

Chapter 2

Literature review

2.1 Map-matching methods

The literature review revealed that a number of studies exist on global positioning system (GPS) point matching to a digital map but most researchers agree on the two most basic categorisations of map-matching algorithms, namely *local*, also called *incremental*, and *global*.

Ying et al. (2014) describe a third approach *multi-track map-matching*, while Lou et al. (2009), although sharing research references with Ying et al. (2014), describe it as *statistical methods*. Li et al. (2013) also describes a *multi-track map-matching* algorithm but argues that all approaches can be categorised as either local or global and that multi-track map-matching is a global method. Miwa et al. (2012) concurs with this categorical breakdown but refers to the two methods as *offline* and *online* due to the nature of their processing frame.

Local methods are also known as online map-matching since they can be used to generate vehicle paths as new samples become available, typically at a rate of 1–30 s per sample. Due to the online nature of these algorithms, they are often used in applications, such as navigation, where a trade-off has to be made between accuracy, robustness and the speed of computation (Lou et al., 2009). While local methods assign portions of a trajectory onto a path based on the local geometry, global algorithms determine the globally optimal path after reading in complete trajectories, i.e. after the object has completed its entire trajectory and therefore is also referred to as *offline map-matching*.

Global methods focus on map-matching accuracy and robustness rather than on speed of execution (Li et al., 2013). Based on the definitions by Li et al. (2013) multi-track map-matching should then be categorised as a subcategory of global methods. Although the definitions of multi-track map-matching of both Ying et al. (2014) and Li et al. (2013), and the statistical method described by Lou et al. (2009), align, perhaps the data-driven nature of the approach might be enough reason for a unique categorisation. Multi-track map-matching is a statistical approach where one matches a large number of possibly sparse trajectories simultaneously to the map by trying to recover regularity among input trajectories Li et al. (2013); (Ying et al., 2014). Based on the more statistical approach of this method reference will be made to multi-track map-matching as a separate category in this dissertation.

2.1.1 Local methods

Local methods try to find a local match of geometries, and try to combine all partial matches into one (Lou et al., 2009). These methods work well on high-sampling rate GPS trajectories and are often used for analysis in a range of intelligent transport systems (ITS) and can provide services such as route guidance, fleet management, road user charging, accident and emergency response, bus arrival information, and other location-based services that require near real-time location information (Quddus et al., 2006). Local methods involves the processing of small amounts of data to identify vehicle route and current position on the road. Although GPS errors may affect the accuracy of the map-matching results, the procedure is concise because the process is repeated each time new data is collected (Miwa et al., 2012).

Quddus et al. (2006) state that these local map-matching algorithms, that utilise road segment connectivity, are generally appropriate for high-frequency positioning data from GPS at least 1 Hz or higher. They state that the use of these algorithms on low-frequency data, such as the data available for this study, especially in urban and sub-urban road networks are not suitable for local methods (Qudus & Washington, 2015). For this reason, this research considered local methods only briefly. Only those that might have similarities to global methods are discussed, while the global methods that were utilised in the study are discussed in detail.

The local method in Transportation Research Board (2002) uses two similarity measures to evaluate the possibility of matching the candidate edges to a GPS point, where a candidate edge is a potential road segment a GPS point can be matched with. They do one similarity measurement for geographic distance and the other for orientation. The sum of the individual measures for a specific edge indicates the possibility of matching the GPS point to the specific road segment the candidate edge represents. Wenk et al. (2006) uses the Dijkstra algorithm to determine the shortest path on a local free space graph constructed of the candidate routes between two partial matching results. The algorithm runs in $O(mn \log(mn))$ time, where m and n are the number of edges in the road network and the number of GPS points respectively.

Chawathe (2007) proposes a segment-based matching method, and assigns confidence values for different sampling points. The algorithm matched high-confidence segments first, and then matches low-confidence segments using previously matched edges. Because multiple sampling points are evaluated at the same time, there are some elements of this method that are found in other global methods. However, local methods only consider a small portion of the trajectory that is close to the current position being matched instead of the entire trajectory as in global methods. When the sampling frequency is very high local methods run fast and perform well. However as the sampling rate decrease these methods have a considerable decrease in accuracy due to *arc-skipping*; when a point cannot be matched to either the last identified link or a link connected to the last identified link and the intermediate links are *skipped* (Lou et al., 2009). Local methods can deal with simple road network structures if the polling frequency is low; however, due to the increase in complexity in road network structures in urban areas, the candidate road segments identified might be crowded, leading to a higher possibility of incorrect links being identified and causing a significant decrease in accuracy (Ying et al., 2014). Qudus & Washington (2015) claims that the use of local methods on sparse GPS data points might lead to a significant decrease in accuracy, in some examples as low as 70% correct link identification compared to the high 90% some methods achieve for high frequency data sets.

It appears that local methods might yield some interesting algorithm principles that

can be applied in sparse **GPS** data sets, such as the one used in this study, but that it will not be as effective as some of the other methods available. The suggested way to improve the accuracy is to not analyse the points continuously, as per the local methods, but rather evaluate all the points of the complete trace at the end of travel using a global method.

It is worth noting that other researchers, such as (Ghiani et al., 2015), who has very specifically researched the use of **GPS** data from waste collection vehicles has not made reference to any of the afore-mentioned methods or classifications. They also refer to their method as *reverse geocoding algorithm*. The principle of their algorithm matches more a local method, since it continuously matches the next data point based on previous points and not the whole set of points at once. Their data polling was also at a much higher frequency, 2s, which makes their method less applicable to this study, even though their area of application is the same.

2.1.2 Global methods

The global methods aim to match the entire trajectory with the road network. Previous studies tried to either search for possible matches between the trajectory and the road segments using the minimum Fréchet distance, a measure of similarity, or to formulate map-matching problems as optimisation problems, trying to find the shortest travel distance between two points matched in the network (Lou et al., 2009). However, as the complexity of the road network structure increases in urban areas, and the low-sampling rate of **GPS** data points becomes more of an issue, a significant decrease in accuracy occurs (Ying et al., 2014). These shortcomings are addressed through the development of more intricate and complex methods.

Alt et al. (2003) and Brakatsoulas et al. (2005) present methods based on Fréchet distance or its variants and are suitable for comparing whole trajectories as they take the continuity of curves into account. In Alt et al. (2003) the minimum Fréchet distance is determined by finding a monotone path in the free space created by two **GPS** points, from the lower left corner to the upper right corner. The algorithm runs in $O(mn \log^2 mn)$ time, where m and n are the number of edges and number of nodes in the road network respectively. Brakatsoulas et al. (2005) continued on the work by Alt et al. (2003) by proposing the use of the average Fréchet distance to reduce the effect of outliers. They also propose the use of the weak Fréchet distance as it reduces the runtime to $O(mn \log mn)$.

Yin & Wolfson (2004) introduce a weight-based map-matching method that can get up to 94% correctness depending on the **GPS** sampling rate. They make use of the edit distance to measure the similarity between trajectory and matched road segments, where the edit distance is the smallest number of insertions, deletions, and substitutions required to change inferred path (**IP**) to true path (**TP**). The accuracy of their method decreases to below 60% when the sampling rates of **GPS** data points increased above 120s. They consequently proposed that their method be used as an online method with high-frequency sampling data (Yin & Wolfson, 2004).

Since the data for this study was at various sampling rates, regularly up to 5 min apart, depending on the vehicle behaviour, the above-mentioned algorithms were not suitable for this study. It was decided to investigate other global methods that specialise in low-sampling **GPS** trajectories.

Lou et al. (2009) proposed a novel method called spatial-temporal (**ST**) map-matching for low sampling rate trajectories. **ST** matching not only considers the spatial geometric and topological structures of the road network but also the speed constraint of the road network. They managed to achieve 70% accuracies up to 5 min interval sampling **GPS**

data, which is unattainable for local methods at such low data polling frequencies. Based on Lou et al. (2009)'s work, Yuan et al. (2010) proposed an interactive voting-based map-matching algorithm to improve accuracy of the results. They were able to improve the matching accuracy using the same data as Lou et al. (2009) with an average of 10% for the same dataset and road network over different data polling frequencies. Even when sampling at 10 min intervals they were able to achieve 67% accuracy.

Yuan et al. (2010) and Ying et al. (2014), further developed the ST algorithm of Lou et al. (2009) breaking the model down into two modules; an offline mining model and an online matching module. The second model is called an online model due to the nature of the analysis being done on the data and not pre-analysing the trajectories compared to the mining model. This should not be confused with the local methods, which matches points continuously as they become available, and is also referred to as an online map-matching method. The mining module contains three submodules:

- The first is a *Spatial-Temporal centrality* estimation, which measures ST betweenness of each road segment based on both the road network structure and the crawled GPS trajectories.
- The second model, called *marginal velocity estimation* model, extracts the possible velocities of each road segment from Google Maps. This is a novel approach compared to the work by Lou et al. (2009) and others, which used the speed restriction of the road segments. They argue that most drivers do not achieve the speed restriction and that this is only useful when the road segments being evaluated have a considerable difference in speed restriction, e.g. highway versus a service road.
- The third model, called *road network decomposition* model, groups the road segments based on their characteristics.

In the online module, Ying et al. (2014) propose a two-phase method to evaluate matching a trajectory onto a path in a road network. The process is broken down into three steps:

- Firstly, calculate the ST score and derive several candidate paths.
- Secondly, the association score of each candidate path is evaluated.
- Finally, a weighted average of geographic score and semantic score for each candidate path is calculated to select the most probable path for matching the trajectory over the road network.

Ying et al. (2014)'s experimentation results show that their method is 100 times faster than Lou et al. (2009)'s ST method and also shows significant improvement in accuracy, from 47% to 52% based on their specific comparative dataset. Figure 2.1 presents a diagram of their setup.

Global methods appear to be a suitable method for this project given their accuracy on sparse GPS trajectories and the effectiveness on big datasets.

2.1.3 Multi-track methods

Multi-track methods are also commonly used for matching GPS observations especially when there is significant uncertainty in the data. Statistical methods can be considered

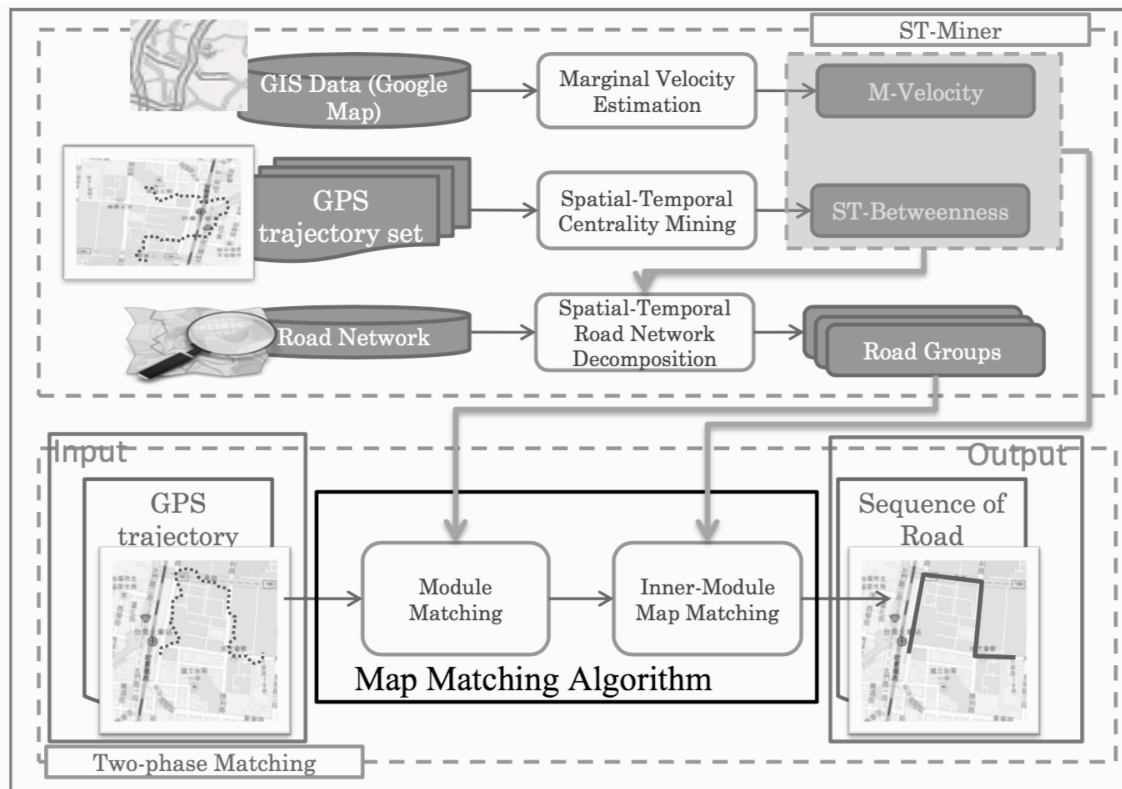


Figure 2.1: Urban map-matching framework (Ying et al., 2014)

a subset of the general max-weight algorithms (Li et al., 2013). In general, a candidate set is determined for each point in the data set, and each candidate path is assigned a weight based on its distance from the observation and topological similarity with candidate routes of neighbouring samples. The focus of multi-track map-matching is to try to match a large number of possibly sparse trajectories simultaneously to the map. Usually, multi-track map-matching adopts data-driven techniques to improve the matching result using historical trajectories. Li et al. (2013) further state that the the advantage of this approach comes from the observation that human generated trajectories, either from human operated vehicles or humans themselves, show a high degree of temporal and spatial regularity.

Javanmard & Zangh (2012) used a data set of sparse trajectories to try to group them based on the same starting and ending positions. These multiple sparse trajectories are used to produce an accurate path on the map by extracting a global order from the partial order of the sample points, then using a single-track map-matching algorithm to produce the matched path. This work focussed on matching different trips on the same route, each with very sparse samples. Although this method is realistic and might even have been applicable to the study dataset, it was unknown at the time of this study whether the trips recorded in the available data set can be grouped by travels on the same route. For example, if the waste collection vehicle data provided for this study was grouped based on areas serviced by the vehicle we could use multiple trajectories to map the routes taken by waste vehicles to service an area.

Pink & Hummel (2008) used a map-matching method based on the *Bayesian Classifier* and incorporated a *hidden Markov model* to recover accurate paths from relatively sparse trajectories. Their method showed an increased robustness compared to other methods by

exploiting vehicular motion constraints in an extended Kalman filter and by interpolating the given road network using cubic splines. Since their work focussed on the inclusion of vehicle orientation in the map-matching algorithm, it could not be applied in this study because only the position of the vehicle was available. Pink & Hummel (2008)’s tests were also only conducted on 1 Hz GPS trajectories, and the available data set has GPS trajectories of varying frequencies.

Wenk et al. (2006) proposed a simplification of the more complex hidden Markov model based method that maintains its capabilities to cope with the noises and sparsity of the raw GPS data. An interesting addition to their algorithm is called *adaptive clipping*, and uses track metadata in the form of error estimates to reduce the road network graph. This is because, in their test data, the road network did not include roads that actually existed and which the vehicle traversed. This method provided them the ability to break trajectories into sub-trajectories and enabled them to produce higher accuracy matches on their dataset. Since the adaptive clipping method was not required for the dataset of the current study, this method was not considered.

Figure 2.2 indicates the difference between single-track map-matching results (Red) and multi-track map-matching results (Blue) and the original trajectory (dotted black) and true path (Green). In most cases single track map-matching prefers cutting the corner while data-driven map-matching tends to follow larger, popular roads, although the distances travelled along the two paths are very similar and thus even the speed analysis would show little deviation from the segment free speed.

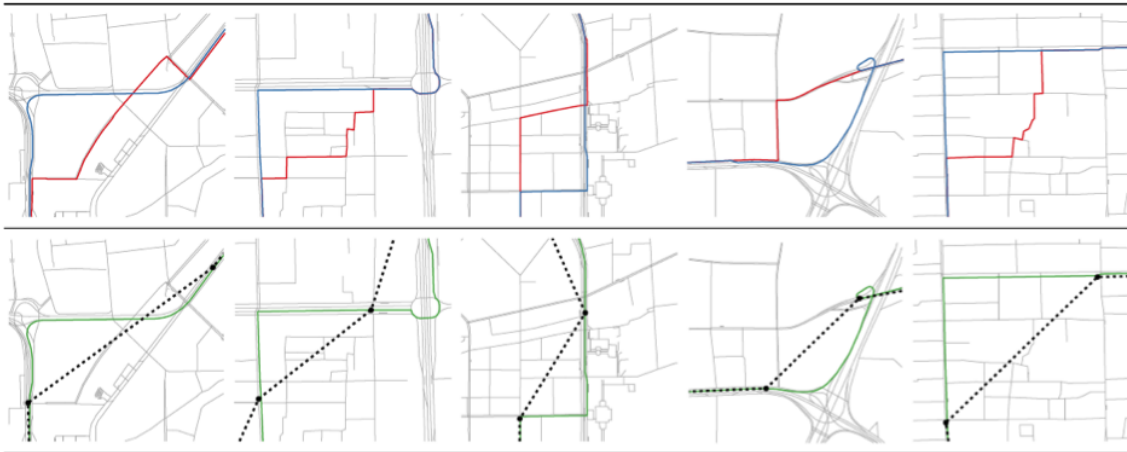


Figure 2.2: Differences between single- and multi-track map-matching (Li et al., 2013)

Multi-track map-matching methods might provide more accurate results in sparse GPS datasets, but because most methods use historical trajectories to find the most likely path for a trajectory, they were not suitable for the dataset available for the current project due to the trajectories not grouped based on service area. In conclusion, the best option for this study was to implement a global method similar to Lou et al. (2009)’s ST method. Future work on this study can further enhance the algorithm based on the work of Yuan et al. (2010) and Ying et al. (2014).

2.2 Evaluation criteria

The two main measures for success of any map-matching algorithm are in terms of efficiency and accuracy, or in other words, run time and inference quality. Many of the

aforementioned research papers and articles list the actual run time of the algorithm, as a function of the number of points and network segments, as the main criteria for the success of an algorithm. The runtime of an algorithm is calculated by means of parameter variation and by fitting a linear model on the output, thus getting the relationship between runtime and network complexity or [GPS](#) points to match.

To determine the accuracy of an algorithm there needs to be a true path, i.e. the correct path of the object, to compare the output of the algorithm to ([Lou et al., 2009](#)). But often the true path is only available when an experimental high-frequency data set is available from which the sparse data set is created, or if the actual route has been logged through manual recording.

[Miwa et al. \(2012\)](#) used two high-frequency data sets with intervals of 5 s and 50 m respectively. They created test data sets from these data sets by increasing the polling intervals to 20 s, 45 s and 90 s for the time-based data set and 100 m, 200 m and 450 m for the distance-based data set. This also allowed them to test the robustness of their algorithm on the sampling frequency, based on time as well as distance. [Miwa et al. \(2012\)](#) still had to determine the correct route by hand, but this is believed to contain very few to no errors due to the high-frequency of the [GPS](#) data. This enabled them to use the following formulas to calculate the accuracy of their algorithm.

$$\text{ARR} = \frac{\text{Length of correctly matched route}}{\text{Total length of correct route}} \quad (2.1)$$

$$\text{IARR} = \frac{\text{Length of incorrectly matched route}}{\text{Total length of matched route}} \quad (2.2)$$

If the matched route includes all of the links of the correct route accuracy ratio of route by length, total matched route ([ARR](#)) will be 1.0 even if incorrect links are included. On the other hand, inaccuracy ratio of route by length ([IARR](#)) will indicate how much of the matched route was incorrectly identified. Therefore, if the correct route is perfectly matched, [ARR](#) is 1.0 and [IARR](#) is 0.0. [Miwa et al. \(2012\)](#) states that this provides better insight into the performance of the algorithm as opposed to the commonly used accuracy ratio of plot matched to correct links ([ARP](#)) index, as used by [Lou et al. \(2009\)](#).

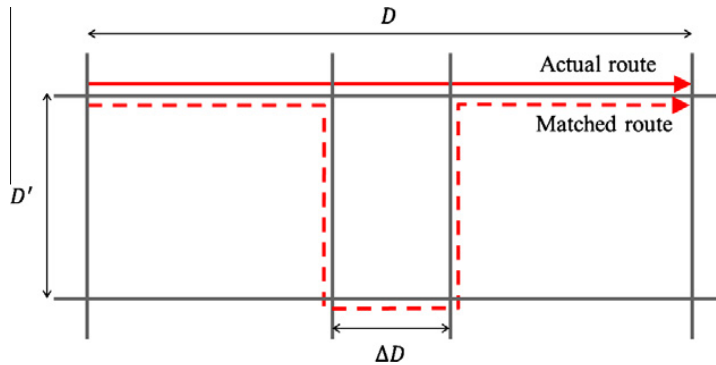


Figure 2.3: Example of [ARR](#) and [IARR](#)

Figure 2.3 shows an example of the relationship between [ARR](#) and [IARR](#). If ΔD decreases to 0 so will [ARR](#) increase to 1.0, which indicates a possible perfect match. However, the [IARR](#) will only approach 0 if D' reduces to 0.

[Lou et al. \(2009\)](#) used a similar data set as [Miwa et al. \(2012\)](#), and were able to synthesise experimental data sets from the human-labelled high-frequency data sets. [Lou](#)

et al. (2009) provided very similar metrics to Miwa et al. (2012), but when comparing the length of the matched route to the actual route, they include the total length of the matched route and not only the correctly matched sections. Their definition is formulated as

$$ARR2 = \frac{\sum \text{The length matched road segments}}{\sum \text{The length of identified road segments in correct route}} \quad (2.3)$$

The drawback with this measurement (2.3), is that if the algorithm identifies incorrect segments of similar length to the correct segments it missed, then the measurement can still provide a ratio of 1. Lou et al. (2009) created a second measurement *Accuracy by Number* (2.4) that would have to be relied on to indicate the number of incorrect segments identified.

$$ARRn = \frac{\# \text{ of correctly matched road segments}}{\# \text{ All road segments in correct route}} \quad (2.4)$$

However, (2.4) will still not provide an indication about the length of the incorrectly matched segments. For this reason this study will not be measuring *Accuracy by length*, but rely on *ARR* and *IARR*. Lou et al. (2009)'s measurement of *Accuracy by Number* (2.4), originally referred to them by the acronym *An*, is incorporated as it provides a useful additional insight to indicate how accurately individual segments were identified. To keep consistency in acronyms this measurement is renamed to accuracy ratio of route by number of links (*ARRn*).

Zheng et al. (2012) calculated the inference quality of the algorithm by measuring the similarity between the inferred path *IP* and the true path *TP* using equation (2.5)

$$AI = \frac{LCR(TP, IP)^{length}}{\text{Max}(TP^{length}, IP^{length})} \quad (2.5)$$

where *LCR* is the longest common road segment of *IP* and *TP*, and *length* indicates the overall length of the applicable path noted.

Although Zheng et al. (2012) referred to their inferred route metric as *RI*, and true path as *Ground Truth* or *RG* this study used the more intuitive acronyms, and appropriate to previously defined terminologies *inferred path*, *IP* and *true path*, *TP* respectively. In conclusion, it was decided for this study to use all of the above measurements, excluding *ARR2*, as it is replaced by *ARR* and *IARR*.

Lou et al. (2009) argued that many current global methods have not tried using true paths from real world generated data to evaluate the actual matching accuracy and thus their true effectiveness is unknown. To some degree the same caveat did exist in this study, as there was no true path available for the waste collection vehicles' trajectories. Fortunately, there was a 5 Hz data set available that could be used to determine the accuracy of the algorithm, albeit by hand. All the experimental data sets, however, were generated from a true path and could be used to assess the accuracy as part of the experimentation.

Li et al. (2013) addresses this concern by constructing a benchmark dataset. They selected 100 random sparse trajectories from their Beijing taxi data set and recruited four volunteers with more than five years of driving experience in Beijing to manually indicate the path they believed the taxis travelled. The purpose was to obtain a true path map for these trajectories based on how a real driver would choose a path that followed these trajectories based on his or her knowledge of the roads in and around the city.

For this study, the initial experimentation was done by creating experimental data sets on a simple grid network and then on a real-world road network. The recommendation is

for future projects to follow an approach similar to that used by [Miwa et al. \(2012\)](#) and [Lou et al. \(2009\)](#), and to use an experimental data set that can be labeled by humans, and then synthesised into an experimental data set based on different observations frequencies.

Chapter 3

Algorithm development

3.1 Definition

In this section we define the algorithm and problem statement and how it was adjusted for this specific study. The first set of definitions explicitly define the common terms used in describing map-matching as used by Lou et al. (2009). The second part stipulates how these definitions relate to the Multi-Agent Transport Simulation (MATSim)-specific objects used for this study. MATSim is a collaborative open source project that can be used as a package in Java. This study use version 0.7.0 of MATSim and Java SE 1.7, for for more info visit <https://github.com/matsim-org> and <https://www.oracle.com/technetwork/java/javase/overview/index.html> respectively.

global positioning system (GPS) log: A collection of GPS points $\mathbf{Log} = \{p_1, p_2, \dots, p_N\}$. Each GPS point $p_i \in \mathbf{Log}$ contains latitude p_i^{lat} , longitude p_i^{long} and timestamp p_i^{t} .

GPS Trajectory: A sequence of GPS points, \mathbf{T} , where the time interval between any consecutive GPS points do not exceed a certain threshold ΔT^* , i.e. $\mathbf{T} = \{p_1, p_2, \dots, p_n\}$, where $p_i \in \mathbf{Log}$, and $0 < p_{i+1}^{\text{t}} - p_i^{\text{t}} < \Delta T | \forall 1 \leq i < n$. ΔT is also referred to as the *sampling interval*.

The enumerated black dots in Figure 1.1 is an example of a GPS trajectory. For this sample dataset a GPS point was generated at the end of the true path when the vehicle was switched off, thus generating an additional point that has a ΔT which is less than the other points in the GPS log. Table 3.1 reflects the collection of the GPS points from Figure 1.1. The GPS coordinates use the standard World Geodetic System 84 (WGS84) as coordinate reference system.

Table 3.1: GPS trajectory from example data

Point	Longitude	Latitude	Time s
p_1	-33.962	18.473	0
p_2	-33.961	18.472	25
p_3	-33.960	18.469	50
p_4	-33.959	18.466	100
p_5	-33.958	18.465	150

Road network: A road network is a directed graph $\mathbf{G}(\mathbf{V}, \mathbf{E})$, where \mathbf{V} is a set of vertices representing the intersections and terminal points of the road segments, and \mathbf{E} is a

set of edges representing road segments. Figure 3.1 shows a portion of the MATSim network of the sample data depicted in Figure 1.1, the circles indicate the vertices and the arrows the edges of the road network. Edges e_1 and e_2 have also been annotated with their start and end points.

Road segment: A road segment, e , is a directed edge of a road network such that $e_i \in \mathbf{E}$. The segment e is associated with an id e^{id} , a typical travel speed e^v , a length value e^l , a starting point e^{start} , an ending point e^{end} and a list of intermediate points that describes the road using a polyline. Figure 3.1 shows several road segments from the example road network with the start and end points of edges e_1 and e_2 annotated.

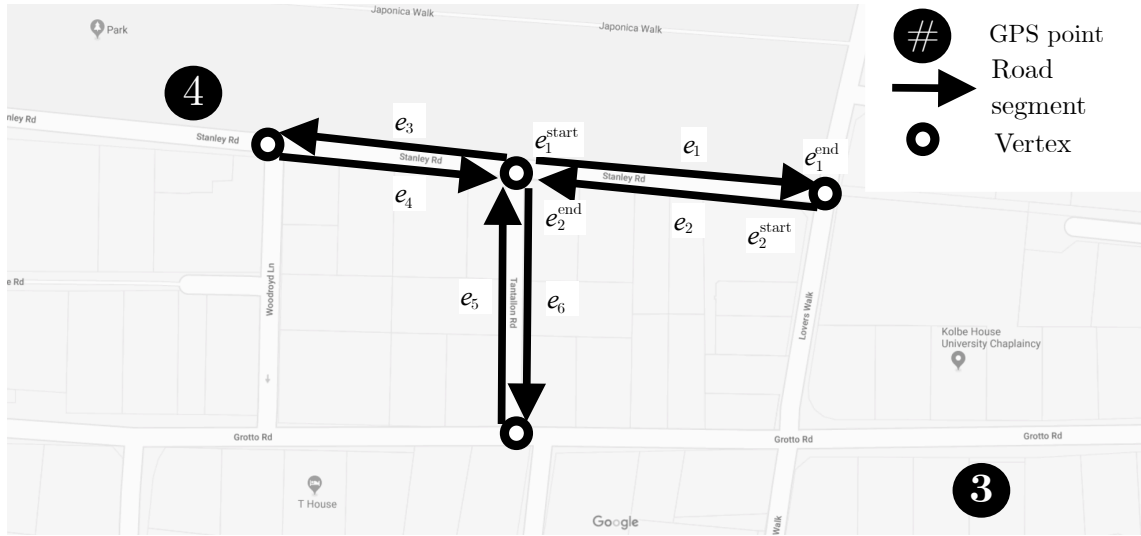


Figure 3.1: Edges on a road network

Path: Given two vertices v_i, v_j in a road network \mathbf{G} , a path \mathbf{P} is a set of connected road segments that start at v_i and end at v_j , i.e. $\mathbf{P} = \{e_1, e_2, \dots, e_n\}$, where $e_1^{\text{start}} = v_i$, $e_n^{\text{end}} = v_j$, $e_k^{\text{start}} = e_{k+1}^{\text{end}}$, $1 \leq k < n$.

Problem statement: Given a raw GPS trajectory, \mathbf{T} , generated by an agent travelling on a path in a road network $\mathbf{G}(\mathbf{V}, \mathbf{E})$, determine the most likely path \mathbf{P} that the agent travelled.

All the previous definitions are general terms used to describe the general format of road networks and GPS trajectories for most map-matching problems. Since the application of the algorithm developed was in MATSim, some specific terminologies need to be defined as well as their relationship to the previous definitions.

GPS point: In order to work with meters, the coordinate system used in MATSim is the SA-Albers projection which is adapted from the Albers equal-area conic projection. Table 3.2 shows the GPS trajectory after the conversion to the SA-Albers projection. A trajectory using the SA-Albers projection is defined as \mathbf{T}_{SAA} .

Link: Similar to the definitions of a road segment, a link L in MATSim represents an edge with only one direction of travel. In order to represent bidirectional traffic between

Table 3.2: [GPS](#) trajectory of example data with SA-Albers projection

Point	x	y	Time s
p_1	-513327.5134	-3716061.097	0
p_2	-513486.9662	-3715934.682	25
p_3	-513722.076	-3715759.626	50
p_4	-514033.4896	-3715661.143	100
p_5	-514115.2881	-3715473.827	150

two points, two links need to be defined, one for each direction of travel. A link has, inter alia, the following fields:

- link identifier, L^{id} - the name of the link;
- a from and to node L^{fn} , L^{tn} , the starting and ending points of a link respectively;
- free speed L^{fs} in m/s . This is the speed an object is allowed to travel on the link, and can vary depending on the time of day; and
- the length of the link L^{length} in meters.

In comparison to the road segment definition, a [MATSim](#) link does not contain intermediate points. To represent a curved road segment multiple links and nodes can be used to accurately represent the road shape in topological format. If the visual aspect of a link is not important a curved road segment can, for example, be represented by just one straight link with the correct length defined programmatically and not graphically. This presents a potential issue in using the spatial and topological analysis part of the map-matching algorithm as it relies on the network to be as topologically and spatially accurate as possible, and [MATSim](#) network can only contain straight lines. As can be seen in Figure 3.2, the [MATSim](#) network available for this study has relatively accurate topological structures. This is achieved by representing curved segments on the road with a number of small, straight links.

Node: A node in [MATSim](#), N , is defined as a topological point where one or more links can be joined and an agent can move from one link to another. A node has, inter alia, the following fields that are used during calculations: node ID N^{id} , list of In links N^{inLinks} (links that end inside this node), list of Out links N^{outLinks} (links that start inside of this node). In Figure 3.2 the nodes are represented by small grey diamonds that connect the links.

Network: A network is a set of links connected to each other via a set of nodes to form and represent a transport network where upon agents can travel and interact with one another. A network is defined as $N_{\text{MATSim}}(\mathbf{L}, \mathbf{N})$, where \mathbf{N} is a set of nodes representing the intersections and terminal points of the road segments, and \mathbf{L} is a set of links representing road segments.

The [MATSim](#) representation of the entire example network is illustrated in Figure 3.2 and includes the true path travelled by the vehicle as well as the [GPS](#) trajectory in the applicable coordinate system.

Path: In [MATSim](#) a path is a list of connected links on which agent can travel on from one node to another, and is defined by the use of \mathbf{P}_l .

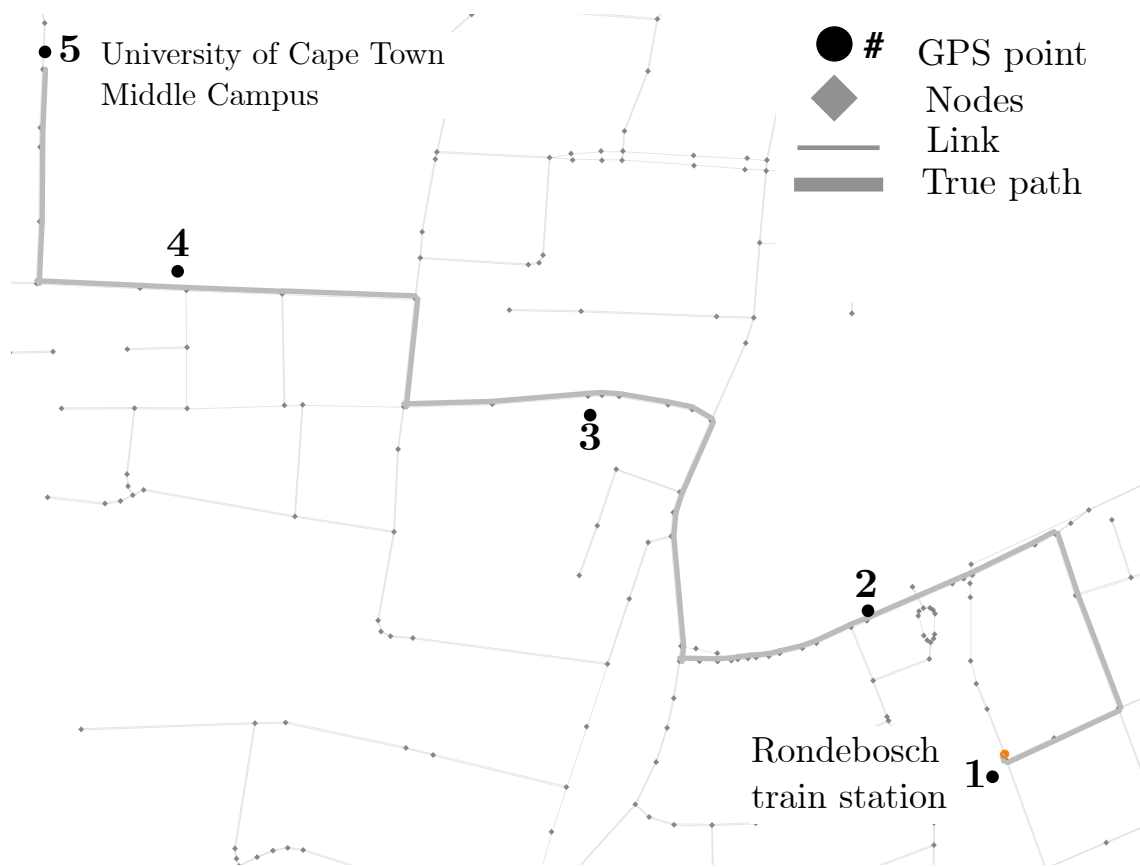


Figure 3.2: MATSim network of sample data

The map-matching problem can now be formulated as follows:

Updated problem statement: Given a raw **GPS** trajectory \mathbf{T}_{SAA} generated by an agent travelling on a path on a road network $N_{MATSim}(\mathbf{L}, \mathbf{N})$, determine the most likely path P_l that the agent has travelled.

3.2 Candidate preparation

For each **GPS** point in a trajectory, there is a potential list of road segments or links on which the agent could have travelled on when generating the **GPS** point. These potential links are called *candidate links*. Given trajectory $\mathbf{T} = \{p_1, p_2, \dots, p_n\}$, there exist a number of candidate links for each point $p_i \in \mathbf{T}$.

Lou et al. (2009) made use of a search radius to identify a number of road segments for each point in the **GPS** trajectory where the road segment is within the radius from the **GPS** point.

There is no native functionality in **MATSim** to efficiently collect all the links that are within a radius of a certain point. However, all the nodes within a radius can be transposed into a specific data structure referred to as a *quadtree* (more detail on the creation of the quadtree refer to section 4.3.2). When all the nodes within a radius have been identified, one can check whether or not the links to and from the nodes are within the specified radius from the **GPS** point by determining the projected distance from the **GPS** point to the nearest point on the link. However, there exists an issue with explicitly defining the radius for collecting all the applicable nodes without taking into consideration the characteristics of the links within the network. The issue is that some nodes might be connected to the link right next to the **GPS** point but the link’s endpoint nodes are outside of the search radius and thus the link will not be evaluated, leading to an incorrect path being inferred for the **GPS** trace. Figure 3.3 illustrates such a possible situation, where **GPS** point 2 is connected to the highlighted line to its left, but the distance to the start and end nodes of this link is 1.7km and 1.5km away respectively. If the search radius for nodes in the quadtree does not include these outlying nodes, the **GPS** point will only have the short residential links to its right evaluated. Thus, the search radius needs to be a function of the longest link within the network to ensure the algorithm does not miss potential candidate links that might possibly be the correct links. Since the search for nodes within a distance is efficiently done using the quadtree the algorithm uses the longest link length within the network as the search radius, which is not an input into the algorithm as per Lou et al. (2009)’s implementation.

Once all the possible nodes have been identified and all the incoming and outgoing links from the nodes have been stored as potential candidate links for the **GPS** point being evaluated, a straight-line distance from the **GPS** point to each potential candidate link is calculated and saved. The list of potential candidate links identified is sorted according to proximity to the **GPS** point. The next step is to use only a certain number of links from the sorted list of potential candidate links for further analysis and result matching. The number of links used greatly affects the efficiency and effectiveness of the algorithm, and a number of experiments was run to determine the influence of these parameters on results.

The next step was to determine a *candidate point* on each candidate link using segment projection, and was defined as follows.

Candidate point: This is the projection of a point p to a link l and is defined as point c on l such that $c = \arg \min_{c_i \in l} \text{dist}(c_i, p)$, where $\text{dist}(c_i, p)$ returns the distance between

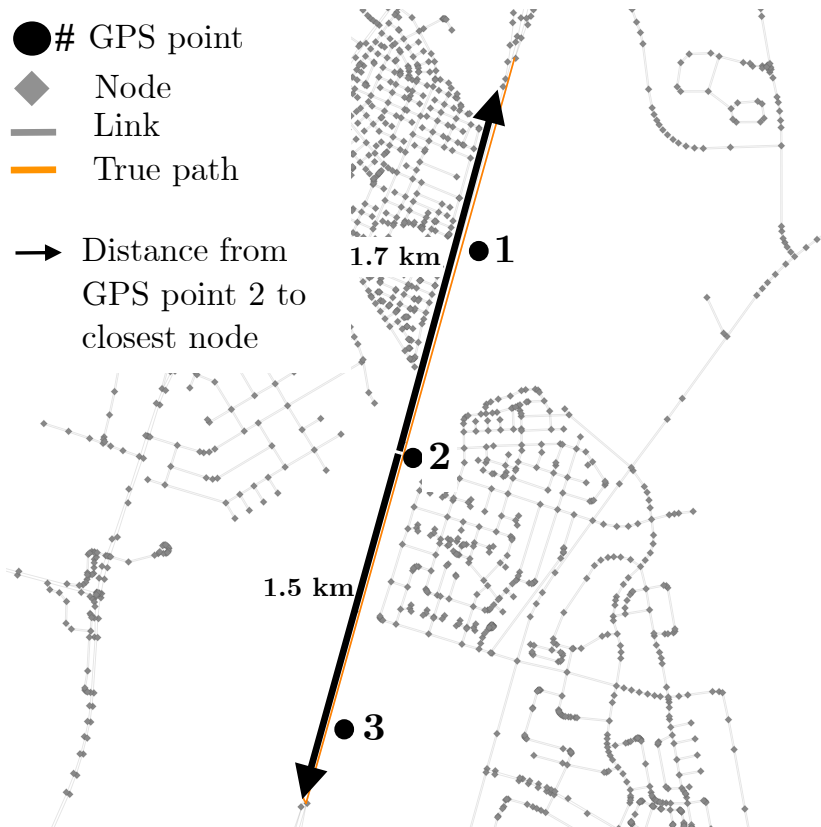


Figure 3.3: Search radius required when searching for candidate links

p and any point c_i on l . The j th candidate link and candidate point of p_i will then respectively be denoted as l_i^j and c_i^j . As shown in Figure 3.4, p_i 's candidate points are c_i^1 , c_i^2 and c_i^3 . The candidate point represents the most likely point on the specific link at which the agent was when recording the GPS point being evaluated.

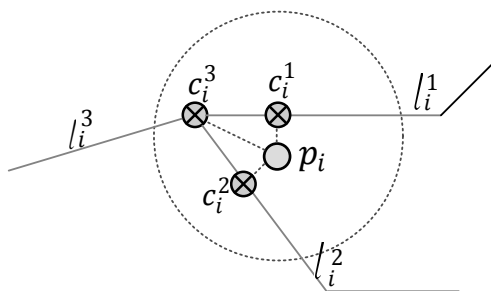


Figure 3.4: Point projection on candidate links

Updated problem statement Once the candidate link sets and their respective candidate points are retrieved for all the sampling points on the trajectory \mathbf{T}_{SAA} , the problem statement becomes how to choose one candidate from each set so that $\mathbf{P}_l = \{c_i^{j_1}, c_i^{j_2}, \dots, c_i^{j_n}\}$ best matches $\mathbf{T}_{SAA} = \{p_1, p_2, \dots, p_n\}$

3.3 Spatial analysis

In the spatial analysis both the geometric and topological information of the road network is used to evaluate the candidate points found in the previous step. The geometric information is used in the *observation probability*, and the topological information in the *transmission probability*.

Observation probability: The observation probability is defined as the likelihood that a **GPS** sampling point p_i matches a candidate point c_i^j and is computed based on the distance between the two points, defined as $dist(c_i^j, p_i)$. The inherent error in a **GPS** measurement causes p_i to be a certain distance from c_i^j and can be approximated using a normal distribution $N(\mu, \sigma^2)$. This distribution can be used to indicate how likely a **GPS** observation p_i can be matched to a candidate point c_i^j on the real road without considering its neighbouring points. Formally the observation probability is defined as: $N(c_i^j)$, of c_i^j with regard to p_i as:

$$N(c_i^j) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(x_i^j - \mu)^2}{2\sigma^2} \right) \quad (3.1)$$

where $x_i^j = dist(c_i^j, p_i)$, the distance between p_i and c_i^j . Similar to [Lou et al. \(2009\)](#) this study makes use of a zero-mean normal distribution with a standard deviation of 20 m for **GPS** error. Using these parameters the observation probability will always output very low probability values for any **GPS** point even if the point is within 0 m of the candidate point. During informal discussions with the original authors they indicated that the final spatial-temporal (**ST**) probability value should rather be viewed as a *normalised score* instead of a probability, per se, due to the low values. In this study the output of the **ST** function was still referred to as a probability.

For the map-matching example used in this section the 8 closest links were used to identify candidate links. The candidate links identified for **GPS** point 2 is shown in [Table 3.3](#) and the candidate link names can be referenced from [Figure 3.5](#).

Table 3.3: **GPS** point candidate links for **GPS** point 2

Candidate link	Distance to link	Observation probability
l_2^1	7.08	0.0187
l_2^2	7.08	0.0187
l_2^3	8.15	0.0184
l_2^4	8.15	0.0184
l_2^5	19.08	0.127
l_2^6	19.08	0.127
l_2^7	19.08	0.127
l_2^8	19.08	0.127

Since the observation probability calculation does not take into account the position context of a **GPS** point, it can sometimes lead to wrong matching results. [Figure 3.6](#) shows such an example. The thick lines represent a highway, and the thin vertical line represents a local road. Although p_i is closer to c_i^1 on the local road, it should match p_i to c_i^2 on the highway if it is already known that p_i 's neighbours p_{i-1} and p_{i+1} are on the highway. This is based on the assumption that a vehicle is unlikely to take a roundabout path [Lou](#)

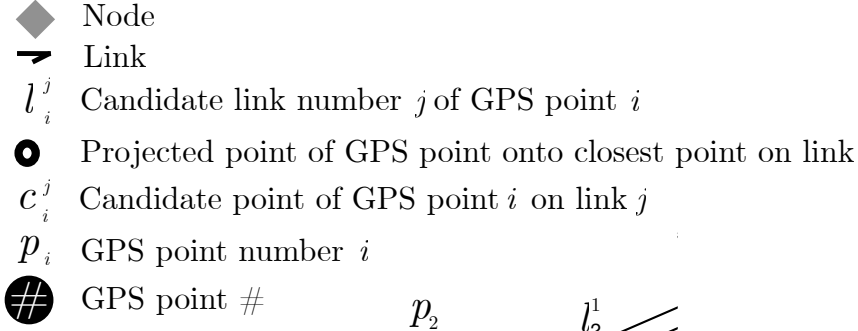


Figure 3.5: Example of candidate links

et al. (2009). To prevent this potential error the algorithm made use of the *transmission probability*.

Transmission probability: Given two candidate points c_{i-1}^t and c_i^s for two neighbouring GPS sampling points p_{i-1} and p_i respectively, the transmission probability from c_{i-1}^t to c_i^s is defined as the likelihood that the true path from p_{i-1} to p_i follows the shortest path from c_{i-1}^t and c_i^s .

The transmission probability is defined as

$$V(c_{i-1}^t \rightarrow c_i^s) = \frac{d_{i-1 \rightarrow i}}{w_{(i-1,t) \rightarrow (i,s)}} t \quad (3.2)$$

where $d_{i-1 \rightarrow i} = \text{dist}(p_i, p_{i-1})$ is the Euclidean distance between p_i and p_{i-1} , and $w_{(i-1,t) \rightarrow (i,s)}$ is the length of shortest path from c_{i-1}^t to c_i^s .

Combining equations (3.1) and (3.2) the spatial analysis function can be defined as: $F_s(c_{i-1}^t \rightarrow c_i^s)$ as the product of observation probability and transmission probability:

$$F_s(c_{i-1}^t \rightarrow c_i^s) = N(c_i^s) * V(c_{i-1}^t \rightarrow c_i^s), 2 \leq i \leq n \quad (3.3)$$

where c_{i-1}^t and c_i^s are any two candidate points for two neighbouring GPS points p_{i-1} and p_i respectively.

Equation (3.3) computes the likelihood that an object moves from c_{i-1}^t to c_i^s using the product of two probability functions; thus, geometric and topological information are both taken into consideration. Note that in practice, it is unlikely for a moving object

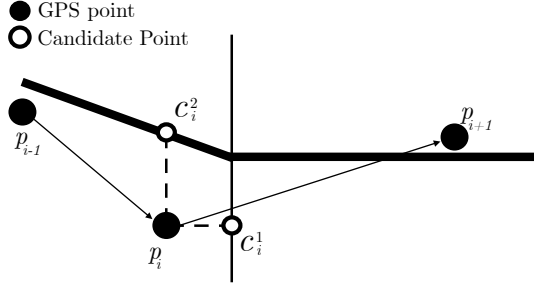


Figure 3.6: Transmission probability example

to always strictly follow the shortest path. Therefore the observation probability $N(c_i^s)$ cannot be omitted from Equation (3.3).

With spatial analysis, for any two neighbouring GPS points p_{i-1} and p_i , a set of candidate paths $c_{i-1}^t \rightarrow c_i^s$ are generated. Each path is assigned a spatial measurement value computed from Equation (3.3).

There exists a shortcoming inherent to low-sampling rate GPS data, because even though the map-matching algorithm determines the most likely link for every GPS point, the assumption is still that for two links found for subsequent GPS points not connected to each other, the links traversed between the two links is the shortest path. In real life very similar paths, in terms of distance, might exist between two links or a driver might take a non-shortest path route and still generate points at similar positions in the network. The temporal analysis aims to address part of this problem but without higher sampling rates there is no accurate way of identifying the correct route if many alternatives exist between two GPS points.

3.4 Temporal analysis

Although the spatial analysis can determine the actual path from the other candidate paths relatively accurately, there are situations where the spatial analysis might choose an unlikely path. When two paths that allow very different travel speeds are situated next to each other, spatial analysis might determine a path that spatially makes sense but when considering the speed at which the agent would travel on the chosen path, this is highly unlikely. The inferred speed might seem unlikely due to either the speed restrictions or the practical speed at which an agent would have to travel to reach the next candidate point within the given time between the subsequent samples. See the example shown in Figure 3.7. The two lines on the left is a highway, and the lines to the right are roads inside a residential area. The spatial analysis function may produce the same value whether two points p_{i-1} and p_i are matched to the highway or to the residential road. However, if one calculates the average speed from p_{i-1} to p_i on a straight line as 100 km/h , the algorithm would match the points to the highway considering the speed limits of the residential road.

More formally, given two candidate points c_{i-1}^t and c_i^s for two neighbouring GPS sampling points p_{i-1} and p_i respectively, the shortest path from c_{i-1}^t to c_i^s is denoted as a list of road segments $[l'_1, l'_2, \dots, l'_k]$. The average speed $\bar{v}_{(i-1,t) \rightarrow (i,s)}$ of the shortest path is

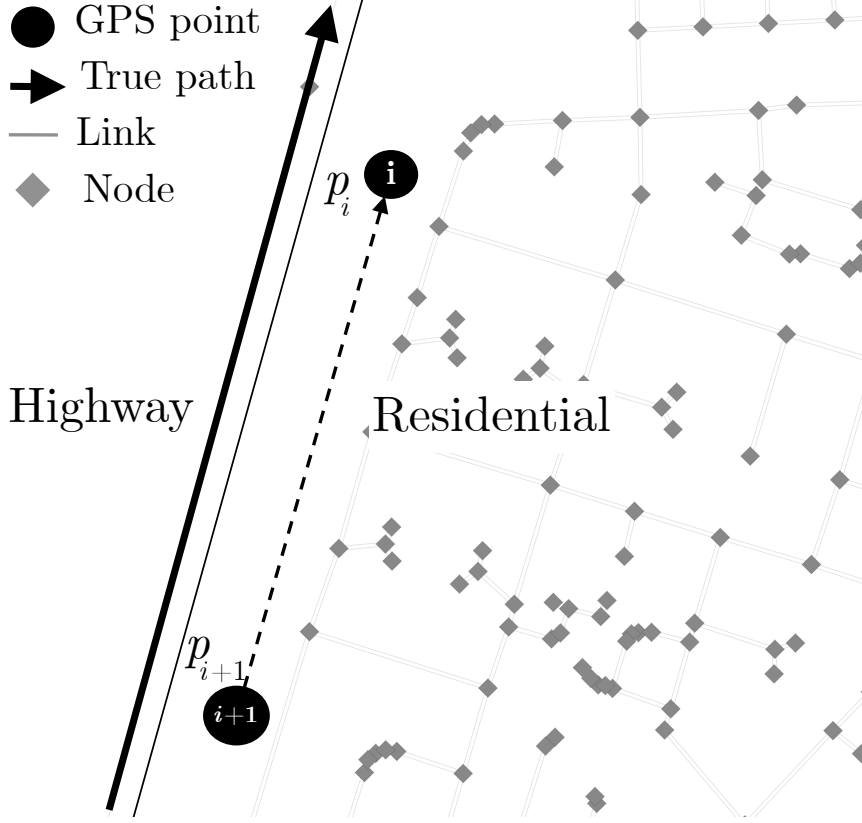


Figure 3.7: Temporal analysis example

computed as follows:

$$\bar{v}_{(i-1,t) \rightarrow (i,s)} = \frac{\sum_{u=1}^k \ell_u}{\Delta t_{i-1 \rightarrow i}} \quad (3.4)$$

where $\ell_u = l'_u{}^{length}$, the length of l'_u , and $\Delta t_{i-1 \rightarrow i} = p_i^t - p_{i-1}^t$, the time interval between two sampling points p_i and p_{i-1} .

Note that each road link l'_u is also associated with a typical speed value $l'_u{}^{fs}$, referred to as the free speed of a link. One of the benefits of using a MATSim network was that the free speed can be a function of the time and as such one is able to incorporate more accurate free speed data based on time of day to evaluate the temporal probability of the algorithm.

The cosine distance is used to test the similarity between the actual average speed from c_{i-1}^t to c_i^s and the speed constraints of the path between the two points. Consider the vector that contains k elements of the same value $\bar{v}_{(i-1,t) \rightarrow (i,s)}$ and the vector $(l'_1{}^{fs}, l'_{21}{}^{fs}, \dots, l'_k{}^{fs})^T$

Temporal analysis: The temporal analysis function is defined as follows.

$$F_t(c_{i-1}^t \rightarrow c_i^s) = \frac{\sum_{u=1}^k (l'_u{}^{fs} \times \bar{v}_{(i-1,t) \rightarrow (i,s)})}{\sqrt{\sum_{u=1}^k (l'_u{}^{fs})^2} \times \sqrt{\sum_{u=1}^k \bar{v}_{(i-1,t) \rightarrow (i,s)}^2}} \quad (3.5)$$

As in the spatial analysis function, c_{i-1}^t and c_i^s are any two candidate points for p_{i-1} and p_i respectively.

3.5 Result matching

Once the spatial and temporal analyses have been completed, a candidate graph \mathbf{G}'_T is generated with every candidate point for every GPS point representing a node in the graph and every link in the graph representing the movement between the neighbouring candidate points.

Formally the candidate graph for trajectory $\mathbf{T} = \{p_1, p_2, \dots, p_n\}$ is $\mathbf{G}'_T(\mathbf{L}'_T, \mathbf{N}'_T)$, where \mathbf{N}'_T is a set of candidate points for each GPS sampling point, and \mathbf{L}'_T is a set of links representing the possibility of moving from one candidate link to another, as depicted in Figure 3.8. Each link also contains attributes for the connection between the two nodes,

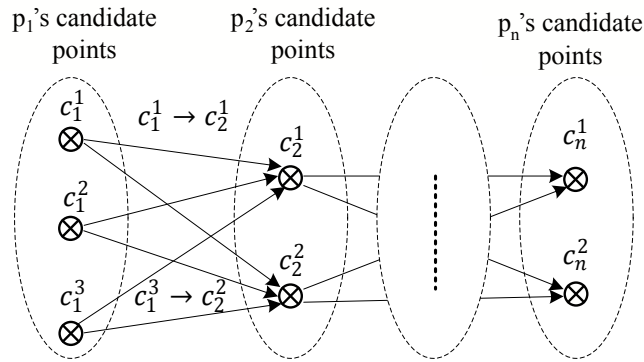


Figure 3.8: Schematic view of candidate graph $\mathbf{G}'_T(\mathbf{L}'_T, \mathbf{N}'_T)$

namely the value of the ST algorithm's probability that the object traversed from the one candidate point to the other, stored as the length of the link. The link also stores some metadata such as the shortest path on the original road network between the candidate links if they are not connected by a common node. The shortest path is used to reconstruct the inferred route once the most probable candidate points have been identified. Each node in \mathbf{G}' also stores the observation probability $N(c_i^s)$ and each link the average speed, $V(c_{i-1}^t \rightarrow c_i^s)$ and temporal analysis $F_t(c_{i-1}^t \rightarrow c_i^s)$ for detailed analysis afterwards.

A candidate path sequence for the entire trajectory \mathbf{T}_{saa} is a path in the candidate graph, denoted as $\mathbf{P}_c = \{c_1^{s1}, c_2^{s2}, \dots, c_n^{sn}\}$. The overall score for such a candidate sequence is $F(\mathbf{P}_c) = \sum_{i=2}^n F(c_{i-1}^{s_{i-1}}, c_i^{s_i})$. From all the candidate sequences the aim is to find the one with the highest overall score as the best matching path for the trajectory. More formally, the best matching path \mathbf{P} for a trajectory \mathbf{T} is selected as:

$$\mathbf{P} = \underset{\mathbf{P}_c \in \mathbf{G}'_T(\mathbf{L}'_T, \mathbf{N}'_T)}{\arg \max} F(\mathbf{P}_c) \quad (3.6)$$

To get the candidate sequence with the highest overall score the entire graph needs to be traversed for all the possible combination of routes from the start to the end point. The score of a route is the the sum of the link lengths in a route where the link length is the ST algorithm's score for the link being the correct link which the agent traversed while generating the GPS points. This process can be executed using existing graphing

algorithms like Dijkstra’s algorithm. But since Dijkstra’s algorithm uses the shortest path as an objective function, instead of using the length of the links as the travel distance utility this implementation uses $1 - l^{\text{length}}$. Since the length of the link is the ST algorithm’s probability that the object traversed from the one candidate point to the other, i.e. the probability that it is the right route, Dijkstra’s algorithm was actually calculating the routes that is least likely to *not* be the wrong route. Changing Dijkstra’s algorithm to essentially calculate the longest route is only possible because the graph is a directed acyclic graph, that is, it is a finite directed graph with no directed cycles, thus no infinite loops will occur when looking for the longest path.

The Dijkstra method requires a single starting and ending node for the algorithm to start at and navigate to. Since the graph has multiple starting points, for point p_1 , and multiple end points, for point p_n , a dummy starting and ending node was created. The starting dummy node was connected to the first tier of nodes with a length equal to the observation probability of the node it was connected to. The ending dummy node was connected to the last tier of nodes with length 1.

Figure 3.9 shows the graphical representation of the map-matching example given in Chapter 1. Every tier in the graph, excluding the starting and ending node, represents the eight possible candidate points for each GPS point, $P_1 \rightarrow P_5$. Once the longest path has been calculated from the graph, the original links from the road network can be inferred from the graph link and node data.

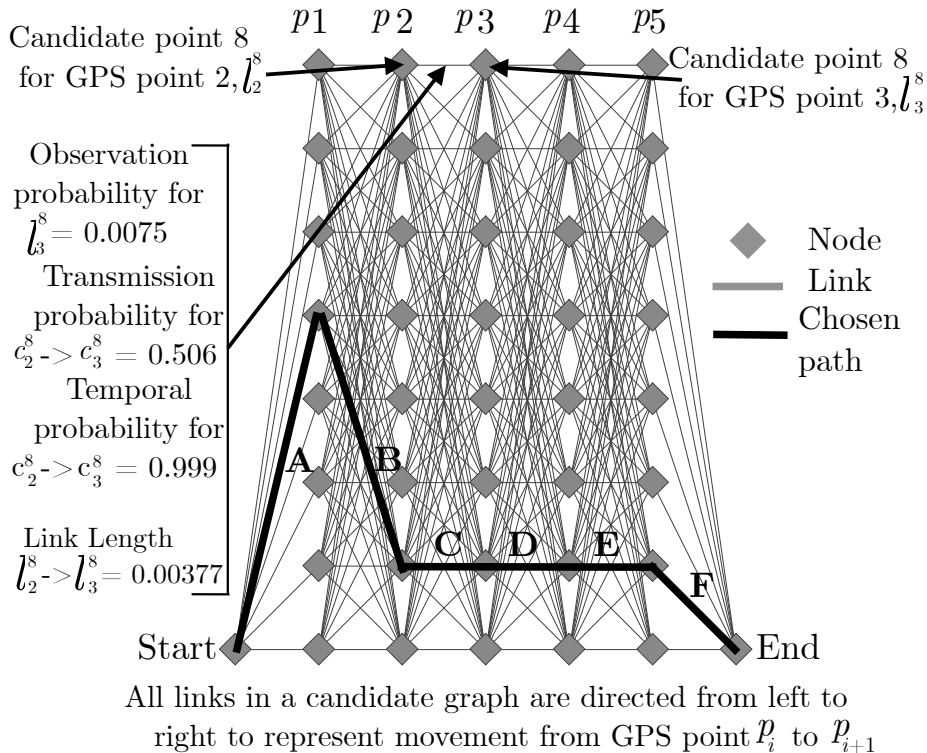


Figure 3.9: Candidate graph for sample data

In Table 3.4, the probability for each chosen link can be seen as well as the average link probability, which was used to calculate the overall probability of the inferred path (IP).

Graph link F did not form part of the calculation as it was only used to create an end point for the Dijkstra algorithm to use and the length of the link is 1.

Table 3.4: Graph link **ST** probability data

Graph link	ST probability
A	0.0159
B	0.0083
C	0.0119
D	0.0153
E	0.0079
Average	0.0119
Min	0.0079
Max	0.0159
Std Dev	0.0033

3.6 Analysing results

Since a true path (**TP**) is available in this example, one can analyse the results of the algorithm using the criteria defined in section 2.2 as per the Table 3.5.

Table 3.5: Example results analyses against **TP**

Measurement	Value
ARR	0.997
IARR	0
ARR_n	0.977
Al	0.997
Average ST value	0.0119

To analyse the results of the algorithm in the absence of a **TP** one can review the inferred speed as well as the **ST** probability value. The inferred speed can be calculated from the **IP** using the time stamps of the **GPS** points while the probabilities can be transferred from the graph links to the associated network links. Since subsequent **GPS** points can be allocated to candidate links that are not connected, some graph links contain a list of multiple network links representing the shortest path between the one candidate link and the subsequent candidate link. An original network link can also be associated with more than one graph link if the link is a candidate link for more than one **GPS** point. In order to assign a probability to the original network links, the probability of each graph link can be assigned to the candidate link and intermediate links with which it was associated. If there are shared original network links between graph links, an average between the probabilities and calculated speed was assigned to these links. The result of this analysis can be seen in Figure 3.10. From Figure 3.10b, we can infer what the waste collection vehicle was doing on each by part of its journey by using the speed as a proxy. It appears from the red lines that the vehicle might have been collecting waste, while the yellow lines imply it was most likely deadheading, i.e. driving to a new service area and not collecting any waste. This simple analysis of inferred speed answers the specific question posed by the use case of this study.

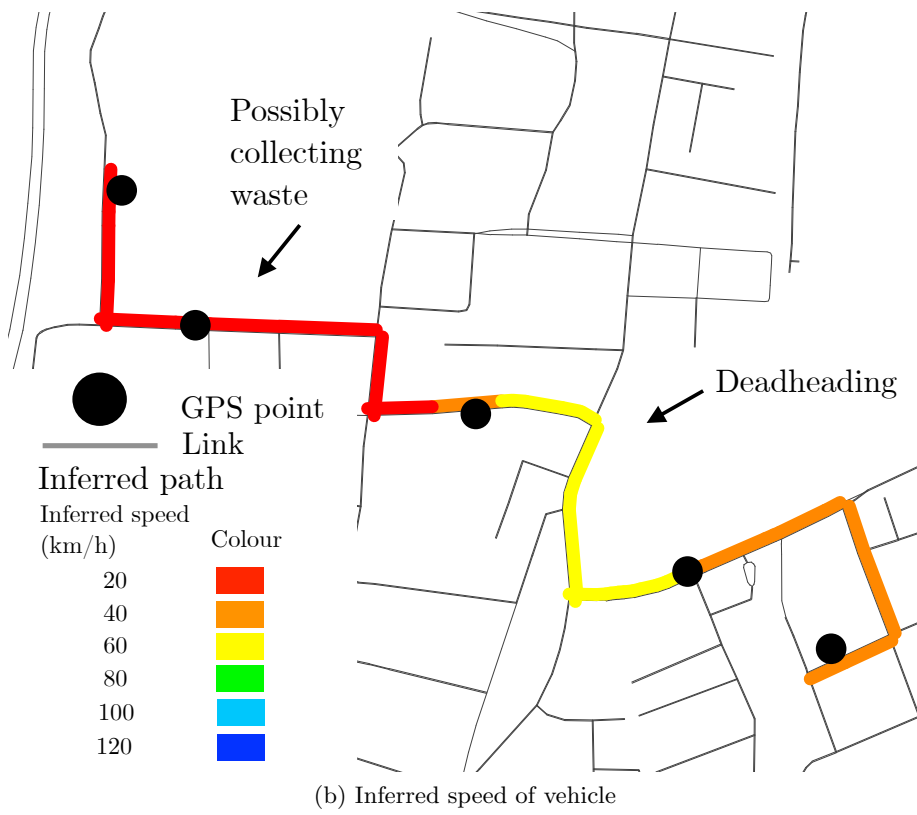
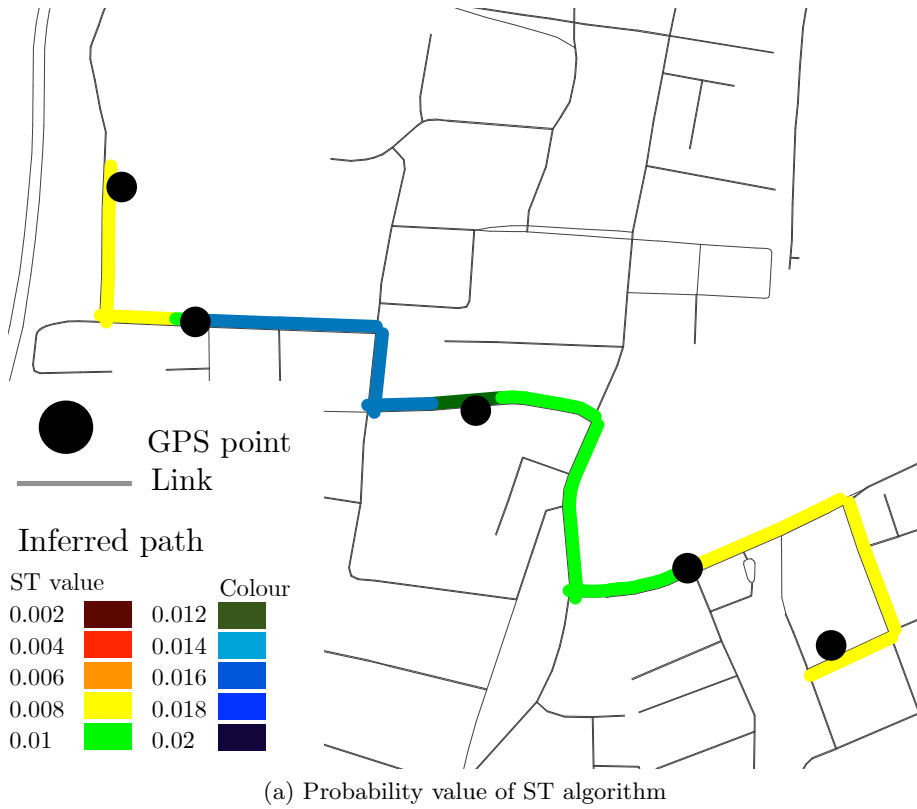


Figure 3.10: Analyses of example map-matching results

To gauge the *confidence level* of the inferred route one can analyse the assigned probabilities from the ST algorithm, as per Figure 3.10a, but without further experimentation to create a benchmark for probabilities versus accuracy, it is difficult to conclude a confidence level in the results purely based on the probabilities of different route sections. It should be noted that low speeds travelled by the vehicle would lead to lower probabilities because the temporal part of the ST algorithm analyses the calculated speed against the free speed of the network and the more it differs, the lower the value. This poses a potentially significant impact on the analysis of waste collection vehicles since they regularly travel at very low speeds while servicing an area.