

Bus Arrival Time Prediction

Computer Science Tripos – Part II

Christ's College

April 5, 2017

Proforma

Name:	Marius Latinis
College:	Christ's College
Project Title:	Bus Arrival Time Prediction
Examination:	Computer Science Tripos – Part II, July 2017
Word Count:	TODO(ml693): figure out
Project Originator:	Dr Richard Mortier
Supervisor:	Dr Richard Mortier

Original Aims of the Project

Implement an algorithm which given the most recent GPS data predicts when the bus will arrive at the future stops. Evaluate the algorithm and show that on average it predicts with error less than 80s.

Work Completed

All that has been completed appears in this dissertation.

Special Difficulties

None

Declaration

I, Marius Latinis of Christ's College, being a candidate for Part II of the Computer Science Tripos, hereby declare that this dissertation and the work described in it are my own work, unaided except as may be specified below, and that the dissertation does not contain material that has already been used to any substantial extent for a comparable purpose.

Signed [signature]

Date April 5, 2017

Contents

1	Introduction	4
2	Preparation	5
3	Implementation	6
3.1	Data Preprocessing	6
3.2	Route Detection Algorithm	7
3.2.1	Why an Algorithm is Needed	7
3.2.2	<i>FollowsPath</i> Predicate Definition	8
3.2.3	<i>FollowsPath</i> Predicate Implementation	9
3.3	Arrival Time Prediction	10
3.3.1	Basic Idea	10
3.3.2	Optimisation Nr. 1	10
3.3.3	Optimisation Nr. 2	10
3.4	Real-time System	11
4	Evaluation	12
4.1	Route Detection Evaluation	12
4.1.1	Sensitivity Score	12
4.2	Arrival Time Prediction Evaluation	14
4.2.1	Evaluation Metrics Used	14
4.2.2	All Routes Evaluation	18
4.2.3	Comparison with the Timetable	19
4.2.4	Evaluating Optimisations	21
5	Conclusion	22
	Bibliography	22
A	Links	23
B	Detailed Evaluation Results	24
B.1	Predictions for Madingley Park	24
C	Project Proposal	27

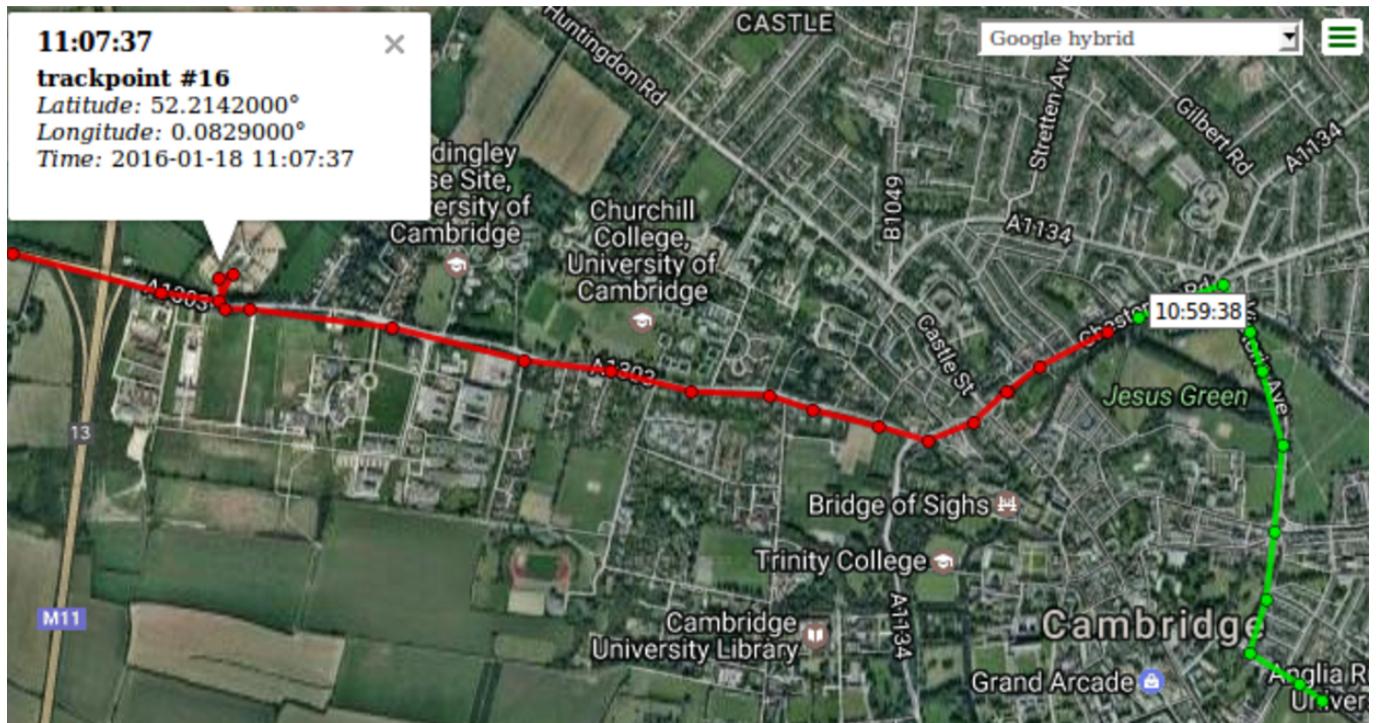
Chapter 1

Introduction

Motivation

In European countries public buses is a popular form of transportation. However, buses often arrive later than the timetable announces. Therefore, the project aims to build an algorithm which accurately predicts the bus arrival times based on the most recent GPS data.

Actual Problem Overview



The most recent GPS data (represented by green in the diagram above) shows where the bus has travelled up to now. The future data (represented by red) is not known during the prediction phase. My goal is to predict the future data given the present data.

In this specific example I want to predict when the bus will reach Madingley Park starting from Chesterton Road. One can see it takes about 8 min. (10:59–11:07) to travel such distance.

Chapter 2

Preparation

Starting Point

A company named Vix is sending real-time GPS data to Cambridge servers. The GPS data covers East-England bus companies **TODO(ml693): which???**. Cambridge converts the binary GPS data into a human readable JSON format. To start with, I was given 2 months JSON data (June and July 2016). That was the starting point data.

I started coding completely from scratch.

Preliminary Work

TODO(ml693): put something here

Chapter 3

Implementation

3.1 Data Preprocessing

Below is the starting point input data example:

```
"entities": [{"received_timestamp":1476477040,"vehicle_id":"122","label":"ATS-3603","latitude":51.91048,"longitude":-0.5006834,"bearing":150.0,"timestamp":1470092134}, {"received_timestamp":1476477040,"vehicle_id":"132","latitude":51.891663,"longitude":-0.45255125,"bearing":6.0,"timestamp":1470091865}, {"received_timestamp":1476477040,"vehicle_id":"187","latitude":51.493637,"longitude":-0.14694783,"bearing":282.0,"timestamp":1470092143}, {"received_timestamp":1476477040,"vehicle_id":"193","latitude":51.909676,"longitude":-0.5108846,"bearing":18.0,"timestamp":1470092134}, {"received_timestamp":1476477040,"vehicle_id":"194","latitude":51.882168,"longitude":-0.41519493,"bearing":114.0,"timestamp":1470092134}, {"received_timestamp":1476477040,"vehicle_id":"195","latitude":51.882168,"longitude":-0.41519493,"bearing":114.0,"timestamp":1470092134}], "received_timestamp":1476477040}
```

One file contains a single snapshot of all buses, with one entry per bus. An entry such as:

```
vehicle_id": "132", "latitude": 51.891663, "longitude": -0.45255125, "bearing": 6.0, "timestamp": 1470091865}
```

means that a vehicle nr. 56 was at geographical position $(51.88010, -0.42051104)$ on the 1st August 2016, GMT time 22:51:05¹.

A new snapshot file arrives once in 30s. The data preprocessing part takes all input files and converts them into another format:

File day15/bus300/trip2

time	latitude	longitude
2016-10-15 10:37:01	52.18485	0.20276
2016-10-15 10:37:32	52.18530	0.19788
2016-10-15 10:38:01	52.18548	0.19515
2016-10-15 10:38:32	52.18530	0.19343
2016-10-15 10:39:02	52.18557	0.19170
2016-10-15 10:39:32	52.18432	0.19070
2016-10-15 10:40:02	52.18271	0.19012
2016-10-15 10:40:32	52.18253	0.18840
2016-10-15 10:41:01	52.18297	0.18423
2016-10-15 10:41:31	52.18584	0.18509
2016-10-15 10:42:30	52.18647	0.18653

File day15/bus400/trip3

time	latitude	longitude
2016-10-15 12:01:35	52.20215	0.12460
2016-10-15 12:01:35	52.20215	0.12460
2016-10-15 12:02:35	52.19937	0.12705
2016-10-15 12:02:35	52.19937	0.12705
2016-10-15 12:03:35	52.19543	0.13064
2016-10-15 12:03:35	52.19543	0.13064
2016-10-15 12:04:34	52.19462	0.13423
2016-10-15 12:05:34	52.19444	0.13581
2016-10-15 12:06:35	52.19274	0.13581
2016-10-15 12:10:35	52.19247	0.13567
2016-10-15 12:11:35	52.19256	0.13409

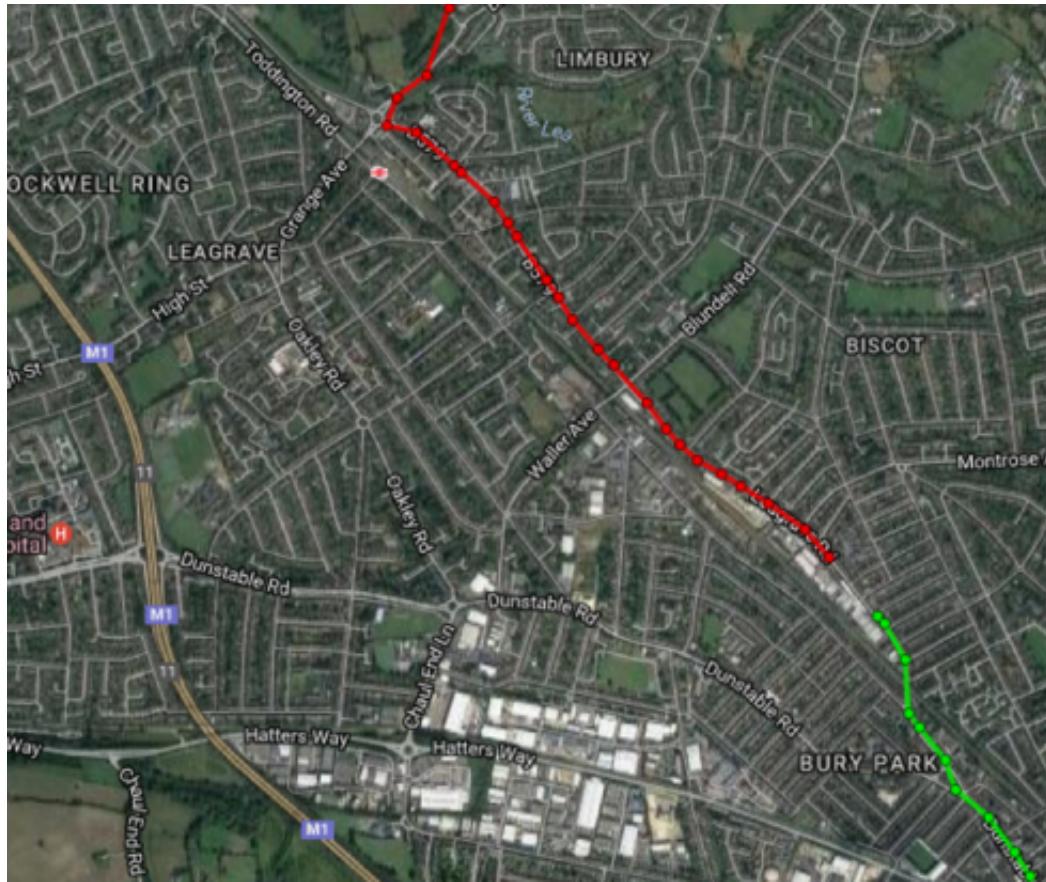
At the end I have a file for each trip made by one bus in one day².

¹calculated from timestamp field

²note that a bus on a particular day could have made multiple trips

3.2 Route Detection Algorithm

3.2.1 Why an Algorithm is Needed



The GPS points indicated in green show how the bus has moved so far. The red points show where the bus will travel in the future, a path I do not know yet. Predicting when the bus will arrive at the stop requires knowing **where** a bus will travel next. How?

Approach Using Static Data

A standard way to find out where the bus will travel next is to know its route in advance. Some GPS entries contain additional fields, such as a **label**:

```
"vehicle_id": "122", "label": "ATS-3603", "latitude": 51.90439, "longitude": -0.5070053,
```

An extra field **label** names the route a bus follows. Auxiliary data might indicate what route the label **ATS-3603** corresponds to. One can then look at the route and determine a sequence of stops a bus will visit next.

Unfortunately, this approach has drawbacks:

- Not all entries contain auxiliary fields. After all, the GPS device on the actual vehicle is responsible only for providing the latitude and longitude coordinates, not the route a bus follows.
- Good static data to infer routes based on additional fields is not straightforward to find. Field values might not match (e.g. a route named **SW-4** in the real-time data is actually **SWX-4**).

Buses are often changing routes³ and data does not always indicate that.

- The static data only shows the sequence of bus stops, not a detailed sequence of streets a bus will follow. It is much better to know an actual path for accurate arrival time prediction.

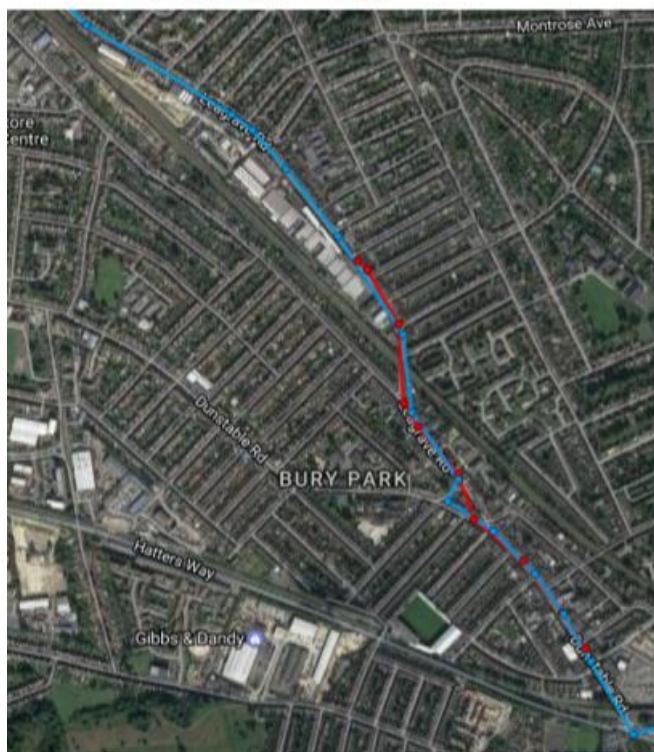
New Approach Overview

To overcome the static data limitations, I propose a new way to infer the route a vehicle is following. I will align a sub-trip how the bus has moved so far with a historical path for which the route is known. The path aligning best to the current sub-trip will be used to predict where a bus travels next. To perform the alignment, I constructed an algorithm⁴ that can tell whether one sequence of GPS points follows another.

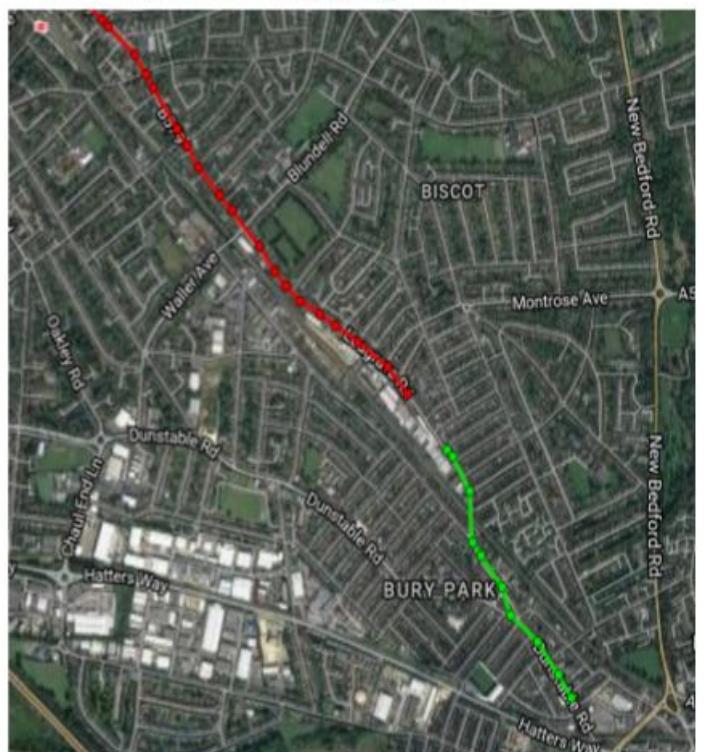
3.2.2 *FollowsPath* Predicate Definition

I define a predicate *FollowsPath*(*sub-trip*, *path*) to be true iff a *sub-trip* follows a *path*:

FollowsPath(red_subtrip, blue_trip) = true



FollowsPath(green_subtrip, red_trip) = false



Two pictures show when I want the predicate to return true and when false. For the first example, the predicate returns true because it is intuitively clear that a red sub-trip follows a blue path. In the second example, the output is false, as the two GPS sequences do not intersect at all.

³a vehicle is unlikely to follow the same route every day

⁴from now on referred to as a boolean predicate

3.2.3 *FollowsPath* Predicate Implementation

Firstly, consider the case when a **sub-trip's** GPS points (P_0, P_1, \dots, P_{n-1}) **exactly** follow another **path**:



Each point P_i is on some segment of another trip's path. For example, P_1 is on the segment AB . Let $x = |AP_1|$ and $y = |P_1B|$. Then the error⁵ function defined as

$$err(P_1, AB) = \frac{x+y}{|AB|} - 1$$

equals 0. For each point P_i , I find a segment S_{P_i} such that $err(P_i, S_{P_i})$ is minimised. The *FollowsPath* predicate returns true iff

$$S = \sum_{i=0}^{n-1} err(P_i, S_{P_i}) < n\epsilon$$

for a suitably chosen threshold value ϵ ⁶.

Note that $S = 0$ when a **sub-trip** follows a **path** exactly. Therefore, in this case, the *FollowsPath* predicate returns true, which is what I want.

The case when a **sub-trip** does follow a **path**, but not exactly (the real world situation), is illustrated in the following picture:



In this example I am aligning a **sub-trip** (R_0, \dots, R_9) with a **path** (B_0, \dots, B_7, \dots). The computation of S proceeds as follows:

$$\begin{aligned} S &= err(R_0, B_0B_1) + err(R_1, B_0B_1) + err(R_2, B_2B_3) + err(R_3, B_2B_3) + err(R_4, B_2B_3) + \\ &\quad err(R_5, B_3B_4) + err(R_6, B_4B_5) + err(R_7, B_5B_6) + err(R_8, B_6B_7) + err(R_9, B_6B_7) \\ &= (1.0294 - 1) + (1.0300 - 1) + (1.0021 - 1) + (1.0730 - 1) + (1.0208 - 1) + \\ &\quad (1.0714 - 1) + (1.0286 - 1) + (1.0158 - 1) + (1.0126 - 1) + (1.0125 - 1) < 10 \cdot 0.3 = n\epsilon \end{aligned}$$

and the predicate again returns true, as desired.

Finally, when a sub-trip does not follow a path, triangle inequality will ensure large values of $err(P_i, S_{P_i})$ ⁷, thus the sum S will exceed $n\epsilon$, and the predicate will return false.

The practical evidence that *FollowsPath* predicate works is shown in the evaluation section.

⁵error function takes a point and a segment as an input

⁶ $\epsilon = 0.3$ is one choice

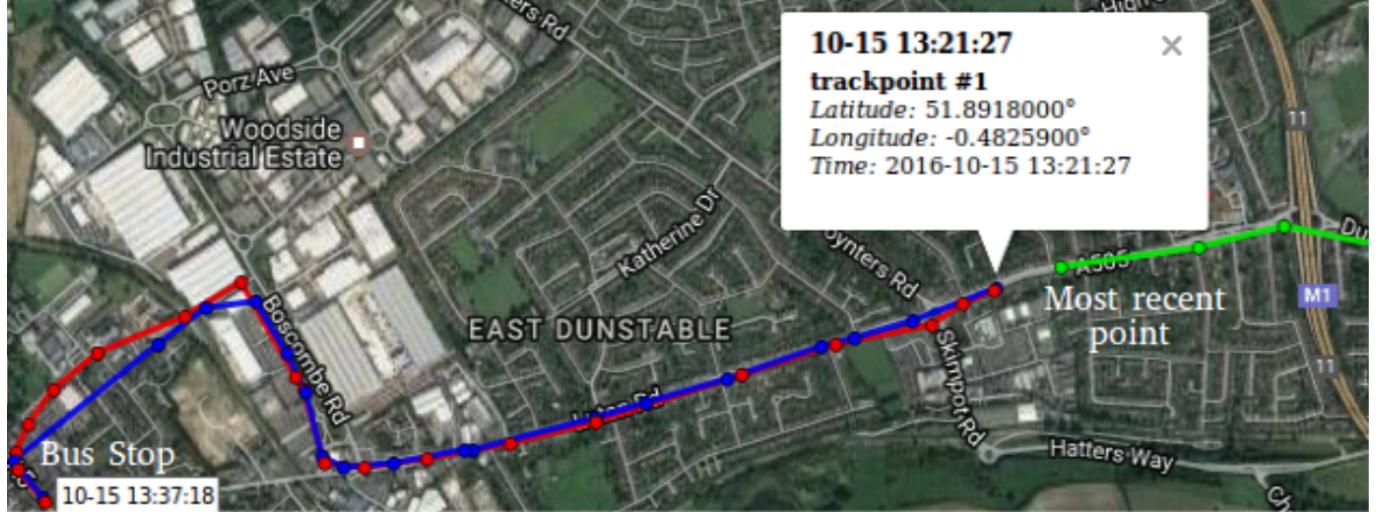
⁷ $err(P_i, S_{P_i}) = \frac{x_i+y_i}{|S_{P_i}|} - 1 \gg (1+\epsilon) - 1 = \epsilon \implies S \gg n\epsilon$

3.3 Arrival Time Prediction

3.3.1 Basic Idea

The *FollowsPath* predicate allows to detect the route a bus is following. The route tells the next stops a bus will visit. This section shows how to predict the arrival time to each future bus stop.

Historical trips (which follow the same route) are used to make a prediction:



In the picture above 2 example historical sub-trips (blue and red) are shown. Both of them start where the vehicle currently located is (indicated by the most recent green point). Both of them represent the history how past vehicles travelled until the bus stop. I compute the duration how long it took for each historical trip to reach the bus stop. Knowing multiple duration values I return the median as the predicted arrival time.

3.3.2 Optimisation Nr. 1

If in the historical trips set certain trips have **just** travelled along the same route (e.g. in the last 30min.), I predict the median time just from those recent trips. Otherwise, the basic approach is used.

3.3.3 Optimisation Nr. 2

Suppose the most recent data



shows the distance a bus has travelled in the last 16 minutes. Certain historical trips



have travelled the same distance in a very different amount of time⁸ (e.g. 4 min.). Predicting arrival time using these historical trips can be misleading. Hence I filter out such trips. I.E. I take only those historical trips into account that have travelled the same distance in roughly the same amount of time⁹.

Optimisations are tested in the evaluation section.

3.4 Real-time System

TODO(ml693): finish implementing it

⁸this could have happened due different traffic congestion level in a different time of the day

⁹up to 30sec. difference is allowed

Chapter 4

Evaluation

Add a few introductory sentences;

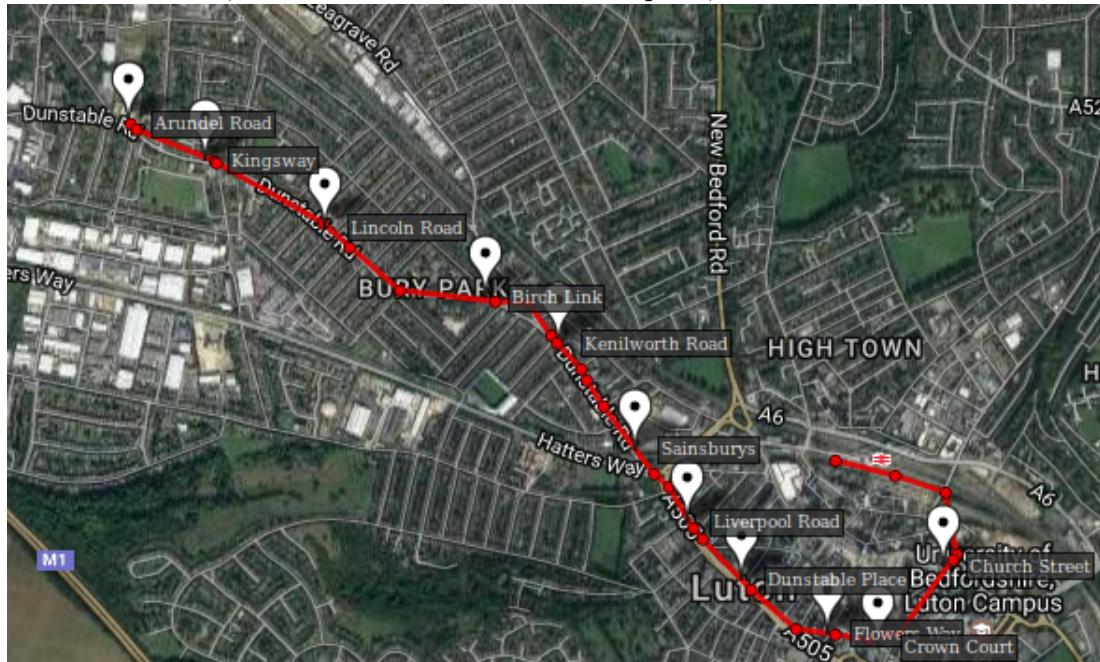
4.1 Route Detection Evaluation

I break the evaluation of route detection algorithm into multiple steps.

4.1.1 Sensitivity Score

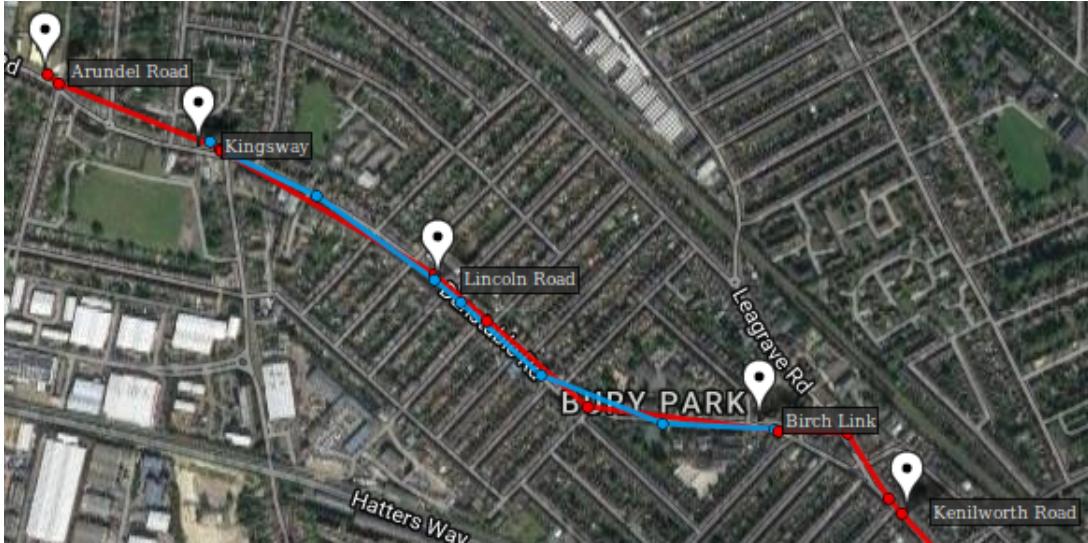
The first step is to check whether the $\text{FollowsPath}(\text{sub-trip}, \text{path})$ predicate returns *true* when the *sub-trip* indeed follows the *path*.

I have 10 routes, each route has an associated path, one of which is shown below¹:



I have a set of trips T_p that actually follow the path p . For each trip $t \in T_p$ I take a *sub-trip*

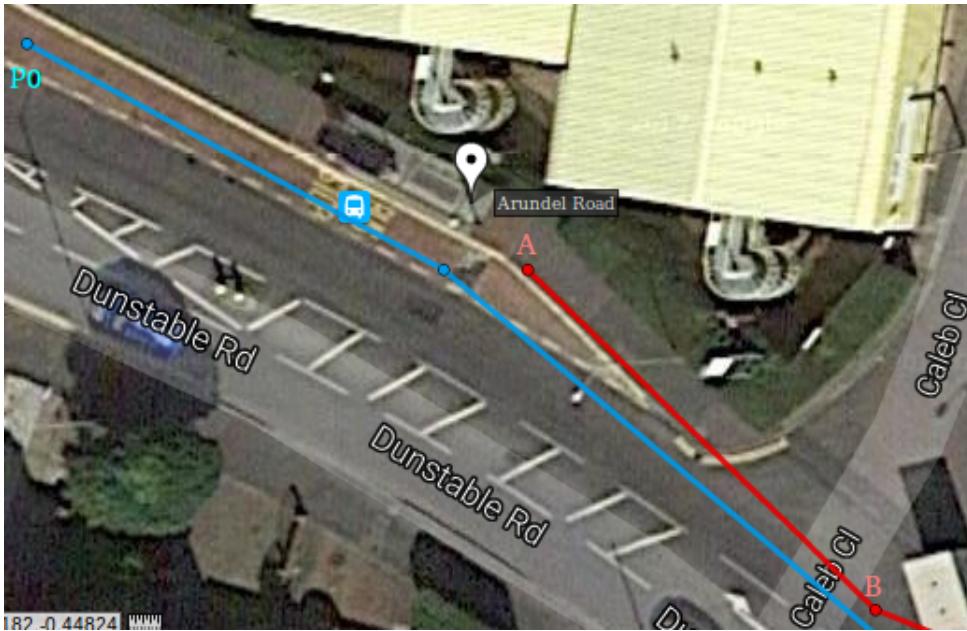
¹Arundel Road on the top left is the first bus stop here



and check whether the *FollowsPath*(*sub-trip*, *path*) predicate returns true or false.

If the *sub-trip* does **not** start from the first stop, the *FollowsPath* predicate has 100% success rate for each route.

If the *sub-trip* starts from the first stop, the success rate drops significantly (and becomes unpredictable) due to the reason illustrated in the following picture:



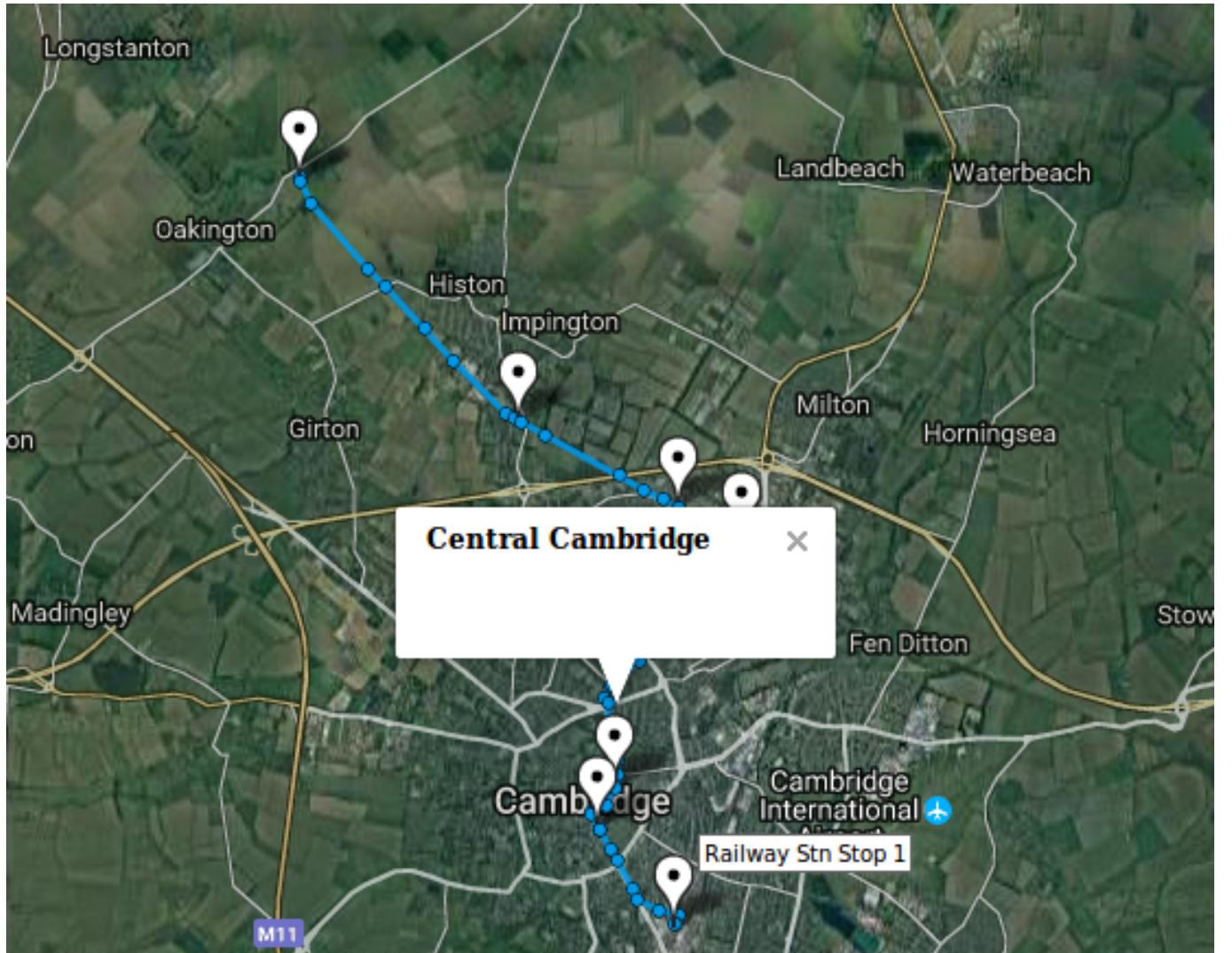
One can see that $\text{err}(P_0, S_{P_0}) = \frac{|P_0A| + |P_0B|}{|AB|} - 1$ might become large. In fact, accumulated errors for the first few points (when the bus is not on any route yet) are so large the total sum exceeds $n\epsilon$ even if errors for the next points are small.

The conclusion is thus to discard the first few GPS points and start using *FollowsPath* predicate after the bus has moved for a while and is following some route.

4.2 Arrival Time Prediction Evaluation

4.2.1 Evaluation Metrics Used

This evaluation section is the most important. It shows how well the algorithm predicts the arrival time. A few evaluation metrics are used repetitively throughout the section. Hence, I want to present them. For the presentation I will use the example route 1046531-20150531-20150830, which begins outside of Cambridge and ends at the Cambridge Railway station:



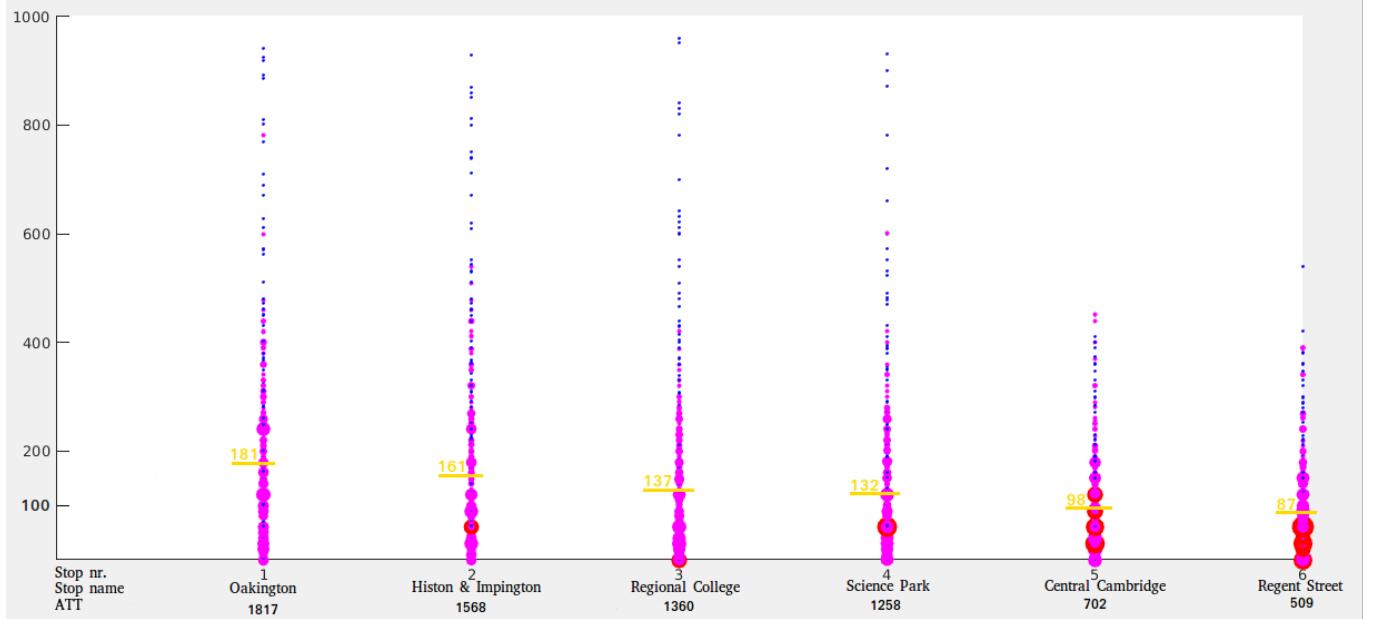
- (i) **The MAE value.** Suppose I have a set of trips T , such that each trip eventually reaches a bus stop of interest. Suppose for trip $t \in T$ the algorithm predicts the bus arrival time at stop to be $t_{predicted}$, but the bus actually arrives at t_{actual} . The absolute prediction error for trip t is $E_t = |t_{actual} - t_{predicted}|$. The **mean** absolute prediction error for the whole set T is

$$MAE = \frac{\sum_{t \in T} E_t}{|T|}$$

For example, suppose I predict when the buses will reach the Railway Station starting from the Central Cambridge stop. I have 558 trips following this route to test my prediction. For each test trip, I assume only its history up to the Central Cambridge stop is known. I predict the arrival time to the Railway Station. After the prediction is made, I then read the rest of the history and check when the trip has actually arrived at the Railway Station. For this example

trip, the resulting MAE value is **98** sec. This number shows the average error the prediction algorithm makes. One can also use the MAE value to compare 2 algorithms. I claim that the prediction algorithm A performs better for some set than the algorithm B , if $MAE_A < MAE_B$.

- (ii) **The scatter plot.** The MAE value gives a first clue how well the algorithm performs. However, just the MAE value alone is misleading. 98 seconds can be either a small or a large error. It depends how far in the future I am predicting. It takes about 10 minutes to reach the Railway Station from the Central Cambridge stop. Thus, 98 seconds is an error **with respect to** the Central Cambridge stop. The scatter plot formalises this relation:



The x-axis labels the stop number i . For each i , I predict when the bus will reach the Railway Station (the last stop) starting from the i -th stop. The y-axis labels the absolute prediction error². A point (i, y) tells that for some trip the prediction from the i -th stop to the last stop has an absolute error equal y . The more repeated errors (i, y) I have, the larger point I plot. For completeness, in the diagram I show the MAE value³ and the average travel time⁴ from the i -th stop to the last stop.

A few things can be observed in this plot:

- Firstly, the MAE value is smaller for each larger index i . The less time it takes for a bus to reach the last stop, the more accurate the prediction becomes.
- In my case, most prediction errors are smaller than the MAE value, except for a few large errors that make the MAE value increase. The large errors occur due to unusual traffic patterns:

²Measured in seconds.

³Labelled in yellow.

⁴For each trip I calculate the time it actually takes to reach the last stop, then return the mean value.



For some bus it took 4 minutes to pass just a few buildings. Since most historical trips usually pass this segment in a much shorter time⁵, the predicted arrival time to the bus stop is much smaller than the actual arrival time. Therefore, the large error value is generated.

- The scatter plot shows the absolute prediction errors only, not whether the bus arrived earlier or later than predicted. There is a simple reason for this: the prediction errors are symmetric, due to how my algorithm constructed is. On average, half of the errors are positive, half negative. Thus, it is enough to plot only the absolute errors and not lose any information.

All prediction points are shown on one plot. Hence, the scatter plot can be used as a representative summary of what the algorithm predicts.

- (iii) **The table of errors.** Traffic is different at a different time of the day. I split the day into 4 disjoint time intervals:

1. 8am–10am. The morning rush hour.
2. 10am–5pm. The working hours.
2. 5pm–8pm. The evening rush hour.
3. 8pm–8am. I call all other hours as the 'night' time.

⁵e.g. 30s

I evaluate the prediction for each of the time intervals separately. The **table of errors** shows the aggregated results:

		The morning rush hour	The working hours	The evening rush hour	The night time
From stop nr. 1	MAE	196	137	193	181
Oakington	ATT	1979	1775	1797	1740
From stop nr. 2	MAE	169	129	152	137
Histon & Impington	ATT	1722	1537	1549	1483
From stop nr. 3	MAE	160	112	145	118
Regional College	ATT	1486	1342	1346	1277
From stop nr. 4	MAE	148	104	125	109
Science Park	ATT	1379	1238	1247	1186
From stop nr. 5	MAE	99	92	93	90
Central Cambridge	ATT	745	684	715	656
From stop nr. 6	MAE	87	79	92	87
Regent Street	ATT	551	500	505	487

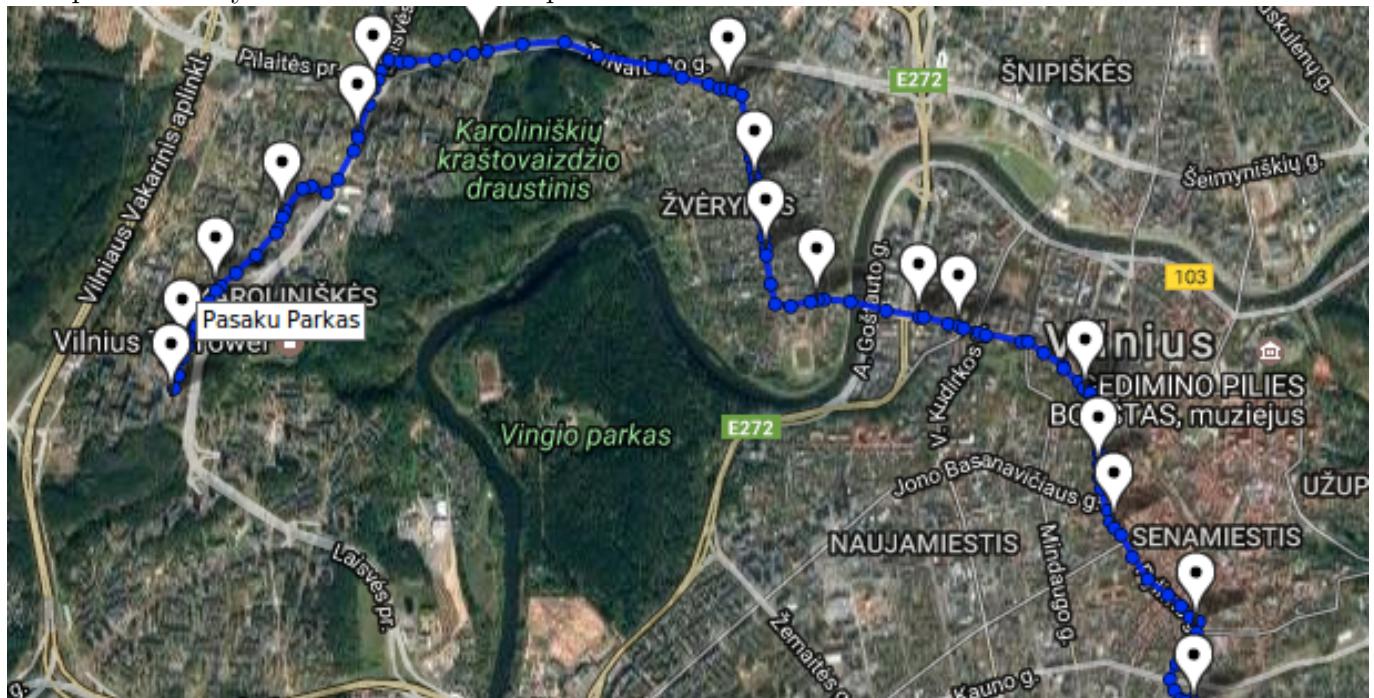
The table of errors shows the different prediction errors for different times of the day. It is useful for those who are interested in a particular route and want to check the extent up to which the prediction system can be trusted for that route. For example, suppose I want to take a morning train to London. I can look at the morning column of this table and see that on average it takes 745s. to reach the train station from the Central Cambridge. I can also see that an algorithm on average predicts the arrival time to the train station with an error of 99s.

Note that the average travel time, as well as the mean absolute prediction error values, are larger for the first and third intervals of time. Hence, for this route the traffic is more congested and less predictable during the rush hours.

4.2.2 All Routes Evaluation

4.2.3 Comparison with the Timetable

The prediction system is useful if for a particular route⁶ it is more accurate than the timetable:



		The morning rush hour	The working hours	The evening rush hour	The night time
From stop nr. 1	MAE	103/99	179/180	103/105	105/102
Stotis	ATT	1905	2035	1870	1877
From stop nr. 2	MAE	105/99	166/180	96/105	105/102
Rudininku	ATT	1722	1860	1716	1714
From stop nr. 3	MAE	86/99	155/180	94/105	95/102
Traku	ATT	1563	1682	1556	1551
From stop nr. 4	MAE	88/99	148/180	94/105	85/102
Reformatu	ATT	1476	1585	1465	1463
From stop nr. 5	MAE	84/99	133/180	93/105	88/102
Islandijos	ATT	1345	1442	1324	1333
From stop nr. 6	MAE	79/99	122/180	88/105	85/102
Pamenkalnio	ATT	1151	1242	1149	1158
From stop nr. 7	MAE	78/99	118/180	81/105	82/102
Jokubo Jasinskio	ATT	1099	1190	1098	1108
From stop nr. 8	MAE	70/99	116/180	76/105	81/102
Liubarto Tiltas	ATT	1023	1112	1021	1029
From stop nr. 9	MAE	66/99	103/180	73/105	77/102
Kestucio	ATT	930	1007	926	937
From stop nr. 10	MAE	63/99	89/180	70/105	74/102
Latviu	ATT	871	934	868	881
From stop nr. 11	MAE	56/99	76/180	69/105	71/102
Seliu	ATT	777	823	782	797
From stop nr. 12	MAE	44/99	57/180	58/105	57/102
Teodoro Narburto	ATT	557	593	575	583
From stop nr. 13	MAE	38/99	48/180	50/105	49/102
Sietyno	ATT	428	435	413	419
From stop nr. 14	MAE	39/99	42/180	47/105	45/102
Laisves Prospektas	ATT	375	380	361	367
From stop nr. 15	MAE	30/99	33/180	39/105	32/102
Karoliniškių Poliklinika	ATT	242	239	233	229
From stop nr. 16	MAE	22/99	26/180	34/105	27/102
Ugniagiesiu	ATT	152	154	152	147
From stop nr. 17	MAE	17/99	22/180	25/105	20/102
Atminties	ATT	73	77	76	72

⁶For variety, I have chosen a route from my home city Vilnius in this evaluation section. The route has 18 stops and I predict the arrival time to the last stop 'Pasaku Parkas' from each of the other stops.

I compare the prediction algorithm performance with the 'performance' of the timetable⁷. Two tables of errors are merged into one and presented here. The left number of each coloured cell shows the average prediction error of my algorithm. The right number of each coloured cell shows the average prediction error of the timetable.

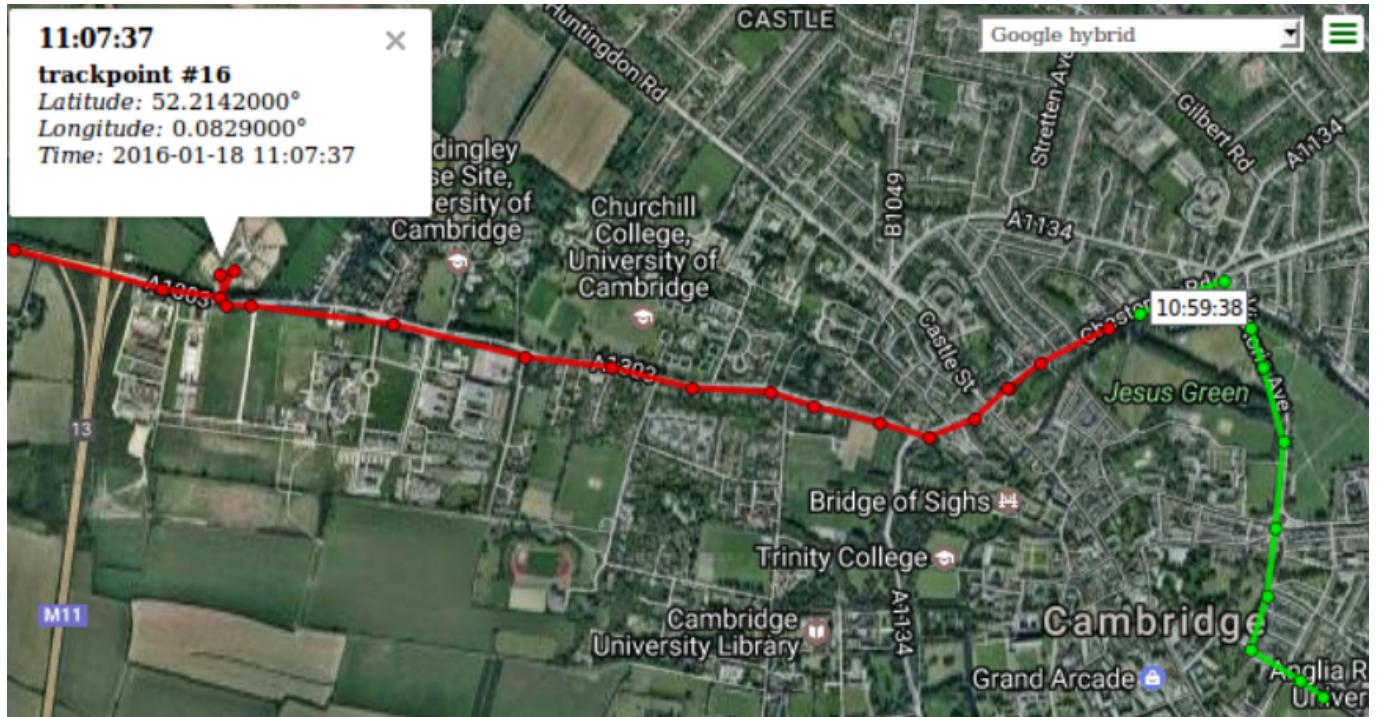
Note that for each column the right numbers are all the same. This is because the arrival time announced by the timetable does not change while the bus is moving.

Only 4 *MAE* values out of 68 are larger for my prediction algorithm. At the beginning of the route, the bus is far away from the last stop, and I have a little amount of recent information how it travels. That is the reason why the prediction errors for the first stops might be worse than what the timetable announces. However, starting from the later stops, the prediction algorithm outperforms the timetable.

⁷Assume the arrival time prediction is what the timetable announces. Thus, one can view the timetable as just another prediction algorithm. Then the *MAE* values for such an algorithm are computed as usual. I compare the *MAE* values of the 2 algorithms.

4.2.4 Evaluating Optimisations

I first run the non-optimised algorithm to predict when the bus will reach the Madingley Park starting from the Chesterton Road:



Training set having 30 trips results in $MAE_0 = 105s^8$. I then apply the 1st, and the 1st + 2nd optimisations on the same training set. The new MAE values are

$$MAE_1 = 88s < MAE_0$$

$$MAE_{1,2} = 69s < MAE_1$$

Hence, I conclude that for this route the optimisations are worth to be applied.

⁸detailed evaluation results are given in the appendix

Chapter 5

Conclusion

I implemented an algorithm which gives sensible prediction results. The algorithm itself is simple, thus anyone coming after me to improve the system should not have trouble read and understand what has been done so far.

If started again, I would have worked in exactly the same way.

Appendix A

Links

1. Dynamic time warping is an algorithm for measuring similarity between two temporal sequences which may vary in speed. For instance, similarities in walking could be detected using DTW, even if one person was walking faster than another. I used DTW algorithm as a core part to implement my *FollowsPath* predicate, with *err* function being used instead of a regular distance function $d(x, y)$. For more info look at

https://en.wikipedia.org/wiki/Dynamic_time_warping

2. Arrival time prediction problem has been studied before. Two papers that I read and compared my project with are:

- Bin Yu, William H.K. Lam, Mei Lam Tam: Bus arrival time prediction at bus stop with multiple routes

<http://www.sciencedirect.com/science/article/pii/S0968090X11000155>

- Pengfei Zhou, Yuanqing Zheng, Mo Li: How Long to Wait? Predicting Bus Arrival Time with Mobile Phone based Participatory Sensing

<http://www.ntu.edu.sg/home/limo/papers/sys012fp.pdf>

3. Project code can be found at:

<https://github.com/ml693/bus>

Appendix B

Detailed Evaluation Results

B.1 Predictions for Madingley Park

Results without optimisations:

day18-bus14365-subtrip0 started at 2016-01-18 13:45:46, arrived at 2016-01-18 13:52:06, prediction error is -80
day18-bus14368-subtrip2 started at 2016-01-18 21:14:38, arrived at 2016-01-18 21:21:18, prediction error is -60
day18-bus14366-subtrip0 started at 2016-01-18 11:46:35, arrived at 2016-01-18 11:55:36, prediction error is 81
day18-bus14369-subtrip1 started at 2016-01-18 17:46:33, arrived at 2016-01-18 17:55:13, prediction error is 59
day18-bus403-subtrip0 started at 2016-01-18 11:23:07, arrived at 2016-01-18 11:28:07, prediction error is -60
day18-bus14374-subtrip0 started at 2016-01-18 16:15:35, arrived at 2016-01-18 16:22:35, prediction error is 19
day18-bus14361-subtrip0 started at 2016-01-18 10:16:14, arrived at 2016-01-18 10:28:34, prediction error is 340
day18-bus14370-subtrip0 started at 2016-01-18 10:30:47, arrived at 2016-01-18 10:44:47, prediction error is **380**
day18-bus14365-subtrip1 started at 2016-01-18 21:44:37, arrived at 2016-01-18 21:50:17, prediction error is -80
day18-bus407-subtrip1 started at 2016-01-18 10:15:37, arrived at 2016-01-18 10:28:40, prediction error is 334
day18-bus14368-subtrip1 started at 2016-01-18 13:14:46, arrived at 2016-01-18 13:20:06, prediction error is -100
day18-bus14363-subtrip1 started at 2016-01-18 15:46:05, arrived at 2016-01-18 15:53:45, prediction error is -20
day18-bus14375-subtrip0 started at 2016-01-18 15:14:50, arrived at 2016-01-18 15:21:10, prediction error is -69
day18-bus401-subtrip2 started at 2016-01-18 11:46:37, arrived at 2016-01-18 11:55:38, prediction error is 81
day18-bus14370-subtrip3 started at 2016-01-18 19:30:41, arrived at 2016-01-18 19:38:42, prediction error is 121
day18-bus410-subtrip0 started at 2016-01-18 16:15:38, arrived at 2016-01-18 16:22:38, prediction error is 20
day18-bus402-subtrip0 started at 2016-01-18 10:30:38, arrived at 2016-01-18 10:44:37, prediction error is **379**
day18-bus14373-subtrip0 started at 2016-01-18 16:46:02, arrived at 2016-01-18 16:53:42, prediction error is -1
day18-bus14372-subtrip1 started at 2016-01-18 14:14:19, arrived at 2016-01-18 14:20:19, prediction error is -120
day18-bus14367-subtrip2 started at 2016-01-19 00:35:49, arrived at 2016-01-19 00:40:49, prediction error is -160
day18-bus14378-subtrip0 started at 2016-01-18 11:00:21, arrived at 2016-01-18 11:07:41, prediction error is 20
day18-bus402-subtrip1 started at 2016-01-18 19:30:08, arrived at 2016-01-18 19:38:37, prediction error is 89
day18-bus400-subtrip0 started at 2016-01-18 10:59:38, arrived at 2016-01-18 11:07:37, prediction error is 29
day18-bus14087-subtrip0 started at 2016-01-18 09:30:34, arrived at 2016-01-18 09:36:34, prediction error is -120
day18-bus409-subtrip0 started at 2016-01-18 08:25:38, arrived at 2016-01-18 08:32:09, prediction error is -29
day18-bus14371-subtrip1 started at 2016-01-18 08:54:47, arrived at 2016-01-18 09:01:07, prediction error is -70
day18-bus14072-subtrip3 started at 2016-01-18 20:51:07, arrived at 2016-01-18 20:57:47, prediction error is -79
day18-bus14072-subtrip1 started at 2016-01-18 12:45:18, arrived at 2016-01-18 12:53:18, prediction error is 20
day18-bus14362-subtrip1 started at 2016-01-18 14:45:16, arrived at 2016-01-18 14:50:36, prediction error is -81
day18-bus405-subtrip2 started at 2016-01-18 13:46:07, arrived at 2016-01-18 13:52:08, prediction error is -59

Results with the optimisation nr. 1:

day18-bus14365-subtrip0 started at 2016-01-18 13:45:46, arrived at 2016-01-18 13:52:06, prediction error is -80
day18-bus14368-subtrip2 started at 2016-01-18 21:14:38, arrived at 2016-01-18 21:21:18, prediction error is 0
day18-bus14366-subtrip0 started at 2016-01-18 11:46:35, arrived at 2016-01-18 11:55:36, prediction error is 81
day18-bus14369-subtrip1 started at 2016-01-18 17:46:33, arrived at 2016-01-18 17:55:13, prediction error is 59
day18-bus403-subtrip0 started at 2016-01-18 11:23:07, arrived at 2016-01-18 11:28:07, prediction error is -60
day18-bus14374-subtrip0 started at 2016-01-18 16:15:35, arrived at 2016-01-18 16:22:35, prediction error is 19
day18-bus14361-subtrip0 started at 2016-01-18 10:16:14, arrived at 2016-01-18 10:28:34, prediction error is 340
day18-bus14370-subtrip0 started at 2016-01-18 10:30:47, arrived at 2016-01-18 10:44:47, prediction error is **40**
day18-bus14365-subtrip1 started at 2016-01-18 21:44:37, arrived at 2016-01-18 21:50:17, prediction error is -60
day18-bus407-subtrip1 started at 2016-01-18 10:15:37, arrived at 2016-01-18 10:28:40, prediction error is 334
day18-bus14368-subtrip1 started at 2016-01-18 13:14:46, arrived at 2016-01-18 13:20:06, prediction error is -100
day18-bus14363-subtrip1 started at 2016-01-18 15:46:05, arrived at 2016-01-18 15:53:45, prediction error is -20
day18-bus14375-subtrip0 started at 2016-01-18 15:14:50, arrived at 2016-01-18 15:21:10, prediction error is -69
day18-bus401-subtrip2 started at 2016-01-18 11:46:37, arrived at 2016-01-18 11:55:38, prediction error is 81
day18-bus14370-subtrip3 started at 2016-01-18 19:30:41, arrived at 2016-01-18 19:38:42, prediction error is 121
day18-bus410-subtrip0 started at 2016-01-18 16:15:38, arrived at 2016-01-18 16:22:38, prediction error is 20
day18-bus402-subtrip0 started at 2016-01-18 10:30:38, arrived at 2016-01-18 10:44:37, prediction error is **39**
day18-bus14373-subtrip0 started at 2016-01-18 16:46:02, arrived at 2016-01-18 16:53:42, prediction error is -1
day18-bus14372-subtrip1 started at 2016-01-18 14:14:19, arrived at 2016-01-18 14:20:19, prediction error is -31
day18-bus14367-subtrip2 started at 2016-01-19 00:35:49, arrived at 2016-01-19 00:40:49, prediction error is -160
day18-bus14378-subtrip0 started at 2016-01-18 11:00:21, arrived at 2016-01-18 11:07:41, prediction error is -200
day18-bus402-subtrip1 started at 2016-01-18 19:30:08, arrived at 2016-01-18 19:38:37, prediction error is 89
day18-bus400-subtrip0 started at 2016-01-18 10:59:38, arrived at 2016-01-18 11:07:37, prediction error is -241
day18-bus14087-subtrip0 started at 2016-01-18 09:30:34, arrived at 2016-01-18 09:36:34, prediction error is -120
day18-bus409-subtrip0 started at 2016-01-18 08:25:38, arrived at 2016-01-18 08:32:09, prediction error is -29
day18-bus14371-subtrip1 started at 2016-01-18 08:54:47, arrived at 2016-01-18 09:01:07, prediction error is -42
day18-bus14072-subtrip3 started at 2016-01-18 20:51:07, arrived at 2016-01-18 20:57:47, prediction error is -79
day18-bus14072-subtrip1 started at 2016-01-18 12:45:18, arrived at 2016-01-18 12:53:18, prediction error is 20
day18-bus14362-subtrip1 started at 2016-01-18 14:45:16, arrived at 2016-01-18 14:50:36, prediction error is -81
day18-bus405-subtrip2 started at 2016-01-18 13:46:07, arrived at 2016-01-18 13:52:08, prediction error is 41

Results with the optimisations nr. 1 and nr. 2:

day18-bus14365-subtrip0 started at 2016-01-18 13:45:46, arrived at 2016-01-18 13:52:06, prediction error is -100
day18-bus14368-subtrip2 started at 2016-01-18 21:14:38, arrived at 2016-01-18 21:21:18, prediction error is -20
day18-bus14366-subtrip0 started at 2016-01-18 11:46:35, arrived at 2016-01-18 11:55:36, prediction error is 41
day18-bus14369-subtrip1 started at 2016-01-18 17:46:33, arrived at 2016-01-18 17:55:13, prediction error is 40
day18-bus403-subtrip0 started at 2016-01-18 11:23:07, arrived at 2016-01-18 11:28:07, prediction error is -59
day18-bus14374-subtrip0 started at 2016-01-18 16:15:35, arrived at 2016-01-18 16:22:35, prediction error is 0
day18-bus14361-subtrip0 started at 2016-01-18 10:16:14, arrived at 2016-01-18 10:28:34, prediction error is 280
day18-bus14370-subtrip0 started at 2016-01-18 10:30:47, arrived at 2016-01-18 10:44:47, prediction error is 40
day18-bus14365-subtrip1 started at 2016-01-18 21:44:37, arrived at 2016-01-18 21:50:17, prediction error is -51
day18-bus407-subtrip1 started at 2016-01-18 10:15:37, arrived at 2016-01-18 10:28:40, prediction error is 302
day18-bus14368-subtrip1 started at 2016-01-18 13:14:46, arrived at 2016-01-18 13:20:06, prediction error is -100
day18-bus14363-subtrip1 started at 2016-01-18 15:46:05, arrived at 2016-01-18 15:53:45, prediction error is -20
day18-bus14375-subtrip0 started at 2016-01-18 15:14:50, arrived at 2016-01-18 15:21:10, prediction error is -42
day18-bus401-subtrip2 started at 2016-01-18 11:46:37, arrived at 2016-01-18 11:55:38, prediction error is 60

day18-bus14370-subtrip3 started at 2016-01-18 19:30:41, arrived at 2016-01-18 19:38:42, prediction error is 122
day18-bus410-subtrip0 started at 2016-01-18 16:15:38, arrived at 2016-01-18 16:22:38, prediction error is 20
day18-bus402-subtrip0 started at 2016-01-18 10:30:38, arrived at 2016-01-18 10:44:37, prediction error is 39
day18-bus14373-subtrip0 started at 2016-01-18 16:46:02, arrived at 2016-01-18 16:53:42, prediction error is -20
day18-bus14372-subtrip1 started at 2016-01-18 14:14:19, arrived at 2016-01-18 14:20:19, prediction error is -31
day18-bus14367-subtrip2 started at 2016-01-19 00:35:49, arrived at 2016-01-19 00:40:49, prediction error is -180
day18-bus14378-subtrip0 started at 2016-01-18 11:00:21, arrived at 2016-01-18 11:07:41, prediction error is 20
day18-bus402-subtrip1 started at 2016-01-18 19:30:08, arrived at 2016-01-18 19:38:37, prediction error is 108
day18-bus400-subtrip0 started at 2016-01-18 10:59:38, arrived at 2016-01-18 11:07:37, prediction error is 29
day18-bus14087-subtrip0 started at 2016-01-18 09:30:34, arrived at 2016-01-18 09:36:34, prediction error is -60
day18-bus409-subtrip0 started at 2016-01-18 08:25:38, arrived at 2016-01-18 08:32:09, prediction error is -9
day18-bus14371-subtrip1 started at 2016-01-18 08:54:47, arrived at 2016-01-18 09:01:07, prediction error is -42
day18-bus14072-subtrip3 started at 2016-01-18 20:51:07, arrived at 2016-01-18 20:57:47, prediction error is -81
day18-bus14072-subtrip1 started at 2016-01-18 12:45:18, arrived at 2016-01-18 12:53:18, prediction error is -1
day18-bus14362-subtrip1 started at 2016-01-18 14:45:16, arrived at 2016-01-18 14:50:36, prediction error is -80
day18-bus405-subtrip2 started at 2016-01-18 13:46:07, arrived at 2016-01-18 13:52:08, prediction error is -99

Appendix C

Project Proposal

(see next page)

Bus Arrival Time Prediction

Computer Science Tripos - Part II - Project Proposal

Marius Latinis (ml693@cam.ac.uk)

Project Originator: Dr. Richard Mortier

Project Supervisor: Dr. Richard Mortier

Director of Studies: Prof. Ian Leslie

Project Overseers: Dr. Markus Kuhn & Prof. Peter Sewell

Document Creation Date: 13 October 2016

Introduction

In European countries public buses is a popular form of transportation. However, buses often arrive later than announced at the stop. This annoys passengers and makes them complain. Wouldn't it be nice to have a good algorithm that predicts the bus arrival times precisely? This project targets such question.

Starting Point

Communication with buses streaming GPS data is already established. A file showing the GPS snapshot of 2 months (June and July) will be given to me as a starting data to work with. Real time bus GPS data will be presented for the optional extension.

According to Dr. Lewis, the prepared and convenient bus GPS data files to work with are written in JSON format. They are roughly in one to one correspondence with the GTFS-realtime files, which are the files actually sent from the buses. Hence, I will work with JSON format as it is more readable and easier to process. That's the starting point data.

Work to be done

The project breaks down into the following parts:

1. Prepare maps for processing. A graph model of the map is considered in this project. The streets will be represented as edges with vertices being their intersections. One needs to **find** a map that contains streets covering the bus stops of consideration.

2. Prepare bus GPS data for processing. Bus GPS data also has to be prepared. The short description of what the bus data looks like according to Dr. Lewis is as follows:

"Buses announce their location data using GTFS-realtime format once every 30s. GTFS-realtime format is then converted to a more readable JSON file. Each file contains information (current GPS location, unique bus identifier, bus route identifier and more) of around 1000 buses."

We hope that **30s** granularity will be enough to estimate the time it takes for a bus to travel from one intersection to the next intersection. Since it is unlikely that the bus will announce its data exactly at the intersection, we will have to interpolate the data.

The **core** project's part is to work only with buses in Cambridge and its neighbourhood. However, there might exist buses in the data that are out of the project's scope (i.e. a bus travelling to London). These will have to be filtered out.

3. Prediction algorithm implementation. Based on the data extracted in paragraphs (1) and (2), the algorithm has to inform when the next bus arrives to which stop. Suppose we are aiming to predict when a bus β being in a current location l_1 will reach the location l_n after having travelled through locations l_2, l_3, \dots, l_{n-1} in between (location is a vertex in a graph). For each i we average the most recent time $t_{i, i+1}$ taken for other buses to travel from l_i to l_{i+1} . Then the predicted time for β to reach l_n will be a sum $\sum_{i=1}^{n-1} t_{i, i+1}$.

Research concerning arrival time prediction has already been done in another paper¹. The new prediction algorithm will differ in a few ways. Firstly, that paper essentially predicts a single value $t_{1, n}$ and outputs it as a result. Instead, we will calculate multiple intermediate values $t_{i, i+1}$, each of which will correspond to a predicted time for a bus to travel across one route's edge. We hypothesize that accurately predicting time for shorter segments and then summing all times up will lead to a better precision (read part 3 how the hypothesis will be tested). Secondly, the paper uses general machine learning models as a tool to predict arrival time. We will use calculations specific to the arrival prediction problem. That will ease the job of refining the algorithm in case it is not working well enough (see the first extension as an example how). It will also make the algorithm's optimisation task easier, as the algorithm itself targets a specific problem, rather than being a general ML tool.

// The pseudocode for the prediction algorithm

For each bus β streaming GPS data:

- a) Extract l_1 - the current location of β
- b) Look at β route to see the next locations l_2, \dots, l_n
// The (c) part will be implemented following the summation idea above
- c) Predict β arrival times t_2, \dots, t_n for each l_2, \dots, l_n and store these predictions
in a multimap data structure $Bus[\beta] \text{Arrives}\{(l_2, t_2), \dots (l_n, t_n)\}$

4. Prediction algorithm evaluation. Mathematical evaluation will be measured based on how well the algorithm predicts new arrival times. We will compare the prediction results with the actual time after we know when the bus has arrived. We want to use an evaluation metric similar to what's used elsewhere² (also known as the MAE_t value):

"In the trace-driven study, we [...] compare the predicted arrival time with the actual arrival time of the campus buses to compute the average of the absolute prediction error."

Various error results are claimed (based on the bus distance to the stop) with **80s** being one of the values in that paper. Hence, we set a goal to build a model that would produce results with an error smaller than **80s**. We will compute the MAE_t with $t_{1,n}$ being replaced by the sum $\sum_{i=1}^{n-1} t_{i,i+1}$. If the new MAE_t value is smaller than **80s**, we claim that the **hypothesis** and hence the project was successful.

However, just assessing the project based on one number does not look convincing. To improve the assessment we propose to do a more in depth **quantitative evaluation**. Rather than computing a one global MAE_t value, we will also compute the MAE_t for each bus route and for each segment. In addition to that, we will count the total number of buses that arrived more than **80s** later, regardless of how late they arrived. All of these numbers will be reported in the progress report and final dissertation. The aim is to provide a few interesting metrics that will allow other researchers to have a choice in case they want to optimise for a particular one.

Finally, one more way to evaluate the project is by analysing how easily the algorithm can be run on a large scale. We will present this analysis in the dissertation by reporting its current complexity (w.r.t. the graph and JSON files size) and giving a clue what would be the sequential algorithm's part if run on multiple machines. That will provide others an idea of whether it is worthful to extend the project. However, **we would like to**

emphasize that scalability is not the main concern of the project, so going into too much detail here would be an overkill.

These four parts together form the core of the project.

Possible extensions

5. Optimise the prediction algorithm from the precision point of view. An advise how to do that was given by the overseer Dr. Markus Kuhn:

"I would start with first estimating average edge times as a function of time of day and time of week, as load on the road network is usually a quite periodic process. I would then look into a textbook on standard multivariate statistical analysis for algorithm ideas on how to exploit statistical dependencies between spatially and temporally related edge traversal times. For example you could determine covariance matrices that describe how recently measured time intervals correlate with time intervals in the near future on either the same edge or on other nearby edges in your route graph. These averages and correlation matrices could then be used to calculate conditional probability distributions for traversal times in the near future, which could be added to make predictions. If you try to independently estimate the probability distribution of traversal time for each of thousands of edges, you may end up with quite noisy predictions, if you have more parameters to estimate than there are data points. This is where dimensionality-reduction techniques (such as principal component analysis) can become useful, where you try to represent the large vectors of edge traversal times that you try to predict as linear combinations of a much smaller number of base vectors (for principal component analysis, these will be eigenvectors of the covariance matrices), which still represent well the overall traffic situation in e.g. a town, but with far fewer parameters to estimate. Fewer parameters are easier to estimate from limited data, and hence allow faster updates."

6. Turning the static GPS data analysis into the real time processing. If our prediction algorithm gives promising results on the static data, extending it to work with the real time data would allow people to see when their desired bus actually arrives at the stop. This will involve writing a batch processing job that takes GPS files as input and continuously reports new arrival times as output. According to Dr. Lewis, the processing will be similar to the processing of stock transactions. Hence, we will split this task into the learning (i.e. surfing the net of how transactions are processed) and implementation parts.

7. Building a webpage to display prediction results. Given that we have predicted results, it would be nice to show them to the public. That's what the webpage will be used for. The aim is to have one stateless page where all the info is placed. Below is the example of the webpage content,

Station's Name	Next buses to arrive and arrival times
Covent Garden	<u>No. 2 in 4min., No. 1 in 7min.</u>
Leicester Square	<u>No. 1 in 4min., No. 8 in 6min., No. 7 in 10min.</u>

Resources required

The main resource required is “Bus GPS position data via Ian Lewis”.

Other resources are my personal laptop running Ubuntu, GitHub, Google Docs, Hermes emailing system, SRCF.

Success Criteria

The project will be a success if I have implemented a working prediction algorithm that achieves **error** ($\text{error} = \text{average absolute error of real time prediction values}$) less than **80s**. In addition to this main error value, we will plot an in depth **quantitative evaluation** statistics. Good results of other numbers will reinforce the main success criterion.

Timetable (planned starting date is 21/10/2011)

- Michaelmas term weeks 3 - 4: find maps open source data and make a graph model of it. A milestone for this part is to have a file containing a nice graph representation of Cambridge bus routes.
- Michaelmas term weeks 5 - 7: write a function which takes a JSON file as an input, extracts every bus with its GPS position, and for each bus finds the closest vertex where the bus currently is. One will also need to filter out those buses that

are present in the file but not in the project's scope (e.g. a bus travelling from Cambridge to London). Furthermore, some data might turn out to be garbage (e.g. bus presence along the wrong route) and this has to be taken into account. As it is important to prepare the main GPS data well, I am therefore giving 3 weeks for this part. This extra time will also help to rethink about the size of the graph once I started using it.

- Michaelmas term week 8: implement an “empty” framework. The framework will take JSON files and Cambridge map as an input and run the fake prediction algorithm on the input. The unimplemented fake part is (roughly) a function that takes a list of numbers as an input (the time it takes for other buses to travel across an edge) and produces a single number as an output (the predicted arrival time). The milestone here is to make sure that the framework itself is running, even if the results it displays are wrong.
- Christmas holidays weeks 1 - 3: replace the fake algorithm by the core prediction function. Evaluate the algorithm by computing the error numbers. If an error is too large, attempt to find a simple refinement of the core algorithm.
- Christmas holidays week 4: gather the best current algorithm statistics and write a progress report. **These results will be given for the midpoint presentation. Note that if everything is OK up to here, that means the core part of the project has been completed.** For the Lent term I am planning to work on extensions and am giving a further timetable:
- Lent weeks 1 - 3: work on optimising the algorithm from the precision point of view. One week can be spent reading about multivariate normal distribution, as it is likely to be used. Second week would be allocated for refining the algorithm. Third week would be spent comparing the new results with the best we have so far.
- Lent week 4: read how the stock data is processed in real time.
- Lent weeks 5 - 7: based on the info read implement the real time processing system.
- Lent week 8: create a webpage to display the results.
- Easter vacation: double check the work done before, write the first dissertation's draft.
- Easter term weeks 1 - 2: create the final dissertation's version.
- Easter term week 3: submit the dissertation.

Literature Appendix

1. Bus arrival time prediction at bus stop with multiple routes - Bin Yu ,William H.K. Lam, Mei Lam Tam
(<http://www.sciencedirect.com/science/article/pii/S0968090X11000155>)
2. How Long to Wait?: Predicting Bus Arrival Time with Mobile Phone based Participatory Sensing - Pengfei Zhou, Yuanqing Zheng, Mo Li
(<http://www.ntu.edu.sg/home/limo/papers/sys012fp.pdf>)
3. Bus Based Probe Urban Travel Time Prediction - Wenjing Pu(<https://books.google.co.uk/books?id=Tdps9w5jHKEC&pg=PA117&lpg=PA117&dq=multivariate+normal+distribution+for+bus+prediction&source=bl&ots=p5Aqeyo2NO&sig=MI2XsklclPBf9ChuR7ZdJLX2Qmc&hl=en&sa=X&ved=0ahUKEwj8soTzuenPAhVFBMAKHcC2CM8Q6AEIMzAD#v=onepage&q&f=false>)
4. Applied Multivariate Statistical Analysis - R. Johnson, D. Wichern
(<https://www.pearsonhighered.com/program/Johnson-Applied-Multivariate-Statistical-Analysis-6th-Edition/PGM274834.html>)
5. Multivariate Normal Distribution - Wikipedia
(https://en.wikipedia.org/wiki/Multivariate_normal_distribution)