories that are very hard to satisfy and, since at  $\gamma = 100\%$ , all of them need to be served,  $MI$  values shoot up for both the algorithms.

**Performance Results:** The running time results portrayed in Fig. 4 show that NetClus is 1-2 orders of magnitude faster than MIF. The high running time of MIF is due to a large number of calls to the shortest path algorithm (once for each site on the chosen trajectory in each of the  $k$  iterations). Moreover, to guard against a poorly selected random initial trajectory, MIF is repeated thrice. On the other hand, NetClus is fast since it uses



(a) Varying  $k$  at  $\gamma = 90\%$  (b) Varying  $\gamma$  at  $k = 5$

Fig. 4: MAX-TIPS: Running time performance.

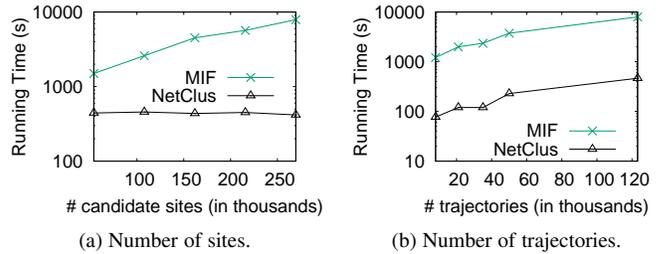


Fig. 5: MAX-TIPS: Scalability results ( $k = 5$  and  $\gamma = 90\%$ ).

pre-computed distances between trajectories and centers of the clusters that they pass through.

**Memory Footprint:** The memory consumption of MIF and NetClus on the BL dataset (at the default values of  $k$  and  $\gamma$ ) are about 8.6GB and 3.2GB respectively. The low memory footprint of NetClus is due to clustering of the site space and consequent compressed representation of trajectories.

**Scalability:** To ascertain the scalability of NetClus and MIF with respect to the number of candidate sites and trajectories, we next took different sized samples from the BL dataset. Fig. 5 shows that NetClus scales better with the number of sites due to its clustered representation. It is faster than MIF by at least an order of magnitude for all the situations.

**Synthetic Datasets:** Fig. 6 shows the evaluation of MAX-TIPS on three synthetic datasets emulating traffic in Atlanta, New York, and Bengaluru. Bengaluru has the smallest sized road network. Therefore, it has the best  $MI$  value and running time. Atlanta and New York have much larger road network, and the trajectories are distributed all over the network. Thus, they are harder to be served, and consequently exhibit high  $MI$  values. Their running times are high owing to large number of candidate sites to be processed.

## 6.3 AVG-TIPS Results

We evaluated the performance of three different algorithms for AVG-TIPS, Opt, HCC, and GREAT on the desired number of service locations  $k$ , varied in the range  $[1, 10]$ , with default as 5. The metrics evaluated are (i) average inconvenience,  $AI(Q) = TI(Q)/|T|$ , and (ii) running time.

**Choice of swap-fraction in HCC algorithm:** Referring to Sec. 5.2, recall that HCC scans only a fraction  $sf$  of the total

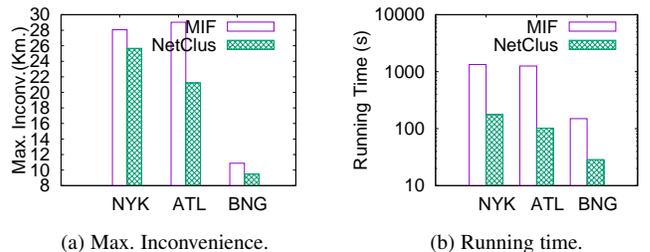
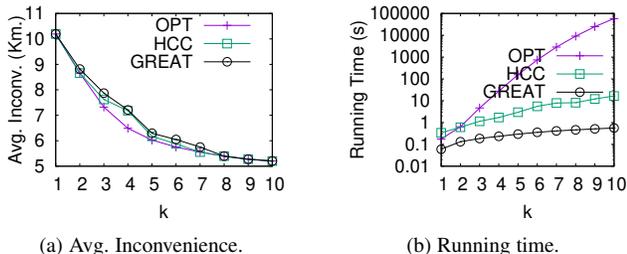


Fig. 6: MAX-TIPS: Synthetic datasets (at  $k = 5$  and  $\gamma = 90\%$ ).



(a) Avg. Inconvenience.

(b) Running time.

Fig. 7: AVG-TIPS: Comparison with optimal.

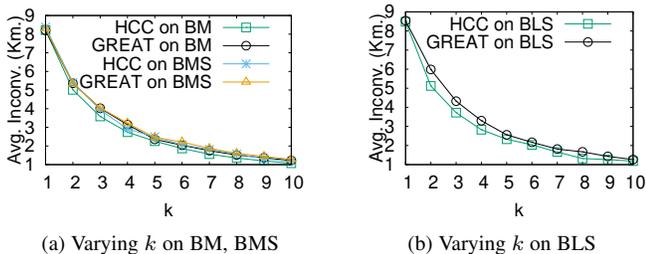
(a) Varying  $k$  on BM, BMS(b) Varying  $k$  on BLS

Fig. 8: AVG-TIPS: Quality results.

number of swaps, referred to as the swap-fraction. Table 3 shows the performance of HCC for different values of  $sf$  (shown as percentage of total number of swaps). The second column lists the relative error in  $AI$  value w.r.t.  $sf = 100\%$ , and the third column indicates the corresponding speed up in running time. For our experimentation, we choose  $sf = 5\%$  as this offers a nice balance with an error of less than 3% and a speed up of more than 4 times.

**Comparison with Optimal:** Since AVG-TIPS is NP-Hard, the ILP-based optimal algorithm given in Sec. 5, can be run only on small datasets. Therefore, as in the case of MAX-TIPS, we evaluate it on the Beijing-Small dataset. Fig. 7 shows that both HCC and GREAT offer almost the same quality as Opt. However, Opt takes hours of query time even on such a small dataset. Consequently, we drop Opt from further consideration.

**Quality Results:** As explained earlier, HCC and GREAT cannot be executed on the BL dataset due to high memory overhead. Hence, to ascertain the effect of site sampling and trajectory sampling on the quality, we use Beijing-Medium (BM), and its sampled counterpart, the Beijing-Medium-Sampled (BMS) datasets. Fig. 8a shows that  $AI$  values achieved by both HCC and GREAT on the sampled dataset BMS is almost as good as the full BM dataset. The average relative error in  $AI$  values due to sampling is only about 10% for HCC and 4% for GREAT. Fig. 8b reports the quality on the BLS dataset. HCC offers roughly 10% better quality than GREAT on an average. Since the error due to sampling is fairly low, it is expected that the  $AI$  values on the BL dataset are likely to be similar to those in the BLS dataset.

**Performance Results:** Fig. 9a shows that sampling offers a speed up of about 2 orders of magnitude on the running times for both HCC and GREAT on BM. Fig. 9b shows that GREAT is about

$sf$ (in %)	with respect to $sf = 100\%$	
	Error	Speed-up factor
1	10.92	7.22
2	4.49	4.97
5	2.86	4.27
10	2.27	3.40
25	1.16	2.37
50	0.00	1.53

TABLE 3: Performance of HCC with varying swap-fraction.

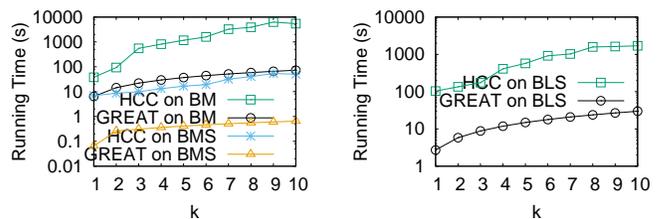
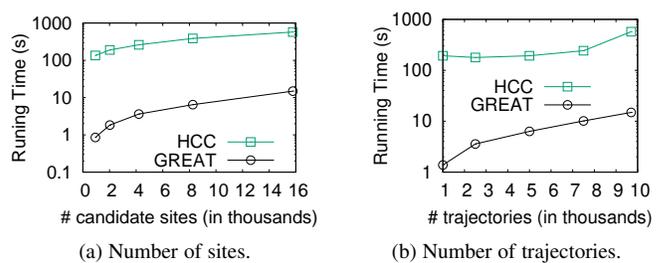
(a) Varying  $k$  on BM, BMS(b) Varying  $k$  on BLS

Fig. 9: AVG-TIPS: Running time performance.



(a) Number of sites.

(b) Number of trajectories.

Fig. 10: AVG-TIPS: Scalability results (at  $k = 5$ ).

2 orders of magnitude faster than HCC when evaluated on BLS. HCC is slower due to its repeated swaps. In addition, it is run thrice to avoid a poor random initial choice.

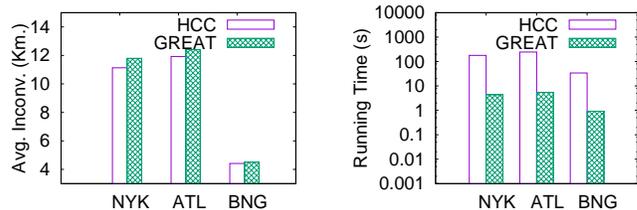
**Memory Footprint:** The memory consumption of HCC and GREAT on BM at  $k = 5$  are roughly 18GB and 11.5GB respectively; and, that on BMS are about 3.8 GB and 0.5GB respectively. The sampling techniques are, thus, highly effective in lowering the memory footprints of both the algorithms.

**Scalability:** We next examine the scalability of HCC and GREAT with respect to the number of candidate sites and trajectories. We use the same setting as in Sec. 6.2 for MAX-TIPS, although the BM dataset is used instead of BL. Fig. 10 shows that both are scalable with GREAT being faster by about 2 orders of magnitude.

**Synthetic Datasets:** The next result shows the effect on synthetic datasets. Using the same setup as discussed in Sec. 6.2, Fig. 11 shows that GREAT is about 1-2 orders of magnitude faster than HCC, while sacrificing no more than 10% in accuracy over all the three synthetic datasets. The  $AI$  values for New York and Atlanta are significantly higher than those of Bengaluru due to their large road network, and even distribution of trajectories. The running times follow the same trend as in the case of MAX-TIPS.

## 6.4 Comparison with Baseline

In this section, we compare the performance of our algorithms with the baseline techniques for static and mobile users. For MAX-TIPS, the baseline techniques used for static and mobile users are AppMinMax [10] and EN (Extended Network method) [12] respectively; and for AVG-TIPS, the baseline techniques for static and mobile users are CLARANS [38] and EN [12] respectively.



(a) Avg. Inconvenience.

(b) Running time.

Fig. 11: AVG-TIPS: Synthetic datasets (at  $k = 5$ ).

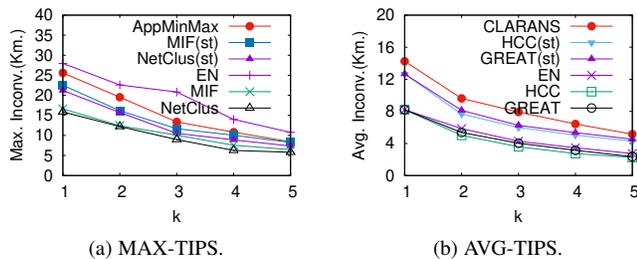


Fig. 12: Comparison with Baseline.

Further, we also show the impact of considering trajectories as opposed to one or more static user locations such as homes and offices. Fig. 12 shows the quality comparisons for MAX-TIPS and AVG-TIPS. For MAX-TIPS, we use the BL dataset, while for AVG-TIPS, the BM dataset was used.

To model static users, we use either of the two end-points of a trajectory as the representative static location of its user. We run AppMinMax over these static user locations. This algorithm reports  $k$  facility locations that minimize the maximum distance of any user to its nearest facility. To study the effect of factoring in two static locations of each user, we consider only the two end-points of each trajectory and ignore the intermediate points. The corresponding MIF and NetClus versions are referred to as MIF(st) and NetClus(st) respectively.

For mobile users, there is no existing work that aims to minimize the maximum user-inconvenience. The nearest baseline technique is the EN method that reports a single facility location that minimizes the average distance of any user to its nearest facility. To generate  $k$  facility locations, we repeat the algorithm  $k$  times, assuming  $k - 1$  existing facilities. This algorithm along with MIF and NetClus are evaluated over *full* trajectories, i.e., without skipping any intermediate point in any of the trajectories.

Fig. 12a shows that considering full trajectories is always much better than considering a single or two static locations per user. While NetClus performs the best, EN performs the worst. This is because EN does not address the objective of MAX-TIPS, but minimizes the average inconvenience.

Next, we evaluate AVG-TIPS over the BM dataset. Here, the CLARANS algorithm (discussed in Sec. 5.2) serves as the baseline technique for static users' scenario. To assess the effect of considering only two static locations per user, we run versions of HCC and GREAT as HCC(st) and GREAT(st) respectively. The EN method works as the baseline techniques for mobile users.

As in the case of MAX-TIPS, Fig. 12b once again validates our claim that using trajectories is always more beneficial than considering one or two static locations per user. While CLARANS running over single static location per user performs the worst, HCC, running on *full* trajectories performs the best. The performance of EN is also good, but marginally lesser than that of HCC.

## 6.5 Existing Facilities

Fig. 13 shows the effect of existing facilities on MAX-TIPS and AVG-TIPS. We choose 2 existing facilities randomly. The corresponding maximum and average inconveniences are 18.51 Km and 8.53 Km respectively. Using the proposed algorithms, we locate  $k'$  new facilities for  $k' = 1, \dots, 5$ . The  $k$  values ( $2 + k'$ ) shown in the figure denotes the total number of facilities including the existing ones. For MAX-TIPS, MIF(ex) and NetClus(ex) represent the  $MI$  values as recorded by MIF and NetClus respectively, while factoring in the *existing facilities*. Similarly, for

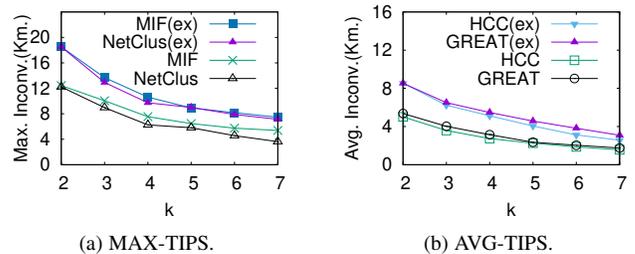


Fig. 13: Effect of Existing Facilities.

AVG-TIPS, HCC(ex) and GREAT(ex) represent the  $AI$  values as recorded by HCC and GREAT respectively, while considering the *existing facilities*. The  $MI$  (respectively,  $AI$ ) values shown by MIF and NetClus (respectively, HCC and GREAT) correspond to the scenario when these algorithms are used to locate new facilities in the *absence* of any existing facilities. Since the existing facilities are chosen randomly, the initial  $MI$  and  $AI$  values (as mentioned above) are quite high. As a result, the  $MI$  values of NetClus(ex) and MIF(ex) continue to be much higher than those of NetClus and MIF respectively, for all values of  $k$ . Following the same argument, the  $AI$  values of HCC(ex) and GREAT(ex) are much higher than those of HCC and GREAT respectively.

## 6.6 Summary of Experiments

We summarize our experimental findings as follows. NetClus offers the best performance for MAX-TIPS on multiple real and synthetic datasets, both in terms of efficiency and quality. It outperforms MIF by more than 20% in quality and is faster by an order of magnitude. For AVG-TIPS, firstly, we observe that it is better to apply the sampling techniques when the dataset is very large, such as Beijing-Large. The error in quality due to sampling is reasonably low for both HCC and GREAT. While GREAT is about 2 orders of magnitude faster than HCC, its quality is about 10% lower than HCC. Thus, if the goal is to achieve low average inconvenience regardless of high but practical running times, one may choose HCC. On the other hand, if a fast query time is desired with reasonably high accuracy in quality, GREAT may be chosen.

## 7 CONCLUSIONS

In this paper, we introduced two facility location problems over user-trajectories, namely, MAX-TIPS and AVG-TIPS, that aim to minimize the maximum and the average user-inconvenience, respectively. We showed that both these problems are NP-hard and proposed one optimal algorithm and two efficient heuristics for each of them. The heuristics can work both in the presence or absence of existing facilities. Empirical evaluation over large-scale real and synthetic datasets show that the proposed solutions are effective in terms of quality and efficient in terms of space and running time.

In future, we will explore other trajectory-based facility location problems.

## REFERENCES

- [1] Z. Drezner, *Facility location: A survey of applications and methods*. Springer, 1995.
- [2] H. W. Hamacher and Z. Drezner, *Facility location: applications and theory*. Springer, 2002.
- [3] Y. Du, D. Zhang, and T. Xia, “The optimal-location query,” in *STD*, 2005, pp. 163–180.
- [4] D. Zhang, Y. Du, T. Xia, and Y. Tao, “Progressive computation of the min-dist optimal-location query,” in *PVLDB*, 2006, pp. 643–654.
- [5] P. Ghaemi, K. Shahabi, J. P. Wilson, and F. Banaei-Kashani, “Optimal network location queries,” in *SIGSPATIAL*, 2010, pp. 478–481.
- [6] X. Xiao, B. Yao, and F. Li, “Optimal location queries in road network databases,” in *ICDE*, 2011, pp. 804–815.
- [7] Z. Chen, Y. Liu, R. C.-W. Wong, J. Xiong, G. Mai, and C. Long, “Efficient algorithms for optimal location queries in road networks,” in *SIGMOD*, 2014, pp. 123–134.
- [8] J. Qi, R. Zhang, L. Kulik, D. Lin, and Y. Xue, “The min-dist location selection query,” in *ICDE*, 2012, pp. 366–377.
- [9] J. Qi, R. Zhang, Y. Wang, A. Y. Xue, G. Yu, and L. Kulik, “The min-dist location selection and facility replacement queries,” *World Wide Web*, vol. 17, no. 6, pp. 1261–1293, 2014.
- [10] R. Liu, A. W.-C. Fu, Z. Chen, S. Huang, and Y. Liu, “Finding multiple new optimal locations in a road network,” in *Advances in Geographic Information Systems*, 2016, p. 36.
- [11] S. Mitra, P. Saraf, R. Sharma, A. Bhattacharya, S. Ranu, and H. Bhandari, “NetClus: A scalable framework for locating top-k sites for placement of trajectory-aware services,” *ICDE*, 2017.
- [12] J. Cui, M. Wang, H. Li, and Y. Cai, “Place your next branch with mile-run: Min-dist location selection over user movement,” *Information Sciences*, vol. 463–464, pp. 1–20, 2018.
- [13] M. Wang, H. Li, J. Cui, K. Deng, S. S. Bhowmick, and Z. Dong, “Pinocchio: Probabilistic influence-based location selection over moving objects,” *TKDE*, vol. 28, no. 11, pp. 3068–3082, Nov 2016.
- [14] Z. Chen, Y. Liu, R. C.-W. Wong, J. Xiong, G. Mai, and C. Long, “Optimal location queries in road networks,” *ACM Trans. Database Syst.*, vol. 40, no. 3, pp. 17:1–17:41, Oct. 2015.
- [15] J. Qi, Z. Xu, Y. Xue, and Z. Wen, “A branch and bound method for min-dist location selection queries,” in *ADC '12*, 2012, pp. 51–60.
- [16] B. Yao, X. Xiao, F. Li, and Y. Wu, “Dynamic monitoring of optimal locations in road network databases,” *VLDB Journal*, vol. 23, no. 5, pp. 697–720, Oct 2014.
- [17] A. K. M. M. R. Khan, L. Kulik, E. Tanin, H. Hua, and T. Hashem, “Efficient computation of the optimal accessible location for a group of mobile agents,” *ACM Trans. Spatial Algorithms Syst.*, vol. 4, no. 4, pp. 10:1–10:32, Sep 2018.
- [18] C. Chen, D. Zhang, L. Wang, X. Ma, X. Han, and E. Sha, “Taxi exp: A novel framework for city-wide package express shipping via taxi crowd sourcing,” in *UIC-ATC-ScalCom*, 2014, pp. 244–251.
- [19] R. C.-W. Wong, M. T. Özsu, P. S. Yu, A. W.-C. Fu, and L. Liu, “Efficient method for maximizing bichromatic reverse nearest neighbor,” *Proc. VLDB Endow.*, vol. 2, no. 1, pp. 1126–1137, Aug 2009.
- [20] D. YanEmail, Z. Zhao, and W. Ng, “Efficient processing of optimal meeting point queries in euclidean space and road networks,” *Knowledge and Information Systems*, vol. 42, no. 2, pp. 1319–351, Feb. 2015.
- [21] F. Chen, H. Lin, J. Qi, P. Li, and Y. Gao, “Collective-k optimal location selection,” *07 2017*, pp. 339–356.
- [22] X. Li, V. Čeikute, C. S. Jensen, and K.-L. Tan, “Trajectory based optimal segment computation in road network databases,” in *Advances in Geographic Information Systems*, 2013, pp. 396–399.
- [23] M. J. Hodgson, “The location of public facilities intermediate to the journey to work,” *European J. of Operational Research*, vol. 6, no. 2, pp. 199–204, 1981.
- [24] O. Berman, R. C. Larson, and N. Fouska, “Optimal location of discretionary service facilities,” *Transportation Science*, vol. 26, no. 3, pp. 201–211, 1992.
- [25] O. Berman, D. Krass, and C. W. Xu, “Locating discretionary service facilities based on probabilistic customer flows,” *Transportation Science*, vol. 29, no. 3, pp. 276–290, 1995.
- [26] O. Berman, D. Bertsimas, and R. C. Larson, “Locating discretionary service facilities, ii: maximizing market size, minimizing inconvenience,” *Operations Research*, vol. 43, no. 4, pp. 623–632, 1995.
- [27] O. Berman, D. Krass, and C. W. Xu, “Locating flow-intercepting facilities: New approaches and results,” *Annals of Operations Research*, vol. 60, no. 1, pp. 121–143, 1995.
- [28] O. Berman and D. Krass, “The generalized maximal covering location problem,” *Computers & Operations Research*, vol. 29, no. 6, pp. 563–581, 2002.
- [29] ———, “Flow intercepting spatial interaction model: a new approach to optimal location of competitive facilities,” *Location Science*, vol. 6, no. 1, pp. 41–65, 1998.
- [30] M. Boccia, A. Sforza, and C. Sterle, “Flow intercepting facility location: Problems, models and heuristics,” *J. Mathematical Modelling and Algorithms*, vol. 8, no. 1, pp. 35–79, 2009.
- [31] T. F. Gonzalez, “Clustering to minimize the maximum intercluster distance,” *Theoretical Computer Science*, vol. 38, pp. 293–306, 1985.
- [32] P. J. Rousseeuw and L. Kaufman, *Finding Groups in Data*. Wiley, 1990.
- [33] M. Charikar, S. Guha, E. Tardos, and D. B. Shmoys, “A constant-factor approximation algorithm for the k-median problem,” in *STOC*, 1999, pp. 1–10.
- [34] K. Jain and V. V. Vazirani, “Approximation algorithms for metric facility location and k-median problems using the primal-dual schema and lagrangian relaxation,” *JACM*, vol. 48, no. 2, pp. 274–296, 1999.
- [35] M. Charikar and S. Guha, “Improved combinatorial algorithms for the facility location and k-median problems and lagrangian relaxation,” *FOCS*, vol. 48, no. 2, pp. 378–388, 1999.
- [36] M. R. Korupolu, C. G. Plaxton, and R. Rajaraman, “Analysis of a local search heuristic for facility location problems,” *J. Algorithms*, vol. 37, no. 1, pp. 146–188, 2000.
- [37] V. Arya, N. Garg, R. Khandekar, A. Meyerson, K. Munagala, and V. Pandit, “Local search heuristics for k-median and facility location problems,” *SIAM J. Computing*, vol. 33, no. 3, pp. 544–562, 2004.
- [38] R. T. Ng and J. Han, “Clarans: A method for clustering objects for spatial data mining,” *TKDE*, vol. 14, no. 5, pp. 1003–1016, 2002.
- [39] Y. Lou, C. Zhang, Y. Zheng, X. Xie, W. Wang, and Y. Huang, “Map-matching for low-sampling-rate GPS trajectories,” in *GIS*, 2009, pp. 352–361.
- [40] H. Bhandari, “Clustering of sites in a road network based on coverage of trajectories,” Master’s thesis, Indian Institute of Technology, Kanpur, 2016.
- [41] G. L. Nemhauser, L. A. Wolsey, and M. L. Fisher, “An analysis of approximations for maximizing submodular set functions,” *Mathematical Programming*, vol. 14, no. 1, pp. 265–294, 1978.
- [42] J. Yuan, Y. Zheng, X. Xie, and G. Sun, “Driving with knowledge from the physical world,” in *KDD*, 2011, pp. 316–324.
- [43] J. Yuan, Y. Zheng, C. Zhang, W. Xie, X. Xie, G. Sun, and Y. Huang, “T-drive: driving directions based on taxi trajectories,” in *SIGSPATIAL GIS*, 2010.
- [44] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms*. MIT Press, 2009.

## APPENDIX A PROOFS OF THEOREMS

### A.1 Proof of Theorem 2 (Quality of MIF)

*Proof.* Let  $\mathcal{Q} = \{s_1, \dots, s_k\}$  be the set of  $k$  sites returned by MIF, with a maximum inconvenience of  $d$ . Let  $\mathcal{RT} = \{T_1, \dots, T_k\}$  be the  $k$  representative trajectories chosen by MIF. Assuming  $\mathcal{S} = V$ , it follows that each node in  $T_i$  is a candidate site. Since  $s_i$  is the nearest candidate site to  $T_i$ , it must be that  $s_i \in T_i$ .

Now, let us consider the case  $k = 1$ . Let  $d$  be the maximum inconvenience due to the selection of the site  $s_1 \in T_1$ . Now suppose  $s_1^*$  is the site reported by an optimal algorithm with optimal maximum inconvenience value  $d^*$ . Hence,  $\forall T_q \in \mathcal{T}$ ,  $d_r(T_q, s_1^*) \leq d^*$ . Therefore, there exists a site  $s_1' \in T_1$ , such that  $d_r(T_q, s_1') \leq 2d^*$  for any trajectory  $T_q \in \mathcal{T}$ , because there exists a path from  $s_1'$  via  $s_1^*$  to trajectory  $T_q$  of distance at most  $2d^*$ . Thus, the maximum inconvenience  $d$  due to the site  $s_1$  is  $d \leq 2d^*$ .

Now, consider the case  $k \geq 2$ . Consider an optimal solution  $\mathcal{Q}^* = \{s_1^*, \dots, s_k^*\}$  with maximum inconvenience  $d^*$ . Let  $Clus(s_i^*)$  denote the set of trajectories served by  $s_i^*$ . Using the pigeon-hole principle, we conclude that at least two representative trajectories  $T_p, T_q \in \mathcal{RT}$ ,  $p < q$ , must belong to one of the  $k$  clusters, say  $Clus(s_i^*)$ . Since the maximum inconvenience of the solution  $\mathcal{Q}$  is  $d$ , therefore,  $d_r(T_q, s_p) \geq d$ . Moreover, since  $T_p, T_q \in Clus(s_i^*)$ , we get  $d_r(T_p, s_i^*) \leq d^*$ ,  $d_r(T_q, s_i^*) \leq d^*$ . From this, we infer that there must be a site  $s_p' \in T_p$  such that  $d_r(s_p', s_i^*) \leq d^*$  and, therefore,  $d_r(T_q, s_p') \leq 2d^*$ . Since the length of any trajectory is at most  $L$ ,  $d_r(s_p, s_p') \leq L$ . From this, we conclude that  $d \leq d_r(T_q, s_p) \leq d_r(T_q, s_p') + d_r(s_p, s_p') \leq 2d^* + L$ .  $\square$

### A.2 Proof of Theorem 3 (Complexity of MIF)

*Proof.* We assume that the number of road segments (edges) is  $O(n)$  as the road networks are roughly planar. Thus, from any node  $v_i \in V$ , the distances to all other nodes in the network can be computed in  $O(n \log n)$  time, using Dijkstra's shortest path algorithm [44].

We analyze the computation cost of any iteration as follows. Recall that  $\mathcal{NF}$  map stores the trajectories in a sorted order based on their distance to the nearest facility in  $\mathcal{Q}$ . Since  $|\mathcal{NF}| = m$ , hence to identify the trajectory  $T_i$  at rank  $\gamma \times |\mathcal{T}|$ , requires  $O(1)$  time, using array implementation of  $\mathcal{NF}$ .

Then, for each node  $v \in T_i$  (that  $T_i$  passes through), the distances are computed to all other nodes in  $V$ . If the maximum number of nodes in any trajectory is  $l$ , i.e.,  $|T_i| \leq l$ , then the above distance computation step takes  $O(l \cdot n \log n)$  time. Thus, identifying the nearest candidate site to  $T_i$ , say  $s_i$ , requires  $O(l \cdot n \log n)$  time.

Following this,  $s_i$  is added to  $\mathcal{Q}$ . Then the distances are computed between  $s_i$  and all other trajectories in  $\mathcal{T}$ . To do this, we first compute distances between  $s_i$  and all nodes in  $V$ , which requires  $O(n \log n)$  time. Assume these distances are indexed on site-ids. The distance between  $s_i$  and any trajectory can be computed in  $O(l^2)$  time, as there are at most  $O(l^2)$  distance look-ups. Therefore, the distance between  $s_i$  and all trajectories in  $\mathcal{T}$  can be computed in  $O(ml^2)$  time. If the distance of any trajectory  $T_j \in \mathcal{T}$  to its nearest facility in  $\mathcal{Q}$  is more than that with  $s_i$ , then this value is updated in the  $\mathcal{NF}$  map. This

update step takes  $O(m)$  time over all the trajectories. Sorting the  $\mathcal{NF}$  map takes  $O(m \log m)$  time. Summing up all these costs, over each of the total  $k$  iterations, the total time complexity is  $O(k \cdot l \cdot n \log n + k \cdot m \cdot l^2 + k \cdot m \log m)$ .

Next, we analyze the space complexity. The road network and candidate sites can be stored in  $O(n)$  space. Storing the trajectories in  $\mathcal{T}$  require  $O(l \cdot m)$  space. Maintaining the  $\mathcal{NF}$  map requires  $O(m)$  space. During each iteration of the algorithm, the node-to-node distances are computed for each node of the previously chosen representative trajectory. Storing these distance values require at most  $O(ln)$  space. As these distance values are no longer used in subsequent iterations, they are discarded at the end of each iteration. Thus, the maintenance overhead of the node-to-node distances over the  $k$  iterations is  $O(l \cdot n)$ . Storing the sets  $\mathcal{RT}$  and  $\mathcal{Q}$  require  $O(k) = O(n)$  space. Hence, the total space complexity is  $O(l(n + m))$ .  $\square$

### A.3 Proof of Theorem 4 (Complexity of NetClus)

*Proof.* Assuming the time required to answer a TOPS query by NetClus to be  $O(t_{TOPS})$ , since  $O(\log_2(\tau_{max}/\tau_{min}))$  TOPS queries are executed, the total time is  $O(\log_2(\tau_{max}/\tau_{min}) \cdot t_{TOPS})$ . The time complexity,  $O(t_{TOPS})$ , of NetClus is analyzed in [11], and is beyond the scope of this paper.

Now, let us analyze the space complexity of NetClus. As discussed above, if the index resolution parameter is  $\epsilon > 0$ , then the total number of index instances is  $t = 1 + \lceil \log_{1+\epsilon}(\tau_{max}/\tau_{min}) \rceil$ . For each index instance, the nodes in  $V$  are divided into clusters which requires  $O(n)$  space. In addition, for each cluster, we store the set of trajectories that pass through it. This cost is  $O(m \cdot l)$  where  $m = |\mathcal{T}|$  and  $l$  is the largest number of nodes in any trajectory. For each node in  $V$ , we store its distance to the cluster-center of the cluster it belongs to. This storage cost is  $O(n)$  across all the nodes. Further, for each trajectory, we store its distance to the cluster-center of the clusters it passes through. Since a trajectory cannot pass through more than  $l$  clusters, this storage cost over all the trajectories is no more than  $O(m \cdot l)$ . Therefore, the total space complexity is  $O(t \cdot (n + m \cdot l))$ .  $\square$

### A.4 Proof of Theorem 5 (Complexity of HCC)

*Proof.* Let  $n', m'$  denote the number of sites and trajectories, that are produced after sampling of the sites and the trajectories, respectively. Thus,  $n' \leq n$  and  $m' \leq m$ . We assume that  $k \leq n'$ .

The algorithm is executed for  $t$  trials, and in each trial, the maximum possible number of iterations is  $\eta$ . In a given iteration, the number of swaps that are scanned is  $sf \cdot k \cdot (n' - k)$ . To evaluate each swap, we need to compute the total inconvenience of all the  $m'$  trajectories with respect to the current set of  $k$  medoids. This requires  $O(k \cdot m')$  time (assuming that the distance between each pair of site and trajectory is pre-computed). Hence, the total running time is  $O(t \cdot \eta \cdot sf \cdot k^2 \cdot (n' - k) \cdot m')$  time.

Next, let us analyze the space complexity. Storing the sampled sites and trajectories require  $O(n')$  and  $O(m' \cdot l)$  space where  $l$  is the maximum number of nodes in any trajectory. Storing the pairwise distance between each pair of site and trajectory in the sampled space requires  $O(m' \cdot n')$  space. To store the  $\mathcal{NF}$  maps, we need  $O(m')$  space. In order to store the current set of medoids, we need  $O(k) = O(n')$  space. Summing up all the costs, we find that the total space complexity is  $O(m' \cdot (n' + l))$ .  $\square$

### A.5 Proof of Theorem 7 (Quality of GREAT)

*Proof.* It is easy to see that GREAT returns the optimal solution for  $k = 1$  as it performs an exhaustive search over all the candidate sites in  $\mathcal{S}$ .

Now, suppose  $k \geq 2$ . For any given set of candidate sites  $\mathcal{Q} \subseteq \mathcal{S}$ , consider a function  $f(\mathcal{Q}) = TI(\mathcal{F}) - TI(\mathcal{Q})$ . Since  $Tl(\mathcal{Q})$  is a non-increasing super-modular function (Th. 6), it follows that  $f(\mathcal{Q})$  is a non-decreasing sub-modular function. Further, if the existing set of facilities  $\mathcal{F} \neq \emptyset$ , then the initial total inconvenience for  $\mathcal{Q} = \emptyset$ , can be written as  $Tl(\emptyset) = TI(\mathcal{F})$ . Thus,  $f(\emptyset) = TI(\mathcal{F}) - TI(\emptyset) = 0$ . It is known that the greedy heuristic offers an approximation bound of  $1 - 1/e$  for any non-decreasing sub-modular function  $f$  with  $f(\emptyset) = 0$  [41]. Since  $OPT$  is a set that minimizes  $Tl$ , from the definition of  $f$ , it follows that  $OPT$  must maximize  $f$  as  $Tl(\mathcal{F})$  is a constant. Let  $\mathcal{Q}$  be the set of sites reported by GREAT. Since GREAT essentially mimics the greedy heuristic for non-decreasing sub-modular functions [41], the same approximation bound is applicable for  $f$ . Therefore,  $f(\mathcal{Q}) \geq (1 - 1/e)f(OPT)$ . From the definition of  $f$ , we can write  $Tl(\mathcal{Q}) \leq$

$$(1 - 1/e)(TI(OPT) - TI(\mathcal{F})) + TI(\mathcal{F}) \leq (1 - 1/e)TI(OPT) + TI(\mathcal{F})/e. \quad \square$$

### A.6 Proof of Theorem 8 (Complexity of GREAT)

*Proof.* The space required to store the trajectories and the nodes on the road network are  $O(ml)$  and  $O(n)$  respectively. Further, storing the distance values for each pair of site  $s_i \in \mathcal{S}$  and trajectory  $T_j \in \mathcal{T}$  requires  $O(m.n)$  space. Further, storing the  $\mathcal{NF}()$  maps require  $O(m)$  space. Since  $|\mathcal{Q}| = k = O(n)$ , the overall space complexity is  $O(m.(n + l))$ .

We next analyze the time complexity. In each iteration of the GREAT algorithm, we compute the total inconvenience  $Tl(\mathcal{Q} \cup \{s_i\})$  for each site  $s_i \in \mathcal{S} \setminus \mathcal{Q}$ . For each site  $s_i$ , this computation requires  $O(m)$  time over all the trajectories. Once a site is added to  $\mathcal{Q}$ , the  $\mathcal{NF}$  maps are updated in  $O(m)$  time. Thus, each iteration requires  $O(m.n)$  time. The overall time complexity of GREAT, running over  $k$  iterations is, therefore,  $O(k.m.n)$ .  $\square$