# Enhanced public transport management employing AI and anonymous BT data collection

Minea Marius
*Telematics and Electronics for Transports Dept.*
University Politehnica of Bucharest
Bucharest, Romania
marius.minea@upb.ro

Dumitrescu Cătălin
*Telematics and Electronics for Transports Dept.*
University Politehnica of Bucharest
Bucharest, Romania
catalindumi@yahoo.com

*Abstract*—**The paper proposes a simple, economic and expandable solution for enhancing the data collection process used in public transport and transport demand management. A non-intrusive and anonymous method is employed to collect an estimative number of passengers in vehicles and public transport stops, along with other, relevant data. Machine learning and specific algorithms are used to improve the data collection process. No specific infrastructure equipment is required.**

*Keywords—Bluetooth data collection, machine learning, k-means algorithm, mixture density, transport demand management, MAC, RSSI*

## I. INTRODUCTION

In the present days, urban areas of highly populated cities are facing traffic congestions, emissions, stress and lack of solutions for improving road traffic experience, if the road network infrastructure does not allow for major enhancements. One of the recommended solutions, also mentioned in several official EU Transport Policies documents is to employ modal shifting, from using the private cars to the public transport [1]. Making public transport attractive, reliable and comfortable are key factors in achieving these goals. This objective can be obtained via a comprehensive set of measures for broad data collection and processing regarding the public transport management, efficiency, time schedulling and satisfying the transport demand with a high degree of accuracy. The solution proposed in this paper refers to a method for anonymously data collection employing BT/Wi-Fi enabled devices, followed by a comprehensive set of statistical filtering, machine learning algorithms and other procedures, set for:

- Obtaining an estimate regarding the evolution of the number of passengers transported along the route in a public transport vehicle;

- Obtaining an estimate regarding the transport demand – passengers in stations;

- Improved vehicle location without satellite navigation support;

- Support for traffic congestion behaviour analysis without infrastructure equipment etc.

## II. STATE OF ART AND LITERATURE SURVEY

### A. State of Art

The modern public transport management systems (PTM) employ on-board and infrastructure equipment for handling vehicle positions, regulatory actions and other specific actions. The transport demand is usually managed via specific sensors installed in stations, buses and other relevant places. Also, the locations of vehicles are collected via onboard GPS-enabled transponders. All that equipment needs a lot of maintenance, power supplying and/or should be developed on wide areas in the infrastructure. There have been tested, however, some alternative techniques to collect various information from Wi-Fi and BT enabled devices, as the number of mobile phones and accessories increases daily. Still, despite the huge potential of this methodology this technique is not yet receiving enough attention.

### B. Literature survey

Several papers in the scientific literature address this subject. In [2], the authors explore the potential of using pedestrian data for evaluation and enhancement of public transportation efficiency. They employ a Wi-Fi/BT tablet and specific software to collect relevant origin-destination information from travelers, with the purpose to improve the public transport management in terminals. In [3], N. Abedi et al. present the benefits and critical challenges on the use of Bluetooth and Wi-Fi for crowd data collection and monitoring. They introduce some new concepts, like discovery time, signal strength analysis, antenna detection range assessment and multirange scanning technique. They conclude that collecting efficient crowd data by scanning MAC addresses can be matched with other crowd data collected by other methods in order to enhance the crowd movement dynamic analysis and monitoring. They also consider that the implementation of scanning approaches in a large scale can deliver significant information from spatio-temporal dynamics of people movements. In [4], Naeim Abedi et al. present the benefits and critical challenges on the use of Bluetooth and Wi-Fi for crowd data collection and monitoring. They mention some challenges that include antenna characteristics, environment's complexity and scanning features. A. Lesani et al. present in [5] the benefits and drawbacks of employing wireless data collection techniques with Wi-Fi and BT. Also, Y. Malinovskiy [6] show the benefits of employing such technologies in public transport. Many authors conclude that this technique represents an attractive method for collecting traffic and movement data.

The remaining of this paper is organized as follows: next section presents the principle of the proposed solution and some experimental data collection, section IV concerns on the proposed algorithms and initial testing, section V the conclusions.

### III. CONCEPT AND EXPERIMENTAL DATA COLLECTION

The solution proposed here is focusing on anonymous data collection and processing for improving the public transport management, including location, transport demand and system usage information. Additional information, such as: traffic congestion, origin-destination patterns of travelers etc. is possible to be obtained via superior data filtering and post-processing.

A BT/Wi-Fi device, capable of discovering and recording MAC addresses, time/position, and RSSI levels is the single equipment needed. Data can be collected online, if the sensor is communication-enabled, or downloaded at the depots, when vehicles end their tours. The global information processing concept is shown in figure 1 below.



Fig. 1.   Data collection process and processing

The process consists in collecting BT and Wi-Fi information regarding discoverable devices (usually not phones, but devices connected to phones, such as smart watches, fit bracelets, headphones, car head units, TV sets etc.) consisting of MAC addresses, RSSI levels, time stamps and location stamps or other, relevant information (such as name of device producer, if available).

The fist filtering phase refers to establishing a dual perimeter of analysis, nominated as "inside" and "outside public transport vehicle. Based on specific RSSI levels received, the discovered MACs are categorized in these two stacks. Of course, there might be a certain number of nodes situated outside the public transport vehicle that will fit in the same perimeter as those inside, but a second phase of data analysis is designed to look at the timestamps and presence consistency of all inside nodes, to further eliminate nodes that do not have a permanent presence. It is expected that those nodes belong to MAC addresses received from pedestrians, or passengers in neighboring vehicles.

Before proceeding to effective algorithm conception, a set of initial data has been collected on selected tram and bus lines in Bucharest, Romania, in several workdays, on the same route – with a length of 1.7 km. The first testing purpose was to see if specific devices, such as TV sets, or computers, that have been discovered in fixed positions in buildings near the route, could serve as pinpoints (or "beacons") for a spatial mapping of the public transport route (shown in Figure 2).



Fig. 2.   The selected test route for collecting data

The first analysis was concerning on finding repetitive MAC addresses on the route, for different time period (days and hours). A sample of collected data is presented in Figure 3 below, where with highlighted colors and text are presented the devices that were discovered for several days in the same locations.



Fig. 3.   Sample data with repetitive devices on the test route

As it can be seen in the above figure, there are several devices, such as BT-enabled TVs, that are discovered in several tests on the same route.

Therefore, it can be assumed that these devices could serve as spatial reference points in a machine learning process, in order to develop the configuration of the route. In case of GPS location function failure, these spatial reference points could serve for a relatively precise location of the public transport vehicle.

### IV. THE PROPOSED SOLUTION AND INITIAL TESTING

#### A. General aspects of technology

Bluetooth is a wireless network standard designed for low power consumption and for communication in a limited personal area (PAN) environment. This technology was not specifically designed to locate objects, but Bluetooth enabled devices are ideal for localization because they contain a mechanism for identifying neighbouring devices and performing communication with those devices. Bluetooth access points are similar to Wi-Fi networks, but unlike them, Bluetooth pointers have a great distance from one another (typically between 10-15m).

The accuracy of the Bluetooth system ranges from 2 to15 m. One of the main advantages of the Bluetooth technology is the variable read distance. This technology is capable of reading at 1/ 10 / 50m, being suitable for locating objects. In addition, it can locate up to 7 objects in a 3m perimeter due to the master's connection capabilities. Frequency or channel jumper for device communication can take up to 10s.

However, it is not possible to use RSSI (signal strength) or the quality of the link parameters to find out the location results with a precise measuremen.

Laboratory tests have been performed intially to determine some physical characteristics of the BT radio signals. The

folowing figures show the analysis of BLE (Bluetooth Low Energy) signals, acquisition made with Aaronia Spectran HF 6065 Spectrum Analyzer.

In figure 4 is presented the spectrum of BLE signals in range 2.110 GHz – 2.169 GHz, and figure 5 shows the spectrogram of the signal. Figure 6 illustrates the histogram analysis of the BLE signals, visualizing the energy fingerprints of the communications.



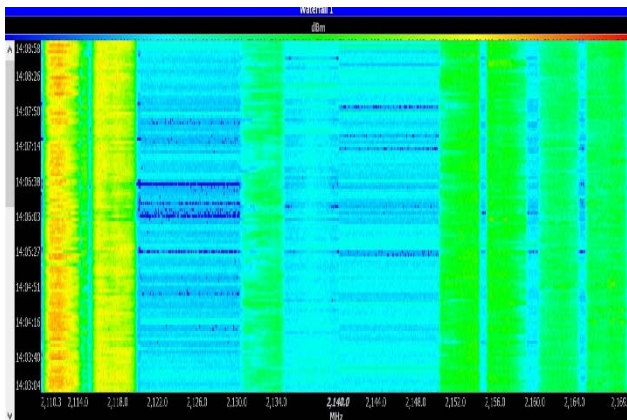Fig. 4.   Spectral representation for the reception of BLE signals.



Fig. 5.   Representing the spectrogram for the received BLE signals.



Fig. 6.   Histogram analysis of the BLE signals.

Figure 7 presents the analysis of de RSSI evolution for the received signals, highlighting the maximum (max hold – red color) and minimum (min hold – purple color) limits for channel power.



Fig. 7.   Channel power analysis for BLE reception (max hold – red color and min hold – purple color).

## B.  Description of algorithms

The proposed (cluster – type) algorithm is employed by a machine learning subsystem for performing an analysis and sorting of received BT/Wi-Fi received MACs. Goals of this approach include:

- discovering and memorizing of MAC addresses that are repeatedly found in same locations on the path of the vehicle, with the purpose to re-use them as milestones along the next route traveling;
- discovering and separating nodes that are located inside the public transport vehicle against the other received nodes; this is achieved via an analysis of RSSI parameters and near-field versus far-field thresholds established by the user;
- performing, if needed, the traceability of specific nodes (this function is used for separating travelers entering or exiting the vehicle - in public transport stops, for example); this might be helpful in achieving information regarding origin-destination patterns of travelers, or in the analysis of service levels;
- performing specific analysis on the outer nodes (in terms of determining the traffic flowing on the section of the road);
- performing a mapping of results, on a specific GIS product.

Grouping (clustering) is an action made with the purpose to partition a set of objects into different groups (clusters), where instances in a group are similar in a certain sense. Clustering is used in many fields, such as: computer learning, form recognition systems, image analysis, bioinformatics, compression, graphics etc.

Amongst other instruments employed in clustering, a proven stable algorithm for this type of application is the *k-means* algorithm. The *k-means* algorithm receives a list of points $X = \{x_i, i = 1: n\}$ as input values. Each point is *d*-dimensional $x_i = (x_{i1}, x_{i2}, ..., x_{id})$. The objective of the *k-means* algorithm is to group the points in the *k* set denoted by $S = \{Sk \mid k = 1: K\}$. The centroid representing the *k* subunit is denoted by $mk$.

Grouping of data should be done in such a way as to minimize the objective function:

$$J(X, S) = \sum_{k=1}^{K} \sum_{x \in S_k} dist(x, m_k) \qquad (1)$$

where $dist$ (.,.) is the Euclidean distance in the d-dimensional space:

$$dist(x, y) = \sqrt{\sum_{i=1}^{d} \left( x_i - y_i \right)^2} \qquad (2)$$

In order to ensure the convergence of the algorithm, more complex initialization techniques have to be applied. In this work the *k-means* method is defined based on the selection of points after a probability distribution that penalizes nearby points using a Gaussian Mixture Regression. This is to ensure a better traceability of results, because RSSI values are better modeled with a Gaussian noise when simulated propagation conditions are employed.

The *k-means* clustering algorithm is a method of determining the clusters that form specific patterns. The procedure is an unsupervised training. The *k-number* of the clusters is known, this being a set of a priori parameters.

Each cluster has a centroid. The algorithm works with $k$ clusters, so $k$ of the points used in the training will be the centers of the $k$ clusters. Since centroid initialization is randomly, there is a possibility that more runs of the algorithm lead to different results. The implementation of the algorithms has been performed employing LabView (Figure 8).



Fig. 8. Software implementation of *k-means* clustering algorithm (upside diagram) and *k-means* algorithm

Each point is associated with the cluster determined by the closest centrid. Distance between point and center can be calculated, for example, as Euclidean distance, but other variants can be chosen.
The flow of the algorithm is:

1. Randomly select $k$ points as the initial centroid.

2. Form $k$ clusters by assigning all points to the closest centroids.

3. Recalculate centroids as following: the new centroid will be the center of gravity determined by cluster points.

4. Steps 2 and 3 resumes until the centroids are no longer changed.



Fig. 9. Realizing the machine learning model using the *k-means* clustering algorithm (right – training model; left – RSSI class based on detection thresholds).
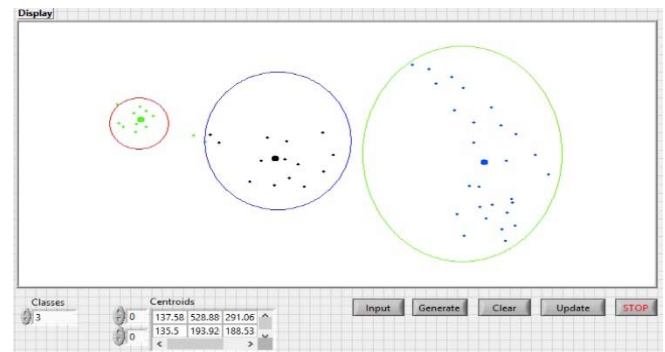


Fig. 10. Detecting RSSI by classes using the clustering algorithm (red – 1 meter; blue – 3 meters; green – 6 meters).

It is possible to correlate these models with density regression function when *Gaussian Mixture models* are employed for determining the common density of the data.

Assume the only unknowns are comprised in the mean vector $\mu_i$, i = 1, 2, ...., n. Thus, the coefficients $\theta_i$ and $\theta$ consist of the elements of $\mu_i$ and $\mu$, respectively. The mixture density is formed as the sum of Gaussian densities, that is, for each class:

$$p(x_k|w_i, \mu_i) = (2\pi)^{\frac{-d}{2}}|\Sigma i|^{\frac{-1}{2}} exp\left\{-\frac{1}{2}\left[(x_k - \mu_i)^T \Sigma_i^{-1}(x_k - \mu_i)\right]\right\}$$

Where:

- $p(x_k|w_i, \mu_i)$ is the probability of occurrence of the event $x_k$, conditioned by the $\mu_i$ conditioning vector;
- $d$ – distance from the centroid to the limit of the class;
- $x_k$ – represents the evolution of data in time.

Pre-multiplying both sides by $\Sigma_i$ yields:

$$\sum_{k=1}^{n} P(w_i|x_k, \hat{\mu})(x_k - \hat{\mu}_i) = 0 \quad i = 1,2, \dots \dots . n \quad (4)$$

Where $\hat{\mu}_i$ is the final vector, resulting from all vectors correspondent to $i = 1,2, \dots \dots . n$.

### C. Analysis of results

The results in (4) illustrates several aspects:

- $\mu_i$ is formed as a weighted summation of the $x_k$, where the weight for each sample is $P(w_i|x_k, \hat{\mu}) / \sum_{k=1}^{n} P(w_i|x_k, \hat{\mu})$ . For the sample where $P(w_i|x_k, \hat{\mu})$ is zero (or small), little is contributed to $\mu_i$ . The term $\mu_i$ may be intuitively chosen, or $w_i$ samples can be employed instead.

This aspect is shown in the following diagrams.



Fig. 11. Results obtained using a gaussian kernel (left – 2D projection; right – 3D projection).

Note that this involves updating the class means by readjustment of the weights on each sample at each iteration. This procedure is similar to the *k-means* clustering algorithm described previously.



Fig. 12. Detecting RSSI by classes using the proposed algorithm (red – 2 meters; blue – 6 meters; green – 4-5 meters).

### D. Interpretation of data leading to extraction of mobility information

Indoor positioning assessment for the development of mobility models can be achieved using several technologies. In this paper, the fingerprinting method is based on the signal strength (RSSI) from a Bluetooth network has been chosen as surveilled element.

Through fingerprinting, it is understood that the surface of the interior space where the location is desired is firstly mapped by measuring the power of the signal from the received Bluetooth nodes and creating a database that will be used later when a node is to be tracked.

Fingerprint clustering is an important step in data pre-processing in order to achieve optimal accuracy, efficiency, and data needs to be collected prior to the actual location process. Clustering of collected data was performed employing the *k-means* algorithm. The results were compared with those obtained from the traditional fingerprinting method.

Data acquisition was manually performed using a spectrum analyzer and a computer. The collection process started with the user locating his position on the floor map, and a map displayed on a computer. Then, the user crossed the entire surface of interest moving in straight lines along the surface of the enclosure. At the end of each straight trajectory, the user had once again marked his position on the map displayed on the computer. The acquisition rate of RSSI was three samples per second. The acquisition was made in two different premises.

The collected data was then divided into a set of learning and a set for testing. Fingerprint clustering can be done in two distinctive ways, using 3D fingerprint coordinates, or using RSSI.

It can be seen that in unexpected cases, the clusters are distributed over several levels of analysis, which is explained by the fact that the level of the analysis stage is less than the maximum allowed horizontal length. Positioning error was dependent on the method used for clustering.
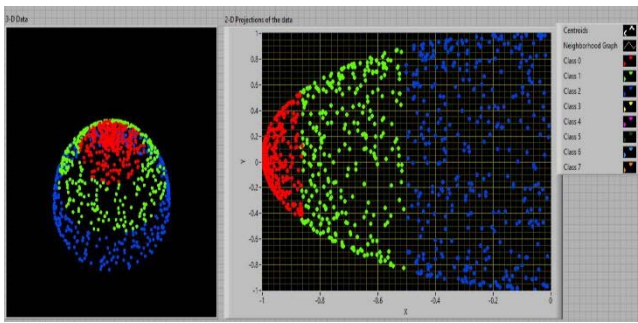
Fig. 13. Emphasizing clustering results in two ways (left – 2D projection; right – 3D projection).

Table 1 lists the errors obtained with the clustering by different methods. The best accuracy both in terms of 2D positioning and position identification was the use of 3D clustering based on *k-means* and the use of a Gaussian Mixture Regression. In addition, the results obtained show an improvement in the positioning time regardless of the method used.

TABLE 1 ASSESSMENT OF CLUSTERING METHODS

| Distance [m] | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| Average positioning error [m] | 4.0 | 3.8 | 4.3 | 4.0 | 3.5 | 3.6 |
| In class node detection probability [%] | 98 | 97 | 98 | 98 | 97 | 98 |

Figure 14 presents the results obtained for the evaluation different methods for clustering. The evaluated methods are RSSI clustering and MGD (Multivariate Gaussian Distance); RSSI clustering and GMR (Gaussian Mixture Regression); 3D clustering (*k-means*) and GMR (proposed).
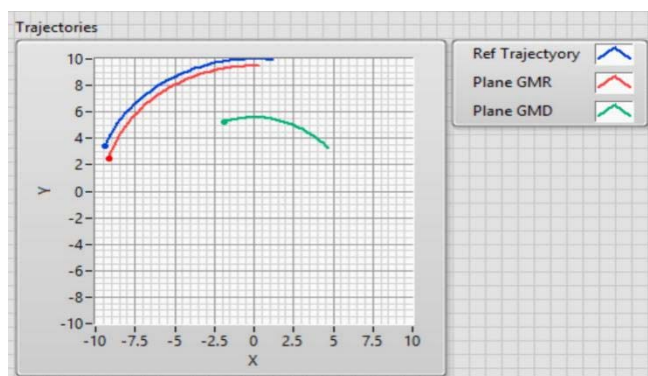


Fig. 14. Estimated trajectory versus real trajectory

The tests showed that this method could serve to collect a strong database regarding the presence and movement of different devices, on two selected "interior" and "exterior" areas of the public transport vehicle. Further field tests will concentrate on the antennas positions and patterns, in order to obtain the best setup for a correct data collection.

## V. CONCLUSION

In this paper a methodology to determine presence and movement patterns of passengers and other relevant elements related to a public transport system has been proposed. The approach is based on anonymous detection of BT (or Wi-Fi) enabled devices inside and outside a public transport vehicle repeatedly traveling on its route. Based on several statistical filtering and specific algorithms, firstly the data is sorted to discover static, repetitive nodes on the path, to re-use them in the next travels as pinpoints, or location references. Secondly, a set comprised of a *k-means* algorithm and a Gaussian Mixture model with regression are employed to perform future selection of data from the samples: discovery of vehicle inside and outside nodes, tracking and clustering of these nodes to further perform origin-destination patterns and advanced public transport efficiency analysis, such as level of service.

Several field tests have been performed to determine the frequency of static nodes detection. The tests showed a permanent presence of above 76%, depending on the testing hours and days of the week. Laboratory tests have been performed to determine the reception parameters of a BT receiver, with a spectrum analyzer, in order to shape the basic elements for the model and to analyze the reception characteristics of BT signals, in terms of RSSI time evolution.

Further, a model was developed in LabView 15, consisting in an unsupervised machine learning, employing clustering method. The obtained results showed the feasibility of the proposed method and possible future development to achieve more function and information from the collected data. The authors consider that the proposed method could be simply implemented in the public transport system with minimal investment, leading to a better transport demand and efficiency management, along with the improvement of the public transport comfort, in the benefit of the passengers. This could contribute in the future to the attractivity of this mode of transportation and drastic reducing of the personal cars' usage.

## REFERENCES

[1] * * * - "White Paper on Transport", Directorate General for Mobility and Transport, EC, ISBN 978-92-79-18270-9, European Union, 2011

[2] Neveen Shlayan, K. Ozbay, A. Kurcku – "Exploring pedestrian Bluetooth and WiFi detection at public transportation terminals", 2016 IEEE 19th International Conference on Intelligent Transportation Systems (ITSC), Windsor Oceanico Hotel, Rio de Janeiro, Brazil, November 1-4, 2016

[3] N. Abedi, A. Bhaskar, E. Chung – "Bluetooth and Wi-Fi MAC Address Based Crowd Data Collection and Monitoring: Benefits, Challenges and Enhancement". Australasian Transport Research Forum 2013 Proceedings, 2 - 4 October 2013, Brisbane, Australia

[4] N. Abedi, A. Bhaskar, E.Chung – " Bluetooth and Wi-Fi MAC Address Based Crowd Data Collection and Monitoring: Benefits, Challenges and Enhancement". 36th Australasian Transport Research Forum (ATRF), At Brisbane, Australia, 2013

[5] A. Lesani, S. Jackson, L.M. Moreno – "Towards a WIFI-Bluetooth system for traffic monitoring in different transportation facilities". Presentation, Dept. of Civil Engineering, McGill University

[6] Y. Malinovskiy, N. Saunier and Y. Wang, "Pedestrian travel analysis using static bluetooth sensors," Transportation Research Record:

Journal of the Transportation Research Board, vol.2299,  pp.137 -149,
2012