

1 **Modelling long-distance route choice using mobile phone call detail**
2 **record data: A case study of Senegal**

3
4 **Andrew Bwambale**
5 Choice Modelling Centre
6 Institute for Transport Studies
7 University of Leeds
8 34-40 University Road, LS2 9JT, Leeds, United Kingdom
9 Email: ts13ab@leeds.ac.uk

10
11 **Charisma F. Choudhury**
12 Choice Modelling Centre
13 Institute for Transport Studies
14 University of Leeds
15 34-40 University Road, LS2 9JT, Leeds, United Kingdom
16 Email: C.F.Choudhury@leeds.ac.uk

17
18 **Stephane Hess**
19 Choice Modelling Centre
20 Institute for Transport Studies
21 University of Leeds
22 34-40 University Road, LS2 9JT, Leeds, United Kingdom
23 Email: S.Hess@its.leeds.ac.uk

24

25 **Abstract**

26 The growing mobile phone penetration rates (world-wide) have led to the emergence of
27 large scale call detail records (CDRs) that could serve as a low-cost data source for travel
28 behaviour modelling compared to the commonly used data sources such as stated
29 preference (SP) data. However, to the best of our knowledge, there is no previous study
30 specifically evaluating the potential of CDR data in the context of route choice behaviour
31 modelling. Being event-driven, the data is discontinuous and only able to yield partial
32 trajectories, thus presenting serious challenges for route identification, especially in
33 highly overlapping networks. This paper proposes techniques for inferring the users'
34 chosen routes or subsets of their likely routes from partial CDR trajectories, as well as
35 data fusion with external sources of information such as route costs and travel times, and
36 then adapts the broad choice framework to the current modelling scenario where the road
37 network is highly overlapped. The model results show that CDR data can capture the
38 expected sensitivities towards route attributes and the behaviour towards overlapping
39 routes. The value of travel time estimates derived from the models are found to be

1 realistic for the study area (Senegal). These findings are timely for developing countries
2 with low budgets for transport studies.

3 **Keywords:** Route choice behaviour, Broad choice, Mobile phone data, Call detail
4 records, Value of travel time

5 **1. Introduction**

6 The modelling of route choice behaviour for long journeys, inter-city and inter-
7 regional trips is an important component of large-scale transport planning models. The
8 emergence of technologies that enable passive collection of mobility trajectories has led
9 to a step change in route choice modelling. Most of these studies are based on Global
10 Positioning System (GPS) data, primarily from navigational devices (Li et al., 2018,
11 Hess et al., 2015, Broach et al., 2012, Bierlaire and Frejinger, 2008), and more recently
12 from smartphones (Bierlaire et al., 2010, Papinski et al., 2009). The navigational
13 devices are still not used widely in all parts of the world (in countries of the Global
14 South in particular). Although smartphone ownership is on the rise throughout the
15 world, the quality of data from these devices strongly rely on the availability of the
16 Assisted-GPS (A-GPS)¹ feature, internet connectivity and data storage capacity, which
17 leads to small sample sizes as seen in most related studies (e.g. Nitsche et al., 2014,
18 Bierlaire et al., 2013, Nitsche et al., 2012, Bierlaire et al., 2010). Thus both sources pose
19 a substantial risk of sampling bias.

20 This problem can be overcome by taking advantage of the large scale
21 anonymous datasets that are already being passively collected by operators for different
22 purposes, and applying these to transport studies. Such datasets have already yielded
23 promising results various mobility studies. Examples include; social media data
24 (Hawelka et al., 2014, Hasan et al., 2013), smart card data (Chakirov and Erath, 2012,
25 Agard et al., 2006), and network-generated mobile phone data such as Call Detail
26 Records (CDRs)² and Global System for Mobile communications (GSM)³ data (Çolak
27 et al., 2015, Jiang et al., 2013, Schlaich et al., 2010). However among these, network-
28 generated mobile phone data is particularly a promising source due to the high mobile
29 phone penetration rate world-wide (GSM Association, 2017).

30 A review of literature shows that there have been a few route identification
31 studies using network-generated mobile phone data (e.g. Nie et al., 2015, Leontiadis et
32 al., 2014, Hoteit et al., 2014, Schlaich et al., 2010), however, most of these studies end
33 on route identification and do not attempt to investigate the factors affecting route
34 choice behaviour. At the moment, only Schlaich (2010) combines GSM trajectories
35 with traffic state information to analyse the influence of variable message signs (VMS)
36 and other factors on route choice. However, the success of Schlaich's study could in
37 part be attributed to the use of GSM data, which is network-driven and semi-continuous
38 in nature as opposed to CDR data, which is event-driven and discontinuous, thus
39 presenting serious trajectory identification challenges. Despite these challenges, the data
40 is more readily available as it is stored longer for billing purposes as opposed to GSM

² CDR data reports the time stamped locations of communication events (i.e. voice calls, text messages, and data calls) as well as the details of the request (i.e. the duration and direction).

³ GSM data reports the IDs of all the GSM cells traversed by an active mobile phone at regular time intervals (irrespective of the calling or texting patterns of the users).

1 data, which is discarded after each location area update operation to save computer
2 memory. Thus, methods developed to harness the potential of such a widely available
3 data source can have wide benefits in both the developed and developing worlds. This
4 motivates this research where we investigate the feasibility of CDR data for modelling
5 route choice behaviour.

6 However, it is important to underscore the practical challenges that stem from
7 the use of CDR data, and how this impacts our work. Since the data is event-driven,
8 there is an increased risk of not capturing any trajectories, especially for very close O-D
9 pairs with short travel times. For this reason, we limit our scope to long-distance inter-
10 regional route choice as there is an increased possibility that a user will use his/her
11 phone at different points during the trip, thus enabling the capture of his/her partial
12 trajectory. The fact that only partial trajectories are observed means that in some cases,
13 it is not possible to precisely infer the chosen routes. In such cases, route choice is
14 observed at a broad sub-group level (e.g. northern, southern etc.), where each sub-group
15 comprises of a small set of possible routes. This prompts us to adapt the ‘Broad Choice
16 Modelling Framework’, developed in the context of vehicle type choice (Wong, 2015)
17 to route choice modelling using noisy CDR data. To the best of our knowledge, the
18 Broad Choice Modelling Framework is the only model structure that can deal with the
19 case when the choice data (dependent variable) has multiple levels of aggregation (i.e.
20 unique choice for some observations while broad sub-group for the other ones) as in this
21 case. We note that this may be problematic in dense inter-urban networks, where it
22 would be difficult to identify a small enough subset of possible routes using a few CDR
23 locations. However, with the increasing trend of mobile internet usage (Gerpott and
24 Thomas, 2014), the frequency of CDR locations is likely to improve significantly in the
25 near future, which adds further to the timeliness of the present paper.

26 Although CDR data is only able to capture the partial trajectories of frequent
27 phone users, the samples are usually large, thus increasing the possibility that they are
28 representative enough to capture rational route choice behaviour. The need to
29 investigate this assertion forms the basis of our validation exercise where we obtain
30 stable results. The developed models are used to estimate the value of travel time (VTT)
31 for Senegal (the study area), yielding reasonable estimates. The study is timely in the
32 sense that it extends the application of CDR data beyond travel pattern visualisation to
33 econometric modelling of travel behaviour. This could motivate reliable and low cost
34 policy formulation in both developed and developing countries.

35 The rest of the paper is arranged as follows; section 2 presents a review of
36 relevant literature, section 3 presents the data description, section 4 presents the data
37 processing conducted, section 5 presents the modelling framework, section 6 discusses
38 the model results, while section 7 presents the summary and conclusions.

39 **2. Literature review**

40 This section briefly reviews literature on the applications of mobile phone data in
41 transport studies, as well as different models of route choice.

42 *2.1. Previous applications of mobile phone data to transportation studies*

43 The last few decades have seen significant research effort aimed at investigating the
44 potential application of mobile phone data to various aspects of transportation studies,
45 so far yielding promising results. To date, the data has been widely applied in human

1 mobility modelling to understand individual's travel patterns (e.g. Ahas et al. 2010,
2 Ahas et al. 2015, Deville et al., 2016, Diao et al., 2016, Xu et al., 2015, Calabrese et al.,
3 2013, Isaacman et al., 2012, Shaw et al. 2016, Song et al., 2010, Yuan et al. 2012, Zhao
4 et al. 2013) and correlating it with phone usage patterns (e.g. Yuan et al. 2012), tourism
5 surveys to understand the attractiveness of different tourist sites (Ahas et al., 2008,
6 Ahas et al., 2007), urban area (land use) classification to facilitate urban planning (Pei et
7 al., 2014, Yuan and Raubal, 2012), analysing commute (Kung et al., 2014, Ahas et al.,
8 2010) and migration patterns (e.g. Wang et al. 2019), estimating trip making rates (e.g.
9 Çolak et al., 2015), developing origin-destination matrices (e.g. Çolak et al., 2015, Iqbal
10 et al., 2014, Calabrese et al., 2011, White and Wells, 2002), detecting travel modes to
11 estimate mode shares (e.g. Qu et al., 2015, Wang et al., 2010, Reddy et al., 2008), and
12 estimating traffic parameters to assess the traffic conditions of key roads (e.g. Bolla et
13 al., 2000). However, we place focus on studies related to route identification. A few of
14 these studies have used GSM data, which reports the complete mobile phone location
15 area sequences of each user, thus enabling the easy identification of routes through
16 sequence matching (e.g. Ma et al. 2013, Tettamanti et al., 2012, Schlaich et al., 2010)
17 and probabilistic methods such locality-sensitive hashing and graph clustering (e.g.
18 Görnerup, 2012). Furthermore, a growing number of studies have focussed on analysing
19 the potential of CDR data, which is more widely available compared to GSM data and
20 yet challenging to use as mentioned in the introduction section. For example, Doyle et
21 al. (2011) use the virtual cell paths technique to extract user trajectories from CDR data,
22 and generate the kernel density paths for different routes to validate their findings.
23 Saravanan et al. (2011) analyse the spatial and temporal information of CDR events
24 over a long period of time to establish the daily routines and routes of the users. Hoteit
25 et al. (2014) join subsequent triangulated CDR locations using linear, cubic, and nearest
26 - neighbour interpolation to model the potential trajectories. Leontiadis et al. (2014)
27 calculate the weights for each road segment within the cell areas linked to a user's
28 communication events and determine the shortest weighted path for a given OD pair.
29 Nie et al. (2015) mark each route with a subset of k optimal cell handover sequences
30 extracted from the full set of possible handover sequences and match these with those
31 observed in the cell phone hand over data (similar to CDR data) based on the degree of
32 similarity. In this study, we follow a slightly similar approach, however, instead of
33 using a similarity index, we pursue the idea of unique and shared location area
34 sequences in the context of broad choice modelling as explained later.

35 **2.2. Existing route choice models**

36 The vast majority of route choice models belong to the family of discrete choice models
37 (see Ben-Akiva and Lerman, 1985 for details), with the multinomial logit (MNL) model
38 (McFadden, 1974) being the most widely used. However, the MNL model is affected by
39 the irrelevance of independent alternatives (IIA) property, which can be problematic for
40 highly overlapping routes (Ramming, 2002). This has motivated the development of more
41 advanced route choice models to address this challenge. Examples include the nested
42 recursive logit model (Mai et al., 2015), the c-logit model (Cascetta et al., 1996), the path
43 size logit model (Ben-Akiva and Ramming, 1998), the link nested (cross-nested) logit
44 model (Vovsha and Bekhor, 1998), and the multinomial probit model with logit kernel
45 (Daganzo et al., 1977). Details of how some of these models overcome the overlapping
46 route problem are discussed in section 5.3 of this paper.

47 An important point to note is that the complexities of route choice modelling go
48 beyond the overlapping route problem. Choice set generation is a key challenge,

1 especially in highly overlapping dense urban networks, where several alternative routes
2 can be possible, and yet individuals do not consider all the alternatives while making
3 choices (Prato, 2009). Several choice set generation methods have been proposed in the
4 literature including, the k-shortest path algorithms (e.g. Shier, 1979, Bellman and Kalaba,
5 1960), the labelling approach (Ben-Akiva et al., 1984), link elimination approaches (e.g.
6 Azevedo et al., 1993, Bellman and Kalaba, 1960), link penalty approaches (e.g. Roupail
7 et al., 1995, De La Barra et al., 1993), simulation approaches (e.g. Sheffi and Powell,
8 1982), doubly stochastic generation functions (e.g. Nielsen, 2000), constrained
9 enumeration methods (e.g. Prato and Bekhor, 2006), and probabilistic methods (e.g.
10 Cascetta and Papola, 2001, Manski, 1977). However, since the focus of this paper is long
11 distance trips, where the alternatives are usually few in number, choice set determination
12 is more straightforward as discussed later in section 4.2 of this paper.

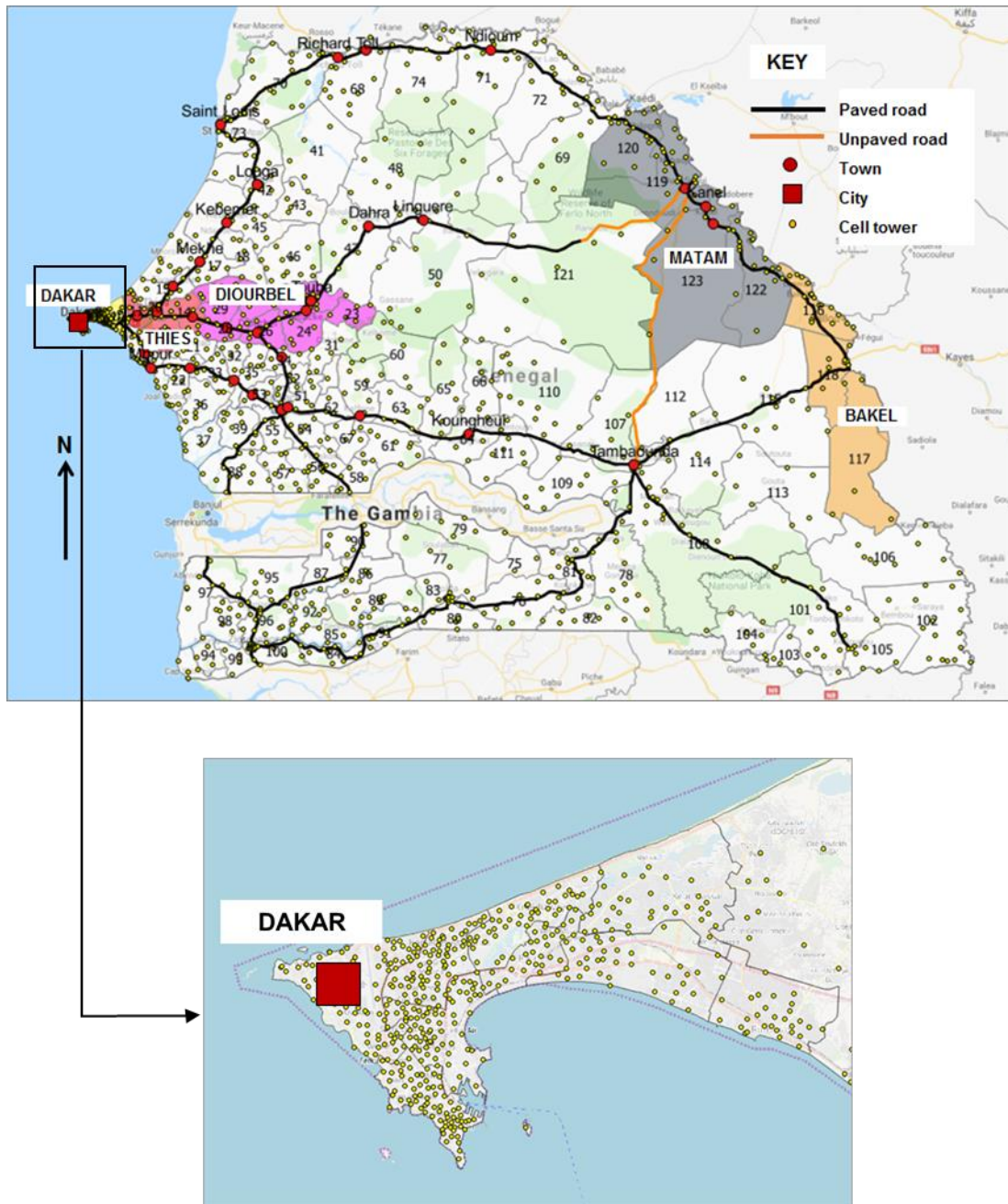
13 As mentioned earlier, the discontinuous nature of CDR data presents serious route
14 identification challenges, resulting in route choice observations at different levels of
15 aggregation (i.e. unique route choices for some travellers while broad sub-groups of the
16 possible routes for the rest). Although there is no specific route choice model dealing
17 with such a scenario, recent advances in other fields of choice modelling have addressed
18 this challenge using the broad choice framework (proposed by Wong, 2015 and later
19 used by Habibi et al. 2017, Lloro and Brownstone 2018, Wong et al. 2018, Yip et
20 al. 2018 in the context of vehicle choice and by Tran et al. 2019 in the context of hotel
21 choice), which can be adapted to deal with the current route choice modelling scenario,
22 discussed later in this paper. It may be noted that according to the choice modelling
23 literature, this framework is the only one that can deal with the case where the choices
24 are uniquely observed only for some respondents and for the rest, only the sub-group of
25 choices are observed.
26

27 **3. Data**

28 This study uses CDR data collected from Senegal as part of the Orange Data for
29 Development (D4D) challenge (de Montjoye et al., 2014).

30 **3.1. Study area**

31 Senegal is located in West Africa with a population of approximately 13.5 million
32 according to the 2013 population census (ANSD, 2016). Road transport accounts for
33 over 99% of all passenger travel (World Bank, 2004). The only long-distance train
34 service (the Dakar-Niger line) was discontinued in May 2010 (Imedia and Calao
35 Production, 2013). The country has a sparse national road network (see Figure 1), and
36 for some O-D pairs, there is only one feasible alternative, making them unsuitable for
37 route choice modelling. This study ignores such O-D pairs and only considers those
38 where alternative routes exist. In total, twelve distant O-D pairs are considered, and
39 these are; Dakar-Bakel, Dakar-Matam, Thies-Bakel, Thies-Matam, Diourbel-Bakel,
40 Diourbel-Matam, and the corresponding O-D pairs for the reverse directions as shown
41 in Figure 1. The long travel times between these regions increase the possibility of
42 capturing the users' partial trajectories as explained earlier.



1

2 Figure 1. Study area (Google Maps 2017, Worldatlas 2017, ArcGIS 2013)

3 **3.2. CDR data**

4 The CDR data was collected between January and December 2013 and aggregated to
 5 the arrondissement (district) level by the data provider. The geographical location of the
 6 arrondissements is presented in Figure 1 where we also show the tower locations for
 7 illustration purposes.

8 The original CDR data comprised of 9 million unique users (67% of the study
 9 area population). This was pre-processed to retain frequent phone users (i.e. those with
 10 interactions on 75% of the days in a year) and randomly split into smaller monthly
 11 rolling sub-samples made available for research (see de Montjoye et al., 2014 for

1 details). The user IDs in each sub-sample are anonymised to prevent possible re-
 2 identification across the different months.

3 The data for each month comprises of about 150,000 users. Assuming that the
 4 most commonly observed arrondissement for each user during the month is their home
 5 district, the monthly population sampling rate ranged from 2.4% in Dakar (the capital)
 6 to 0.4% in the rural regions. On average, these users together generated over 40 million
 7 records per month (see excerpt of the CDR data in Table 1a). The data is reduced to
 8 remove duplicate records resulting in the processed arrondissement visitation data
 9 presented in Table 1b.

10 Table 1a. Excerpt of the raw CDR data

Anonymised User ID	Timestamp	Arrondissement ID ⁴	
130599	13-01-02 20:10	25	
130599	13-01-13 13:10	7	
130599	13-01-19 23:50	19	
130599	13-01-19 23:50	19	
130599	13-01-22 01:30	2	
130599	13-01-22 01:30	2	
130599	13-01-28 20:20	4	
130599	13-01-28 20:20	4	Discarded from the data
130599	13-01-29 19:40	4	
130599	13-01-29 19:50	4	
130599	13-01-29 20:00	4	
130599	13-01-29 20:40	4	
130599	13-01-29 21:20	4	
130599	13-01-29 21:50	4	
130599	13-01-29 21:50	4	
130599	13-01-29 21:50	4	
130599	13-01-29 21:50	4	

11
 12

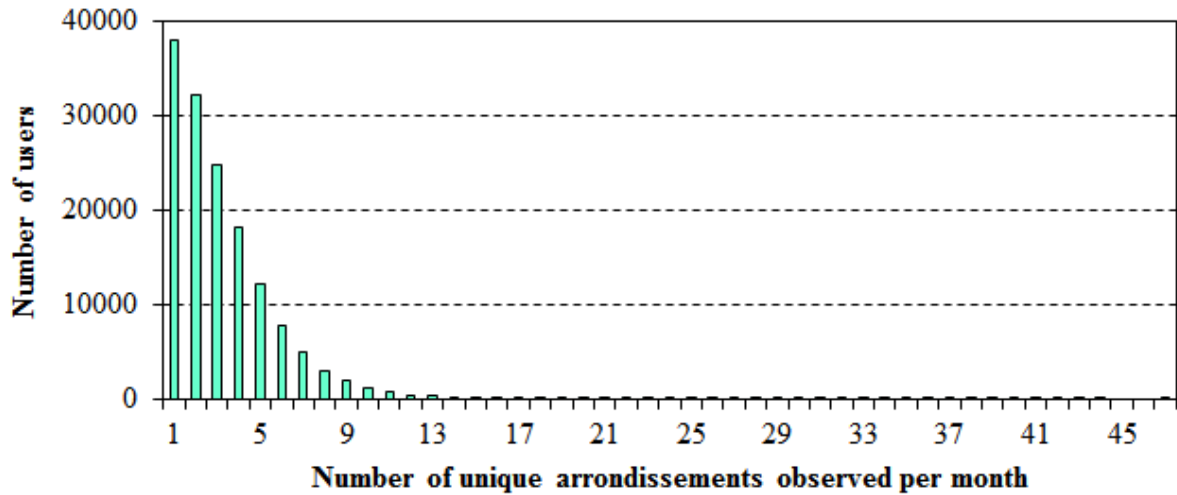
13 Table 1b. Excerpt of the processed arrondissement visitation data

Anonymised User ID with monthly identifier (e.g. January)	Arrondissement ID	1 st observation	Last observation
130599.01	25	13-01-02 20:10	13-01-02 20:10
130599.01	7	13-01-13 13:10	13-01-13 13:10
130599.01	19	13-01-19 23:50	13-01-19 23:50
130599.01	2	13-01-22 01:30	13-01-22 01:30
130599.01	4	13-01-28 20:20	13-01-29 21:50

14
 15
 16
 17
 18
 19
 20

The overall level of user mobility is illustrated in Figure 2. As shown, most users visited less than three unique arrondissements per month. The low levels of inter-arrondissement mobility led to the capture of few trajectories as reflected in the final sample size (see Section 4.2).

⁴ The geographical locations of the arrondissement IDs are presented in Figure 1



1

2 Figure 2. Average monthly arrondissement observation frequency distribution

3

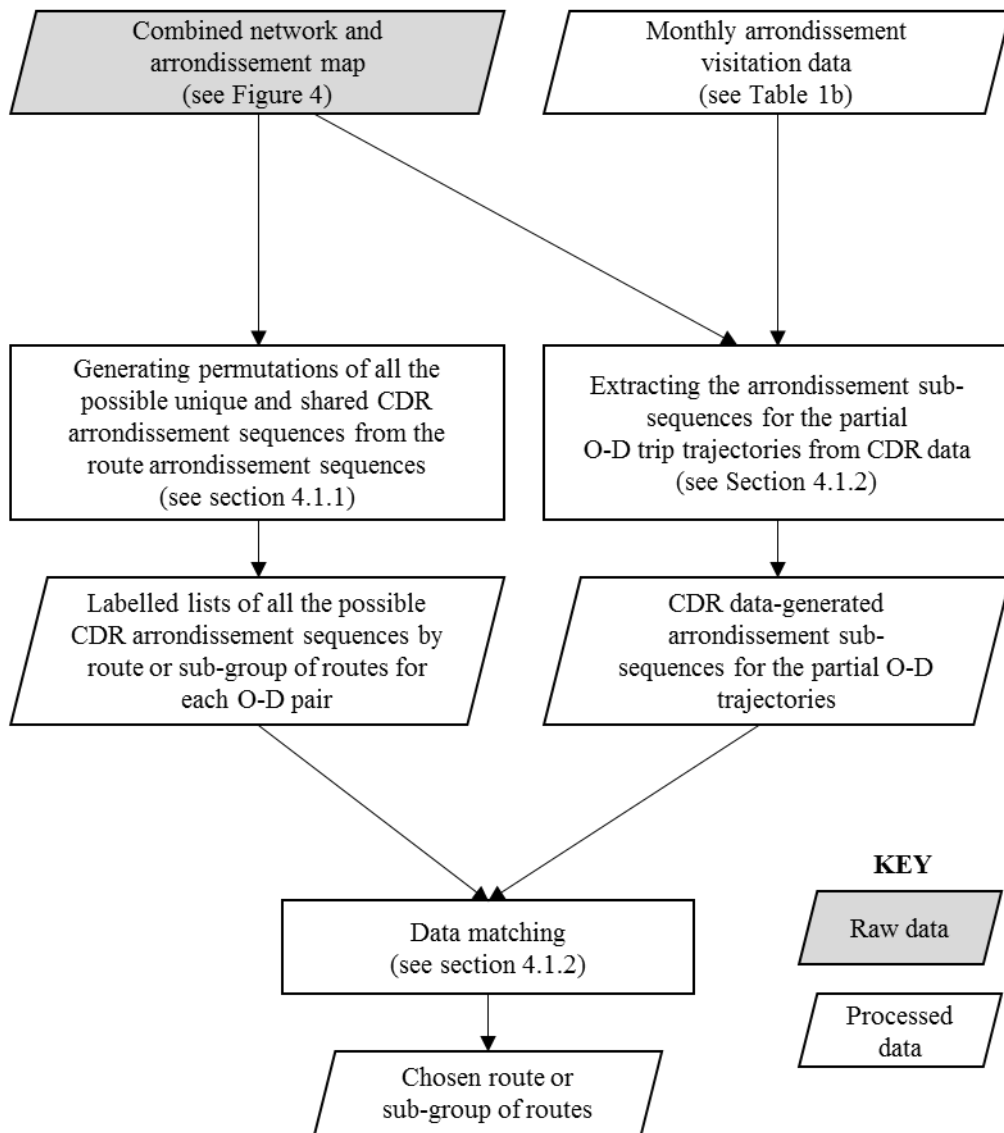
4 Although it is difficult to detect false tower jump movements in aggregate CDR
 5 data (Çolak et al., 2015, Iqbal et al., 2014), this is not a big factor for distant O-D pairs
 6 where the origin and destination arrondissements have been grouped into regions.

7 **4. Data preparation for analysis**

8 This section describes the analysis carried out on the processed arrondissement
 9 visitation data in Table 1b to identify the routes followed, as well as the processes of
 10 estimating the route attributes.

11 **4.1. Route identification**

12 The route identification process is summarised in Figure 3. This is divided into two
 13 main stages as described in the subsequent sections.



1

2 Figure 3. Summary of the route identification process

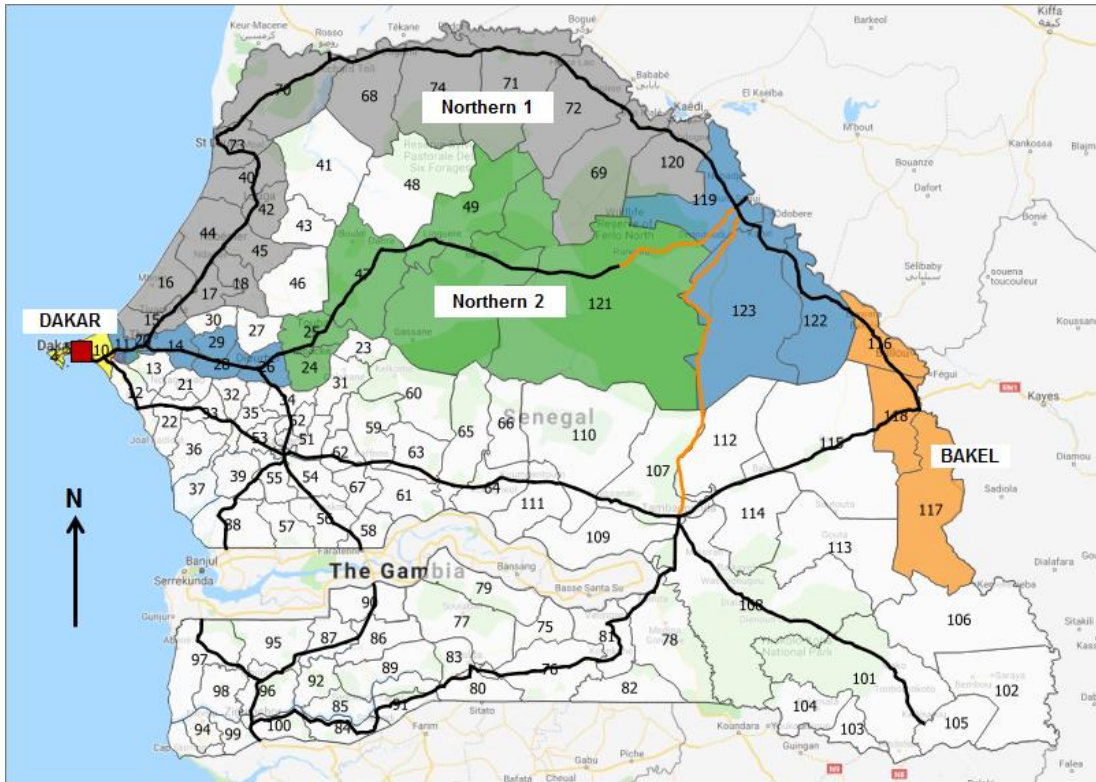
3

4 *4.1.1. Generation of unique and shared CDR arrondissement sequences*

5 Route arrondissement sequences (which are extracted from maps) show the order of all
 6 the arrondissements traversed by a particular route between a given O-D pair. On the
 7 other hand, CDR arrondissement sequences (which are extracted from the CDR data)
 8 show the order of the arrondissements in which a user used his/her phone during the
 9 trip, and are subsets of the route arrondissement sequences.

10 For any given trip along a particular route, several possible CDR arrondissement
 11 sequences can be observed depending on the number and the location of the CDR
 12 events. These can be obtained by generating permutations of different sizes based on the
 13 route arrondissement sequence (in which order matters and no repetitions are allowed).
 14 However, since most of the O-D pairs have overlapping routes, some of the CDR
 15 arrondissement sequences can be linked to more than one route if all the intermediate
 16 CDR events occurred along the shared sections. In this case, it would only be possible

1 to observe the subset of routes that were potentially followed (see illustration in Figure
 2 4 where the blue areas indicate the shared arrondissements, while the grey and green
 3 areas indicate the unique arrondissements for the Northern 1 and 2 routes respectively).



4

5 Figure 4. Arrondissement paths (Dakar-Bakel O-D pair as an example)

6

7 The generated CDR arrondissement sequences linked to each route were cross-
 8 referenced to identify the permutations linked to unique routes (*i.e. unique CDR*
 9 *arrondissement sequences*) and those shared across multiple routes (*i.e. shared CDR*
 10 *arrondissement sequences*). The outcome of this analysis was a list of all the possible
 11 CDR arrondissement sequences labelled with the associated routes or sub-groups of the
 12 possible routes (see example in Table 2).

13 *4.1.2. Extraction the O-D pair trip trajectories from CDR data*

14 The processed arrondissement visitation data for each user (see excerpt in Table 1b) was
 15 analysed to extract sub-sequences linked with possible trips between the regions of
 16 interest following the criteria below;

- 17 • The first and the last arrondissements in the sub-sequence must be located
 18 within different regions of interest, and the user must not be observed in an
 19 upstream or a downstream region of interest within the same day for origins and
 20 destinations respectively;
- 21 • The dwell time in the origin and the destination regions of interest must be
 22 longer than that required to directly traverse each of them to increase the
 23 possibility that these are the trip start and end locations. A user needs to use
 24 his/her phone at least twice in each of these regions to calculate the dwell time;

- The intermediate arrondissements in the sub-sequence must all be associated with one of the defined arrondissement/corridors paths (see Figure 4) to ensure only direct trips are retained; and
- The timestamp difference between the origin and the destination must not exceed 24 hours, which is used as an upper limit to distinguish between users with direct trips but delay to use their phones on arrival, and those with intermediate destinations, thereby arriving late.

The extracted sub-sequences meeting all the above criteria were either assigned to unique routes or sub-groups of the possible routes by cross-referencing with the labelled lists generated in Section 4.1.1. Table 2 presents an excerpt of the route assignment data.

Table 2. Excerpt of the route assignment data

Anonymised User ID with monthly identifier (e.g. January)	CDR Trajectory	Route/ Broad sub-group
131891.01	Dakar-11-C1-Bakel	Northern 1
131891.01	Dakar-C1-Bakel	Northern 1
132801.01	Dakar-122-Bakel	Northern (Northern 1/ Northern 2)*
132801.01	Dakar-28-C2-123-Bakel	Northern 2
132801.01	Dakar-C1-123-Bakel	Northern 1

* CDR trajectory can belong to both the Northern 1 and Northern 2 routes

In this data, 70% comprised of unique assignments, while 30% comprised of broad assignments. Since this is a scenario where for some users and/or trips, we know the chosen route at a more disaggregate level than for others, we use the Broad Choice Modelling Framework (Wong, 2015) to analyse route choice behaviour.

4.2. Estimation of route attributes

For model estimation, it is critical to determine the choice set and the attributes of the alternatives. We assumed that the choice set is comprised of the routes that have ever been chosen by the different users. Routes not chosen by any user for the whole year were excluded from the choice set. On average, each origin-destination pair had four alternatives. Given that the total dataset is comprised of 9,453 records from 6,497 users, this is not a very restrictive assumption. Given the choice sets, we reviewed previous studies to identify the attributes typically used in route choice models and their availability status for Senegal as summarised in Table 3. Although data on six explanatory variables was available, the final model specification contains three explanatory variables only as the inclusion of the other variables led to correlation problems and/or illogical model results. A detailed explanation of the variable specification process is presented in Section 6.1. The subsequent sections summarise the processes of estimating some of these attributes.

Table 3. Attributes typically used in route choice models

Attribute	Sample references	Data availability	Remarks
------------------	--------------------------	--------------------------	----------------

Individual socio-demographics	(Ramming 2002, Zhang and Levinson 2008)	x	Mobile phone data is usually anonymous
Travel time	(Hess et al. 2015, Ramming 2002, Ben-Elia and Shiftan 2010)	✓	Can be derived from traditional data or maps
Travel cost	(Hess et al. 2015)	✓	Can be estimated using the vehicle operating costs
Distance	(Hamerslag 1981, Bitzios and Ferreira 1993)	✓	Can be calculated from maps
Scenic characteristics	(Ben-Akiva et al. 1984, Zhang and Levinson 2008)	✓	Can be derived from maps
Safety (e.g. presence of black spots)	(Ben-Akiva et al. 1984, Ben-Elia and Shiftan 2010)	x	Data could not be obtained
Urban developments along the route	(Ben-Akiva et al. 1984, Zhang and Levinson 2008)	✓	Data can be obtained from maps
Time or distance on uninterrupted flow facilities (e.g. freeways)	(Bierlaire and Frejinger 2008, Ramming 2002)	x	No such facilities between the regions of interest at the time of data collection
Traffic congestion	(Bitzios and Ferreira 1993)	x	Data could not be obtained
Road quality (e.g. road surface conditions)	(Ben-Akiva et al. 1984)	✓	Data from the national roads agency is available
Road signs (e.g. direction signs)	(Wootton, Ness, and Burton 1981)	x	Data could not be obtained

1
2

3 4.2.1. *Link length and surface attributes*

4 Data on the link lengths and surface attributes (i.e. paved or unpaved) was derived from
5 the Senegal roads GIS layer (ArcGIS, 2013). This was updated to reflect the situation in
6 2013 relying on road condition reports sourced from government and other relevant
7 websites (Ageroute Senegal, 2017, ANSD, 2017, Logistics Cluster, 2013).

8 4.2.2. *Travel time*

9 Travel time cannot be reliably estimated from the CDR data as users do not necessarily
10 use their phones at the moment of departure or arrival. The typical travel times for most
11 links in 2013 were obtained from the website of Logistics Capacity Assessment
12 (Logistics Cluster, 2013). For links not covered by this website, we relied on Google
13 Maps (Google Maps, 2017a).

1 4.2.3. Travel cost

2 Travel cost was estimated in terms of the vehicle operating costs (VOCs) per user (i.e.
 3 fuel and non-fuel costs). After a review of several VOC estimation techniques, we
 4 settled for the HDM-III model (Watanatada et al., 1987) due to its applicability to
 5 developing countries and input data availability. The HDM-III model is an earlier
 6 version of the more advanced HDM-4 (Kerali, 2000), which we could not use due to
 7 input data constraints.

8 The HDM-III model relies on vehicle calibration data (where we used default
 9 values) and other basic input data (see Table 4) to estimate the VOCs for each vehicle
 10 type. The model works by defining link-specific relationships between the International
 11 Roughness Index - IRI (Sayers et al., 1986) and speed, and using the IRI values at the
 12 respective average link speeds (derived from Sections 4.2.1 and 4.2.2) to estimate the
 13 link-specific VOCs, which are summed to estimate the route VOCs.

14 Table 4. HDM-III basic input data

Input	Measure	Unit	Car	Inter-urban taxi	Minibuses	Bus	Source
Terrain type	Rise & fall Horizontal curvature	m/km deg/km	Estimated directly for each link using Google Earth and Auto CAD Civil 3D				(Autodesk 2017, Google Earth 7.1.8.3036 2016)
Desired max speed	Desired max speed	km/hr	90	90	90	90	(WHO 2016)
Economic unit costs (Excluding taxes)*	Vehicle cost price	\$	19,452	27,098	49,474	67,465	(ADF 2011)
	Fuel type	NA	Petrol	Diesel	Diesel	Diesel	(ADF 2011)
	Fuel costs	\$/litre	1.017	0.727	0.727	0.727	(ADF 2011)
	Lubricants	\$/litre	4.688	6.466	6.466	6.466	(ADF 2011)
	Tyres + tubes	\$	83.36	83.36	83.36	187.53	(ADF 2011)
	Maintenance costs	\$/hour	0.284	0.310	0.310	0.541	(ADF 2011)
	Crew costs	\$/hour	0	0.511	0.511	0.511	(ADF 2011)
	Interest rate	%	5.4	5.4	5.4	5.4	(World Bank 2017b)
Utilisation	Mileage per year	km	25,000	50,000	50,000	60,000	(ADF 2011)
	Hour driven per year	hours	350	750	750	1250	(ADF 2011)
Service life	Service life	years	12	12	12	12	(ADF 2011)
Gross vehicle weight	Gross vehicle weight	tons	1.2	2.0	3.0	11	(Watanatada et al. 1987)

15 * Prices adjusted for 2010-2013 inflation and the 2013 USD exchange rate (World Bank 2017a, c)

1 The estimated route VOCs need to be converted to person costs. Given the
 2 anonymous nature of CDR data, we use information on the typical occupancy rates and
 3 mode shares (see Table 5) to estimate the weighted average VOCs per user for each
 4 route.

5 In this research, we have used the same average travel cost for all the users along
 6 a particular route between a given O-D pair. An improved approach would have been to
 7 estimate user-specific travel costs based on the corresponding travel speeds, however,
 8 this was not possible due to difficulties in obtaining the user-specific travel times as
 9 discussed in the previous section.

10 Table 5. Typical occupancy rates and mode shares in Senegal (World Bank 2004)

Vehicle type	Average number of passengers	Passengers per km	Mode share (%)
Cars	3	1746.7	0.183
Interurban Taxi	7	1830.4	0.191
Minibus	14	2149.6	0.225
Buses	25	3838.6	0.401

11
 12

13 The other attributes considered were the scenic characteristics and the urban
 14 developments along each route. Scenic variables were estimated in terms of route
 15 lengths traversed through nature reserves (Google Maps, 2017b), while urban
 16 developments were reflected as the number of towns along each route as shown in
 17 Figure 1 (Worldatlas, 2017).

18 5. Modelling framework

19 We use discrete choice models in this study since route choices are discrete and
 20 mutually exclusive. To develop these models, we apply the random utility theory
 21 (Marschak, 1960), a well-established approach for estimating discrete choice models.

22 5.1. Basic model

23 Suppose U_{nr} is the utility of choosing route r by individual n . This can be expressed as;

$$24 \quad U_{nr} = V_{nr} + \varepsilon_{nr} \quad (1)$$

25 Where V_{nr} and ε_{nr} are the systematic and the random parts utility respectively.
 26 The systematic utility is a function of the observed route attributes, and may be
 27 expressed as $V_{nr} = \beta' X_{nr}$, where X_{nr} is a vector of the attributes of route r for
 28 individual n and β is a vector the model parameters. We assume that the random term
 29 ε_{nr} is independently and identically distributed across the alternatives following a type I
 30 extreme value distribution, and use the Multinomial Logit (MNL) model to estimate the
 31 route choice probabilities as follows (see McFadden, 1974 for details);

$$P_n(r) = \frac{\exp(V_{nr})}{\sum_{r^* \in C_n} \exp(V_{nr^*})} \quad (2)$$

Where, $P_n(r)$ is the probability of individual n choosing route r , and C_n is the choice set. Given the route choice probabilities, the model parameters can be estimated by maximising the log-likelihood function below;

$$LL = \sum_n \sum_r [Z_{nr} \cdot \ln(P_n(r))] \quad (3)$$

Where Z_{nr} is a dummy variable, which is equal to 1 if and only if user n chooses route r .

5.2. Accounting for broad choices

The log-likelihood function in Equation (3) assumes that all the route choices are uniquely observed, and is inadequate for the current scenario, where we also have broad sub-group choices. Therefore, we use the broad choice modelling framework proposed by Wong (2015) to account for this situation.

In the broad choice framework, the choice probabilities of the broad sub-groups are expressed as a sum of the choice probabilities of the member alternatives. For example, the choice probability of the ‘Northern’ broad sub-group is the sum of the ‘Northern 1’ and the ‘Northern 2’ route choice probabilities (see Figure 4 and Table 2). The joint probabilities of the broad sub-groups capture the aggregate shares at the unique route choice level using the relative probabilities of the constituent routes.

The goal of model estimation is to maximise the probabilities of both the observed routes and the broad sub-groups for users with unique and broad choices respectively. The log-likelihood function is specified as follows (Wong, 2015);

$$LL = \sum_n \sum_b \left[Z_{nb} \cdot \ln \left(\sum_{r \in S_b} P_n(r) \right) \right] \quad (4)$$

Where S_b is a set comprising of the routes in broad category b . For uniquely assigned trips, set S_b comprises of only one alternative. Z_{nb} is a dummy variable, which is equal to 1 if and only if user n is associated with category b .

5.3. Accounting for overlap

A major weakness of the MNL model (Equation 2) is the IIA property, which could lead to illogical route choice probabilities for highly overlapping routes as is the case in this study. This is illustrated using the overlapping route problem (Ramming, 2002, Cascetta et al., 1996) in Figure 5.

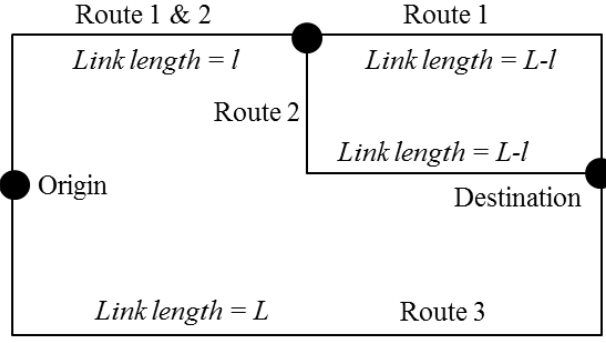


Figure 5. Overlapping route problem (Cascetta et al. 1996, Ramming 2002)

Here, all three routes have the same total length L , however, routes 1 and 2 follow the same alignment for length l followed by distinct sections, each of length $L - l$. The MNL model predicts equal shares for each of the routes irrespective of the overlap length. However, as the overlap increases (as l tends to L), it becomes difficult to distinguish between routes 1 and 2, and it is expected that route 3 will take a share of 50%, while routes 1 and 2 will each take a share of 25%.

Various models accounting for overlap were presented in Section 2.2, however, this study only uses the c-logit (Ramming, 2002, Cascetta et al., 1996) and the path size logit (Ben-Akiva and Ramming, 1998) models for illustration purposes as the modelling scenario did not require complex formulations. In general, these models are modifications of the MNL model, where the systematic route utilities are adjusted using certain correction factors as follows;

$$P_n(r) = \frac{\exp(V_{nr} + \tau_{nr})}{\sum_{r^* \in C_n} \exp(V_{nr^*} + \tau_{nr^*})} \quad (5)$$

Where τ_{nr} is the systematic utility correction factor for route r .

5.3.1. C-logit model

For the c-logit model, the correction factor τ_{nr} is a commonality factor (Cascetta et al., 1996). Different possible specifications have been proposed, however, this study uses the following common specification (Ramming, 2002, Cascetta et al., 1996);

$$\tau_{nr} = \beta_{CF} \ln \left[\sum_{r^* \in C_n} \left(\frac{L_{rr^*}}{\sqrt{L_r L_{r^*}}} \right)^{\gamma_{CF}} \right] \quad (6)$$

Where L_{rr^*} is the overlap length between routes r and r^* , L_r and L_{r^*} are the total lengths of routes r and r^* respectively, β_{CF} and γ_{CF} are the unknown parameters to be estimated. From Equation 6, the ratio in the brackets is proportional to the degree of overlap, while the corresponding logarithm has an inverse negative relationship. Thus β_{CF} and γ_{CF} are expected to have negative and positive signs respectively to allow for the positive adjustment of route utility with decreasing overlap (Ramming, 2002).

1 5.3.2. Path size logit model

2 For the path size logit model, the correction factor τ_{nr} is a path size term, which is
3 computed as the weighted average of the constituent link sizes. The specification
4 adopted for this study is as follows (Ben-Akiva and Ramming, 1998);

$$5 \quad \tau_{nr} = \beta_{ps} \ln \left[\sum_{a \in \Gamma_r} \left(\frac{l_a}{L_r} * \frac{1}{N_{ar}} \right) \right] \quad (7)$$

6 Where $1/N_{ar}$ is the inverse of the number of routes sharing the link a (the link
7 size), and l_a/L_r is a weight representing the proportion contributed by link a to the
8 overall route size. From Equation 7, it is observed that route size is inversely
9 proportional to the degree of overlap, while the corresponding logarithm has a negative
10 proportional relationship. Thus, the path size parameter is expected to be positive to
11 allow for negative adjustment of route utility with increasing overlap.

12 The models accounting for overlap are generally expected to have better fit than
13 the MNL model. Model fit is evaluated using the adjusted-rho square and the likelihood
14 ratio tests (see Ben-Akiva and Lerman, 1985 for details).

15 6. Model results

16 This section presents the modelling results. We start by discussing the variable
17 specification, followed by the model estimation and validation results.

18 6.1. Variable specification

19 The attributes available for possible inclusion in the model are summarised in Table 3.
20 However, these could not all be specified together or in the same way for various
21 reasons as explained below.

22 For travel time and cost, after initial tests using the linear specification, we used
23 the log-transforms of the variables to allow for utility damping with respect to
24 increasing time and cost (see Daly, 2010 for details). Various interactions of these
25 variables with others (such as surface type) were tested, however, this led to correlation
26 problems, and hence generic variables were specified.

27 The urban developments along the alternative routes were incorporated in terms
28 of the average distance between towns rather than the number of towns to avoid
29 situations where longer routes also have more towns. Again, this was specified using the
30 log-transform of the variable for similar reasons as the travel time and cost. This being a
31 largely rural road network with no traffic signals, the average distance between towns is
32 the only variable we could use to capture traffic flow interruptions.

33 An attempt was made to incorporate scenic beauty into the model using either
34 the length or the proportion of route length traversed through nature reserves, however,
35 we could not achieve intuitive model results potentially due to our lack of detailed
36 knowledge about the characteristics of these reserves, and the security levels of the
37 corresponding routes. The final systematic utility specification is as follows;

$$38 \quad V_{nr} = \beta_{l-cost} \ln(C_{nr}) + \beta_{l-time} \ln(T_{nr}) + \beta_{l-dtown} \ln(Dt_{nr}) \quad (8)$$

1 Where C_{nr} , T_{nr} , and Dt_{nr} give the travel cost, the travel time, and the average
2 distance between towns respectively of route r for individual n , and the β s are the
3 corresponding model parameters to be estimated.

4 **6.2. Estimation results**

5 We present the estimation results of the MNL model, the c-logit model, and the path
6 size logit based on the full sample in Table 6 for comparison purposes.

7 7. Table 6. Estimation results

Variable	MNL model		C-logit model		Path size logit model	
	Parameter	t-stat	Parameter	t-stat	Parameter	t-stat
Route variables						
Natural log of travel cost (US Dollars)	-4.2221	-11.60	-4.1644	-11.50	-4.4117	-18.99
Natural log of travel time (Hours)	-1.6668	-13.78	-1.5987	-12.62	-1.1018	-17.87
Natural log of av. distance between towns (Km)	0.5081	2.54	0.1705	1.10	1.5264	2.00
Commonality factor						
Beta			-0.0063	-1.48		
Gamma			0.6563	2.59		
Path size						
Path size parameter					1.6068	7.27
Measures of fit in estimation						
No. of observations	9453		9453		9453	
No. of decision makers	6497		6497		6497	
LL(C)	-11,135.41		-11,135.41		-11,135.41	
LL(F)	-7,758.83		-7,752.91		-7,547.60	
Number of parameters	3		5		4	
ρ_{adj}^2 w.r.t LL(C)	0.3030		0.3033		0.3218	
LR w.r.t LL(C)	6,753.17		6,764.99		7,175.63	
p-value of LR	0.0000		0.0000		0.0000	

8 **7.1.1. Statistical performance of the models**

9 The estimated parameters in each of the three models are statically significant at the
10 95% level of confidence. The only exceptions are the parameters associated with the
11 average distance between towns and the beta coefficient of the commonality term in the
12 c-logit model, however, these were retained in the model as they are important.

13 To evaluate the collective statistical significance of the models, we used the
14 likelihood ratio test (see Ben-Akiva and Lerman, 1985 for details), whose values are
15 reported in the last two rows of Table 6. From these, it is noted that the p-values for all

1 the models are less than 0.01. Therefore, the hypothesis that the respective model
 2 parameters are collectively equal to zero is rejected at the 99% level of confidence.

3 *7.1.2. Route variables*

4 The parameter signs for the travel cost variable are consistent with a priori expectations
 5 in each of the three models. In general, an increase in the cost of an alternative is
 6 expected to have a negative impact on its utility, hence the negative parameter sign. The
 7 same explanation holds for the travel time parameters as individuals generally prefer
 8 shorter travel times.

9 On the other hand, the average distance between towns has a positive parameter
 10 sign. As earlier mentioned, this variable gives an indication of the amount of
 11 uninterrupted flow. Although no traffic congestion problems have been reported in
 12 these towns, traffic generally slows down due to speed control measures leading to
 13 delays. An increase in the average distance between towns therefore indicates more
 14 uninterrupted flow, hence the positive parameter sign.

15 *7.1.3. Overlap correction parameters*

16 For the c-logit model, it is observed that the beta parameter of the commonality term is
 17 negative while the gamma parameter is positive. Similarly, the path size parameter in
 18 the path size logit model has a positive parameter sign. These results are in line with
 19 behavioural expectations as discussed earlier under Equations 6 and 7, an indication that
 20 CDR data is able to capture the behaviour towards overlapping routes.

21 *7.1.4. Model comparison*

22 A comparison of the adjusted rho-square values in Table 6 shows that the models
 23 accounting for overlap (i.e. the c-logit and the path size logit models) perform better
 24 than the MNL model. This is as expected given that the national road network of
 25 Senegal is highly overlapping (see Figure 1 and discussion under Figure 5). It is also
 26 worth noting that the path size logit model outperforms the c-logit model because the
 27 behavioural underpinning of the systematic utility adjustment process in the path size
 28 logit model is stronger than that in the c-logit model (Ramming, 2002).

29 The statistical significance of the improvements associated with accounting for
 30 overlap are evaluated using the likelihood ratios of the c-logit and the path size logit
 31 models with respect to the MNL model (see Table 7).

32 Table 7. Statistical comparison of the models

MNL formulation	C-logit formulation			Path size logit formulation		
LL(F)	LL(F)	LR w.r.t MNL model	p-value	LL(F)	LR w.r.t MNL model	p-value
-7758.83	-7752.91	11.83	0.0027	-7547.60	422.46	0.0000

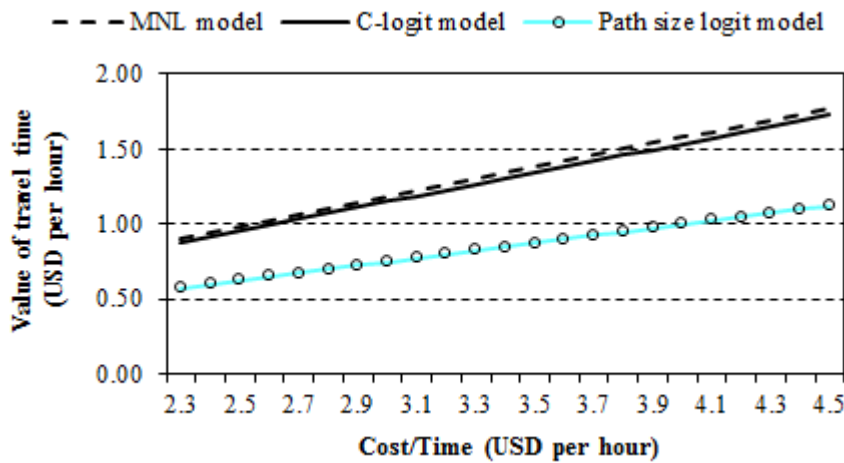
33
 34 From Table 7, it is noted that the p-values for the c-logit and the path size logit
 35 models are all less than 0.01, an indication that accounting for overlap has a statistically
 36 significant effect (at the 99% confidence level) beyond the improvements contributed
 37 by the additional degrees of freedom resulting from the extra parameters (see Ben-
 38 Akiva and Lerman, 1985 for details).

1 7.1.5. Policy insights

2 This section highlights the policy implications of the reported results in terms of the
 3 value of travel time (VTT). This metric quantifies the benefits derived from reduced
 4 travel time in monetary terms, and is useful in transportation cost-benefit analysis
 5 (Mackie et al., 2001). The value of travel time is calculated by taking the ratio of the
 6 partial derivatives of the systematic utility function (V) with respect to the travel time
 7 (T_{nr}) and cost (C_{nr}) as follows;

$$8 \quad VTT = \frac{\partial V_{nr} / \partial T_{nr}}{\partial V_{nr} / \partial C_{nr}} = \frac{\beta_{l-time} C_{nr}}{\beta_{l-cost} T_{nr}} \quad (9)$$

9 Figure 6 shows the variations of VTT with respect to the cost per hour of the
 10 alternatives for the range in the estimation data, and it is observed that the values
 11 increase as the alternatives become more expensive, which is expected. It is noted that
 12 the path size logit model gives the lowest values, which are deemed to be more accurate
 13 due to the superior performance (see Tables 6 and 7).



14

15 Figure 6. Variations in VTT across the alternatives and models

16

17 However, to assess how realistic these estimates are, we computed the average
 18 values for each model using the estimation data, and compared these values with those
 19 derived from other studies and other relevant statistics as summarised in Table 8.

20 Table 8. Comparison of the VTT estimates with other sources

Model	Values in USD/hr and 2013 prices			
	VTT current study	VTT Teye et al. (2017) meta-analysis	Dakar–Diamniadio road toll (Gainer and Chan 2016)	Median hourly wage (Tijdens et al. 2012)
MNL	1.0822			
C-logit model	1.0524	4.3213	2.3411	0.6767
Path size logit model	0.6846			

In the Africa-wide meta-analysis by Teye et al. (2017), VTT is estimated as a function of the GDP per capita. However, the reported mean value (4.3213 USD/hr) seems high when compared to the toll being charged on the new Dakar–Diamniadio toll highway for a time saving of one hour (2.3411 USD/hr), a value that was highly criticised by the Senegalese media as being extremely high (Gainer and Chan, 2016).

Although the median hourly wages do not necessarily translate into the value of travel time, they give a good indication of the range in which these values should fall, and as observed in Table 8, the average VTT estimate for the path size logit model is very close to the Senegalese median hourly wage. We consider this VTT estimate to be more reasonable for Senegal.

7.2. Validation results

The models based on the full sample provide intuitive results in terms of the parameter signs and the relative model performance. To assess the stability and the predictive performance of the models, the dataset was randomly split into five parts at the individual level. Five rolling subsets, each comprising of 80% of the users were generated for model estimation purposes. For each of these, a complementary subset comprising of 20% of the users was generated for validation purposes. The models were re-estimated on each of the 80% subsets, and the parameter estimates applied to the corresponding 20% hold-out subsets to estimate the predictive measures of fit. Table 9 presents the summary outputs from this process.

Table 9. Validation results

Subset	MNL model		C-logit model		Path size logit model	
	LL(F)	Adjusted rho-square	LL(F)	Adjusted rho-square	LL(F)	Adjusted rho-square
<i>Estimation subsets (comprising of 80% of the users)</i>						
Subset 1	-6201.66	0.3179	-6196.67	0.3182	-6045.15	0.3350
Subset 2	-6112.08	0.2999	-6106.49	0.3003	-5952.67	0.3180
Subset 3	-6306.93	0.2952	-6297.40	0.2961	-6138.29	0.3140
Subset 4	-6272.84	0.2993	-6263.20	0.3001	-6125.94	0.3156
Subset 5	-6184.90	0.2970	-6173.96	0.2980	-6018.68	0.3157
<i>Validation subsets (comprising of 20% of the users)</i>						
Subset 1	-1573.97	0.2267	-1569.71	0.2278	-1531.79	0.2469
Subset 2	-1661.15	0.3069	-1657.24	0.3077	-1621.68	0.3229
Subset 3	-1466.52	0.3265	-1466.01	0.3259	-1436.50	0.3398
Subset 4	-1502.27	0.3093	-1502.25	0.3084	-1450.65	0.3325
Subset 5	-1592.25	0.3165	-1590.68	0.3163	-1558.17	0.3306

23
24
25
26
27

1
2
3 The general interpretation of the parameter signs and the relative model
4 performance in each of the 80% estimation subsets remained the same as in the full
5 sample, an indication that the data is representative (detailed results available on
6 request). A comparison of the measures of fit in estimation and validation shows that
7 there is no significant loss in model fit, an indication that the performance of the models
8 during estimation is not due to overfitting, rather it is due to the strong explanatory
9 power of the variables.

10 A comparison of the predictive measures of fit shows that the relative model
11 performance during estimation is mirrored on the holdout samples, with the path size
12 logit model still giving the best model performance due to its behavioural superiority. It
13 would have been interesting to further validate the above results with outputs of route
14 choice models based on traditional data or GSM data, but this was not possible due to
15 lack of data in Senegal.

16 **8. Summary and conclusions**

17 This paper has successfully demonstrated the potential of CDR data to capture rational
18 route choice behaviour for long-distance inter-regional O-D pairs. The broad choice
19 framework was used to leverage the limitations of CDR data where unique route
20 choices could not be observed for some users, and only the broad sub-groups of the
21 possible routes were identifiable. This study is unique in the sense that it adapts the
22 broad choice framework to the context of route choice modelling using noisy CDR data.

23 An examination of the parameter signs shows that CDR data is able to capture
24 the expected sensitivities towards particular route attributes. A review of different
25 models accounting for overlap was conducted, and among these, the c-logit and the path
26 size logit models were considered. A comparison of these models against the
27 multinomial logit model (which does not account for overlap) showed significant
28 improvements in model fit, with the path size logit model giving the best performance.
29 The validation runs based on the 20% holdout samples largely showed the same
30 advantages in prediction, especially for the path size logit model. These results show
31 that CDR is able to capture the expected behaviour towards overlapping routes.

32 This study is timely as it extends the application of CDR data beyond travel
33 pattern analyses to econometric modelling of route choice that can be used for
34 forecasting. The proposed framework can help in the assessment of different policy
35 implications at a low cost compared to traditional approaches, which involve expensive
36 data collection. For example, the models developed in this study can be used to reliably
37 estimate the value of travel time (VTT) as we have demonstrated. This study is thus
38 beneficial to developing countries where budget constraints on transport studies are
39 common and traditional data for transport studies is scarce.

40 We conclude that the study findings serve as a proof-of-concept that CDR data
41 can be successfully used to model route choice behaviour for long-distance inter-
42 regional trips, where there is a strong possibility that a user will use his/her phone
43 during the trip, thereby enabling the capture of their partial trajectories needed for route
44 identification. It may be noted that with the increasing trend of mobile internet usage
45 (Gerpott and Thomas, 2014), the temporal resolution of CDR locations is likely to
46 improve significantly in the near future, and this could make CDR data suitable for
47 evaluating route choice behaviour for short trips.

1 A comparison of the study findings with those based on traditional data from
2 Senegal would have been insightful, however, this was not possible due to data
3 unavailability. Investigating the performance of the proposed approach in urban or
4 intra-city scenarios would be an interesting direction for future research.

5 **Acknowledgements**

6 We would like to thank the Economic and Social Research Council (ESRC) of the
7 UK and Institute for Transport Studies, University of Leeds for funding this research. The
8 authors also acknowledge the financial support by the European Research Council
9 through the consolidator grant 615596-DECISIONS. The research in this paper used the
10 Data for Development (D4D) Challenge dataset from Senegal made available by Sonatel
11 Group and Orange Group.

12 **References**

- 13 Agard, B., Morency, C. & Trépanier, M. 2006. Mining public transport user behaviour
14 from smart card data. *IFAC Proceedings Volumes*, 39, 399-404.
- 15 Ageroute Senegal. 2017. *Programmes & Projets* [Online]. Ageroute Senegal. Available:
16 http://www.ageroute.sn/index.php/rapports-d-activites-annuels/cat_view.html
17 [Accessed 12 April 2017].
- 18 Ahas, R., Aasa, A., Mark, Ü., Pae, T. & Kull, A. 2007. Seasonal tourism spaces in
19 Estonia: Case study with mobile positioning data. *Tourism management*, 28,
20 898-910.
- 21 Ahas, R., Aasa, A., Roose, A., Mark, Ü. & Silm, S. 2008. Evaluating passive mobile
22 positioning data for tourism surveys: An Estonian case study. *Tourism*
23 *Management*, 29, 469-486.
- 24 Ahas, R., Aasa, A., Silm, S. & Tiru, M. 2010. Daily rhythms of suburban commuters'
25 movements in the Tallinn metropolitan area: Case study with mobile positioning
26 data. *Transportation Research Part C: Emerging Technologies*, 18, 45-54.
- 27 Ahas, R., Aasa, A., Yuan, Y., Raubal, M., Smoreda, Z., Liu, Y., Ziemlicki, C., Tiru, M.
28 and Zook, M., 2015. Everyday space-time geographies: using mobile phone-
29 based sensor data to monitor urban activity in Harbin, Paris, and Tallinn.
30 *International Journal of Geographical Information Science*, 29(11), pp.2017-
31 2039.
- 32 Ansd 2016. Rapport projection de la population du Senegal. Dakar, Senegal: Agence
33 Nationale de la Statistique et de la Démographie.
- 34 Ansd. 2017. *Situation Economique et Sociale* [Online]. Agence Nationale de la
35 Statistique et de la Démographie. Available:
36 http://www.ansd.sn/index.php?option=com_sess&view=sess&Itemid=398
37 [Accessed 11 April 2017].
- 38 Arcgis. 2013. *Senegal Roads* [Online]. Available:
39 [https://services1.arcgis.com/4AWkjgqSzd8pqxQA/arcgis/rest/services/Senegal_](https://services1.arcgis.com/4AWkjgqSzd8pqxQA/arcgis/rest/services/Senegal_Roads/FeatureServer/0)
40 [Roads/FeatureServer/0](https://services1.arcgis.com/4AWkjgqSzd8pqxQA/arcgis/rest/services/Senegal_Roads/FeatureServer/0) [Accessed 01 March 2017].
- 41 Azevedo, J., Costa, M. E. O. S., Madeira, J. J. E. S. & Martins, E. Q. V. 1993. An
42 algorithm for the ranking of shortest paths. *European Journal of Operational*
43 *Research*, 69, 97-106.
- 44 Bellman, R. & Kalaba, R. 1960. On k th best policies. *Journal of the Society for*
45 *Industrial and Applied Mathematics*, 8, 582-588.
- 46 Ben-Akiva, M., Bergman, M., Daly, A. J. & Ramaswamy, R. Modeling inter-urban
47 route choice behaviour. Proceedings of the 9th International Symposium on

- 1 Transportation and Traffic Theory, 1984. VNU Science Press Utrecht, The
2 Netherlands, 299-330.
- 3 Ben-Akiva, M. & Ramming, S. 1998. Lecture notes: Discrete choice models of traveler
4 behavior in networks. *Prepared for Advanced Methods for Planning and*
5 *Management of Transportation Networks. Capri, Italy, 25.*
- 6 Ben-Akiva, M. E. & Lerman, S. R. 1985. *Discrete choice analysis: theory and*
7 *application to travel demand*, MIT press.
- 8 Bierlaire, M., Chen, J. & Newman, J. 2010. Modeling route choice behavior from
9 smartphone GPS data.
- 10 Bierlaire, M., Chen, J. & Newman, J. 2013. A probabilistic map matching method for
11 smartphone GPS data. *Transportation Research Part C: Emerging*
12 *Technologies*, 26, 78-98.
- 13 Bierlaire, M. & Frejinger, E. 2008. Route choice modeling with network-free data.
14 *Transportation Research Part C: Emerging Technologies*, 16, 187-198.
- 15 Bolla, R., Davoli, F. & Giordano, F. Estimating road traffic parameters from mobile
16 communications. Proceedings 7th World Congress on ITS, Turin, Italy, 2000.
- 17 Broach, J., Dill, J. & Gliebe, J. 2012. Where do cyclists ride? A route choice model
18 developed with revealed preference GPS data. *Transportation Research Part A:*
19 *Policy and Practice*, 46, 1730-1740.
- 20 Calabrese, F., Di Lorenzo, G., Liu, L. & Ratti, C. 2011. Estimating Origin-Destination
21 flows using opportunistically collected mobile phone location data from one
22 million users in Boston Metropolitan Area. *IEEE Pervasive Computing*, 10, 36-
23 44.
- 24 Calabrese, F., Diao, M., Di Lorenzo, G., Ferreira, J. & Ratti, C. 2013. Understanding
25 individual mobility patterns from urban sensing data: A mobile phone trace
26 example. *Transportation research part C: emerging technologies*, 26, 301-313.
- 27 Cascetta, E., Nuzzolo, A., Russo, F. & Vitetta, A. A modified logit route choice model
28 overcoming path overlapping problems: specification and some calibration
29 results for interurban networks. Proceedings of the 13th International
30 Symposium on Transportation and Traffic Theory, 1996. Pergamon Lyon,
31 France, 697-711.
- 32 Cascetta, E. & Papola, A. 2001. Random utility models with implicit
33 availability/perception of choice alternatives for the simulation of travel
34 demand. *Transportation Research Part C: Emerging Technologies*, 9, 249-263.
- 35 Chakirov, A. & Erath, A. 2012. Activity identification and primary location modelling
36 based on smart card payment data for public transport.
- 37 Çolak, S., Alexander, L. P., Alvim, B. G., Mehndiretta, S. R. & González, M. C.
38 Analyzing Cell Phone Location Data for Urban Travel: Current Methods,
39 Limitations and Opportunities. Transportation Research Board 94th Annual
40 Meeting, 2015.
- 41 Daganzo, C. F., Bouthelier, F. & Sheffi, Y. 1977. Multinomial probit and qualitative
42 choice: A computationally efficient algorithm. *Transportation Science*, 11, 338-
43 358.
- 44 Daly, A. 2010. Cost damping in travel demand models: Report of a study for the
45 Department for Transport. United Kingdom: RAND Corporation.
- 46 De La Barra, T., Perez, B. & Anez, J. Multidimensional path search and assignment.
47 PTRC Summer Annual Meeting, 21st, 1993, University of Manchester, United
48 Kingdom, 1993.

- 1 De Montjoye, Y.-A., Smoreda, Z., Trinquart, R., Ziemlicki, C. & Blondel, V. D. 2014.
2 D4D-Senegal: the second mobile phone data for development challenge. *arXiv*
3 *preprint arXiv:1407.4885*.
- 4 Deville, P., Song, C., Eagle, N., Blondel, V. D., Barabási, A.-L. & Wang, D. 2016.
5 Scaling identity connects human mobility and social interactions. *Proceedings of*
6 *the National Academy of Sciences*, 201525443.
- 7 Diao, M., Zhu, Y., Ferreira Jr, J. & Ratti, C. 2016. Inferring individual daily activities
8 from mobile phone traces: A Boston example. *Environment and Planning B:*
9 *Planning and Design*, 43, 920-940.
- 10 Doyle, J., Hung, P., Kelly, D., Mcloone, S. & Farrell, R. 2011. Utilising mobile phone
11 billing records for travel mode discovery.
- 12 Gainer, M. & Chan, S. 2016. A NEW ROUTE TO DEVELOPMENT: SENEGAL'S
13 TOLL HIGHWAY PUBLIC-PRIVATE PARTNERSHIP, 2003 – 2013. New
14 Jersey, USA: Innovations for Successful Societies, Princeton University.
- 15 Gerpott, T. J. & Thomas, S. 2014. Empirical research on mobile Internet usage: A meta-
16 analysis of the literature. *Telecommunications Policy*, 38, 291-310.
- 17 Google Maps. 2017a. *Google map directions* [Online]. Google. Available:
18 <https://www.google.co.uk/maps/dir> [Accessed 13 April 2017].
- 19 Google Maps. 2017b. *Senegal nature reserves* [Online]. Google. Available:
20 <https://www.google.co.uk/maps/@15.1978209,-15.0824015,8.67z> [Accessed 28
21 December 2017].
- 22 Görnerup, O. Scalable Mining of Common Routes in Mobile Communication Network
23 Traffic Data. *Pervasive*, 2012. Springer, 99-106.
- 24 Groves, R. M. 2006. Nonresponse rates and nonresponse bias in household surveys.
25 *Public opinion quarterly*, 646-675.
- 26 Gsm Association. 2017. *The Mobile Economy 2017* [Online]. Available:
27 [https://www.gsmainelligence.com/research/?file=9e927fd6896724e7b26f33f61](https://www.gsmainelligence.com/research/?file=9e927fd6896724e7b26f33f61db5b9d5&download)
28 [db5b9d5&download](https://www.gsmainelligence.com/research/?file=9e927fd6896724e7b26f33f61db5b9d5&download) [Accessed 04 November 2017].
- 29 Habibi, S., Frejinger, E. and Sundberg, M., 2017. An empirical study on aggregation of
30 alternatives and its influence on prediction in car type choice models. *Transportation*,
31 pp.1-20.
- 32 Hasan, S., Zhan, X. & Ukkusuri, S. V. Understanding urban human activity and
33 mobility patterns using large-scale location-based data from online social media.
34 *Proceedings of the 2nd ACM SIGKDD international workshop on urban*
35 *computing*, 2013. ACM, 6.
- 36 Hawelka, B., Sitko, I., Beinat, E., Sobolevsky, S., Kazakopoulos, P. & Ratti, C. 2014.
37 Geo-located Twitter as proxy for global mobility patterns. *Cartography and*
38 *Geographic Information Science*, 41, 260-271.
- 39 Hess, S., Quddus, M., Rieser-Schüssler, N. & Daly, A. 2015. Developing advanced
40 route choice models for heavy goods vehicles using GPS data. *Transportation*
41 *Research Part E: Logistics and Transportation Review*, 77, 29-44.
- 42 Hoteit, S., Secci, S., Sobolevsky, S., Ratti, C. & Pujolle, G. 2014. Estimating human
43 trajectories and hotspots through mobile phone data. *Computer Networks*, 64,
44 296-307.
- 45 Imedia & Calao Production. 2013. *Le chemin de fer sénégalais* [Online]. AU-
46 SENEGAL.COM. Available: [http://www.au-senegal.com/le-chemin-de-](http://www.au-senegal.com/le-chemin-de-fer,345?lang=fr)
47 [fer,345?lang=fr](http://www.au-senegal.com/le-chemin-de-fer,345?lang=fr) [Accessed 06 August 2017].

- 1 Iqbal, M. S., Choudhury, C. F., Wang, P. & González, M. C. 2014. Development of
2 origin–destination matrices using mobile phone call data. *Transportation*
3 *Research Part C: Emerging Technologies*, 40, 63-74.
- 4 Isaacman, S., Becker, R., Cáceres, R., Martonosi, M., Rowland, J., Varshavsky, A. &
5 Willinger, W. Human mobility modeling at metropolitan scales. Proceedings of
6 the 10th international conference on Mobile systems, applications, and services,
7 2012. Acm, 239-252.
- 8 Jiang, S., Fiore, G. A., Yang, Y., Ferreira Jr, J., Frazzoli, E. & González, M. C. A
9 review of urban computing for mobile phone traces: current methods, challenges
10 and opportunities. Proceedings of the 2nd ACM SIGKDD International
11 Workshop on Urban Computing, 2013. ACM, 2.
- 12 Kerali, H. G. R. 2000. *Overview of HDM-4*, Paris, The World Road Association
13 (PIARC), Paris and The World Bank, Washington, DC.
- 14 Kung, K. S., Greco, K., Sobolevsky, S. & Ratti, C. 2014. Exploring universal patterns in
15 human home-work commuting from mobile phone data. *PloS one*, 9, e96180.
- 16 Leontiadis, I., Lima, A., Kwak, H., Stanojevic, R., Wetherall, D. & Papagiannaki, K.
17 From cells to streets: Estimating mobile paths with cellular-side data.
18 Proceedings of the 10th ACM International on Conference on emerging
19 Networking Experiments and Technologies, 2014. ACM, 121-132.
- 20 Li, L., Wang, S. & Wang, F.-Y. 2018. An Analysis of Taxi Driver’s Route Choice
21 Behavior Using the Trace Records. *IEEE Transactions on Computational Social*
22 *Systems*, 5, 576-582.
- 23 Lloro, A. and Brownstone, D., 2018. Vehicle choice and utilization: Improving
24 estimation with partially observed choices and hybrid pairs. *Journal of choice*
25 *modelling*, 28, pp.137-152.
- 26 Logistics Cluster. 2013. 2.3 *Senegal Road Assessment* [Online]. Logistics Cluster and
27 World Food Programme. Available:
28 <http://dlca.logcluster.org/display/public/DLCA/2.3+Senegal+Road+Assessment>
29 [Accessed 11 May 2017].
- 30 Mackie, P., Jara-Díaz, S. & Fowkes, A. 2001. The value of travel time savings in
31 evaluation. *Transportation Research Part E: Logistics and Transportation*
32 *Review*, 37, 91-106.
- 33 Ma, J., Li, H., Yuan, F. and Bauer, T., 2013. Deriving operational origin-destination
34 matrices from large scale mobile phone data. *International Journal of*
35 *Transportation Science and Technology*, 2(3), pp.183-204.
- 36 Mai, T., Fosgerau, M. & Frejinger, E. 2015. A nested recursive logit model for route
37 choice analysis. *Transportation Research Part B: Methodological*, 75, 100-112.
- 38 Manski, C. F. 1977. The structure of random utility models. *Theory and decision*, 8,
39 229-254.
- 40 Marschak, J. 1960. Binary Choice Constraints on Random Utility Indications. In:
41 ARROW, K. (ed.) *Stanford Symposium on Mathematical Methods in the Social*
42 *Science*. Stanford, California: Stanford University Press.
- 43 Mcfadden, D. 1974. Conditional logit analysis of qualitative choice behavior. *Frontiers*
44 *in Econometrics*, 105-142.
- 45 Nie, J., Zhang, J., Zhong, G. & Hu, Y. 2015. A Novel Approach to Road Matching
46 Based on Cell Phone Handover. *CICTP 2015*.
- 47 Nielsen, O. A. 2000. A stochastic transit assignment model considering differences in
48 passengers utility functions. *Transportation Research Part B: Methodological*,
49 34, 377-402.

- 1 Nitsche, P., Widhalm, P., Breuss, S., Brändle, N. & Maurer, P. 2014. Supporting large-
2 scale travel surveys with smartphones—a practical approach. *Transportation*
3 *Research Part C: Emerging Technologies*, 43, 212-221.
- 4 Nitsche, P., Widhalm, P., Breuss, S. & Maurer, P. 2012. A strategy on how to utilize
5 smartphones for automatically reconstructing trips in travel surveys. *Procedia-*
6 *Social and Behavioral Sciences*, 48, 1033-1046.
- 7 Papinski, D., Scott, D. M. & Doherty, S. T. 2009. Exploring the route choice decision-
8 making process: A comparison of planned and observed routes obtained using
9 person-based GPS. *Transportation research part F: traffic psychology and*
10 *behaviour*, 12, 347-358.
- 11 Pei, T., Sobolevsky, S., Ratti, C., Shaw, S.-L., Li, T. & Zhou, C. 2014. A new insight
12 into land use classification based on aggregated mobile phone data.
13 *International Journal of Geographical Information Science*, 28, 1988-2007.
- 14 Prato, C. & Bekhor, S. 2006. Applying branch-and-bound technique to route choice set
15 generation. *Transportation Research Record: Journal of the Transportation*
16 *Research Board*, 19-28.
- 17 Prato, C. G. 2009. Route choice modeling: past, present and future research directions.
18 *Journal of choice modelling*, 2, 65-100.
- 19 Qu, Y., Gong, H. & Wang, P. Transportation mode split with mobile phone data.
20 Intelligent Transportation Systems (ITSC), 2015 IEEE 18th International
21 Conference on, 2015. IEEE, 285-289.
- 22 Ramming, M. S. 2002. *Network knowledge and route choice*. Ph.D. Thesis,
23 Massachusetts Institute of Technology, Cambridge, USA.
- 24 Reddy, S., Burke, J., Estrin, D., Hansen, M. & Srivastava, M. Determining
25 transportation mode on mobile phones. *Wearable Computers*, 2008. ISWC
26 2008. 12th IEEE International Symposium on, 2008. IEEE, 25-28.
- 27 Roupail, N. M., Ranjithan, S. R., El Dessouki, W., Smith, T. & Brill, E. D. A decision
28 support system for dynamic pre-trip route planning. *Applications of Advanced*
29 *Technologies in Transportation Engineering*, 1995. ASCE, 325-329.
- 30 Saravanan, M., Pravinth, S. V. & Holla, P. Route detection and mobility based
31 clustering. *Internet Multimedia Systems Architecture and Application*
32 *(IMSAA)*, IEEE 5th International Conference, 2011. IEEE, 1-7.
- 33 Sayers, M. W., Gillespie, T. D. & Queiroz, C. a. V. 1986. The international road
34 roughness experiment: establishing correlation and a calibration standard for
35 measurements. *World Bank Technical Paper No. 45*.
- 36 Schlaich, J. 2010. Analyzing route choice behavior with mobile phone trajectories.
37 *Transportation Research Record: Journal of the Transportation Research*
38 *Board*, 78-85.
- 39 Schlaich, J., Otterstätter, T. & Friedrich, M. Generating trajectories from mobile phone
40 data. *Proceedings of the 89th annual meeting compendium of papers,*
41 *transportation research board of the national academies*, 2010.
- 42 Shaw, S.L., Tsou, M.H. and Ye, X., 2016. Human dynamics in the mobile and big data
43 era. *International Journal of Geographical Information Science*, 30(9), pp.1687-
44 1693.
- 45 Sheffi, Y. & Powell, W. B. 1982. An algorithm for the equilibrium assignment problem
46 with random link times. *Networks*, 12, 191-207.
- 47 Shier, D. R. 1979. On algorithms for finding the k shortest paths in a network.
48 *Networks*, 9, 195-214.
- 49 Song, C., Koren, T., Wang, P. & Barabási, A.-L. 2010. Modelling the scaling properties
50 of human mobility. *Nature Physics*, 6, 818-823.

- 1 Tettamanti, T., Demeter, H. & Varga, I. 2012. Route choice estimation based on cellular
2 signaling data. *Acta Polytechnica Hungarica*, 9, 207-220.
- 3 Tran, L.T.T., Ly, P.T.M. and Le, L.T., 2019. Hotel choice: A closer look at
4 demographics and online ratings. *International Journal of Hospitality Management*, 82,
5 pp.13-21.
- 6 Teye, C., Davidson, P., Porter, H. & Bell, M. G. H. 2017. Meta-analysis on the value of
7 travel time savings in Africa. *International Choice Modelling Conference 2017*.
8 Cape Town, South Africa.
- 9 Vovsha, P. & Bekhor, S. 1998. Link-nested logit model of route choice: overcoming
10 route overlapping problem. *Transportation Research Record: Journal of the*
11 *Transportation Research Board*, 133-142.
- 12 Vrtic, M., Schuessler, N., Erath, A., Axhausen, K., Frejinger, E., Stojanovic, J.,
13 Bierlaire, M., Rudel, R. & Maggi, R. 2006. Including travelling costs in the
14 modelling of mobility behaviour. Final report for SVI research program
15 Mobility Pricing: Project B1, on behalf of the Swiss Federal Department of the
16 Environment, Transport, Energy and Communications. IVT ETH Zurich, ROSO
17 EPF Lausanne and USI Lugano.
- 18 Wang, H., Calabrese, F., Di Lorenzo, G. & Ratti, C. Transportation mode inference
19 from anonymized and aggregated mobile phone call detail records. *Intelligent*
20 *Transportation Systems (ITSC)*, 2010 13th International IEEE Conference on,
21 2010. IEEE, 318-323.
- 22 Watanatada, T., Herral, C., Paterson, W., Dhareshwar, A., Bhandari, A. & Tsunokawa,
23 K. 1987. The Highway Design and Maintenance Standards Model. *Volume 1*
24 *Description of the HDM-III Model*. Baltimore and London: The John Hopkins
25 University Press.
- 26 White, J. & Wells, I. 2002. Extracting origin destination information from mobile phone
27 data.
- 28 Wong, T. C. J. 2015. *Econometric Models in Transportation*. Ph.D. Thesis, University
29 of California, Irvine.
- 30 Wong, T., Brownstone, D. and Bunch, D.S., 2018. Aggregation biases in discrete choice
31 models. *Journal of Choice Modelling*.
- 32 World Bank 2004. Performance and impact indicators for transport in Senegal: Detailed
33 statistics - June 2004. World Bank.
- 34 Worldatlas. 2017. *Senegal Facts* [Online]. Worldatlas. Available:
35 <https://www.worldatlas.com/webimage/countrys/africa/senegal/snfacts.htm>
36 [Accessed 22 November 2017].
- 37 Xu, Y., Shaw, S.-L., Zhao, Z., Yin, L., Fang, Z. & Li, Q. 2015. Understanding
38 aggregate human mobility patterns using passive mobile phone location data: a
39 home-based approach. *Transportation*, 42, 625-646.
- 40 Yip, A.H., Michalek, J.J. and Whitefoot, K.S., 2018. On the implications of using
41 composite vehicles in choice model prediction. *Transportation Research Part B:*
42 *Methodological*, 116, pp.163-188.
- 43 Yuan, Y. & Raubal, M. Extracting dynamic urban mobility patterns from mobile phone
44 data. *International Conference on Geographic Information Science*, 2012.
45 Springer, 354-367.
- 46 Yuan, Y., Raubal, M. and Liu, Y., 2012. Correlating mobile phone usage and travel
47 behavior—A case study of Harbin, China. *Computers, Environment and Urban*
48 *Systems*, 36(2), pp.118-130.

- 1 Wang, Y., Dong, L., Liu, Y., Huang, Z. and Liu, Y., 2019. Migration patterns in China
- 2 extracted from mobile positioning data. *Habitat International*, 86, pp.71-80.
- 3 Zhao, Z., Shaw, S.L., Xu, Y., Lu, F., Chen, J. and Yin, L., 2016. Understanding the bias
- 4 of call detail records in human mobility research. *International Journal of*
- 5 *Geographical Information Science*, 30(9), pp.1738-1762.

6