

Genetic Algorithms for Estimating Longest Path from Inherently Fuzzy Data Acquired with GPS

José Villar, Adolfo Otero, José Otero, and Luciano Sánchez *

Computer Science Department, Universidad de Oviedo, Edificio Departamental 1,
Campus de Viesques s/n Gijón (SPAIN)
otero@uniovi.es, jotero@uniovi.es, luciano@uniovi.es,
villarjose@uniovi.es

Abstract. Measuring the length of a path that a taxi must fare is an obvious task: when driving lower than certain speed threshold the fare is time dependent, but at higher speeds the length of the path is measured, and the fare depends on such measure. When passing an indoor MOT test, the taximeter is calibrated simulating a cab run, while the taxi is placed on a device equipped with four rotating steel cylinders in touch with the drive wheels. This indoor measure might be inaccurate, as the information given by the cylinders is affected by tires inflating pressure, and only straight trajectories are tested. Moreover, modern vehicles with driving aids such as ABS, ESP or TCS might have their electronics damaged in the test, since two wheels are spinning while the others are not. To surpass these problems, we have designed a small, portable GPS sensor that periodically logs the coordinates of the vehicle and computes the length of a discretionary circuit. We will show that all the legal issues with the tolerance of such a procedure (GPS data are inherently imprecise) can be overcome if genetic and fuzzy techniques are used to process and analyze the raw data.

1 Introduction

One of the tasks to be performed in the Spanish VTSS is the test and control of the taximeters in the taxicabs. This supervision must be performed every year because the taxicabs' fares are revised and published by the authorities every year. The process a taxicab owner must follow includes driving the taxicab to a specialized garage to change the fares in the taximeter. When the fares are changed, a MOT test must be done. In this MOT test, the tester engineer verifies if both the distance traveled and the waiting time fares lie between the limits imposed.

The verification of the fares can be done in two ways. The simplest way consists in doing a cab run in a previously measured circuit, manually computing the fare. More over, one person from the MOT agency must do it. One second approach is to use a machine capable of the recovering of the speed of the cab to

* This work was funded by Spanish M. of Education, under the grant TIN2005-08386-C05.

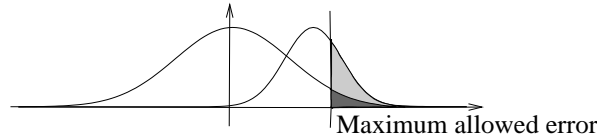


Fig. 1. If the owners of the taxis calibrated their taximeters in good faith, the density of the errors in the measures of taximeters should be centered in 0. Field measures show that the density is centered near 9% (the legal cut point is 10%). A small deviation in the tolerance of our measure, which would be unnoticed under theoretical circumstances (dark gray area,) will cause a high percentage of rejections (light gray area).

select the waiting fare or the traveled fare and to compute the time elapsed and the distance. Currently, such device is used, but fails when active safety systems nowadays present in cars trigger, moreover these systems may be damaged.

In this situation, a new method of testing taximeters must be developed. This system should be designed taking into account that it is not desirable to block one MOT test engineer when testing a taximeter. We have decided to use GPS technology to track the position of a vehicle in an actual road, and process this information on-line [13]. Moreover, the taxi driver can be sent alone to cover a distance, and no personal of MOT agency is needed, making the process cheaper.

There are some drawbacks, though. GPS generates imprecise data, and the degree of imprecision of every sample is different. The differences in tolerances must be taken into account in the algorithm that analyzes the data. The significance of this step is crucial for our system to compute the upper bound of the length of the trajectory, which must be provided in the case that a taximeter is rejected. The legal margin of error of a taximeter in Spain is 10%. We can not reject a taxi with a deviation of 7% if we can not warrant a tolerance lower than 3%, say. This could seem a minor problem, and it would be, if the density of the errors in the taxis resembled the left Gaussian in Fig. 1. Unfortunately, our study revealed that the calibration of taximeters is far from unbiased. Small changes in the tolerance produce important changes in the number of rejections. Therefore, it is needed a procedure to determine the bounds of the measure with high accuracy and it is also needed that all the tolerance errors benefit the owner of the taxi. In other words, we need to compute the lowest upper bound (LUB) of the trajectories compatible with the (imprecise) GPS measures.

In this paper we will explain a new method for estimating the LUB of the trajectory from imprecise data. Through multiobjective genetic algorithms, the measures are filtered to obtain the smallest set of samples that define a multi polygonal covering the input data. The LUB of the path is found by means of a deterministic algorithm that processes this multi polygonal.

The structure of this work follows: In next section, how GPS measures are obtained is detailed. Then, a description of the proposal is done in Sect. 3. The genetic algorithms are detailed in Sect. 3.1, while the deterministic algorithm for estimating the maximum length is detailed in Sect. 3.2. In Sect. 4 experiment and results are shown. Finally, conclusions and future work are presented.

2 GPS-based measures are fuzzy data

The term Global Positioning System (GPS) refers to a set of devices (satellites and receiver) working together to get a fix (the position) of the receiver. The receiver can get some signals from the satellites and compute a set of measures: longitude, latitude, altitude, number of satellites in use, time, etc. Each signal received from a satellite contains information about the time that the signal lasts from the satellite to the receiver.

The higher the number of satellites, the better the accuracy. But even with a high number of satellites in use (12 to 16) the geometry or constellation of the satellites must be taken into account to estimate the fix accuracy. This is done using DOP (Dilution of Precision), a measure of the probability of the effects of the constellation on the fix accuracy; a higher value of DOP indicates a weaker geometry of satellites. In the case of GPS longitude and latitude accuracy, the HDOP (latitude and longitude DOP) value must be taken into account. Related with HDOP is the CEP (Circular Error Probable), a given value of CEP at probability P means that the receiver is inside a circle of radius CEP, centered at the measured fix with that probability. When using consumer-grade receivers, it is very common to obtain accuracies like 3 meter CEP (50%) and 7 meters (90%). Given the number of satellites n used for the measure and an accuracy probability P , the CEP is computed by means of equation Eq. 1. Constants A, B, C and D are device dependant [16].

$$CEP = \left(-\left(A \cdot \left(\frac{C}{n^2} + D \right) \right)^2 + B^2 \right) \cdot \ln(1 - P(Err \leq CEP | HDOP))^{0.5} \quad (1)$$

2.1 Fuzzy interpretation of GPS-values

Under the imprecise probabilities framework, it makes sense to understand a fuzzy set as a set of tolerances, each one of them is assigned a confidence degree, being the lower degree the narrower tolerance [9]. In particular, it has stated that, given an incomplete set of confidence intervals for a random variable, we can build a fuzzy random variable, whose α -cuts are confidence intervals with degree $1 - \alpha$, that contains all the information we know about the unknown random variable [4]. In our case, the GPS sensor provides two confidence intervals at 50% and 90% (the mentioned circle of radius CEP,) and therefore the fuzzy representation of GPS coordinates is immediate.

3 Determining the length of trajectories using fuzzy data

GPS data is recorded at regular time intervals. Each sample is a fuzzy set, as mentioned, whose α -cuts are circles. In turn, every circle is a confidence interval for the coordinates of the taxicab at that moment. It is remarked that taking the centers of these circles is not a valid estimation. We need to compute the LUB of the paths whose extremes are contained in the circles, and this length will always be higher than the value obtained from the centers.

The answer to the problem is not easy, though. If we try to compute the maximum length of all compatible piecewise linear paths that are contained in the circles it is obvious that, the shorter the sampling period, the longer the estimation. This is not correct, and we wish the estimation of the length not to be too influenced by the sampling period [12]. We have decided to process the fuzzy data and remove all redundant information with the help of a genetic algorithm, as we will show in the section that follows.

When using crisp data, the geometric problem of simplifying polygonal lines has been studied in [7]. The most similar approach to ours, up to our best knowledge, uses fuzzy data from a geographical database for reconstruction of 3D images by means of B-splines[1], where a fuzzy point is said to be covered by the fuzzy B-spline if the fuzzy set induced by the latter completely contains the former, we use this concept next.

3.1 Multiobjective fuzzy fitness genetic algorithm for filtering the fuzzy input data

The fuzzy GPS measures are filtered using a multiobjective genetic algorithm. The output is the minimum set of fuzzy input data that defines a fuzzy trajectory covering as many points as possible. Using those fuzzy points, and for each α -cut, a distance value is computed by means of a deterministic algorithm, which will be detailed later.

Every candidate solution is evaluated as follows: we first build a polygonal chain for each α -cut of the selected data, using the tangent surfaces to the selected fuzzy data set ¹. We wish that this chain contain as many data as possible, while having the minimum area.

Both objectives are fuzzy numbers and define a multicriteria problem [3], and two different approaches had been used for solving the problem. The first one is using the NSGA-II algorithm [5,6]. The second approach is using the multi-objective genetic operators simulated annealing (MOSA) [14]. Further details of those algorithms follow.

Coding of individuals Each individual is a boolean vector, marking the corresponding fuzzy input data to be or not part of the hypothesis: those marked with true are used to define the polygonal chain. To generate an individual, a probability value p is given, and for each fuzzy point in the vector of input fuzzy data, it is included in the hypothesis with probability less or equal than p . The origin and the end of the ride must be always included.

Genetic operators The definitions of crossover and mutation must reduce the number of vertexes in the population, and therefore they are unbiased.

¹ This chain might include some extra points not covered by the input data, but this fact always would benefit the taxi, thus it is legally correct.

Given two parents A and B , the offspring are two new chains C and D such that $A \cap B \subseteq C$ and $A \cap B \subset D$; a vertex $v \in A - B$ has a probability p^+ of being in C , and a vertex in $B - A$ has a probability p^- of being in C , where p^- is much lower than p^+ . The chain D is built the same way. Mutation is defined as the random removing of a point of the chain, different from the first or last one. The operation named *toggle* is very similar to mutation, but it can alter the state of inclusion in the hypothesis of a randomly selected fuzzy data. Toggle is used as genetic operation for MOSA. When generating a neighborhood of current individual a random number of toggle operations are done. The number of operations is temperature dependent, and so the neighborhood of new individual, as well.

Multiobjective fuzzy fitness As stated before, two criteria are to be reached: the minimization of polygonal chain area and the maximization of the percentage of data covered. Both of them are fuzzy numbers. This means that it is needed an operator *less than* and an operator *less or equal than*, both defined for fuzzy numbers, so dominance could be evaluated. Some work has been done in evaluating Pareto dominance with fuzzy fitness. In [17] the Pareto dominance concept is extended to fuzzy dominance, and different levels of α -cut are used for each decision making process, using the concept known as α -dominance. In [11] it is proposed a fuzzy rule to determine the *degree of dominance of x over y* , and another fuzzy rule to determine the *degree of been dominated of x by y* . Then, aggregating those rules by means of the max t-conorm, a crisp rank of dominance is obtained for each individual x . In [10] a totally different approach is used. It defines a comparison between fuzzy numbers, so Pareto dominance could be used as stated in its definition. In [8] a generalization of the Pareto dominance concept is proposed. In that work, instead of using especial operators *less than* and *less or equal than*, fuzzy Pareto dominance is defined so the result of such redefinition is that decision surface is obtained. For the purposes of this work, α -dominance approach is used.

3.2 Deterministic longest path estimation

Once the data is preprocessed by means of genetic algorithms, LUB is computed. For each α -cut of the fuzzy b-splines that contains the taxi trajectory, we get a polygonal set constructed with trapezoids, as it can be seen in right side of Fig. 2. The motion direction is indicated by the thin dashed arrow. Each trapezoid vertex is denoted with a pair of integers, those at the left of the arrow have zero at first and those at the right have one at first. The other number is the step in motion sequence. The longest path at each step i goes through $(0, i)$ vertex or $(1, i)$. The set of vertexes that defines the longest path, can be computed by exhaustive exploration of all possible combinations, but this is very expensive in terms of computational cost and proved impracticable in a realistic trajectory with 700 points, for instance. This problem has been studied in the area of Computational Geometry and is related with Longest Path with

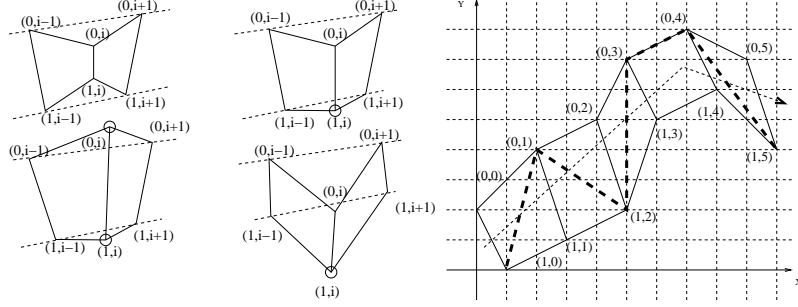


Fig. 2. Left: Possible relative positions of vertex and lines between prior and next vertex. Right: Example of longest path estimation.

Forbidden Pairs [2], that is NPO PB-complete. Because of this and given that in a realistic trajectory the changes of direction and the changes in distance between left and right vertex are limited due to the dynamics of the taxi, the geometry of the road and GPS behavior, we use a heuristic that is lineal in time with the number of vertex. The heuristic is based in the selection of convex vertexes: when a vehicle turns, the longest path goes through the exterior of the trajectory curvature. The convexity of a vertex is analyzed using the straight lines that rely on previous and next vertexes, the possible relative positions of the central vertex can be seen in right side of Fig. 2, where convex vertex are marked with a small circle and the lines that pass through vertex $(0, i - 1)$, $(0, i + 1)$ and $(1, i - 1)$, $(1, i + 1)$ are drawn. From left to right and up to bottom, if both vertexes lie between the lines, both are concave. If only one is outside of the lines, it must be convex. If both are out of the lines, either both are convex (left) or one is concave and the other one convex. In both cases, if the farthest one from the nearest line is chosen, then it is convex.

The heuristic is as follows: the first segment of the longest path goes from a convex vertex in step 1 to the vertex at step 0 that gives the maximum segment length. From vertex 1 to the one before the last, the path goes through this vertex if there is only a convex vertex, through the farthest one if there are two convex vertexes or there is not any convex vertex. Last segment ends in the farthest vertex from the previous one. In right side of Fig. 2 the path computed with this heuristic is marked with a thick dashed line. The first segment goes from $(1, 0)$ to $(0, 1)$ because $(0, 1)$ is convex and the distance to $(0, 0)$ is shorter. Then the longest path continues to $(1, 2)$ because is the only convex. The same situation happens with $(0, 3)$ and $(0, 4)$. Finally, the path ends in $(1, 5)$ because it is farther from $(0, 4)$ than $(0, 5)$.

4 Experiments and results

In the experiments presented here, the parameters of the NSGA-II algorithm are: 4000 generations, 15 individuals in the population, 0.1 and 0.7 of mutation

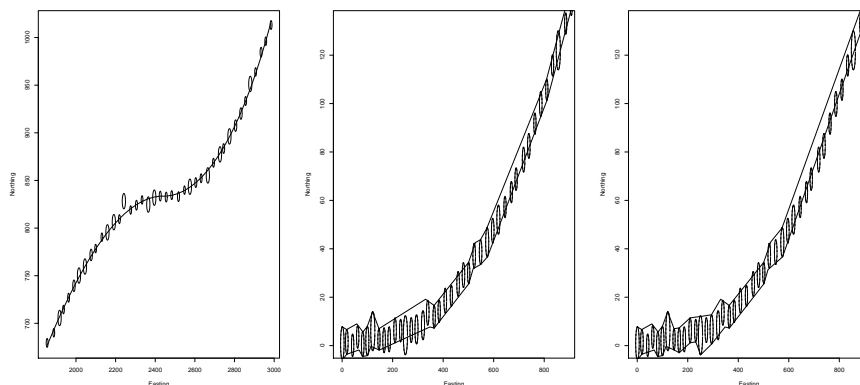


Fig. 3. Left: Example of GPS generated data along with the real trajectory. Center: Part of the first trajectory simplified by NSGA-II algorithm. Right: Same data simplified by MOSA.

and crossover probabilities, $p^+ = 0.7$ and $p^- = 0.01$. Each individual must cover a minimum of 85 percent of input data to be included in the Pareto front. When using MOSA, delta is $1/4000$, T_0 is 1.0 and T_1 is 0.0, while the rest of parameters are the same to those of NSGA-II. We have decided to evaluate our algorithm in a realistic path that covers the situations usually found when the MOT test of a taxi is done, and computing HDOP, CEP, and projecting earth measures adequately [16, 15].

The trajectory is sampled each second, obtaining 1000 points, the total length of the trajectory is 21273.21 meters. At each location, we take a random number from 4 to 9 as the number of available satellites, that we found representative for real data. From this data, we build a dataset of GPS measures, sampled at each second. Each measurement is simulated using the following procedure, with a probability of 0.95, a point is selected that is closer in distance to the real one less than the CEP at that probability. With 0.05 probability the point is selected further than the corresponding CEP from the original data. This resembles the uncertainty that occurs using GPS, and the obtained data can be used to test how tight the bounds obtained with our algorithm are. The reader must remember that the goal is to obtain a multi polygonal chain that covers most of the GPS fixes with minimum number of vertexes and with the minimum area. In left side of Fig. 3 is shown part of the generated data. GPS measures are represented with circles (actually ellipsoids due to scaling issues) with radius equal to 95% CEP and the original trajectory with a continuous line. As it can be seen, most of the circles intersect the trajectory, that is, most of the points of the real trajectory (in fact 95%) are inside the circles with CEP radius, centered in GPS fixes.

We perform two experiments with two subset of the complete dataset with 120 points each. The true length of the first trajectory is 3228.574 meters. The

estimated length of the longest path compatible with the 85 % of the points of the first processed trajectory polygonal chain using NSGA-II is 3471.75, and 3555.34 using MOSA. If the taximeter reports a distance longer more than 10% than this upper bound, it should be rejected because even in the worst case the taximeter is out of tolerance. The distance through the GPS fixes is 3238.521, that is much closer to the real data, but the taxi owner can argue about the uncertainty of the procedure saying that it is inaccurate, if we compute an upper bound of the length compatible with GPS data there is no chance for this.

The length of the second trajectory is 2741.306 meters. The estimated length of the longest path compatible with the 85 % of the points of the corresponding processed trajectory polygonal chain using NSGA-II is 3059.1, while using MOSA is 3130.48. In this case the bound is less tight since the trajectory has stronger turns and this leads to longest path compatible with the data. In center and right of Fig. 3 can be see how the simplification of the trajectory works with NSGA-II and MOSA algorithms showing part of the data from the first trajectory. The data correspond to the individuals with less total length. Both algorithms cover most of the data, but differ in which data must be preserved.

5 Conclusions and future work

During the development of this application we found that if we report directly the data obtained with GPS equipment, there were legality issues about the uncertainty of the measures. Taxi owners could easily gain in courts any reclamation where the uncertainty of the GPS measures were revealed. As result, the upper bound of the trajectory length compatible with GPS data is computed. In this way there is no doubt to reject a taximeter with reported length above of this measure. Additionally, this alternative is less restrictive with the real data given the biased error detected in the taximeters. In the experiments, MOSA has shown to be almost as accurate as NSGA-II but much more faster. We have found that our algorithm performs worst when the trajectory includes more and stronger turns, this issue must be solved in future modifications with an additional heuristic that includes the dynamic behavior of a real driver using the time information in GPS measures.

Future work includes also using different fuzzy dominance approaches that should be tested to better fit the longest path better.

References

1. A.M. Anile, B. Falcidieno, G. Gallo, M. Spagnuolo and S. Spinello, *Modeling uncertain data with fuzzy B-splines*, Fuzzy Sets and Systems 113, 397–410, 2000.
2. Berman, P., and Schnitger, G. (1992), “On the complexity of approximating the independent set problem”, Inform. and Comput. 96, 77-94.
3. C. A. Coello, *An Updated Survey of Evolutionary Multiobjective Optimization Techniques : State of the Art and Future Trends*, 1999 Congress on Evolutionary Computation, IEEE Service Center, 1999.

4. I. Couso, S. Montes, P. Gil The necessity of the strong alpha-cuts of a fuzzy set International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems 9-2, 249-262, 2001
5. K. Deb, S. Agrawal, A. Pratab, and T. Meyarivan, *A Fast Elitist Non-Dominated Sorting Genetic Algorithm for Multi-Objective Optimization: NSGA-II*, In Marc Schoenauer, Kalyanmoy Deb, Günter Rudolph, Xin Yao, Evelyne Lutton, Juan Julian Merelo, and Hans-Paul Schwefel, editors, Proceedings of the Parallel Problem Solving from Nature VI Conference, 849-858, Springer. Lecture Notes in Computer Science, 2000.
6. K. Deb and T. Goel, *Controlled Elitist Non-dominated Sorting Genetic Algorithms for Better Convergence*, In Eckart Zitzler, Kalyanmoy Deb, Lothar Thiele, Carlos A. Coello Coello, and David Corne, editors, First International Conference on Evolutionary Multi-Criterion Optimization, 67-81. Springer-Verlag. Lecture Notes in Computer Science No. 1993, 2001.
7. R. Estkowski and J. S. B. Mitchell, *Simplifying a polygonal subdivision while keeping it simple*, SCG '01: Proceedings of the seventeenth annual symposium on Computational geometry, ISBN 1-58113-357-X, 40-49, ACM Press, New York, NY, USA, 2002.
8. M. Farina and P. Amato, *Fuzzy Optimality and Evolutionary Multiobjective Optimization*, in Carlos M. Fonseca, Peter J. Fleming, Eckart Zitzler, Kalyanmoy Deb and Lothar Thiele (editors), Evolutionary Multi-Criterion Optimization. Second International Conference, EMO 2003, pp. 58-72, Springer. Lecture Notes in Computer Science. Volume 2632, Faro, Portugal, April 2003.
9. Goodman, Nguyen. *Uncertainty Models for Knowledge-based Systems*. North-Holland. 1985
10. M. Hapke, A. Jaszkievicz and R. Slowinski, *Pareto Simulated Annealing for Fuzzy Multi-Objective Combinatorial Optimization*, Journal of Heuristics, 6(3), 329-345, August 2000.
11. M. Köppen, K. Franke and B. Nickolay, *Fuzzy-Pareto Dominance Driven Multiobjective Genetic Algorithm*, In Proceedings of the 10th IFSAWorld Congress (IFSA 2003), pages 450-453, Istanbul, Turkey, June, 2003.
12. N. Meratnia and R. A. de By, *Trajectory representation in location-based services : problems and solutions*, in Proceedings of the 3rd IEEE Workshop on Web and Wireless Geographical Systems (W2GIS 2003) in conjunction with the Fourth International Conference on Web Information Systems Engineering (WISE), Rome, Italy, 2003.
13. A. Otero and J. Otero and L. Sánchez and J. R. Villar, , *Longest path estimation from inherently fuzzy data acquired with GPS using genetic algorithms*, 2nd International Symposium on Evolving Fuzzy Systems, University of Lancaster, UK, 2006
14. L. Sánchez, J. Otero and J. R. Villar, *Boosting of fuzzy models for high-dimensional imprecise datasets*, in Proceedings of the Information Processing and Management of Uncertainty in Knowledge-Based Systems IPMU06, Paris, 2006.
15. Snyder, J. P., "Map Projections Used by the U. S. Geological Survey", 2nd edition, Geol. Survey Bulletin 1532, 313 p., U. S. Government Printing Office, Washington, D. C., 1982.
16. Wilson, D., "David L. Wilson's GPS Accuracy Web Page", <http://users.erols.com/dlwilson/gps.html>.
17. J. Zhang, B. Pham, and P. Chen, *Fuzzy Genetic Algorithms Based on Level Interval Algorithm*, In Kazmierczak, E, Eds. Proceedings The 10th IEEE International Conference on Fuzzy Systems, pages pp. 1424-1427, Melbourne, Australia, 2001.