# Crowdedness Spot Acquisition by Using Mobility Based Clustering

## Sujit A. Navale<sup>1</sup>, Yogesh B Gurav<sup>2</sup>

<sup>1</sup>Savitribai Phule Pune University, Pune, Maharashtra, India

<sup>2</sup>Professor, Savitribai Phule Pune University, Pune, Maharashtra, India

Abstract: Detecting hot spots of moving vehicles in an urban area is absolutely required to many smart city applications. Crowdedness spot is a crowded area with a number of irregular objects. The practical investigation on hot spots in smart city offerings many unique features, such as highly mobile environments, the non-uniform biased samples, and supremely limited size of sample objects. The traditional density-based clustering algorithms flop to capture the actual clustering property of objects, making the outputs meaningless. In this paper we recommend a novel, called mobility-based clustering which is non-density-based approach. The basic idea is that sample objects are hired as "sensors" to recognize the vehicle crowdedness in nearby areas using their instant mobility, rather than the "object representatives". As such the mobility of samples is certainly incorporated. Several important factors beyond the vehicle crowdedness have been identified and techniques to remunerate these effects are proposed. Furthermore, taking the identified crowdedness spots as a label of the taxi, we can identify one individual taxi to be a crowdedness taxi that crosses a number of different crowdedness spots. We estimate the presentation of our methods and baseline approaches based on real traffic situations and real-life data sets.

Keywords: Data mining, Mobility-based clustering, traffic detection, vehicle, crowdedness, intelligent transportation systems, vehicular and wireless technologies

## 1. Introduction

Many metropolitan cities are facing a number of serious problems, such as frequent traffic jams, unexpected emergency events, and even disasters. Many of these problems are relative to crowded moving objects such as vehicles, trains, etc. Detecting hot spots of moving vehicles in an urban area is completely necessary to many smart city applications. Informally, areas of high crowdedness of vehicles can be described as hot spots of vehicles. The hot spots with especially high crowdedness are usually the sites of traffic congestions. An immediate application for hot spot study is that we can anticipate vehicle speeds based on the crowdedness distribution. Indeed, hot spots are often the potential sites of interests due to the higher likelihood of the events and opportunities (e.g. traffic jam, exhibitions, and commercial promotions). However, because of the privacy issues or localization equipment limitations it is hard to collect the location information of all the vehicles in the city.

In this work, we study the hot spot related issues using taxi's statistics as the sample data. A better understanding of the city traffic via a quantitative research on hot spots is the ultimate goal of this research. To do so we have several key tasks to achieve: 1) to define and quantify the vehicle crowdedness of an area; 2) to picture the crowdedness distribution of the city and detect the hot spots; and 3) to investigate the evolution of hot spots.

Given the dynamic temporal and spatial information of moving vehicles, hot spots can be considered as a general case of object clustering in mobile environments. In recent years, web related clustering, evolutionary clustering in low mobility environments and uncertain data streams have also drawn a lot awareness. In our application structure, however, some new unique features make previous well-designed algorithms fail to express the real clustering property of moving vehicles.

The first major challenge is incomplete information. Existing algorithms (for static or mobile) are all densitybased approaches that use inter-node distances as a critical measure. They depend on the location information of target objects to utilize the clustering property. However, in many practical applications, it is unlikely to obtain such information from the total population of vehicles. It greatly reduces the effectiveness of density-based algorithms. Second, the sample object set is a specific type of vehicles. It has very limited abstraction to represent general vehicles. Besides these, there are also some practical concerns such as extreme dynamism, and high object mobility.

To deal with these challenges, we suggest a novel, nondensity based approach called mobility-based clustering. Mobility-based clustering is based on a simple observation that usually vehicles are deliberate to have high mobility. A vehicle of high mobility can largely designate a low crowdedness and vice versa. By this, the sample vehicles are not simply used as objects but appoint as "sensors" to recognize the vehicle crowdedness in nearby areas. The main advantages of mobility-based clustering are several folds. First, mobility-based clustering is less sensitive to the size of the sample object set, though a larger sample set can produce more accurate readings of the crowdedness sensing. Second, mobility-based clustering does not require exact location information and thus is durable to the location inaccuracy. Third, mobility-based clustering naturally includes the mobility of vehicles. It is therefore particularly suitable for high mobility environments. Mobility-based clustering greatly outperforms existing density-based clustering algorithms in terms of forecast accuracy of vehicle density because of these advantages.

The density-based clustering using taxis as samples will generate a quite deviated result. Such a deviation, which is mainly due to the intrinsic limitation of density-based approaches, is our main motivation for this work. In this paper, we make improvement along following conditions. First, to quantify the crowdedness of certain areas, fully taking the mobility and object dynamism as an advantage we propose a novel mobility-based model. This is, to the best of our knowledge, the first attempt of non-density-based object clustering in literature. Second, several key factors, which have the appreciable effect on the accuracy of the vehicle crowdedness measurements, are identified and investigated. Effective techniques to remunerate the negative effects have been developed. Third, we study the hot spots of top vehicle crowdedness values and analysis of results show that several top hot spots are completely locality consistent over time, while more hot spots present more locality difference.

# 2. Preliminaries

In this section, we first present characteristics of the raw dataset used in our work. Moreover, we introduce road network grid. At last, we present the main observations and design principles of mobility-based clustering.

## 2.1 Raw Dataset Characteristics

The GPS receivers frequently report their current states to a data center via GPRS links. The reports include the instant speed, the geographic location and the status of occupied or unoccupied (by guest) of the taxi. In the remainder of this paper, we use the terms "sample" and "taxi" correspondently, as well as the terms "location" and "spot" if not otherwise stated. The term "vehicle" is used to represent the general vehicles including taxi samples, taxis not sampled, private cars and buses.

The GPS system that is installed for civic applications. Due to the low cost of these applications, the data reports mainly have the following limitations. First, the data set is incomplete. Noticeable amount of reports were missing due to weak GPRS signals (via which taxis are connected to the system) or limited bandwidth of GPRS wireless channels. If we use this insignificant sample to represent the large number of general vehicles the error will be important. Moreover, all sample objects are taxis which are only one definite type of vehicles. Taxis are highly enticementoriented that have strong preferences on some desired locations. They would like to accumulate on sites of high customer flows, such as business areas, train stations, and traffic reconnections. Such preferences make it a bad option to employ this one type of vehicles as the representative of others.

Second, due to blocked GPS signals (e.g., taxis in tunnel or surrounded by high buildings) the reported GPS data may not be exact. Since GPRS is a paid communication service, it is costly to periodically report their current status information. In the city, taxis are allowed to report their data at an inconsistent time, with a desired 5 second period. In 5 seconds a vehicle can drive 150 meters at 100 kmph speed. Concerning all these factors, the location errors of vehicles are on the order of hundreds of meters. It becomes impossible to apply the traditional density-based approaches which are desperately relied on the exact locations.

Third, the data is biased in temporal and spatial spaces. For example, 90% roads have no data for more than 80% of the time in a day, and 50% have no data in 12 continuous hours. To the opposite, 80% of the reports are collected from 20% of roads. How to mine meaningful information from the biased samples is another great challenge. Motivated by these new challenges, we propose a novel, mobility-based clustering method.

## 2.2 Road Network Grid

For ease of calculation, we discretize the time dimension in the unit of second. The time instance t = 0 is the initial time. The time instances t and t+1 are consecutive time instances. We discretize the physical space by dynamically partition the whole area into a number of rectangle grids. The size of each grid depends on the intensity of the reports at the area, ranging from 10 meters for report rich areas to 90 meters for report scarce areas. We believe this granularity is sufficient for most applications. Each grid is represented by its center location so that all spots in the grid will be evaluated the same as that of the grid center. A 2-tuple i = (x, y) represents one grid where x is the index of the grid along the longitude and y is that along the latitude. We chiefly introduce domain knowledge and build road network to enhance the grids and fetch much more exact spot locations. Because the road topology and type will impact the vehicle, the speed as well as the drive pattern, hence we study the following problems based on road network grid.

## 2.3 Observations and Design Principles

Different from traditional density-based approaches, mobility based approach is placed on two simple conclusions. The first one is that vehicles prefer high mobility in a infrequent region. To the opposite, for security concerns vehicles will drive slowly when the nearby area is crowded. Motivated by it, we apply vehicles as sensors using their instant speed to sense the vehicle crowdedness of proximity. The second one is that the reported locations can be inaccurate, while the reported speeds are directly obtained from the speedometers installed on taxis so they are usually quite accurate. For safety concerns sudden changes of speeds are uncommon. Therefore the speed errors coming from the unsynchronized reports are also small.

Practically speaking, in mobility-based clustering we collect statistics of taxi speeds at each spot. The spot crowdedness is then a relative measurement regarding the instant speed, the maximum speed, and the minimum speed. Though a higher crowdedness usually leads to a smaller mobility, by high crowdedness a smaller mobility is not always caused. Besides the spot crowdedness, there are many other factors having similar effects on taxi mobility.

Firstly, one fact is that drivers may have various driving styles and nature. In particular, due to enticement-oriented nature, employed taxis (by guests) often have higher speeds than un employed taxis which may be looking for guests. Profiling these different drivers will help to describe taxi motility more accurately.

Secondly, mobility of vehicles is environment dependent. Some roads are designed for high speed traffic, while others are mainly for connection purposes. Traffic lights clearly slow down vehicles, which is not due to the high crowdedness of the spots. We should characterize spots so that to diminish these negative effects.

Thirdly, spot crowdedness may have spatial and temporal correlations. Contiguous spots may have strong connections in between. A crowded spot is very likely to be crowded again in the next time stamp. Hot spots may derive over both time and spatial dimensions. To well capture the crowdedness of spots, we should take all these factors into account so that the derived crowdedness values can appropriately reflect the real crowdedness of spots.

# 3. Crowdedness Spot Acquisition

The crowdedness spot can be considered as a higher level of feature retrieved from the taxi. Hence, we can moreover operate the crowdedness spot to study the taxi. For example, the taxis always cross crowdedness spots may be have more chances to detention the crowded areas' information or pick up passengers; at the same time, these taxis' behavior may help us give more analysis of the city transportation. In this section, we build the support vector machine (SVM)-based intelligent search to categorize the taxis.

In crowdedness taxi intelligent search process. First, a domain expert makes the directed taxi features, uses them to create the learning data sets, and exploits the data sets to train and build the predictive model. Second, the directed features are published to the users. Third, a user selects a feature of interest to regain the relevant list of crowdedness taxis from a search engine. Fourth, the retrieved taxis are analyzed and categorized by the predictive model. Finally, only the taxis that are scored as significant are sent back to the user.

# 4. Field Study Evaluation

In this segment, in the first place, we assess the crowdedness work by the genuine information sets (taxi information and transport information), and second, to confirm the adequacy of our portability based grouping, we contrast it and the current system on various field studies (field cam records) and observational information sets. At long last, we assess the insightful SVM-based crowdedness taxi classifier.

## 4.1 Crowdedness Function Validation

Our errand of approval is to assess which of the direct and factual crowdedness capacities delivers a finer fit crowdedness circulation. We help out the approval through two methodologies. One is by leading the expectation about vehicle speed. By this we search for an ideal design for portability based bunching. The second acceptance is to think about the ideal setting portability based bunching with conventional calculations.

#### 4.2 Mobility-Based Clustering Validation

So as to confirm the adequacy of our versatility based grouping, we lead various field studies. We setup camcorders at foreordained locales in the city, record the genuine activity circumstances in fields, and afterward measure the crowdedness of these territories through a logged off way. The results are indicated as "genuine circumstance". For mobility-based clustering, the created crowdedness is a relative quality. It needs a suitable scalar to create unquestionably the vehicle densities. To acquire this scalar, we have to gather the genuine beginning condition of the spot for adjustment. In practice, this alignment normally acquires generally high cost. We contend, by and by, that this adjustment is a solitary run operation such that the high cost can be amortized over a long operation time.

# 5. Related Work

Clustering: Object bunching is an overall concentrated on issue with a lot of examination endeavors being dedicated in. As of late, bunching moving articles is turning into a hot exploration issue. Superb moving microclusters are alertly kept up which prompts quick and focused grouping results. at the point when the thickness or amount of the items is not that adequate for uncommon application situations, they will fizzle. In our application situation, thickness based systems don't work because of the new emerging peculiarities, for example, great less examples and prominent information point area slips. The clustering algorithm DBSCAN depends on a density-based notion of clusters and is intended to find bunches of subjective shape and additionally to recognize clamor. The generalized algorithm - called GDBSCAN - can bunch point questions and spatially stretched out articles as indicated by both, their spatial and their non-spatial qualities. What's more, four applications utilizing 2D points (astronomy), 3D points (biology), 5D points (earth science) and 2D polygons (geography are displayed, showing the pertinence of GDBSCAN to true issues.

## Traffic data analysis

The discovery of exceptions in spatio-temporal traffic data is a vital exploration issue in the information mining and learning revelation group. However to the best of our insight, the disclosure of connections, particularly causal associations, among located activity anomalies has not been examined in the recent past. We propose calculations which build exception causality trees focused around worldly and spatial properties of discovered exceptions. Incessant substructures of these causality trees uncover repeating communications among spatio-temporal outliers, as well as potential defects in the outline of existing movement systems. The adequacy and quality of our calculations are accepted by probes a huge volume of genuine taxi trajectories in a urban street system. Another classification of related work concentrates on the examination of versatile movement object information. They are primarily intrigued by the discovering ranges of high activity load. The current work expected that the committed sensor gadgets had been conveyed so that the gathering of vehicle crowdedness gets to be direct. In our work, we don't have committed sensors however utilize versatile protests as "sensors" to see the crowdedness.

In this paper, we don't have committed sensors yet utilize portable questions as "sensors" to see the crowdedness. The distinction between our methodology and the skimming auto in intelligent transportation system (ITS) is that the gliding auto in ITS concentrates on the rate of vehicles, and for crowdedness spot issues, they have the normal presumptions that each one buoy auto speaks to various genuine vehicles. In our issue, then again, taxis are bad agents of different vehicles, and thusly, such methodologies will fizzle.

# 6. Conclusion and Future Work

In this paper, we have proposed mobility-based clustering, a novel methodology to distinguish crowdedness spots in an exceptionally versatile environment with to a great degree constrained and one-sided item inspects. The remarkable mobility-based clustering is to utilize speed data to induce the crowdedness of moving articles. Besides, we consider the crowdedness spot classifications and the crowdedness taxi securing from the located crowdedness spots. We assessed the execution of mobility-based clustering based with respect to genuine taxi information gathered in the city through field studies. Future work can be directed along taking after headings. First and foremost, in mobility-based clustering, the velocity data is discriminating. Because of the little example information set, we utilized a basic methodology to gauge the portability of vehicles at the spot of no information. Better portability estimation can create better crowdedness values. Second, there are numerous variables other than spot crowdedness that will have affect on vehicle versatility, for example, activity lights and fender benders. We abandon them for future work. Third, we require more field studies, despite the fact that work escalated, to further confirm the adequacy of the versatility based methodology. Fourth, better street griding strategy is required for recovering a great deal all the more valuable areas. At last, contingent upon different qualities of moving articles, other non-thickness based grouping may be worth further examinations.

# References

- [1] S. Liu, Y. Liu, L. Ni, J. Fan, and M. Li, "Towards mobility-based clustering," in *Proc. ACM SIGKDD*, 2010, pp. 919–928.
- [2] Y. Li, J. Han, and J. Yang, "Clustering moving objects," in *Proc. ACM SIGKDD*, 2004, pp. 617–622.
- [3] W. Liu, Y. Zheng, S. Chawla, J. Yuan, and X. Xing, "Discovering spatio-temporal causal interactions in traffic data streams," in *Proc. ACMSIGKDD*, 2011, pp. 1010–1018.
- [4] R. T. Ng and J. Han, "CLARANS: A method for clustering objects for spatial data mining," *IEEE Trans. Knowl. Data Eng.*, vol. 14, no. 5, pp. 1003–1016, Sep./Oct. 2002.
- [5] K. Sirvio and J. Hollmén, "Spatio-temporal road condition forecasting with Markov chains and artificial neural networks," in *Proc. HAIS*, 2008, pp. 204–211.
- [6] P. S. Castro, D. Zhang, and S. Li, "Urban traffic modelling and prediction using large scale taxi GPS traces," in *Proc. Pervasive*, 2012, pp. 57–72.