

# Monitoring Moving Objects using Uncertain Web Data

Mouhamadou Lamine Ba  
Talel Abdessalem

Sébastien Montenez  
Pierre Senellart

Institut Mines–Télécom; Télécom ParisTech; CNRS LTCI  
Paris, France  
first.last@telecom-paristech.fr

## ABSTRACT

A number of applications deal with monitoring *moving objects*: cars, aircrafts, ships, persons, etc. Traditionally, this requires capturing data from sensor networks, image or video analysis, or using other application-specific resources. We show in this demonstration paper how Web content can be exploited instead to gather information (trajectories, metadata) about moving objects. As this content is marred with uncertainty and inconsistency, we develop a methodology for estimating uncertainty and filtering the resulting data. We present as an application a demonstration of a system that constructs trajectories of ships from social networking data, presenting to a user inferred trajectories, meta-information, as well as uncertainty levels on extracted information and trustworthiness of data providers.

## Categories and Subject Descriptors

H.2.8 [Database Applications]: Spatial databases and GIS; H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*Retrieval models, Information filtering, Relevance feedback*

## Keywords

Moving objects, social Web data, trajectory, uncertainty estimation

## 1. INTRODUCTION

*Moving objects and the Web.* Consider the problem of tracking real-world *objects* such as cars, trains, aircrafts, ships, persons (e.g., celebrities), or, more broadly, populations or groups of humans, natural phenomena such as cyclones, epidemics. Such moving objects are characterized by timestamped location data and other meta-information such as name, size, maximum reachable speed, acceleration patterns, etc. The analysis and mining of spatio-temporal information about moving objects is common in a variety of applications, e.g., for pattern discovery [3, 6, 8, 12] or prediction of trajectories and locations [2, 7]. The overall goal may be to better understand certain natural phenomena, to improve city services, to regulate route traffic, etc. Currently used methods for tracking moving objects are often complex, mostly rely on application-specific resources and costly equipment (e.g., satellite or radar tracking of ships and aircrafts).

The World Wide Web, on the other hand, is a huge source of public information about various real-world moving objects. Timestamped

geographical data about the position of moving objects are disseminated on the Web, notably on location-based social networks or media sharing platforms. Social networking sites like Twitter and Facebook have the ability of recording the real-time location of the user posting a message, thanks to data from the GPS system, or mobile or wireless networks. In addition, these messages may also contain location information as free text. Thus, it is theoretically possible to obtain information about the user herself, or any moving object that the user is referring to in her message. Media on sharing platforms like Flickr and Instagram may be annotated with automatically acquired spatio-temporal information, such as the timestamp and location of a picture as added by modern digital cameras.

In addition, the Web also provides in a variety of online databases more general information about moving objects. For instance, data such as the usual residence of a given individual can often be found online, e.g., in Wikipedia or Yellow Pages services. Characteristics of particular flights and ships are available on specialized Web platforms. In this demonstration paper, we illustrate how to extrapolate on the information extracted from multiple Web sources in order to infer the locations of certain moving objects at given times, and to obtain general information about these. We then visualize these locations, together with hypothetical trajectories, on a map-based representation. This information is uncertain, however, and exhibits many inconsistencies. One of the challenges to overcome is to estimate the inherent reliability of that information.

We claim Web information can be used in a variety of settings where one would like to track moving objects. We illustrate with the applications of celebrity spotting and ship monitoring; the sources, data, and scenario of the demonstration will focus on the latter.

*Celebrity spotting.* Journalists, paparazzi, fans, detectives, intelligence services, are routinely interested in gathering data and following the travels of given individuals, usually celebrities. These individuals may be active in social networks such as Twitter and Instagram, where they share geolocated data (or data tagged with location information); they may also be spotted in real life by users of media sharing platforms, who will upload geolocated pictures. Various Web sites, news articles, etc., provide additional meta-information. Exploiting this mass of information would provide a cost-effective and legal manner to reconstruct a meaningful trajectory.

*Ship monitoring.* Researchers in maritime traffic investigate the routes followed by different kinds of ships to propose traffic optimization methods, to predict pollution levels, or to prevent pirating actions. Though ships do broadcast information about their position using the AIS (Automatic Identification System) [10], this information is not made publicly available, and no historical log is kept. Information about the timestamped location of cruise ships, military vessels, tankers, etc., is common on Web sharing platforms such as Flickr and on other specialized Web sources (see Section 4).

Integrating ship data from multiple Web sources also helps obtaining more complete and certain information about their characteristics.

**Related work and contribution.** Discovering user routines based on geographical data from social networks has been studied in [8], focusing on a specific type of moving object, and does not consider visualization aspects; in addition, the authors do not deal with uncertainty of used information. Web information is uncertain because of imprecise and incomplete data. Moreover, according to [11] location data are inherently uncertain since one is never sure whether those locations are approximate or really precise. In this paper, we extract data about moving objects through keyword search over a set of Web sources. We estimate the amount of uncertainty in each location for all kinds of moving object based on two main criteria: *outliers* and *far-fetched trajectories*. We introduce another criterion, namely *on-land locations*, pertaining to the specific maritime traffic demonstration application. For each non-geographical piece of information, we consider and integrate multiple possible values from different Web sources. The most probable attribute value can be computed using truth discovering algorithms [1, 5].

**Outline.** We present in Section 2 our Web extraction approach. Section 3 describes a method for evaluating the precision of obtained locations, trustworthiness of users and for integrating uncertain attribute values. Section 4 introduces the maritime traffic application and our implementation, while Section 5 details the demonstration scenario. **A video accompanying this demonstration paper is available at <http://dbweb.enst.fr/ships.mpg>.**

## 2. DATA EXTRACTION

We distinguish two types of Web information about a moving object: *location data* and *general information*. We extract object information from Web sources through keyword search. That is, we suppose given the name of the moving object (a key phrase) and crawl data we obtain from a set of Web sources (see Section 4 for the specific sources used for ship monitoring). We focus on location-based platforms and social networks, providing geolocated data items, for object locations.

**Gathering general information.** We collect general information about moving objects based on a *supervised* extraction over a fixed set of Web sources. The main intuition is that for many moving objects, e.g., ships, general information provided by a number of Web sources is structured into Web templates. This is particularly true for domain-specific resources. Inside this template, each particular characteristic has a meaningful label, with a value associated to it. We implement, based on such observation, an extraction process over these Web sources by using source-specific functionality for keyword search, and then crawling and parsing obtained HTML pages. Through hand-written schema mapping rules, we return data items in a global schema as a collection of attributes and corresponding values. There may be some conflicting attribute values from different sources.

**Location extraction.** We extract locations of moving objects by searching Web data items such as pictures, posts, and tweets that have geographical information attached to them (either directly as semantic geolocation information, or as can be extracted by a gazetteer on tags and free text). This type of Web data can be found on the majority of popular location-based networks and social Web platforms like Flickr, Instagram, Facebook, etc. A geolocated Web data item comes with geographical data (latitude and longitude), a date, and additional meta-information such as title, description, set of tags, user name, etc. As an example, picture geolocation is sometimes available as *Exif* tags, automatically recorded by a digital camera at the time the picture was captured. The extraction

proceeds as follows. Given a key phrase and a set of social Web sources, we first look for relevant geolocated data items regarding the input keyword. Then, for each data item we extract geographical data, dates and meta-information as sketched above.

## 3. UNCERTAINTY ESTIMATION

We evaluate in this section the amount of uncertainty in moving object data extracted from the Web. We first estimate the precision of spatio-temporal data according to three criteria and compute from that a trust score for each item provider. Then, we sketch our integration of general information from different sources.

### 3.1 Precision of Spatio-temporal Information

As already mentioned, we harness geolocated Web data items for computing the different locations, thereby hypothetical trajectories, of moving objects. Geographical information associated to their geolocated data items come with imprecisions, however. First, for various reasons a keyword extraction approach is imprecise. As a result, the search may return wrong or irrelevant results regarding the moving object of interest. For instance, when searching geolocated Web data items about a moving object  $O$ , one can get from a given Web source results related to another type of real-world object, e.g., a street with a similar name. Second, even if the results obtained really describe the object  $O$ , either the timestamp or the spatial information may be wrong (because of poorly configured software, purportedly introduced errors, ambiguous location names for gazetteers, etc.). We need an automated manner to detect these potential errors, and estimate the uncertainty of the data.

As a general framework, we estimate the precision of locations related to any  $O$  against two criteria. First, we detect outliers, that is, isolated locations, which represent locations with high probabilities to be impossible compared to other ones, that form a more consistent set. Second, we evaluate the amount of imprecision in geographical data by analyzing whether two successive locations in a chronological sense form a realizable trajectory of  $O$  with respect to its maximum speed. For the purpose of our demonstration application, dealing with ships, we also consider a third criterion which determines whether a location is in (or near) a water area. We next explain how we measure precision in each case.

Let  $I_1, \dots, I_j$  be a chronological sequence of distinct geolocated Web data items about the specific moving object  $O$ . We use the simple point-location model of [11] and represent formally a specific location  $I_j$  of the moving object  $O$  as a couple  $(gd(I_j), dat(I_j))$  where  $gd(I_j)$  is geographical coordinates (latitude and longitude) and  $dat(I_j)$  is a date. Fix two locations  $(gd(I_i), dat(I_i))$  and  $(gd(I_j), dat(I_j))$ . Necessarily,  $i \leq j$  if and only if  $dat(I_i) \leq dat(I_j)$ . Given the roundness of the Earth, the distance  $d_{ij}$  between these two locations of  $O$  is computed via the Haversine formula [9].

**Detecting possible outliers.** An outlier is a location far away from a set of locations, that are all within a given time interval and a maximum distance. Fix a given time window (say, one week for the maritime traffic application). We say that a point is an outlier if it falls in the middle of an interval where it is at distance  $Kd$  of the centroid of all other points of the interval, where  $d$  is the maximum distance between these points and  $K > 1$  is an application-specific constant, e.g.,  $K = 5$ .

**Far-fetched trajectories regarding reference speed.** A possible itinerary of the moving object  $O$  is a connected set of chronologically ordered locations. A trajectory may be far-fetched, i.e., unreasonable, if reaching one location from a previous one is impossible when we consider the reference speed of  $O$ . Let  $V$  be the reference maximum speed of  $O$  induced from gathered general information. Given two consecutive locations  $(gd(I_i), dat(I_i))$  and

$(gd(I_j), dat(I_j))$  with  $j = i + 1$ , we verify whether the following inequality holds:  $d_{ij} \leq V \times (dat(I_j) - dat(I_i))$ . If not, we are in the presence of an impossible trajectory. At least one of these two locations is wrong. As we do not know in advance which one, both are marked as potentially uncertain.

**On-land locations.** This measure pertains to our maritime traffic application in which we are mostly interested in locations falling in water areas. Other applications have similar application-specific ways of detecting impossible points.

A location on land is defined as a point that is out of water areas. All the water regions on Earth can be found on the Web platform Natural Earth<sup>1</sup> in the form of *multi-polygons* for lakes, seas and oceans, and *polylines* for rivers. Based on these shapes and the *ray-casting algorithm*, we check whether a given location  $(gd(I_j), dat(I_j))$  of the moving object  $O$  (here typically a ship) falls within one of the considered polygons or polylines. We estimate all on-land locations for  $O$  in this manner. Observe that some of these locations on land could be relevant for our application. In particular, locations on land, e.g., ports, that are close to water areas. To account for those kinds of interesting locations on land, we introduce a *tolerance factor* by considering the disc with radius of  $x$  and centered on a location. In the demonstration application, we set  $x$  to 0.1 degree of latitude/longitude. Given that, we find on-land locations whose disc, w.r.t. the tolerance factor, intersects one of the polygons or polylines of a water area. The overall set of on-land locations are finally those that do not satisfy this condition.

### 3.2 Computing User Trust Score

Uncertainty about moving objects in social Web can be inherent to the users instead of used imperfect sensors or an error-prone extraction process. In fact, some users can be *out-of scope* while others may intentionally spread false information in order to pollute social platforms. Consequently, one would want to have an indication about the trustworthiness of users sharing items: a trustworthy provider is more likely to share relevant items than an untrusted one. Given a snapshot of Web items about a particular moving object, we estimate hereafter the trustworthiness of their associated users based on the precision level of those items.

Consider again the set  $I_1, \dots, I_j$  of distinct geo-located items about  $O$ . We also suppose known the users which shared these items. Each user can provide more than one item about the same moving object. Suppose a particular user  $U$  posting a subset of items  $O(U)$  of  $O$ . We put items in  $O(U)$  into two groups w.r.t. the precision level of their locations, as previously determined: *less precise items* (outliers, participating in a far-fetched trajectory, on-land location in our application) and *more precise items* (items whose locations do not exhibit anomalies w.r.t. precision measures).

Even though the more precise items seem to be good candidate to reliably describe the mobility of the considered object, we cannot be certain of their correctness. Instead, we assign more precise items with a probability value of  $\alpha$ . In the same spirit, less precise items may still be correct. Their level of imprecision varies, however, in function of their specific types and can be application-driven. For instance, it seems reasonable to give less chance of being correct to on-land items than outliers on water areas for our maritime application. We define our belief about the level of relevancy of on-land items, outlier items, and those implied in far-fetched trajectories as  $\beta$ ,  $\gamma$  and  $\theta$  respectively. We constrain  $\beta$  to be lower than  $\gamma$  and  $\theta$ . In addition, the probability of relevance of any less precise item cannot be greater than that of any more precise item.

Let us fix the set of more precise items  $MP(I)$  and the set of less precise items  $LP(I)$  in  $O(U)$  with  $MP(I) \cap LP(I) = \emptyset$ . Consider

<sup>1</sup><http://www.naturalearthdata.com/>

the set of on-land items  $Land(I)$ , outliers  $Out(I)$ , and far-fetched items  $Far(I)$  in  $LP(I)$  with pair-wise intersections not necessary empty. We define the level of trustworthiness  $Trust(U)$  of user  $U$  as the average of the probabilities of relevance of its provided items. That is, we set:

$$Trust(U) = \frac{\alpha \times |MP(I)| + \beta \times |Land(I)| + \gamma \times |Out(I)| + \theta \times |Far(I)|}{|MP(I)| + |Land(I)| + |Out(I)| + |Far(I)|}.$$

The trust score of each user can be recomputed when external feedbacks or knowledge are available. This can help lower or increase the probability of relevance of some items.

### 3.3 Integrating Uncertain Attribute Values

We do not only extract from the Web geographical information. We also collect general information about the moving object  $O$  in the form of attributes and corresponding values. Attributes are distinguished by meaningful labels specific to the individual sources. In general, Web sources have different level of completeness in terms of the data they provide. In addition, some of them can provide conflicting information, i.e., there can be multiple possible values for a given attribute, coming from different Web sources.

As general information comes from multiple sources, we need to integrate them in order to provide to the user a unique global view. In this integration process, we have to deal with the uncertainty that is inherent to the Web, but also that results from contradictions. We integrate general information about  $O$  from multiple Web sources by first matching values of the same attribute provided by distinct sources, using a manually constructed schema mapping across sources, and then by merging identical values. When a conflict occurs, we consider the value provided by the majority of sources as the most reliable one, but we keep all different values, as will be clear in the demonstration. This process for choosing the most probable values among conflicting ones corresponds to a voting approach. More elaborate voting strategies can be used, such as those given in [5].

## 4. MARITIME TRAFFIC APPLICATION

*Use case.* The use case of our demonstration is the monitoring of ships. We rely on Flickr for collecting a large amount of geographical information about the different locations of a given ship. Flickr provides an easy-to-use API<sup>2</sup>, with a set of predefined functions, e.g., `flickr.photos.search`, for extracting all pictures (each with a unique identifier), together with necessary meta-data including geographical coordinates and dates, whose title, description, or tags contains a certain keyword given in input. The extraction process on top of the Flickr platform is automated using API Blender [4], an open-source library facilitating interactions with the Flickr API. As for the general information on ships, in particular details about their specifications, we integrated information from GrossTonnage<sup>3</sup>, Marinetraffic<sup>4</sup>, ShippingExplorer<sup>5</sup>, ShipSpotting<sup>6</sup>, and Wikipedia<sup>7</sup>. These sources contain general information about various types of ships. The purpose of the first three is to gather data about objects in the maritime domain, especially vessels. However, excepting Marinetraffic that provides partial information under an API, these Web sources do not provide a way to extract specific information from their platforms, and need to be crawled.

<sup>2</sup><https://www.flickr.com/services/api/>

<sup>3</sup><http://grosstonnage.com/>

<sup>4</sup><http://www.marinetraffic.com/fr/ais/home/>

<sup>5</sup><http://www.shippingexplorer.net/en>

<sup>6</sup><http://shipspotting.com/>

<sup>7</sup>[http://en.wikipedia.org/wiki/Main\\_Page](http://en.wikipedia.org/wiki/Main_Page)

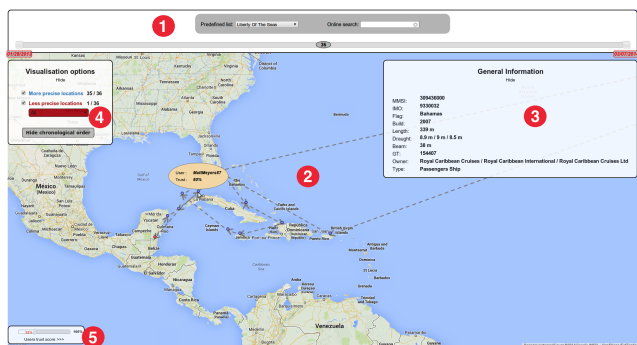


Figure 1: Main interface of our system

**Implementation.** Our demonstration system is a Web application with a map displaying ship locations. We implemented the full system using HTML, CSS, and JavaScript on the client-side, and Python on the server-side. The projection of raw geographical data onto a map uses the popular Google Maps JavaScript API. Finally, features such as filtering options are performed using the jQuery JavaScript library.

## 5. DEMONSTRATION SCENARIO

**Interface.** To interact with the system, a given user can either choose a ship name in a predefined list with locally saved data, or trigger an on-line search over the considered Web sources by providing a keyword (see Region 1 in Figure 1). For live Web search, the user can restrict the proposed set of sources for the extraction of the general information about the requested ship. The system will integrate obtained information when multiple sources are involved. Once information is obtained in local or from the Web, the different locations are displayed on the map and the general information is shown (Regions 2 and 3).

Ship positions are divided by default into two categories with different colors. Blue points on the map correspond to more precise locations, red points to less precise. As for the general information, we only show the most probable value for each attribute. The user has, however, the possibility to see details about possible other values by hovering the mouse over each attribute label. The user can restrict the visualization to ship positions in a given time interval with the slider at the top of the interface. Over this slider, we have the total number of mapping locations. More advanced visualization options are available, as shown in Regions 4 and 5 in Figure 1: The user can filter locations with high or low precision. For low precision locations, she can focus either on those on the land, outliers, or locations leading to impossible trajectories. Finally, the user can visualize hypothetical itineraries, filter users according to given trust scores, or restrict to specific users.

**Example interaction.** An expert in the maritime domain would like to acquire new ships with specific characteristics and history for business purposes. A company sells vessels which may correspond to her needs. Two particular passenger ships “Liberty of the seas” and “Costa Serena” are of interest. Thus, she decides to verify from various Web sources whether the details given by the seller are correct before making a definitive choice. To do so, she uses our maritime traffic application which already holds information about “Liberty of the seas” and “Costa Serena” in local. The user selects the first one and obtains the map view of locations. She primarily overviews the general information about the ship, and observes conflicting values for its draught and its owner. The user thus checks the values of these two attributes, as given by her most trusted Web sources. These values seem to be consistent with the seller’s data after verification. The user remembers that she is very

interested in positions of “Liberty of the seas” at some periods of the year (January to March and August to November). She filters positions corresponding to these date intervals with the slider. Surprisingly, the user remarks that the ship was near the Caribbean Sea and the Mediterranean Sea in these times. These information contradict the seller who had stated that the ship has never left Europe. To obtain more insight about the journeys, the user triggers the view of hypothetical trajectories. She examines the choice list of less precise locations to understand why some points are incorrect. She concludes that all among them are on-land and indeed invalid. Finally, she removes these kinds of locations from the map, which confirms that ship routes mostly cover two main regions.

The user pursues explorations by considering “Costa Serena” now. She only focuses on its positions and past destinations. She notes that less precise locations make the visualization cumbersome with no clear overview on routes. Therefore, the user filters one by one each type of less precise locations for explanations. For instance, she picks on-land locations and notices that all of them are located on *Corsica*, an island which contains a region named “Costa Serena” – she learns this information by clicking on an on-land location and reading the corresponding Flickr page. Observing that providers of less precise locations have trust scores below 100%, the user sets the minimum trust to 80%. Finally, she refines remaining locations w.r.t. given intervals of dates, comparing with data from the selling company, and can therefore make an informed purchase decision.

## 6. ACKNOWLEDGEMENTS

This work was supported in part by the NormAtis project.

## 7. REFERENCES

- [1] M. L. Ba, S. Montenez, R. Tang, and T. Abdesslem. Integration of web sources under uncertainty and dependencies using probabilistic XML. In *UnCrowd*, 2014.
- [2] G. K. D. De Vries and M. Van Someren. Machine Learning for Vessel Trajectories Using Compression, Alignments and Domain Knowledge. *Expert Syst. Appl.*, 39(18), 2012.
- [3] M. Dimond, G. Smith, and J. Goulding. Improving route prediction through user journey detection. In *Proc. SIGSPATIAL*, pages 476–479, 2013.
- [4] G. Gouriten and P. Senellart. API Blender: A uniform interface to social platform APIs. In *WWW*, 2012. Dev. track.
- [5] X. Li, X. L. Dong, K. Lyons, W. Meng, and D. Srivastava. Truth finding on the deep Web: is the problem solved? *PVLDB*, 6(2), 2012.
- [6] Z. Li, B. Ding, J. Han, R. Kays, and P. Nye. Mining Periodic Behaviors for Moving Objects. In *KDD*, 2010.
- [7] M. Morzy. Mining Frequent Trajectories of Moving Objects for Location Prediction. In *MLDM*, 2007.
- [8] F. Pianese, X. An, F. Kawsar, and H. Ishizuka. Discovering and predicting user routines by differential analysis of social network traces. In *WoWMoM*, pages 1–9, 2013.
- [9] R. W. Sinnott. Virtues of the Haversine. *Sky and Telescope*, 68(2), 1984.
- [10] H. Taka, D. Shibata, M. Wada, H. Matsumoto, and K. Hatanaka. Construction of a marine traffic monitoring system around the world. In *OCEANS*, 2013.
- [11] O. Wolfson. Moving Objects Information Management: The Database Challenge. In *NGITS*, 2002.
- [12] H. Yuan, Y. Qian, R. Yang, and M. Ren. Human mobility discovering and movement intention detection with GPS trajectories. *Decision Support Systems*, 2013.