

Mining GPS Data for Traffic Congestion Detection and Prediction

Suhas Prakash Kaklij

Department of Information Technology, Siddhant College of Engineering, Sudumbare, Pune University, Pune 412109, India

Abstract: GPS data is available in the large amount, also for the devices having GPS a large amount data is being collected over time. The mining of this huge data is endorsed in discovery of the areas which face regular traffic congestion. User will have prior awareness of such locations which guide in deciding whether or not to go for that route. Avoidance of such routes will also assist in reduction of congestion of such locations. Also detected that the work which has been carried out till now in this field do not provide very precise and relevant results. The reason behind this is the no proper algorithm are selected and distinguished between on road and off road traffic. To deal with all this we proposed this system. This system will be structured and applied over GPS data i.e. data coming from devices like mobile phones, tablets, on board units etc. In the technique used in this system, these GPS data will be first cauterized using the K-means clustering algorithm. The clusters obtained are filtered out. On further processing these clusters a mining method of Naive bayes algorithm is used for mining for traffic Congestion detection and prediction

Keywords: Traffic Congestion Detection, Traffic Jam Prediction, Traffic Tracking and Tracing.

1. Introduction

Road network is biggest network widely used for Transportation. Each city has its Road network. Roads are used for daily transport not only for the people but for goods and many other things. The biggest problem now a days people facing is Traffic Congestion. The most of the congestion occur early morning or late afternoon because students and employees are going to their works and colleges so they also be late at traffic spot. People are not able to reach their work due to this traffic problem. As per the observation traffic congestion is dynamic in nature, it is not static. Means traffic congestion is variable as time passes and resources provided by current infrastructure are limited. In current emerging IT world we have lot of traffic data available with us in the different formats. With the use of this data we can get the flow of traffic information with respect to the location and time. This traffic information is important not only for current status of traffic but it can helps to analyze and predict upcoming traffic patterns. We can collect such information by processing GPS data. With use of 2G and 3G enabled GPS devices, huge set of data is collected with an average error of 2-15m [2]. These errors can further be decreased using some of the correction strategies such as map-based correction given in [2]. It is real time data which gives convenience to mine the traffic patterns of particular area. We can evaluate such data to get the traffic congestion patterns which in turn helps to identify the location where traffic congestion is possible. Prediction is also possible for traffic congestion of relative routes with respect to time.

2. Related Work

Substantial amount of efforts experimented in the field of analyzing traffic patterns. H. Inose et al. In 1967, as given in [3], proposed how traffic signals are work systematically. It works for the minimization of delay in time of vehicles and providing appropriate and preferential offsets to the optimal graph in the road network. In 2002, Ashbrook et al., as given

in [4], projected user consequential locations and end user activities using GPS data. As projected the city is divided in to clusters using K-means clustering further classified into a Markov Model. Thus, their work targeted on analyzing user GPS data to mine user momentous locations. As per the year 2010, Lipan et al. in [5], mined traffic patterns from GPS data concentrated from public transport. Their work focuses on observing bus schedules. Association guidelines are built on clusters where individual cluster has its own moderate speed. In 2011 [6] Mandal K and his team used probe vehicle technique for traffic congestion monitoring, Traffic information from probe vehicles has great potential for improving the estimation accuracy of traffic situations System tries to monitor the traffic flow pattern and then detect the congestion. As given in [7], Yao et al. proposed a speed pattern model which provide assumptions for traffic conditions and speed pattern with the help of machine learning. In 2013[1] Anand Gupta and his team proposed a groundwork for traffic congestion detection concentrating more on innovative algorithm which in turn cut down conflict of data for traffic Jam and Traffic signal. These efforts have given significant and helpful results. On the other hand, to the best of the authors' discoveries, not much attention has been given to detection and prediction of traffic congestion with properly manipulating of on road & off road data, As well as the Conflict between the Traffic signal and Traffic Jam. Also no appropriate selection of mining and clustering algorithm.

3. Motivation

Detecting traffic jam based on simple rules, such as using a probe vehicle technique for improving the estimation accuracy of traffic situations, velocity-based approach, and fuzzy logic might not handle the problem stated previously with great effect due to the following reasons

Section headings come in several varieties:

1. Selection of no proper Clustering algorithms.
2. Use of less relevant mining methods.

3. Lack of segregation between road and off road traffic
 These are basic motivation circumstances for proposing and developing such TCD –framework.

4. Contribution

To accomplish the points indicated in the motivation section, we have proposed a framework called TDC- where in T stand for “Traffic”, C for “Congestion”, and D for “Detection”. This framework is hybrid strategy of two distinct finest methods namely K-means and Naive Bayes which together helps to give more accurate result for Traffic congestion detection. Whereas framework proposed by Anand Gupta and it team more emphasis given on algorithm which reduces conflict of data for traffic Jam and Traffic signal

5. Organization

To explain the framework - TCD, We organized paper as follows: Section-II explains the structure of the framework and the algorithms associated with it. Section-III Shows Algorithm used in framework Section IV Shows how Experiment is to be carried out Section V shows the expected result as obtained. Section-VI concludes the paper describing

6. Proposed Framework

Description of DTC Framework (Refer to Fig 2): The logs files of GPS that are being collected overtime comprise of GPS data coming from Mobiles phones, Notepad devices or GPS enable devices .Data might be in different forms like for users who might be in a motionless state like sitting or in motion such as walking. The initial GPS data collected from mobile devices have log with information of Device ID, Location details like <Lat, Long>, Time related data like Time and date is sent through first stage Data generation module to Data importing module where we load this attributes to database .Many time GPS data available in two different form one having speed information as attribute and one without speed attribute. In our data generation module we do not have speed attribute generated, so we can calculate it by using Haversine Formula in fig 1

In this formula Φ_1 , Φ_2 are the latitudes of two given points and λ_1 , λ_2 are their corresponding longitudes. Now we have stored the attributes to data base as <Device Id, Latitude, Longitude, Time, Speed, and Day>

$$d = 2r \arcsin \left(\sqrt{\sin^2 \left(\frac{\Phi_2 - \Phi_1}{2} \right) + \cos(\Phi_1) \cos(\Phi_2) \sin^2 \left(\frac{\lambda_2 - \lambda_1}{2} \right)} \right)$$

Figure 1: Haversine formula to calculate distance between two points

All this information stored is used to detect the location of the particular device id. Once Location is detected for device id then by using location information and the speed of the device id we can categories it into On Road and Off Road. The logic to get On Road and Off Road data is simple .We have Device ID with its location including the <latitude, longitude> ,We also have the <latitude, longitude> of the Road . So coordinates which are not in the data set of Road

data coordinates all such devices are categorized as Off Road data and Device Id which are in the range of the Road <latitude, longitude> data is categorized as on Road data to avoid the confusion between data set . Now proceed with On Road data and ignored the Off Road data. The on road data is further inputted to Filtering process in which data is process in 2 parts .In part 1, analysis of data is done where data is analyzed from different perspective so that it helps to determine the transportation medium of each Device ID .In this process first threshold average speed of each device is decided so according to data available along with the speed and distance information average speed is calculated and the transportation medium is decided for each device id.

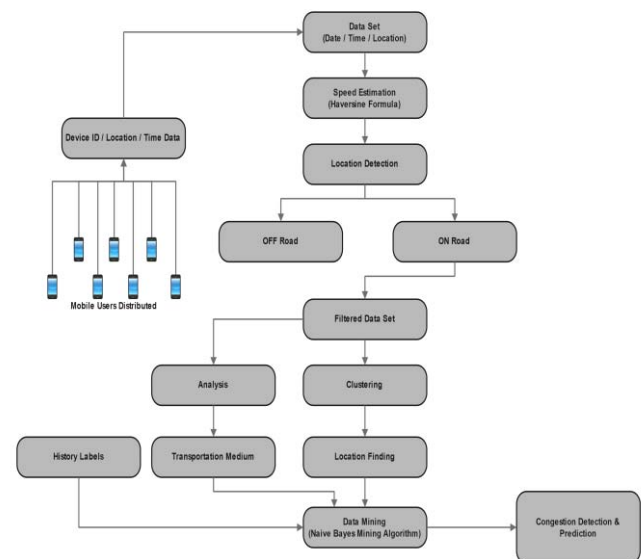


Figure 2: TCD Framework

In part 2, with the help of latitude, longitude, average speed and a unique Cluster-ID divide a city map into clusters of different sizes. We have chosen the most efficient, faster K Means clustering algorithm. The GPS raw data available with all the above parameters are applied to this K-Means clustering .The one of the main property for selection of the K-Means is whatever clusters are resulted those are non-hierarchical and they do not overlap, also the K-Means produces tighter clusters than hierarchical clustering [10]. Once clustering is do with help of the clusters location of each cluster is detected [8] The output of both part 1 & 2 is inputted to mining process where data is process by Naive Bayes Algorithm [9] for traffic congestion detection. Which is more advance over r the J48 decision tree classifier [11]. We also have the Historical data available in addition to part 1 & 2 data so that efficient traffic congestion prediction is done

7. Algorithm

Clustering K-Means Algorithm:

K-means is used to solve the familiar clustering problem. In this a given data set is classify in to specified number of clusters. Main focus is to allocate k number of centroid for each cluster. As different location causes different result, these centroids should be placed in a crafty way. The good choice is to place each centroid as much far as possible. Now

choose each point associated to a given data set and position it to the closest centroid. When all points are covered, the first step is completed. Now we need to re-calculate k new centroids as center of the clusters resulting from the last step. Once we have these k new centroids, again previous data set points and closest new centroid needs to tie together. We may notice that the k centroids keep changing their location one by one all changes are done. Its internally uses Euclidian algorithm to find the distance in case of distance between two points is too short. This algorithm focuses on minimizing an objective function, the objective function [8]

$$J = \sum_{j=1}^k \sum_{i=1}^n \|x_i^{(j)} - c_j\|^2 \quad (1)$$

Where $\|x_i^{(j)} - c_j\|$ = distance measure between a data

point X_i and the cluster center C_j , n is an indicator of the distance of the data points from their respective cluster centers.

Algorithm steps:

- 1: Assign K points considering objects which are going to clustered .These are the first points of centroid
- 2: Check the group which has closest centroid and assign object to that group for checking the latest centroid
- 3: Check all objects are considered for allocation or not, once done then every time recalculate K-centroid Position.
- 4: Repeat Steps 2 and 3 until no more changes are done. This yields a separation of the objects into groups from which the metric to be minimized can be calculated.

Consider that we have n sample feature x_1, x_2, \dots, x_n all from the same class, and we are aware of that they fall into k compact clusters, $k < n$. Assume m_i be the mean of the vectors in cluster i. If the clusters are well detached, we can use a minimum-distance classifier to separate them. We can say that x is in cluster i if $\|x - m_i\|$ is the minimum of all the k distances. This suggests the below method for finding the k means:

Make initial prediction for the means m_1, m_2, \dots, m_k

- 1: Continue till there are no changes in any mean
- 2: Use the projected means to classify the samples into clusters
- 3: For i from 1 to k
- 4: Swap m_i with the mean of all of samples for cluster i
- 5: end_for
- 6: end_until

2. Naive Bayes Algorithm:

A naive Bayes classifier consider that the absence (or presence) of a specific feature of a class is unrelated to the presence (or absence) of any other feature, given the class variable. The Naive Bayes classifier is created on Bayes' theorem with independence assumptions between predictors. A Naive Bayesian model is easy to construct, with no complex iterative parameter estimation which makes it mostly useful for very large datasets [9]

The naive Bayesian classifier Algorithm Steps:

Step 1. Consider D be a training set of tuples and their associated class labels. Every tuple is characterized by an n-dimensional attribute, $X=(x_1, x_2, \dots, x_n)$, describing n measurements made on the tuple from n attributes, respectively, A_1, A_2, \dots, A_n .

Step 2. Assume that there are m classes, C_1, C_2, \dots, C_m . Given a tuple, X, the classifier will expect that X belongs to the class having the greatest posterior probability, conditioned on X. That is, the Naïve Bayesian classifier expects that tuple x belongs to the class C_i if and only if

$$P(C_i|X) > P(C_j|X) \text{ for every } 1 \leq j \leq m \text{ and } j \neq i$$

Thus we try to maximize $P(C_i|X)$. The class C_i for which $P(C_i|X)$ is maximized is called the maximum posteriori hypothesis. By Bayes' theorem

$$P(C_i|X) = \frac{P(X|C_i)P(C_i)}{P(X)} \quad (2)$$

Step 3: Now $P(X)$ is constant for all classes, only $P(X|C_i)P(C_i)$ need be maximized. If the class prior probabilities are not known upfront, then it is commonly supposed that the classes are equally likely, that is, $P(C_1) = P(C_2) = \dots = P(C_m)$, and we would therefore maximize $P(X|C_i)$. Else, we maximize $P(X|C_i)P(C_i)$.

Step 4: Known data sets with many attributes, it would be extremely expensive to compute $P(X|C_i)$. In order to decrease computation in evaluating $P(X|C_i)$, the naïve assumption of class conditional independence is made. This supposes that the values of the attributes are conditionally independent of each another shows (3) and (4), given the class label of the tuple (i.e., that there are no dependence relationships among the attributes).

Thus,

$$P(X|C_i) = \prod_{k=1}^N P(x_k|C_i) \quad (3)$$

$$= P(X_1|C_i) \times P(X_2|C_i) \times \dots \quad (4)$$

So the predicted class label is the class C_i for which $P(X|C_i)P(C_i)$ is the maximum.

(a) Consider that if A_k is categorical, then $P(X_k|C_i)$ is the number of tuples of class C_i in D showing the value x_k for A_k , divided by $|C_i, D|$, where $|C_i, D|$ is the number of tuples of class C_i in D. Now using the values of the C_i and X_k we can check for relation between then with appropriate selection of parameters

(b) A continuous-valued attribute is typically supposed to have a Gaussian distribution with a mean μ and standard deviation σ as per (5) and (6)

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (5)$$

So

$$P(x_k|C_i) = g(x_k, \mu_{ci}, \sigma_{ci}) \quad (6)$$

Step 5: So to predict the class label of X, $P(X|C_i)P(C_i)$ is calculated for each class C_i . The classifier expects that the class label of tuple X is the class C_i if and only if

$$P(X|C_i)P(C_i) > P(X|C_j)P(C_j) \quad (7)$$

So the predicted class label is the class C_i for which we result $P(X|C_i)P(C_i)$ is the maximum.

8. Experiment

The experiment will be carried out on Lenovo Intel® Core™ i5 Processor, 6 GB RAM, Windows 7 64 bit OS. Complete Experiment will be carried out on Single node Hadoop framework. Database used is MongoDB .Planning to consider 300+ dataset from the different GPS enables phones. Data set which don't have speed attributes for them speed is calculated using the Haversine formula as shown in Table I.

Table 1: Table After Processing With Haversine Formula

Longitude	Latitude	Distance (km)	Time	Speed(km/hr)
57.45879126	22.4157892	0.05563	11:21:46	28.6162
57.45924545	22.4161793	----	11:21:53	----

This data is then inputted to location detection module to detect the exact location of the device. Now the on road and off road traffic is detected. After that filtering of the data set carried out which further provides input to the K-means clustering algorithm which results in to different clusters. Meanwhile the same dataset is processed through the Transport medium segregation module .Once transport medium is identified ,the clustered data and Transport medium segregated data is provide as input to the Naïve bayes mining algorithm . It also uses the historical data as input to do the prediction. Hadoop Distributed File System (HDFS) will be used in the experiment, so this framework can handle the huge amount of data. Also the map-reduced programming method used which help to reduce the time of execution.

9. Expected Result

The experimentation on system will give the competent outcomes. We will investigate all real-time data received from GPS and the prediction done by the framework. The outcome achieved by the framework will be plotted on graph. In this the Actuals are plotted on x-axis and Predications are plotted on Y-axis. To avoid the obscuration and to spread the data we will add the 65% jitter in figure. We are expecting the prediction up to 95%.

There is chance of miss predictions due to the few scenarios like traffic Signal and persistent change in the mode of transport that is to even 6%. So the comprehensive expected relevancy is 89%. It has been observed that usually miss prediction arise due to the traffic signal. So it is expected that Framework will predict and forecast the location where congestion is happening with maximum certainty.

10. Conclusion

In the present paper, innovative framework is proposed to identify constant traffic congestion locations adopting the data coming in from various types of GPS enabled devices

like mobiles, iPad, OBU etc. With expected accuracy up to 89%. Different methodology to detect & predict congestion is mention. Various different modules are used to isolate on road and off road traffic. Framework has adaptability to segregate traffic medium. All this features benefit to scale down the deficiency of preceding frameworks. The finest clustering and mining algorithm used to obtain maximum accuracy. The framework is flexible to different cities by changing the city-dependent thresholds. Further use of Hadoop framework provides capability to handle the traffic big data with improved performance.

11. Acknowledgment

Author would like to take this occasion to express our profound gratitude and sincere regard to my Project Guide Prof. Sonali Rangdale, for her guidance, valuable feedback and constant help throughout the duration of the project. Her valuable suggestions were of immense help throughout my project work. Her perceptive criticism kept me working to make this project in a much better way. Working under her was an extremely knowledgeable experience for me.

References

- [1] Anand Gupta, Sajal Choudhary, Shachi Paul; "DTC: A Framework to Detect Traffic Congestion by Mining versatile GPS data" In the 1st international Conference Emerging Trends & Application in Computer Science (ICETACS 01), pp.97-103, Sept 2013
- [2] M. Modsching, R. Kramer, K.T. Hagen; "Field trial on GPS Accuracy in a medium size city: The influence of built up", In the Proceedings of 3rd Workshop on Positioning, Navigation and Communication (WPNC 06), pp. 209-218, Hannover, March 16, 2006.
- [3] H. Inose, H. Fujisaki, T. Hamada; "Theory of road-traffic control based on macroscopic traffic model", Electronics Letters, Volume 3 , Issue 8, pp. 385- 386, 1967.
- [4] D. Ashbrook, T. Starner; "Learning significant locations and predicting user movement with GPS", In the Proceedings of Sixth IEEE Interna- tional Symposium on Wearable Computers (ISWC 02), pp. 101-108, Seattle, WA, 2002.
- [5] F. Lipan and A. Groza; "Mining traffic patterns from public transportation GPS data", In the Proceedings of the 6th International Conference on Intelligent Computer Communication pp. 123- 126, Cluj-Napoca, Romania, August 26 - 28, 2010.
- [6] Mandal,K."Road Traffic Congestion Monitoring and Measurement using Active RFID and GSM Technology" In Intelligent Transportation Systems (ITSC), 14th International IEEE Conference, pp. 1375 - 1379, Oct 2011
- [7] Y. H. Ho, Y. C. Wu, M. C. Chen, T.J. Wen, Y.S. Sun ; "GPS Data Based Urban Guidance", In the Proceedings of International Conference of Advances in Social Networks, pp. 703-708, Kaohsiung, Taiwan, July 25-27, 2011
- [8] Shi Na , Liu Xumin , Guan Yong , " Research on k-means Clustering Algorithm: An Improved k-means

- Clustering Algorithm” Intelligent Information Technology and Security Informatics (IITSI), Third International Symposium on April 2010
- [9] Duan Wei ,Lu Xiang-yang “Weighted Naive Bayesian Classifier Model Based on Information Gain “ Intelligent System Design and Engineering Application (ISDEA), International Conference on Oct-2010
- [10] G. Nathiya, S. C. Punitha, M. Punithavalli; ”An Analytical Study on Behavior of Clusters Using K Means, EM and K* Means Algorithm”, In International Journal of Computer Science and Information Security, Volume 7, Issue 3, 2010.
- [11] V.P. Bresfelean; ”Analysis and Predictions on Students’ Behavior Using Decision Trees in Weka Environment”, In the Proceedings of Information Technology Interfaces (ITI ’07), Cavtat / Dubrovnik, Croatia, pp.51-56, June 25-28, 2007

