

# Report for Urban Computing project2

Xiang Xinye  
SCSE

Nanyang Technological University  
Singapore  
XI0001YE@e.ntu.edu.sg

Tan Jie Heng Alfred  
SCSE

Nanyang Technological University  
Singapore  
TANJ0307@e.ntu.edu.sg

Yin Wenqi  
SCSE

Nanyang Technological University  
Singapore  
S220157@e.ntu.edu.sg

**Abstract**—In this study we investigate the taxi trajectory patterns within Porto, Portugal, using the city's intricate road network as a foundation. The research consists of six pivotal stages: sourcing the road network and trajectory data, visualizing GPS points, overlaying trajectories onto the road network, enhance the visualization of the route, in-depth route analysis, and refining map matching algorithms. Through these steps, we seek to unravel the complex interplay between taxi movements and urban road layouts. This exploration not only unveils transportation trends but also augments map matching methodologies, paving the way for advanced urban mobility studies in the future.

## I. INTRODUCTION

Urban road networks are a intricate web of roads, merging together for urban traffic to flow seamlessly from one destination to the next. These networks tend to be a complex structure, making map matching a challenging task.

Map matching is the process of aligning a sequence of observed user positions to the road network on the map. These positions are typically received as GPS signals. However, because of the complex nature of maps, there could be multiple possible matches for a single observed position. In aggregation, this could lead to multiple possible trajectories for a given sequence of observed positions. In this project, we attempt to preprocess the position data obtained from ??, before matching these positions to the map using the open-source library *Fast Map Matching*. Finally, with the routes obtained from the map matching, we shall do further analysis to glean some insights into the traffic network within Porto city.

## II. TASK 1: DATA PREPARATION

In this section, we discuss the data sources utilized in this study and elucidate the process of preparing the data for rigorous analysis.

### A. Harnessing OpenStreetMap Data

OpenStreetMap (OSM), a free-to-use geographic database, serves as the primary source of geographic information in this project. In the context of this study, where we focus on Porto City in Portugal, the OSMnx library facilitates extraction and manipulation of data, such as the map data of Porto City. The following code snippet shows the steps to download and visualize the road network of Porto City:

Listing 1: Acquiring and visualizing the map.

```
1 place = "Porto, Portugal"
```

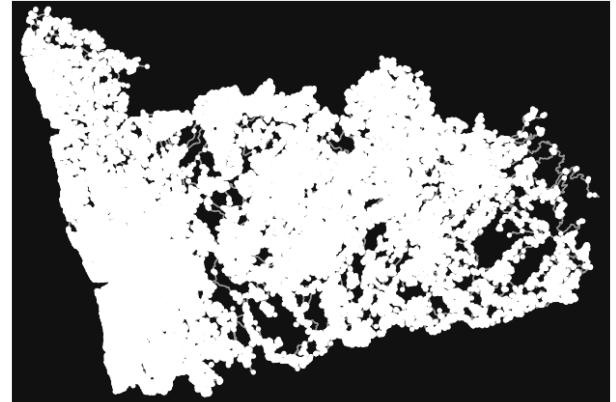


Fig. 1: Visualization of Porto City's road network.

```
2 # Procure the road network data
3 start_time = time.time()
4 G = ox.graph_from_place(place, network_type
   ='drive', which_result=2)
5 print("Time taken:", time.time() -
      start_time, "seconds")
6 # Render and archive the road network
7 ox.plot_graph(G)
8 ox.save_graphml(G)
```

Through OSMnx, the road network data was easily retrieved and stored, allowing us to perform in-depth network analysis. Fig. 1 shows the visualization of the Porto city's road network. The road network has a graphical structure: Each vertex represents a junction or intersection in a road network, while each edge represents a path that connect these intersections. There could therefore be multiple edges between each vertex. From the visualization in 1, we can see that the road network is very dense, with many intersections within the entire Porto City.

### B. Taxi Trajectory Data on Kaggle

In this subsection, we elaborate on the taxi trajectory dataset sourced from Kaggle. The tracjectory data used in this research is curated from a Kaggle competition [1], which is only accessible after the creation of a Kaggle account. Within the dataset, there are taxi trajectories, accompanied with other metadata such as whether any GPS location has been dropped (i.e., missing\_data). For the purposes of this project, we only

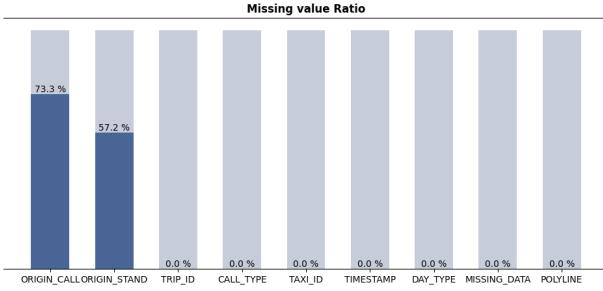


Fig. 2: Analysis missing data.

work with a subset of the first 1500 trips from the train.csv of the dataset. To this end, we extracted the first 1500 taxi trajectories and their associated metadata, placing them into a new CSV file, which we subsequently named as *train-1500.csv*. In the ensuing sections of this study, we will analyze and work on the smaller train-1500 dataset.

### III. TASK 2: GPS POINT VISUALIZATION

Before further analysis, it will be helpful to first visualize the taxi trajectories and gain a better understanding of our dataset. In this section, we overlay the GPS points from the taxi trajectory dataset onto the road network of Porto to visualize the some taxi trajectories. We used the OSMnx library, a powerful tool tailored for OpenStreetMap data, to visualize only the first 15 taxi trips, so as to prevent cluttering in the visualization.

#### A. Data Preparation for Trajectory Boundaries

Prior to visualization, it is essential to preprocess the data. From our dataset, we extract the first 15 trajectories, before computing their encompassing geographical boundaries (delineated by the area these trajectories cover). The subsequent code gives a glimpse into this data processing:

Listing 2: Extraction and boundary calculation for the first 15 trajectories.

```

1 df_15 = pd.read_csv('data/train-1500.csv',
2                     nrows=15)
3 trajectories = []
4 trajectory_count = 0
5 for polyline in df_15['POLYLINE']:
6     polyline = polyline[2:-2]
7     trajectory = []
8     if len(polyline) <= 2:
9         print(f'Skip empty row')
10    else:
11        for coordinate in polyline.split(
12            ','):
13            longitude, latitude = map(float,
14                coordinate.strip().split(
15                    ',')))
16            trajectory.append((longitude,
17                latitude))
18 x_max, y_max, x_min, y_min,
19             is_first =
20             calculate_boundaries(x_max,
21                 y_max, x_min, y_min,
22                 is_first,
23                 longitude,
24                 latitude)
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
279
280
281
282
283
284
285
286
287
288
289
289
290
291
292
293
294
295
296
297
297
298
299
299
300
301
302
303
304
305
306
307
308
309
309
310
311
312
313
314
315
316
317
317
318
319
319
320
321
322
323
324
325
326
327
327
328
329
329
330
331
332
333
334
335
336
337
337
338
339
339
340
341
342
343
344
345
346
347
347
348
349
349
350
351
352
353
354
355
356
357
357
358
359
359
360
361
362
363
364
365
366
367
367
368
369
369
370
371
372
373
374
375
376
377
377
378
379
379
380
381
382
383
384
385
386
387
387
388
389
389
390
391
392
393
394
395
396
397
397
398
399
399
400
401
402
403
404
405
406
407
407
408
409
409
410
411
412
413
414
415
416
416
417
418
418
419
420
421
422
423
424
425
426
426
427
428
428
429
430
431
432
433
434
435
436
436
437
438
438
439
440
441
442
443
444
445
446
446
447
448
448
449
450
451
452
453
454
455
456
456
457
458
458
459
460
461
462
463
464
465
466
466
467
468
468
469
470
471
472
473
474
475
476
476
477
478
478
479
480
481
482
483
484
485
486
486
487
488
488
489
490
491
492
493
494
495
496
496
497
498
498
499
500
501
502
503
504
505
506
506
507
508
508
509
510
511
512
513
514
515
515
516
517
517
518
519
519
520
521
522
523
524
525
526
526
527
528
528
529
530
531
532
533
534
535
536
536
537
538
538
539
540
541
542
543
544
545
545
546
547
547
548
549
549
550
551
552
553
554
555
556
556
557
558
558
559
559
560
561
562
563
564
565
565
566
567
567
568
569
569
570
571
572
573
574
575
575
576
577
577
578
579
579
580
581
582
583
584
585
585
586
587
587
588
589
589
590
591
592
593
594
595
595
596
597
597
598
599
599
600
601
602
603
604
605
605
606
607
607
608
609
609
610
611
612
613
614
614
615
616
616
617
618
618
619
620
621
622
623
623
624
625
625
626
627
627
628
629
629
630
631
632
633
634
635
635
636
637
637
638
639
639
640
641
642
643
644
644
645
646
646
647
648
648
649
650
651
652
653
654
654
655
656
656
657
658
658
659
660
661
662
663
664
664
665
666
666
667
668
668
669
670
671
672
673
673
674
675
675
676
677
677
678
679
679
680
681
682
683
684
684
685
686
686
687
688
688
689
690
691
692
693
693
694
695
695
696
697
697
698
699
699
700
701
702
703
703
704
705
705
706
707
707
708
709
709
710
711
711
712
713
713
714
715
715
716
717
717
718
719
719
720
721
721
722
723
723
724
725
725
726
727
727
728
729
729
730
731
731
732
733
733
734
735
735
736
737
737
738
739
739
740
741
741
742
743
743
744
745
745
746
747
747
748
749
749
750
751
751
752
753
753
754
755
755
756
757
757
758
759
759
760
761
761
762
763
763
764
765
765
766
767
767
768
769
769
770
771
771
772
773
773
774
775
775
776
777
777
778
779
779
780
781
781
782
783
783
784
785
785
786
787
787
788
789
789
790
791
791
792
793
793
794
795
795
796
797
797
798
799
799
800
801
801
802
803
803
804
805
805
806
807
807
808
809
809
810
811
811
812
813
813
814
815
815
816
817
817
818
819
819
820
821
821
822
823
823
824
825
825
826
827
827
828
829
829
830
831
831
832
833
833
834
835
835
836
837
837
838
839
839
840
841
841
842
843
843
844
845
845
846
847
847
848
849
849
850
851
851
852
853
853
854
855
855
856
857
857
858
859
859
860
861
861
862
863
863
864
865
865
866
867
867
868
869
869
870
871
871
872
873
873
874
875
875
876
877
877
878
879
879
880
881
881
882
883
883
884
885
885
886
887
887
888
889
889
890
891
891
892
893
893
894
895
895
896
897
897
898
899
899
900
901
901
902
903
903
904
905
905
906
907
907
908
909
909
910
911
911
912
913
913
914
915
915
916
917
917
918
919
919
920
921
921
922
923
923
924
925
925
926
927
927
928
929
929
930
931
931
932
933
933
934
935
935
936
937
937
938
939
939
940
941
941
942
943
943
944
945
945
946
947
947
948
949
949
950
951
951
952
953
953
954
955
955
956
957
957
958
959
959
960
961
961
962
963
963
964
965
965
966
967
967
968
969
969
970
971
971
972
973
973
974
975
975
976
977
977
978
979
979
980
981
981
982
983
983
984
985
985
986
987
987
988
989
989
990
991
991
992
993
993
994
995
995
996
997
997
998
999
999
1000
1000
1001
1001
1002
1002
1003
1003
1004
1004
1005
1005
1006
1006
1007
1007
1008
1008
1009
1009
1010
1010
1011
1011
1012
1012
1013
1013
1014
1014
1015
1015
1016
1016
1017
1017
1018
1018
1019
1019
1020
1020
1021
1021
1022
1022
1023
1023
1024
1024
1025
1025
1026
1026
1027
1027
1028
1028
1029
1029
1030
1030
1031
1031
1032
1032
1033
1033
1034
1034
1035
1035
1036
1036
1037
1037
1038
1038
1039
1039
1040
1040
1041
1041
1042
1042
1043
1043
1044
1044
1045
1045
1046
1046
1047
1047
1048
1048
1049
1049
1050
1050
1051
1051
1052
1052
1053
1053
1054
1054
1055
1055
1056
1056
1057
1057
1058
1058
1059
1059
1060
1060
1061
1061
1062
1062
1063
1063
1064
1064
1065
1065
1066
1066
1067
1067
1068
1068
1069
1069
1070
1070
1071
1071
1072
1072
1073
1073
1074
1074
1075
1075
1076
1076
1077
1077
1078
1078
1079
1079
1080
1080
1081
1081
1082
1082
1083
1083
1084
1084
1085
1085
1086
1086
1087
1087
1088
1088
1089
1089
1090
1090
1091
1091
1092
1092
1093
1093
1094
1094
1095
1095
1096
1096
1097
1097
1098
1098
1099
1099
1100
1100
1101
1101
1102
1102
1103
1103
1104
1104
1105
1105
1106
1106
1107
1107
1108
1108
1109
1109
1110
1110
1111
1111
1112
1112
1113
1113
1114
1114
1115
1115
1116
1116
1117
1117
1118
1118
1119
1119
1120
1120
1121
1121
1122
1122
1123
1123
1124
1124
1125
1125
1126
1126
1127
1127
1128
1128
1129
1129
1130
1130
1131
1131
1132
1132
1133
1133
1134
1134
1135
1135
1136
1136
1137
1137
1138
1138
1139
1139
1140
1140
1141
1141
1142
1142
1143
1143
1144
1144
1145
1145
1146
1146
1147
1147
1148
1148
1149
1149
1150
1150
1151
1151
1152
1152
1153
1153
1154
1154
1155
1155
1156
1156
1157
1157
1158
1158
1159
1159
1160
1160
1161
1161
1162
1162
1163
1163
1164
1164
1165
1165
1166
1166
1167
1167
1168
1168
1169
1169
1170
1170
1171
1171
1172
1172
1173
1173
1174
1174
1175
1175
1176
1176
1177
1177
1178
1178
1179
1179
1180
1180
1181
1181
1182
1182
1183
1183
1184
1184
1185
1185
1186
1186
1187
1187
1188
1188
1189
1189
1190
1190
1191
1191
1192
1192
1193
1193
1194
1194
1195
1195
1196
1196
1197
1197
1198
1198
1199
1199
1200
1200
1201
1201
1202
1202
1203
1203
1204
1204
1205
1205
1206
1206
1207
1207
1208
1208
1209
1209
1210
1210
1211
1211
1212
1212
1213
1213
1214
1214
1215
1215
1216
1216
1217
1217
1218
1218
1219
1219
1220
1220
1221
1221
1222
1222
1223
1223
1224
1224
1225
1225
1226
1226
1227
1227
1228
1228
1229
1229
1230
1230
1231
1231
1232
1232
1233
1233
1234
1234
1235
1235
1236
1236
1237
1237
1238
1238
1239
1239
1240
1240
1241
1241
1242
1242
1243
1243
1244
1244
1245
1245
1246
1246
1247
1247
1248
1248
1249
1249
1250
1250
1251
1251
1252
1252
1253
1253
1254
1254
1255
1255
1256
1256
1257
1257
1258
1258
1259
1259
1260
1260
1261
1261
1262
1262
1263
1263
1264
1264
1265
1265
1266
1266
1267
1267
1268
1268
1269
1269
1270
1270
1271
1271
1272
1272
1273
1273
1274
1274
1275
1275
1276
1276
1277
1277
1278
1278
1279
1279
1280
1280
1281
1281
1282
1282
1283
1283
1284
1284
1285
1285
1286
1286
1287
1287
1288
1288
1289
1289
1290
1290
1291
1291
1292
1292
1293
1293
1294
1294
1295
1295
1296
1296
1297
1297
1298
1298
1299
1299
1300
1300
1301
1301
1302
1302
1303
1303
1304
1304
1305
1305
1306
1306
1307
1307
1308
1308
1309
1309
1310
1310
1311
1311
1312
1312
1313
1313
1314
1314
1315
1315
1316
1316
1317
1317
1318
1318
1319
1319
1320
1320
1321
1321
1322
1322
1323
1323
1324
1324
1325
1325
1326
1326
1327
1327
1328
1328
1329
1329
1330
1330
1331
1331
1332
1332
1333
1333
1334
1334
1335
1335
1336
1336
1337
1337
1338
1338
1339
1339
1340
1340
1341
1341
1342
1342
1343
1343
1344
1344
1345
1345
1346
1346
1347
1347
1348
1348
1349
1349
1350
1350
1351
1351
1352
1352
1353
1353
1354
1354
1355
1355
1356
1356
1357
1357
1358
1358
1359
1359
1360
1360
1361
1361
1362
1362
1363
1363
1364
1364
1365
1365
1366
1366
1367
1367
1368
1368
1369
1369
1370
1370
1371
1371
1372
1372
1373
1373
1374
1374
1375
1375
1376
1376
1377
1377
1378
1378
1379
1379
1380
1380
1381
1381
1382
1382
1383
1383
1384
1384
1385
1385
1386
1386
1387
1387
1388
1388
1389
1389
1390
1390
1391
1391
1392
1392
1393
1393
1394
1394
1395
1395
1396
1396
1397
1397
1398
1398
1399
1399
1400
1400
1401
1401
1402
1402
1403
1403
1404
1404
1405
1405
1406
1406
1407
1407
1408
1408
1409
1409
1410
1410
1411
1411
1412
1412
1413
1413
1414
1414
1415
1415
1416
1416
1417
1417
1418
1418
1419
1419
1420
1420
1421
1421
1422
1422
1423
1423
1424
1424
1425
1425
1426
1426
1427
1427
1428
1428
1429
1429
1430
1430
1431
1431
1432
1432
1433
1433
1434
1434
1435
1435
1436
1436
1437
1437
1438
1438
1439
1439
1440
1440
1441
1441
1442
1442
1443
1443
1444
1444
1445
1445
1446
1446
1447
1447
1448
1448
1449
1449
1450
1450
1451
1451
1452
1452
1453
1453
1454
1454
1455
1455
1456
1456
1457
1457
1458
1458
1459
1459
1460
1460
1461
1461
1462
1462
1463
1463
1464
1464
1465
1465
1466
1466
1467
1467
1468
1468
1469
1469
1470
1470
1471
1471
1472
1472
1473
1473
1474
1474
1475
1475
1476
1476
1477
1477
1478
1478
1479
1479
1480
1480
1481
1481
1482
1482
1483
1483
1484
1484
1485
1485
1486
1486
1487
1487
1488
1488
1489
1489
1490
1490
1491
1491
1492
1492
1493
1493
1494
1494
1495
1495
1496
1496
1497
1497
1498
1498
1499
1499
1500
1500
1501
1501
1502
1502
1503
1503
1504
1504
1505
1505
1506
1506
1507
1507
1508
1508
1509
1509
1510
1510
1511
1511
1512
1512
1513
1513
1514
1514
1515
1515
1516
1516
1517
1517
1518
1518
1519
1519
1520
1520
1521
1521
1522
1522
1523
1523
1524
1524
1525
1525
1526
1526
1527
1527
1528
1528
1529
1529
1530
1530
1531
1531
1532
1532
1533
1533
1534
1534
1535
1535
1536
1536
1537
1537
1538
1538
1539
1539
1540
1540
1541
1541
1542
1542
1543
15
```

TABLE I: Sample Taxi Trajectory Data

TRIP ID	CALL TYPE	ORIGIN CALL	ORIGIN STAND	TAXI_ID	TIMESTAMP	DAY TYPE	MISSING DATA	POLYLINE
1372636858620000589	C	NaN	NaN	20000589	1372636858	A	False	[-8.618643,41.141412],[-8.618499,41.141376],...
1372637303620000596	B	NaN	7.0	20000596	1372637303	A	False	[-8.639847,41.159826],[-8.640351,41.159871],...
1372636951620000320	C	NaN	NaN	20000320	1372636951	A	False	[-8.612964,41.140359],[-8.613378,41.14035],...
1372636854620000520	C	NaN	NaN	20000520	1372636854	A	False	[-8.574678,41.151951],[-8.574705,41.151942],...
1372637091620000337	C	NaN	NaN	20000337	1372637091	A	False	[-8.645994,41.18049],[-8.645949,41.180517],...

Note: This is a sample table with a long title that will automatically wrap to the next line when the width is not sufficient.

```

11     x_coords, y_coords = zip(*trajectory)
12     ax.scatter(x_coords, y_coords, s=30,
13                 label=f'Trip {i+1}', color=
14                 color_map[i], edgecolor='k',
15                 linewidth=0.5)

```

### C. GPS Point Integration with Road Network

With the foundational road network in place, supplemented by the GPS points of the 15 trips, we combine them for a unified visualization. Furthermore, in our implementation, we assigned distinct colors to each trip for clear reference. The result is an informative display that highlights the positions of the taxis, contextualized by Porto’s road grid (Fig. 3).

In the figure, the white vertices are the vertices for the road network as described above. The other coloured vertices are the GPS locations of the 15 trips (colour coded with the key on the top right of the image). We can see that most of the GPS positions lie on the edges of the road network, which is expected, since the trips are taken on the road (i.e., the edge). However, other GPS readings are noisier and they do not lie directly on any edge. Because of the dense nature of the road network, it is not straightforward when deciding exactly which road the GPS location is referring to. On top of this, some of the GPS locations of the same trip are very close to one another spatially (e.g., trip 14). This appears on the map as just clutter of points and, without more context, it is again not easy to visualize the actual trajectory of the driver. Nonetheless, this rough visualization of the GPS locations allow us to have some insights into the data that we have (i.e., slightly noisy) and how we could further improve the analysis. This paves the way for deeper analysis and map matching in subsequent sections.

## IV. TASK 3: MAP MATCHING

After gaining some insights into the data we have procured, we can now perform map matching. Map matching is the process of aligning a sequence of observed user positions (i.e., map trajectory data) with the road network on a digital map to produce a sequence of road segments (termed as ‘routes’) representing the trajectory. This is essential in various applications, such as determining the exact path a vehicle has traveled, even if the GPS data has inaccuracies. For this

purpose, the open-source tool Fast Map Matching (FMM) was employed.

### A. Background of FMM

In this subsection, we define the problem we hope to solve. Firstly, the road network can be conceptualized as a graph,  $G = (V, E)$ , with edge set  $E$  and vertex set  $V$ . The goal of map matching is to find a path, denoted by  $route_t = \text{Path}(\hat{p}_1^{(t)}, \hat{p}_2^{(t)}, \dots, \hat{p}_{N_t}^{(t)})$ , where  $\hat{p}_i$  is a projected point of the GPS point  $p_i$ , onto the road network, given a fixed and ordered tuple of GPS points, denoted by  $trip_t = (p_1^{(t)}, p_2^{(t)}, \dots, p_{N_t}^{(t)})$ . Here,  $N_t$  refers to the number of GPS points in trip  $t$ , and the  $\text{Path}(\cdot)$  is a function that takes into sequences of points to produce an actual trajectory route on the road network.

FMM allows us to efficiently match each sequence of GPS points to a possible route in the road network. This is done by creating an upper bounded origin-destination (hash) table (UBODT) [2]. This table stores information about all pairs of shortest paths, between every two nodes, below an upper bound,  $\Delta$ . This is done by iteratively applying the single-source Dijkstra algorithm, which is a classical algorithm for finding the shortest path between nodes in a weighted graph. Map matching is therefore done by referencing to the UBODT instead of the more expensive spatial queries.

With the UBODT, we can perform the map matching using a Hidden Markov Model (HMM) based algorithm. In HMM, we assume the following:

- **Hidden states:** These are states, that affect the model, that are not privy to us. In our case, this refers to the road segments.
- **A Markov process:** This is a (strong) assumption that the current state only depends on the previous state.
- **Observations:** These are observable outcomes produced by the model. In this case, this refers to the GPS locations of each trip.

In order to reduce the space of possibilities, we search for at most  $k$  candidate points on the road network around each GPS point. In other word, for each  $p^{(t)}$ , we have  $\hat{p}^{(t)} \in C(p, k, r) = \{C_1, C_2, \dots, C_s\}$ , where  $s \leq k$  and  $C_j$  is the  $j^{\text{th}}$  candidate. Here, each  $C_j = (p'_j, e, dist, \lambda)$ , where  $e$  is an edge on the network,  $p'_j$  is the projection of point  $p_j^{(t)}$  onto the edge,  $dist$  is the distance between  $p_j^{(t)}$  and  $p'_j$  and  $\lambda$  is the distance

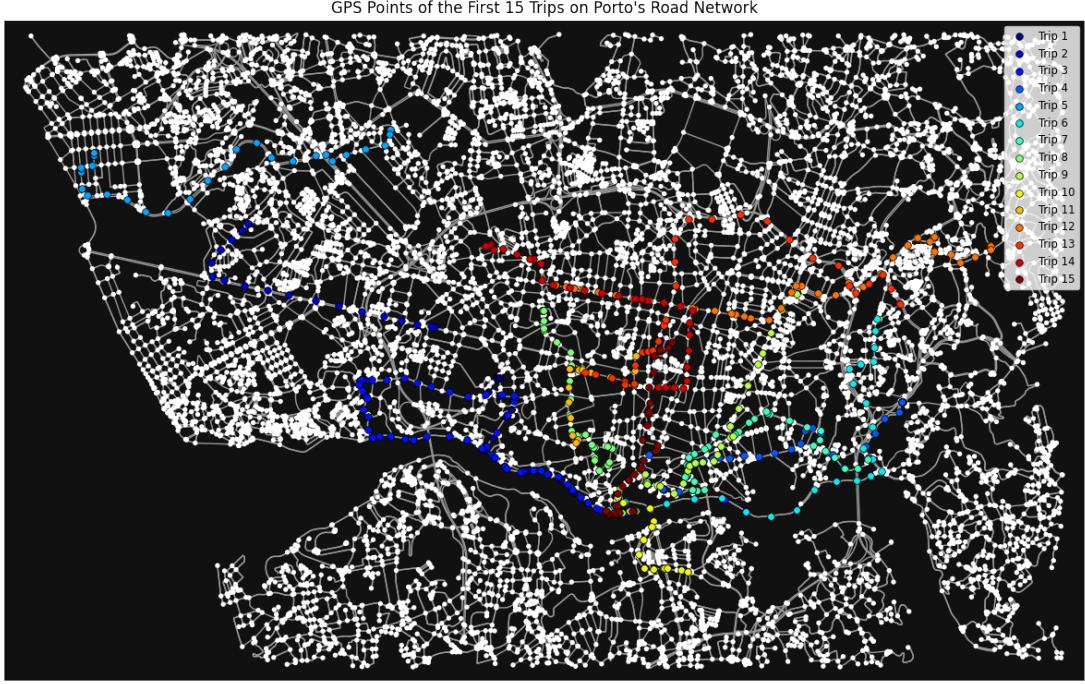


Fig. 3: Visualization of 15 Taxi Trajectories.

of  $p'_j$  from the starting vertex of edge  $e$ . We also limit the search of these candidates to be at most  $r$  units away from  $p$ . This means that for each  $C_j$ , we have  $\|C_j - p\| \leq r$ .

With this, we have a list of candidate sets  $\{C(p_i^{(t)}, k, r)\}_{i=1}^{N_t}$  for each trip  $t$ . This list of candidate sets are therefore the possible sequences hidden states in our HMM, and we are solving for the most probable sequence of hidden states given the GPS points observed.

In order to do this, we have to first find the shortest path (SP) distance between  $C_i$  and  $C_k$ , for all  $C_i \in C(p_{n-1}^{(t)}, k, r)$  and  $C_k \in C(p_n^{(t)}, k, r)$ . In other words, we have to find all pairwise distance between the candidates of the  $(n-1)^{\text{th}}$  GPS point and the candidates of the  $n^{\text{th}}$  GPS point. These values can easily retrieved from the UBODT. Now, denote  $SP_{n-1,n}(i, k)$  as the shortest path distance between  $C_i$  and  $C_k$  as defined above. Further, we let  $d_{n-1,n}$  be the Euclidean distance between  $p_{n-1}$  and  $p_n$ . Then, following [2], the transition probability is defined as:

$$tp(C_i, C_k) = \frac{\min(d_{n-1,n}, SP_{n-1,n})}{\max(d_{n-1,n}, SP_{n-1,n})} \quad (1)$$

Generally, transition probability is the likelihood of going from a given hidden state to the next given hidden state. On the other hand, the emission probability is the likelihood of observing a particular observation given the current state. Here, the emission probability,  $ep(C_k)$ , is given to be

$$ep(C_k) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(C_k.dist)^2/2\sigma^2} \quad (2)$$

where  $C_k.dist$  is the  $dist$  entry in  $C_k \in C(p_n^{(t)}, k, r)$ . By modelling this way, we are assuming that the GPS points follows a zero mean Gaussian distribution (i.e., noisy data) and  $\sigma$  is the standard deviation of the GPS error.

With these probabilities and assuming a Markov process, we can calculate the best matched path as,

$$\text{Path}(\hat{p}_1^{(t)}, \dots, \hat{p}_{N_t}^{(t)}) = \text{Path}[\arg \max_{i \in \{1, \dots, N_t\}} \sum_{n=1}^N tp(C_{n-1}, C_n) \times ep(C_n)] \quad (3)$$

where  $C_{n-1}$  and  $C_n$  are some candidate points in the candidate sets of  $p_{n-1}$  and  $p_n$  respectively. In other words, we find the sequence  $(\hat{p}_1^{(t)}, \dots, \hat{p}_{N_t}^{(t)})$  by optimizing for both the emission and transition probabilities. This is done by the Viterbi algorithm [3].

Notice, however, that in equation 3 we consider the transition probability between  $C_{n-1}$  and  $C_n$ . This is afforded by assuming a Markov process, making the computation much less complex. However, in reality,  $C_n$  would depend on all  $C_i$  for  $1 \leq i < n$  (albeit to varying degree), not just  $C_{n-1}$ . Specifically, the current point on the road network cannot be entirely independent from all the points that are not the previous point. Therefore, the assumption of Markov process leads to a coarser modelling of reality, as we disregard dependencies further in space and time.

Finally, penalty was implemented to resolve issues of reverse movements (e.g., moving back and forth). Although there are some fringe cases where vehicles do that, such routes are more often than not poorly matched. With this penalty, it enforces that a route is a directed path, instead of a circuit or a cycle.

## B. Methodology

With the above background, we use the FMM library in this task.

- 1) **Data Preparation:** We have to input both the trajectory data and the road network. To ensure compatibility with FMM, we have to indicate the start and end nodes of each edge in the road network. The trajectory data of each trip is made up of a sequence of timestamped GPS coordinates, extracted from the columns of train\_1500.
- 2) **Execution:**
  - FMM was installed and set up according to the documentation provided in its GitHub repository.
  - The trajectory data was input to FMM alongside the road network data.
  - FMM processed the data and mapped each trajectory to the most probable sequence of road segments, forming a route.
- 3) **Output Verification:** The resulting routes were verified by overlaying them on the road network to ensure they align with actual roads. In the next task, the first 15 routes were visually inspected to further validate the accuracy of the map matching process.

**Listing 4:** Extraction and boundary calculation for the first 15 trajectories.

```

1 # Load the network
2 network = Network(network_path, "fid", "u",
3                   "v")
4 print(f"Nodes {network.get_node_count()}")
5 edges {network.get_edge_count()}")
6 graph = NetworkGraph(network)
7
8 # Generate UBODT if not exists or if
9 # regeneration is forced
10 if not os.path.exists(ubodt_path) or
11     regenerate_ubodt:
12     ubodt_gen = UBODTGenAlgorithm(network,
13                                     graph)
14     status = ubodt_gen.generate_ubodt(
15         ubodt_path, threshold, binary=False
16         , use_omp=True)
17     if not status:
18         print("Error generating UBODT.")
19         return
20     print("UBODT generated successfully.")
21
22 # Load UBODT
23 ubodt = UBODT.read_ubodt_csv(ubodt_path)
24
25 # Create FMM model
26 model = FastMapMatch(network, graph, ubodt)
27 fmm_config = FastMapMatchConfig(k, radius,
28                                 gps_error)

```

## C. Results

The map matching process was executed through a tailored Python function leveraging the FMM (Fast Map Matching) library. This function, titled `fmm_map_matching`, serves the purpose of mapping raw GPS trajectories to the nearest

road segments in a predefined road network. The function therefore solves for the optimal sequence of hidden states, given a sequence of observations. In order to produce sound and accurate results, the function takes in the following arguments (i.e., parameters):

- **Candidates,  $k$ :** This parameter determines the maximum number of road segments candidates to consider for each GPS point. In our implementation, we take  $k = 6$ .
- **radius,  $r$ :** The value determines the range within which road segment candidates for a given GPS point are searched. In other words, this is the radius of the circle with a GPS point as the centre. The algorithm then performs range query within this circle. Here, we take radius = 0.05.
- **gps\_error:** This corresponds to  $\sigma$  seen in equation 2, which is from the assumption that the GPS signals are distributed as a zero-mean Gaussian distribution.
- **threshold:** This is the upper bound of UBODT. With a value of 0.02, it serves to prune the search space during the construction of UBODT.

For the input, the road network data was sourced from the Porto road network, and the trajectory data was retrieved from train-1500. If necessary, the function was also configured to regenerate the UBODT file. Each trajectory from the train-1500 was processed to match it with the road network. The resulting mapped routes, along with details about the matching process, were written to `matched_routines.csv`. In Table II, we show the first five rows in the `matched_routines.csv`.

We exported the following outputs from the `fmm_map_matching` function

- 1) **Index:** This is a standard index column representing the matched routes.
- 2) **cpath:** An abbreviation for "complete path", each entry in this column is constructed from opath (see below). In particular, each entry is a list of road segment identifiers that connect the edges in opath such that it forms a continuous path. This represents  $route_t$  in terms of edges on the road network.
- 3) **mgeom:** Short for "matched geometry", this is the geometric representation (in LINESTRING format) of the matched paths on the road network. Each LINESTRING is composed of a series of points (longitude, latitude) that represent where the GPS point or trajectory was matched to on the road segment.
- 4) **opath:** Abbreviation for "optimal path", each entry stores a list of road segment identifiers for the corresponding projected GPS trajectory data (i.e., list of  $\hat{p}$  for the trip). This need not form a continuous path.
- 5) **offset:** This is a list of values representing the relative position of the GPS point, on its corresponding road segment, to the start of the road segment. Therefore, a value close to 0 indicates that the GPS point is near the beginning of its matched segment, whereas a value close to 1 denotes that it is near the end of the segment.
- 6) **length:** Each entry in this column is a list of values

TABLE II: First five rows of `matched_routines.csv`

Index	cpath	mgeom	opath	offset	length	spdist
0	[1054, 4247, ...]	LINESTRING(-8.6186233, 41.141456,...)	[1054, 1054,...]	[0.005765244498 603186,...]	[0.00707063699 2543766,...]	[0.0, 0.0001460268816041218,...]
1	[37924, 179951, ...]	LINESTRING(-8.6398592, 41.159752...)	[37924, 37924,...]	[0.0016796833 541194568,...]	[0.0023412979 739263237,...]	[0.0, 0.0005046527883385151,...]
2	[]	LINESTRING()	[]	[]	[]	[]
3	[]	LINESTRING()	[]	[]	[]	[]
4	[109634, 157753,...]	LINESTRING(-8.6458216, 41.180406,...)	[109634, 109634,...]	[0.00044984402 869269397,...]	[0.0005863222 76888652,...]	[0.0, 4.397033822113258e-05,...]

which represent the length of the corresponding road segment in `cpath`

- 7) **spdist:** Short for "shortest path distances", this is the  $SP_{n-1,n}$  outlined above. For instance, the  $i^{\text{th}}$  element in an entry, is the distance between  $\mathbf{p}'_i$  and  $\mathbf{p}'_{i-1}$ , where  $\mathbf{p}'_j$  is the matched point on the road network for the  $j^{\text{th}}$  GPS point  $\mathbf{p}_j$  in this trip. Note that all the first value of the column is defined to be 0.0.

After map matching, we realized some of the entries are empty. For instance, index 2 in Table II have empty entries for the `cpath`, `mgeom`, and `opath` columns. This is because the has map matching failed to produce a sequence of road segment, given the parameters listed above. For example, if there is just one GPS point within a trajectory which could not be paired with a plausible edge (perhaps  $r$  is too restrictive), then the entire trajectory cannot be matched to a route. This could be mitigated by relaxing the constraints, but it might come at the expense of accuracy, as routes produced will be of low confidence (i.e., increased uncertainty threshold). We will further investigate this in the subsequent section.

Nonetheless, our `fmm_map_matching` function, built on the FMM library, successfully and efficiently matched most of the trips to a route on the road network. In order to have a more in-depth understanding, however, we would have to visualize these information on the road network; ultimately, these are temporal data.

## V. TASK 4: ROUTE VISUALIZATION

This section delves into the visualization process of the routes that were generated from the map matching task of the first 15 trips. Visualization is crucial not only for verifying the results but also for a more comprehensive understanding of the spatial relationship between the trajectories and the road network. Here, we aim to produce a clear and accurate representation of the first 15 routes on the road network, distinguishing each trip either by different colors.

### A. Methodology

- 1) **Tool Selection:** The OSMnx tool was primarily chosen for this task due to its compatibility with OpenStreetMap data and its capacity to generate high-quality visualizations. OSMnx facilitates the extraction and visualization of urban street networks from OpenStreetMap, making it suitable for our use-case.

- 2) **Road segments:** The sequences of road segments obtained from the previous map matching task serve as our primary input. These routes were already mapped to the road network and therefore ready for visualization.

- 3) **Visualization Process:** We then visualize the first 15 routes in the following manner:

- Using OSMnx, we extract the road network of the relevant region, similar to what we did in section III.
- The first 15 routes were then plotted, with different colours, on this road network. The colour-coding allows for a clearer visualization, especially in areas where the routes overlap or intersect.

- 4) **Validation:** The visualized routes are then compared against the original trajectories and the digital map to ensure they accurately represent the path taken.

### B. Results

The visualization can be seen in Fig. 4. Because some of the trips in the first 15 trips of train-1500 were not matched by the function, we also include the subsequent few trips of train-1500. This means that, some of the routes in Fig. 4 do not belong to any of the trajectories shown in Fig. 3.

Looking at the matched trajectories (e.g., the dark blue route in upper right corner of Fig. 4 and the light blue route in upper right corner of Fig. 3), we can see that they are well-matched. From the visualization, we can see that most of the trips went through the city-centre, with only a few going to the fringes of the city. This can be seen from the fact that many of these routes intersect and at times overlap in some segments. These overlapping segments indicate that the road segments are commonly used across many trips. Some routes seem to be exiting the network, however this is because we used the Porto, Portugal road network, and these routes crossed out of the city.

The visualization allows us to have a much better understanding about how the FMM performed, as well as constructing a better picture of how the trajectories were taken. Furthermore, the color-coding distinguishes the different trips, allowing for easier identification of each route. However, it will become infeasible if we were to perform analysis on the visualization, especially if lots of data come into play. For that, we will have to look back at the `matched_routines.csv` file.

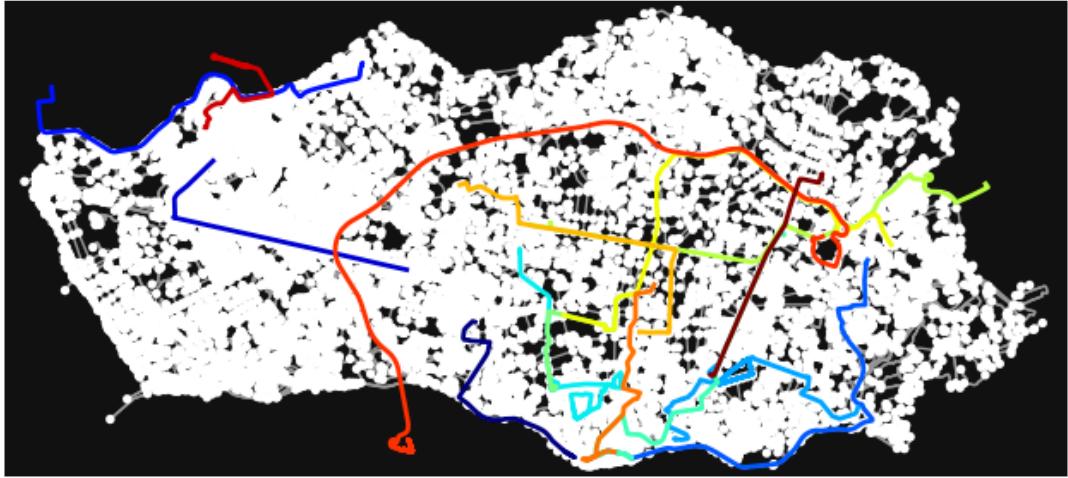


Fig. 4: Visualization of 15 Matched Taxi Routes. The colour coding scheme here differs from Fig. 3.

## VI. TASK 5: ROUTE ANALYSIS

In this section, we intended to answer two questions:

- 1) What are 10 road segments that are traversed the most often?
- 2) What are the 10 road segments that have the longest average travelling time, ignoring those that are not traversed?

### A. Most frequently traversed road segments

The idea behind this is simple: We count the number of routes that a road segment contains. This can be done by identifying the road segment within each route. We can identify the road segment by using the entries in `opath` and `cpath` which are identifiers for each road segment. However, each road segment need not have a unique FID – for example, one of the route has  $FID_1 = 1100$  and  $FID_2 = 4231$ . If we were to simply count the appearance of all unique FIDs, then the number of traversal for this road will be lower than in reality.

On the other hand, each road segment has a OpenStreetMap ID (OSMID) as well. Unfortunately, it suffers the same issue as that of the FID – some road segments have multiple OSMID. Additionally, some road segments share the same OSMID. In order to properly count the number of traversal per road segment, we have to address these issues. To that end, we produce a new unique identification for each road segment, by mixing in some other information in each road segment.

*1) Methodology:* Here, we outline how we overcame the issues outlined above. In order to produce a unique ID for each road segment, we perform the following:

- 1) Given the FID of a road segment, we obtain its OSMID, start node,  $u$ , end node,  $v$ , and its length,  $d$ .
- 2) We then define a variable  $uvd$  as a tuple  $(\min(u, v), \max(u, v), d)$ . This serves as to make the ID unique for each road segment.

With this new, unique ID, we have a one-to-one mapping of each road segment on the road network to an identity. This

is a dictionary named `road_fid2id_info`. Now, we can count the number of times a road segment has been traversed as follows:

- 1) **Create a new dictionary:** We create a dictionary `road_id2frequency`, in order to count the frequency of each unique road segments
- 2) **Iterate over all the routes:** This counts the frequency of each road segment. This takes the FID from the `cpath` then convert these FIDs to the unique IDs as in `road_fid2id_info`.
- 3) **Removal of consecutive duplicates:** Since the GPS updates are taken every 15 seconds, there will be routes with repeated road segments in the `cpath`. However, this should only count as one traversal, instead of multiple traversal. Therefore, we have to remove them.
- 4) **Sorting the frequency of roads:** After iterating over all the routes, we just have to sort the frequency of the road segments in descending order.
- 5) **Obtain results:** We then return the FID of the top 10 road segments and their counts.

The algorithm to do this is outlined below in VI-A1

Listing 5: 10 most frequently traversed road segments

```

1 road_id2frequency = defaultdict(int)
2 for route in results:
3     repeat_c_ids = [road_fid2id_info[c][0]
4                      for c in route['cpath']]
5     c_ids = []
6     for r in repeat_c_ids:
7         if len(c_ids) == 0 or r != c_ids[-1]:
8             c_ids.append(r)
9     for ID in c_ids:
10        road_id2frequency[ID] += 1
11    # Do the counting
12 sorted_id = sorted(road_id2frequency.keys()
13                   , key=lambda x: road_id2frequency[x],
14                   reverse=True)
15    # Sort the list by frequency

```

```

13 result = []
14 for i in range(5):
15     ID = sorted_id[i]
16     fids = tuple(road_id2fids[ID])
17     road_name = gdf_edges.loc[fids[0], "name"]
18     result.append((fids, road_name,
                    road_id2frequency[ID]))

```

Finally, we have the following ranking. With these, we could visualize the road segments with the highest traversal frequency as in Fig. 7

TABLE III: 10 FIDs road segments that are traversed the most often.

Rank	Road ID	Use Count
1	100535	171
2	4239	164
3	42	147
4	1461	141
5	158673	135
6	1559	134
7	1547	133
8	47	131
9	188463	124
10	50	121

### B. Road segments with longest average travelling time

To solve this problem, we have to be able to find the time spent on a road segment. However, we note that all the information given by the FMM algorithm were spatial in nature. Fortunately, we know that the GPS points are updated periodically at every 15 seconds. With some assumptions, we could make some inferences from the information we have.

We know that consecutive projected GPS points (i.e., `opatch`) are 15 seconds away. Therefore, we can find the time taken on a road segment by considering pairwise neighbours of the projected GPS point. Note that there are two possible cases: (1) Two consecutive GPS points are matched to the same road segment; (2) Two consecutive GPS points are on different road segments. For the former case, we simply take the time spent on the particular road segment to be 15 seconds.

For the latter case, where two consecutive projected GPS points are on different road segments, we have to find a way to "distribute" the 15 seconds to all these road segments. In order to do this, we made an assumption that the speed of the vehicle is constant between these two updates. Suppose the road segments that are in between these two projected GPS points are  $e_1, \dots, e_k$  (these can be found within the entries of `cpath`). Then, we can distribute the travelling time,  $t_{e_r}$ , on each road segment as,

$$t_{e_r} = 15 \times \frac{l_{e_r}}{\sum_{i=1}^k l_{e_i}}, \quad (4)$$

where  $l_{e_r}$  is the length of the road segment that the vehicle covers and can be calculated as,

$$l_{e_r} = \|e_r\| - \text{offset}, \quad (5)$$

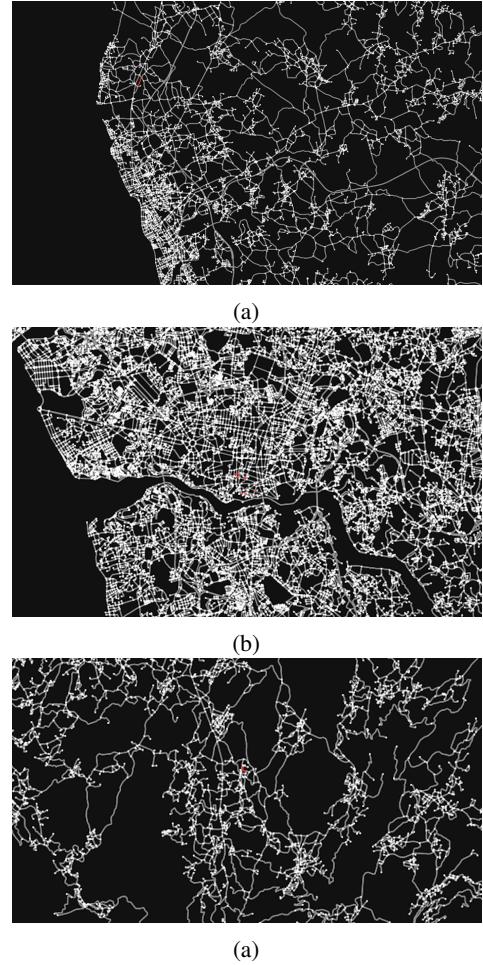


Fig. 7: Visualization of the most frequent traversed road segments in red

where  $\|e_r\|$  is the length of the road segment and offset is as described in IV.

Suppose that, of all the trips,  $N_r$  of them passes through  $t_{e_r}$ , then there will be  $N_r$  of such  $t_{e_r}$ . Therefore, we can calculate the average travelling time,  $T_r$ , on  $e_r$  to be,

$$T_r = \frac{1}{N_r} \sum_{i=1}^{N_r} t_{e_r}^{(i)} \quad (6)$$

1) *Methodology:* Here, we outline the steps taken to perform what was outlined above:

- 1) We first create a dictionary called `road_id2travel_patch`. This allows us to store both the time spent on a road segment and the distance between it and its preceding projected GPS point.
- 2) For each route, we calculate the travelling time for each road segment in this route as outlined above. The travelling time is then updated in `road_id2travel_patch`

In our implementation, we added some constraints in order to clean the data. For instance, we enforced that the `spdist` has to be more than  $10^{-7}$  for every projected GPS point. This is to ensure that the taxis, which could be idling, are travelling at a reasonable speed for it to be considered in our calculation. Furthermore, we removed road segments that have less than 10 routes pass through them, as they are segments that are least travelled, resulting in significantly different time difference. Finally, we also excluded road segments whose total travelled distance by all the taxis is less than 1.5 times the segments' length. Intuitively, this means that the entire road segment is not well covered by all the taxis' routes. For example, the taxis could only be travelling through one specific part of the road segment. Therefore, this time may not be representative.

Listing 6: 10 largest average travelling time road segments

```

1   mean_pass_times = dict() # (log_number
2     , mean_pass_time)
3   for ID, logs in road_id2travel_patch.
4     items():
5     log_number = len(logs)
6     if log_number < ignore_log_number:
7       # if the number of logs
8       smaller than ignore_log_number,
9       ignore this road ID
10      continue
11      road_length = road_fid2id_info[
12        road_id2fids[ID][0]][4]
13      total_pass_distance = sum([log[1]
14        for log in logs])
15      log_rate = total_pass_distance /
16        road_length
17      if log_rate < ignore_log_rate: #
18        if the total distance travelled
19        on this road by cars is
20        smaller than 1.5 length of road
21        , ignore this road ID.
22        continue
23      total_pass_time = sum([log[0] for
24        log in logs])
25      mean_pass_time = total_pass_time /
26        total_pass_distance *
27        road_length
28      mean_pass_times[ID] = (log_number,
29        mean_pass_time, log_rate)
30
31      sorted_pass_time = sorted(
32        mean_pass_times.items(), key=lambda
33          x: x[1][1], reverse=True)[:5]
34      fid_time_lognumber = [(tuple(
35        road_id2fids[s[0]]), s[1][1], s
36        [1][0], s[1][2]) for s in
37        sorted_pass_time]

```

A snippet of the algorithm can be seen in list:task 5-2. Here, the `log_number` ( $= 10$ ) is the minimum number of routes that each road segment must have, and `log_rate` ( $= 1.5$ ) is the ratio of the total distance travelled to road length as mentioned above. Finally, the function returns `fid_time_lognumber` which gives a sequence of road segments' FID with decreasing time travelled on it. We also visualized this result in 8

is the visualization.

Listing 7: 10 largest average travelling time road segments

```

1   road_id2travel_patch = defaultdict(list)
2
3   for route in results:
4     c_index = 0
5     repeat_c_ids = [road_fid2id_info[c][0]
6       for c in route['cpath']]
7     c_ids = []
8     for r in repeat_c_ids:
9       if len(c_ids) == 0 or r != c_ids
10         [-1]:
11         c_ids.append(r)
12     last_o_id = road_fid2id_info[route['
13       opath'][0]][0]
14     for o_index in range(1, len(route['
15       opath'])):
16       cur_o_id = road_fid2id_info[route['
17       opath'][o_index]][0]
18     if last_o_id == cur_o_id:
19       # if two consecutive GPS points are in the
20       same road
21       if route['spdist'][o_index] > 1
22         e-7:
23         road_id2travel_patch[
24           cur_o_id].append((15,
25             route['spdist'][o_index
26             ]))
27     else:
28       # if two consecutive GPS passed many
29       different roads
30     # Omitted: calculate on each of the road
31       respectively before combining
32
33     mean_pass_times = dict()
34     # Omitted: filter out the road segments
35       that is unqualified
36
37     sorted_pass_time = sorted(mean_pass_times.
38       items(), key=lambda x: x[1][1], reverse
39       =True)[:10]
40     # Sort and get the top 10
41     fid_time_lognumber = [(tuple(road_id2fids[s
42       [0]]), s[1][1], s[1][0], s[1][2]) for s
43       in sorted_pass_time]

```

These two analysis provide useful insight, especially in a dense road network, where some road segments face significant bottleneck issues. For example, after finding out those routes most travelled, in terms of count and time, authorities could use these information to plan future development of the road network to move some vehicles away from these road segments. However, this would require further analysis that is beyond the scope of this project. Nevertheless, these simple analyses provide valuable insights.

## VII. TASK 6: IMPROVEMENTS TO MAP MATCHING

As discussed in section IV, some of the trajectory data are not being matched, due to the lack of plausible road segments for some of the GPS points in the trip. Put another way, the GPS points are likely noisy data and "too far off" from the actual road network. If we could remove these uninformative data points from the trajectory data, we could potentially match these routes. Furthermore, existing, matched routes, could

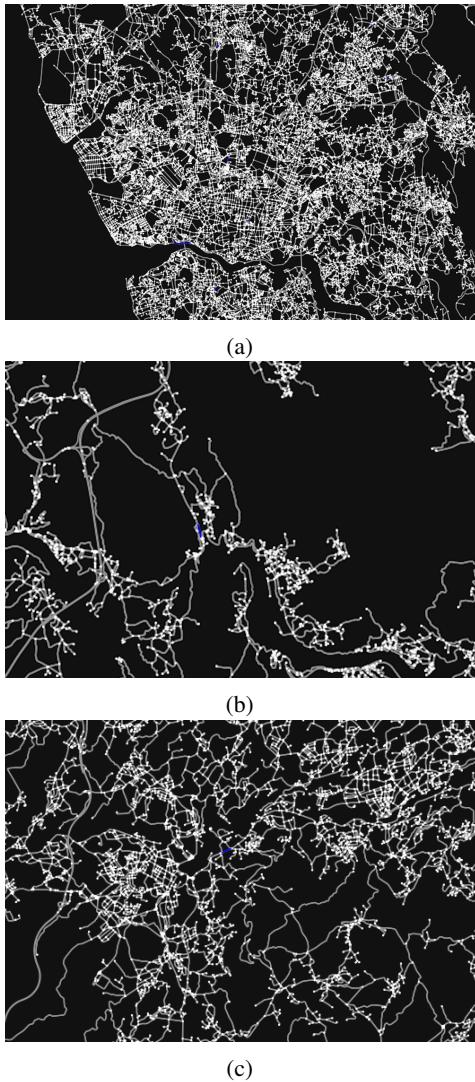


Fig. 8: Visualization of road segments with the largest average travelling time in blue

theoretically be less noisy as well. For example, some routes could suffer from some jittering, but by removing the points causing these jitters (some degree of being outliers), the routes could be smoothed out.

#### A. Outlier detection

We start first by defining what constitutes as an outlier, and have a rough understanding of how many of such outliers there are. There are multiple ways to define an outlier – an example would be to consider the distance of a raw GPS point from an existing road segment. Although feasible, this will be slightly computationally expensive, since we have to store both the road segments and GPS points in memory, while doing pairwise distance comparison.

We devised a different definition which involves the Euclidean distance between consecutive GPS points. In particular, suppose the  $i^{\text{th}}$  GPS point,  $p_i$ , is an outlier, then  $\|p_{i-1} - p_i\| > t$ . Intuitively, if it is "too far off" from its previous neighbour,

the GPS point is likely to be an outlier, supposing the previous neighbour is not an outlier. In order to decide on this threshold  $t$ , our group considered two heuristics:

- 1) The **distribution of the distances**. This is to ensure that the threshold chosen can both weed out the outliers and keep most of the GPS points.
- 2) The **frequency of the GPS updates**, which is 15 seconds. This provides a gauge on how plausible certain GPS can be logged.

For the first point, we considered plotting a variogram cloud. However, with the sheer amount of data, this quickly proves infeasible as the visualization will be too dense to be understood. Instead, since we want to understand the distribution, we plotted a histogram (Fig. 9), where the horizontal axis represents the Euclidean distance, in kilometres, between two consecutive points, and the vertical axis represents the number of such points.

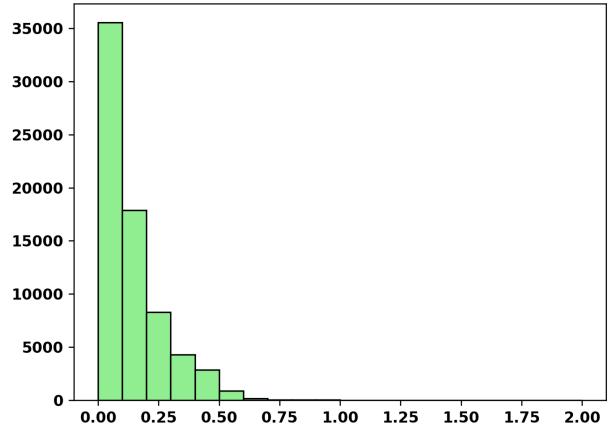


Fig. 9: Distribution of Euclidean distances between two consecutive points. Horizontal axis are in kilometres.

From the histogram, we can see that the Euclidean distance follows a left-tailed distribution – majority of the GPS points is less than 0.5 km from its neighbours. This is realistic, as it translates to speeds less than 120 km/h, with an update frequency of 15 seconds. We could therefore set  $t = 0.5$ .

#### B. Methodology

In order to remove the outliers, we first define a  $\text{Dist}( \cdot )$  function, which calculates the Euclidean distance between two points,  $p$  and  $p'$ :

$$\text{Dist}(p, p') = 100 \sqrt{(p_1 - p'_1)^2 + (p_2 - p'_2)^2} \quad (7)$$

The distance is scaled by 100 as the GPS data were initially at the scale of  $10^{-2}$  km. With this, we have the following algorithm:

Listing 8: Outlier removal.

```

2 #CSV file
3 train1500 = pd.read_csv('./data/train-1500.
4
5 for t in range(1500):
6     gps_points = eval(train1500['POLYLINE'
7         ][t])
8     assert isinstance(train1500['POLYLINE'
9         ][t], str)
10
11     #Checking the distances between
12     #consecutive GPS points
13     if len(gps_points) < 2:
14         modified_points = gps_points
15     else:
16         modified_points = gps_points[:1]
17         last_point = gps_points[0]
18         threshold = threshold_min
19         for i in range(1, len(gps_points)):
20             next_point = gps_points[i]
21             if distance(last_point,
22                 next_point) < threshold: # a valid point
23                 modified_points.append(
24                     next_point)
25                 last_point = next_point
26                 threshold = threshold_min
27             else:
28                 threshold += threshold_step
29                 threshold = min(threshold,
30                     threshold_max)
31         train1500['POLYLINE'][t] =
32             list2string(modified_points)
33
34     #Output trajectories without outlier GPS
35     #points as modified_train1500.csv
36 train1500.to_csv(folder+/
37     'modified_train1500.csv', index=0)

```

Here, instead of using one threshold,  $t$ , we have three parameters:  $\text{threshold}_{\min}$ ,  $\text{threshold}_{\text{step}}$ ,  $\text{threshold}_{\max}$ . By default, each GPS point,  $p_i$ , will be checked against its previous neighbour,  $p_{i-1}$ : If  $\|p_i - p_{i-1}\| < \text{threshold}_{\min}$ , then we consider  $p_i$  as an inlier. Conversely, if  $\|p_i - p_{i-1}\| \geq \text{threshold}_{\min}$ , then we consider  $p_i$  as an outlier. When this happens, when it comes to determining if  $p_{i+1}$  is an outlier, it will have to be compared with  $p_{i-1}$ , since  $p_i$  has been removed. But the difference between these two updates is now 30 seconds instead of 15 seconds. Naturally, the threshold should be larger to account for this increased time difference. Therefore, when this happens (i.e.,  $p_i$  is being removed as an outlier), we add a  $\text{threshold}_{\text{step}}$  to  $\text{threshold}_{\min}$ , in order to 'loosen' the criterion between these two points with larger time difference. Finally,  $\text{threshold}_{\max}$  ensures that the threshold used will always be less than or equal to  $\text{threshold}_{\max}$ .

In our experiment, we used the values  $\text{threshold}_{\min} = 0.5$ ,  $\text{threshold}_{\max} = 2$ , and  $\text{threshold}_{\text{step}} = 0.2$ . The choice of  $\text{threshold}_{\min}$  comes from the reasoning that most of the GPS points have less than or equal to that distance, as well as that it is heuristically a reasonable

threshold. On the other hand, the  $\text{threshold}_{\text{step}}$  is chosen to be smaller than  $\text{threshold}_{\min}$  and roughly the mean of the Euclidean distance of all the points. This prevents an explosion of the threshold used, while making a reasonable estimate. With these parameters, we obtained a `modified_train1500.csv` file. Subsequently, we generated a `modified_matched_routines.csv` by the map matching algorithm used in section IV. Finally, we generate the visualization of the routes as in the steps used in section V. This can be seen in Fig. 10.

Comparing Fig. 10 to Fig. 3 and Fig. 4, we can see that some trajectories that were previously unmatched (e.g., trip 2), are now able to be matched. By removing these outliers, we managed to provide a candidate road segment to all the the inlier GPS points. This in turns allow us to have a nonzero probability of obtaining the best sequence of road segments for all these routes.

However, looking at the `modified_matched_routines.csv` file, we noticed, again, that some of the trips were once again not matched. For instance, trip 9 does not have any entry in `cpath`, suggesting no routes were assigned to it. This once again means that there exists a GPS point in that trip that is not assigned a road segment. To circumvent this, a naive way would be to lower  $\text{threshold}_{\min}$  once again, to weed out these GPS points. However, a potential problem with this is that some trips may no longer have sufficient GPS points for us to confidently interpolate a route for it. In this case, we suffer some accuracy issues once again, as the sparser the GPS points, the more inaccurate the routes generated will be. A remedy for this would be to remove modified trips with  $\leq N$  GPS points entirely.

Regardless, the removal of outliers allow for us to visualize routes that were previously unmatched. However, some routes still remain to be unmatched and other methods of cleaning the data may be required.

## VIII. CONCLUSION

In this project, we worked with visualizing GPS points on the road network of Porto city in Portugal. Furthermore, we used an open-source library called FMM to perform map matching, giving us not only the practical insights of the algorithm, but also the theoretical part of it. Finally, with the matched routes, we were able to perform some detailed analysis on the routes, for example highlighting routes that were most traversed, and improve on the map matching results.

## IX. GROUP MEMBER CONTRIBUTIONS

### A. Contributions by Xiang Xinye

I have written the task 5 including the code and the report.

### B. Contributions by Yin Wenqi

I have written the task 1 and 2 as well as the related code and report.

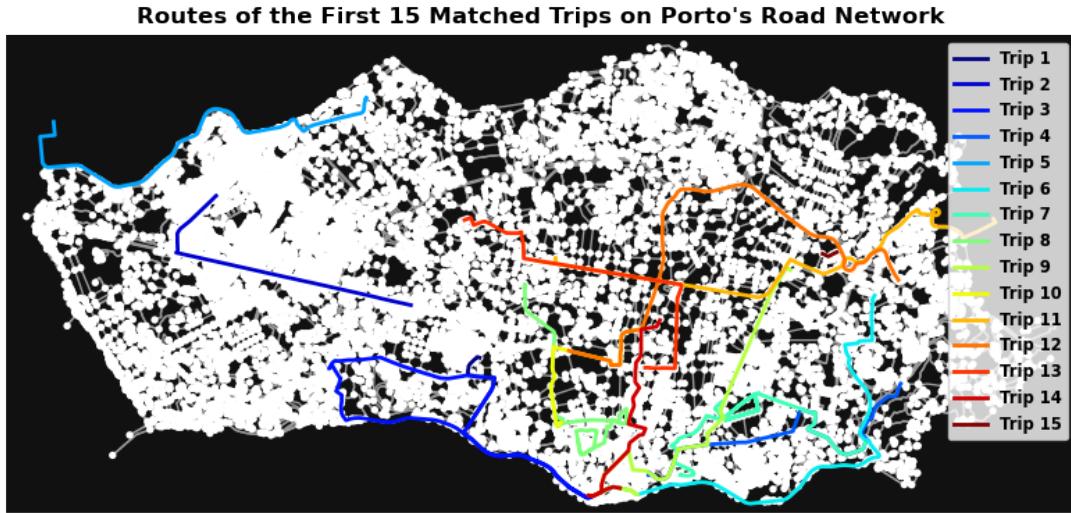


Fig. 10: Visualization of first 15 routes after outlier removal.

### C. Contributions by Tan Jie Heng Alfred

I have written most of section IV, V, and all of section VII, including the codes and the report. I have also done analysis for the entire report.

## APPENDIX

### REFERENCES

- [1] K. Inc., “Kaggle: Your home for data science,” 2023. Accessed on April 20, 2023.
- [2] C. Yang and G. Gidófalvi, “Fast map matching, an algorithm integrating hidden markov model with precomputation,” *International Journal of Geographical Information Science*, vol. 32, no. 3, pp. 547–570, 2018.
- [3] J. G. Forney, “The viterbi algorithm,” 1973.