

Technical University of Munich

Development of a Destination Choice Model for Ontario

by

Joseph Molloy

A thesis submitted in partial fulfillment for the
degree of Master of Science in Transport Systems

in the

Chair of Modeling Spatial Mobility
Department of Civil, Geo and Environmental Engineering

November 2016

Declaration of Authorship

I, AUTHOR NAME, declare that this thesis titled, 'THESIS TITLE' and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

Date:

“Write a funny quote here.”

If the quote is taken from someone, their name goes here

Technical University of Munich

Abstract

Chair of Modeling Spatial Mobility

Department of Civil, Geo and Environmental Engineering

Master of Science in Transport Systems

by Joseph Molloy

The Thesis Abstract is written here (and usually kept to just this page). The page is kept centered vertically so can expand into the blank space above the title too...

Acknowledgements

The acknowledgements and the people to thank go here, don't forget to include your project advisor...

Contents

Declaration of Authorship	i
Abstract	iii
Acknowledgements	iv
List of Figures	vii
List of Tables	viii
Abbreviations	ix
1 Literature Review	1
1.1 Intercity Transport Models	2
1.2 Discrete Destination Choice Models	2
1.3 Trip Chaining	3
1.4 Recent Intercity Transport Models	4
1.5 Mnlogit R Package	5
1.5.1 Data format	5
1.5.2 Formula Specification	5
1.6 Location based social networks	6
1.7 Objectives	7
2 Data Acquisition and Analysis	9
2.1 Travel Survey of Residents of Canada	9
2.1.1 Introduction	9
2.1.2 Method	10
2.1.3 Data	10
2.1.3.1 Spatial Resolution	10
2.1.3.2 Error Detection and Imputation	11
2.1.3.3 Weighting	11
2.1.4 Microfile Format	12
2.1.4.1 Trip Datafile	13
2.1.4.2 Visit Datafile	14
2.2 Filtering of Trip Records	14

2.3	Zone System	17
2.4	Aggregating of Zonal Data	20
2.5	Foursquare	21
3	Model Design	22
4	Gravity Model	23
4.1	Design	23
4.2	Calibration	23
4.3	Results	23
5	Destination Choice Model	24
5.1	Estimation	24
5.1.1	Socio-economic Variables	24
5.1.2	Origin-Destination Interactions	24
5.1.3	Incorporating LBSN Data	24
5.1.4	Model Subsets	24
	Seasons	24
	Income Strata	24
5.1.5	Final Model	24
6	Implementation	25
6.1	MTO Model Integration	25
6.2	Algorithm	25
6.3	Disaggregation to TAZs	25
7	Validation	26
7.1	Case study of a new Ski Resort	26
8	Epilogue	27
8.1	Conclusions	27
8.2	Discussion	27
8.3	Future Work	27
A	An Appendix	28
	Bibliography	30

List of Figures

2.1	Dividing external zones into east and west	15
2.2	An example of an external origin-destination pair that passes through Ontario	16
2.3	Estimated vs. Observed Trip Distance	17
2.4	Different methods of aggregating internal zones to match the TSRC spa- tial resolution.	19

List of Tables

Abbreviations

LAH List Abbreviations **Here**

Chapter 1

Literature Review

Also known as Intercity models, long-distance transport models were first proposed in the 1960s, with two of the earliest being developed in the United States and Canada respectively (Canadian Transport Commission [1971](#)). These models were comparatively basic, with the demand component of the model only incorporating zonal population and income as attraction measures, and trip time, cost and convenience as impedance measures. More recent demand models include attributes such as auto ownership and household size.

Long-distance models are commonly defined to contain trips of certain length or longer, as opposed to the much more common urban model. TRB's NCHRP Report 735 notes that current state-wide models and travel surveys in the United States have used a range of thresholds to define long-distance trip-making, "either 50, 75, or 100 miles as the minimum threshold for trips to be considered long-distance." (Schiffer [2012](#))

According to Miller ([2004](#)), a distinct class of intercity travel demand models exist, which have unique characteristics when compared to urban models. "An intercity travel demand model is designed to forecast travel demand between two or more urban areas ... rather than travel within a given urban region". He also highlights two main features of such models. Firstly, he argues that an intercity travel demand model should apply to a well-defined travel corridor, containing a small number of major cities. Secondly he suggests that such models are almost always designed to model the impact of new travel modes such as high speed rail, or other policy initiatives.

Miller also notes that while urban models and methods are well documented in open literature, and applied in published policy analysis, intercity models are often the intellectual property of the consultants involved. The models are infrequently published in the scientific literature, and the travel data on private travel modes is often closely

guarded, meaning that the models are often hard to replicate, if they are published. It follows that “intercity travel demand models tend to be a less attractive/feasible application area for academic researchers than the more data-rich urban field.” (Miller 2004)

1.1 Intercity Transport Models

Until the 1980s, intercity transport models were exclusively designed as aggregate models, which distributed trips between origin-destination pairs and modes using the gravity model proposed by (Casey 1955). However, as early as 1962, the deficiencies in this approach were identified by numerous researchers (1962; 1962). In 1967, (Wilson 1967) first proved the theoretical validity of the gravity model, following two decades of its use in practice. These theoretical foundations encouraged further use of the gravity model in the development of transport demand models.

1.2 Discrete Destination Choice Models

Despite their widespread acceptance and use in practice, the aggregate gravity model has some fundamental flaws as a modelling methodology. As an alternative, McFadden (1973) and M. E. Ben-Akiva (1974) proposed the use of the logit model as a disaggregate method to model travel demand. McFadden, in his pioneering paper, noted that “When the model of choice behaviour under examination depends on unobserved characteristics in the population, the testable implications of the individual choice model are obscured.” (McFadden 1973) Further research focused on the use of disaggregate models for the trip distribution step of the classic four step model. These came to be known as destination choice models, the focus of this thesis. A thorough investigation of the suitability of discrete choice models as opposed to aggregate methods for transportation modelling was conducted by Spear (1977). Spear noted that the

- Individual choice models are more data efficient than conventional (i.e. gravity) models.
- They can utilise the variation in socio-economic data much better, avoiding ecological fallacies.
- The probabilistic nature of the dependent variable allows for the modelling of interdependent choices, such as mode choice and trip chaining decisions.

Since then, the application of disaggregate models in transport demand modelling has been continually refined, with important research done in both modelling destination choice and mode choice in this manner. Daly (1982) focused on representing the attractiveness of a destination, in a destination choice model, while further work was done by M. E. Ben-Akiva (1974) and Anas (1983) in defining the structure and application of such models. Train (2009) comments that “discrete choice models cannot be calibrated using a simple curve fitting, as since the dependent variable, as a probability cannot be observed”. Instead, maximum likelihood estimation is used. Since the utility of every alternative must be calculated, this technique was prohibitive before the advent of modern computers for large scale problems. This may go some way to explaining the persistent popularity of aggregate models, due to their simplified computational requirements. In chapter 9 of *Discrete Choice Analysis*, Akiva and Lerman (1985) present a comprehensive discussion of destination choice models. They note that “destination choice is characterised by a very large number of alternatives”, and that the selection of resolution of the choice set is a very important consideration. They further discuss the challenge of data availability for destination attractiveness. Since the attractiveness of data is not always available at a destination level, “the alternatives in a destination choice model must be based on aggregate alternatives”. Even with the modern GPS and social data available to the modern modeller, this is clearly still an issue. This is an important point that Ben-Akiva and Lerman make, that while destination choice models can model the decisions of individual travellers, they still need to rely on some level of aggregation for modelling the utility of each destination.

Discrete choice models have been repeatedly shown to provide better results than aggregate methods when modelling travel behaviour (1984; 2013) compared the gravity model and multinomial logit destination choice model when integrated into a model for Maryland. They found that the destination choice model performed much better than the gravity model for a state-wide travel demand model.

1.3 Trip Chaining

Adler and M. Ben-Akiva (1979) were one of the first to model the interdependencies between links in a trip chain. They defined a theoretical and empirical model of trip chaining behavior to do so, based on utility theory, and accounting for the tradeoffs involved in multistep chain trips. They, like most researchers in the area, focus on daily travel patterns within urban models. However, they note that “It is important that the determinants of non-work travel patterns that include multiple-sojourn tours be better understood”. To do this, they model the utility to a given household of a particular travel

pattern as a function of scheduling convenience, activity duration, income, destination attributes and socio-economic characteristics of the household. One of the significant advantages of a disaggregate destination model is the ability to model tours. Due to the nature of the data, trip chaining commonly isn't included in long distance travel models. Moeckel, Fussell, and Donnelly (2015) considered its inclusion, however the proportion of multi-link trips was found to be too small, and the trip lengths between stops were not recorded in the National Household Travel Survey (NHTS) data they used.

Kitamura (1984) incorporated trip chaining directly into an analysis of destination choice. He used an approach called Prospective Utility that “represents the expected utility of the visit to that zone and also those visits that may be made”. In essence, this theory postulates that with two destinations A and B of equal utility, opportunity B will be more attractive than A to a trip maker when it is surrounded by destinations supporting other opportunities that the trip maker wishes to pursue.

1.4 Recent Intercity Transport Models

(M. Outwater et al. 2010) developed a state-wide model to model high speed rail. They combined both stated and revealed preference data in their attributes. For destination choice, they looked at destination attraction, employment and household characteristics, the region and area type, trip purpose, distance class, and party size. While not a combined destination-mode choice model, they combine some network data to calculate auto and non-auto accessibility, for peak and off-peak respectively. Destination was estimated using a simple multinomial logit model. The authors also note that their modelling shows “that an individual may value different trip characteristics for different distance-categories of travel.”. Destination zones were found to be more attractive if better assessable at the 95% confidence level, underlying the importance of including the accessibility of the zone as an attribute to the model, not just the travel time. They also modelled the area type of a zone, as one of rural, suburban or urban. Interaction terms were also created between zones, under the assumption that urban to urban trips are much more common.

More recently, models are also being designed on a larger, more ambitious scale. One such example is the new national model of long-distance travel in the United States (M. L. Outwater et al. 2015). This model focuses includes multiple advancements on previous models, including modelling at an individual household level, a high level of spatial detail for destination choice, and the vertical integration of all 4 steps of the classical model. Unlike activity based models, it uses a temporal resolution of months and weeks, not days, and jointly predicts destination and mode choice together.

1.5 Mnlogit R Package

In this thesis, the R package `mnlogit` (CITE) is used to estimate the multinomial model of destination choice. This package provides improvements over the classic `mlogit` (CITE) package, by reducing memory usage, allowing for more model parameters, and performing the maximum likelihood estimation in parallel for significant decrease in estimation runtime.

1.5.1 Data format

The model input must be provided in a long format. The format of the input can be described as followed: Let S be the set of input trips, and R be the total choice set of possible destination alternatives.

For each trip $s \in S$, an arbitrary choice set R_s is required, where

$$R_s = \{a | R_s \subseteq R \wedge \exists t \in S : destination(t) = a\}$$

The model input is then constructed by adding a row for each trip s and alternative $a \in R_s$, giving a total number of rows $\sum_{s \in S} |R_s|$ for the model input. A boolean column is also added that indicates whether a particular choice was chosen for that trip or not.

1.5.2 Formula Specification

The `mnlogit` package accepts model formulas structured using the R formula package. A `mlogit` formula consists of 4 parts: $choice \sim Y|X|Z$, where X, Y, Z are as followed:

- Choice: the LHS of the equation, the column that indicates if an alternative was chosen or not.
- X : Individual i specific variables with alternative k specific coefficients $\vec{X}_i \vec{\beta}_k$.
- Y : Alternative specific variables with alternative independent coefficients $\vec{Y}_i \vec{\alpha}$.
- Z : Alternative specific variables with alternative specific coefficients $\vec{Z}_i \vec{\gamma}_k$.

In the context of destination choice, individual variables are those such as income, gender and education level that pertain to the traveler. A coefficient must be estimated for every destination, as the value of the individual vary does not vary across the choice set for each trip, only between trips. Alternative specific variables can have coefficients

independent or dependent of the choice set. The coefficient only needs to be dependent on the alternative Z when the parameter has a different meaning across the choice set. This is more commonly used in mode choice modelling; for example, the calculation of cost might vary between car and train journeys.

Parts of the equation can also be excluded by specifying 0 or -1 in the respective section. Intersects can also be removed by adding $+0$ to the respective section. The estimation will return an error if the equation cannot be solved, and the most common reason for this is high multi-colinearity between parameters specified in the model.

1.6 Location based social networks

The ubiquity of mobile GPS transceivers, especially in the smart phone market, has enabled a new category of social networks, called location based social networks (LBSN), which associate social networking data with a geo-referenced location. Different social networks have taken advantage of this opportunity in different ways. Facebook enables a user to mark themselves as safe during a natural disaster, flicker can show a map of where your images were taken, and google maps can provide accurate travel times by identifying areas of congestion.

Most location based social networks, such as facebook, tripadvisor and foursquare enable users to 'checkin' to a 'venue', such as a shop, tourist attraction or airport, and provide tips, ratings and reviews. When these services are used by millions of people around the world, in different countries and cities, a enourmous amount of data is collected, which can be used in a multitude of ways to explore mobility patterns.

Lindqvist et al. (2011) looked at how and why people use location sharing services such as foursquare, and discussed how users manage their privacy when using such services. (Cheng et al. 2011) collected 22 million checkins across 220,000 users to quantitatively assess human mobility patterns. 53% of their checkins came from foursquare, highlighting the dominance of foursquare in the LSBN space. They sampled location sharing status updates from the public Twitter feed, identifying users with geo-referenced tweets, and then collected the most recent 2000 geo-labeled tweets.

Noulas et al. (2012) used foursquare data to design gravity model based on Stouffer's theory of intervening opportunities. They found that while no universal law exists between mobility and distance, a universal behavior in all cities when measured with their rank-distance variable exists. Regarding the potential applications of LBSNs in future research, they note that the scale of data collected by foursquare provides the means to analyse and compare mobility patterns in different parts of the world, surpassing

cultural, geographical and political borders. They also warn “there may be a strong demographic bias in the community of Foursquare users”, before noting that “it is encouraging that the analysis and models developed in the context of the present work demonstrate strong similarities across multiple urban centers and different countries.”

Abdulazim et al. (2015) introduces a framework for inferring activity travel given nearby land use information gathered from LSBNs. Their results suggest that daily activity travel can be automatically inferred from LSBN data, and they present a generic method for acquiring land use data from LSBN services such as foursquare. The authors also present a case study for the greater Toronto and Hamilton area, Ontario, Canada, a subset of the study area for this thesis. SA et al. (2015) investigated the potential for cell phone and foursquare data to replace the use of Travel Surveys in calculating and Origin Destination Demand matrix. They found that The cell phone and Foursquare data were consistent with OD pairs expected to have higher trip volumes, but that some differences existed. Jin et al. (2014) proposed a doubly constrained gravity model based on LBSNs. They were able to achieve significant reductions in O-D estimation errors caused by sampling bias when compared to a singly constrained model.

1.7 Objectives

Disaggregate models provide clear advantages over aggregate methods in modelling trip distribution. While aggregate methods still hold sway in the modelling of long distance travel, due to the availability of data and more powerful computers, the modelling of destination choice using logit models is becoming more popular. Destination choice models provide more flexibility in attribute selection, and more efficient use of data. Most models include the basic socio-economic variables and a description of zone attractiveness.

It is often difficult or simply not appropriate to take a model that has been designed for another geographical region and apply it in study area of concern. Firstly, the data available for the study area will most likely be different to those available for other study areas. The data may provide more variables that weren't available to modellers working on other regions, or be more restricted, forcing the modeler to be creative in designing parameters that can represent the travel behavior in the region. Secondly, it is very difficult to design accurate models that work effectively when transposed to new study areas. This is due not just to obvious geographical differences, but variations in policy and culture that are difficult to reflect in a destination choice model. If every possible parameter reflecting was added to the model, not only would it be computationally infeasible, but there would be a high risk of overfitting in the model. The fact that the destination choice models already presented in the literature are individually unique

supports this notion that modelling is both a science and an art, and that there is no “one size fits all” model.

For these reasons, the design of a destination choice model for Ontario in itself reflects a new contribution to the field. The field of transport modelling is advanced every time the process of a designing a new model is performed. Future researchers can then look at the body of previous models, and use statistical analysis, their experience and intuition to select variables that best suit the requirements and use cases for which their model will be designed.

The second contribution of this thesis is the application of LSBNs to improve the utility modelling of destination choice. While work has been done on investigating mobility patterns, and the generation of OD matrices using LSBNs, their application to disaggregate models in transport has not been considered. This thesis will explore how foursquare check-in data can be used in the calculation of destination utility, and whether it can effectively replace other socio-economic variables such as employment. This is a significant contribution, as socio-economic data is not always available, especially at the high modelling resolutions made possible by the advances in microsimulation and computing technologies. check-ins data also provides an opportunity to model important traits of destination utility, such as the presence of national parks, that are not reflected in standard socio-economic variables.

Chapter 2

Data Acquisition and Analysis

2.1 Travel Survey of Residents of Canada

2.1.1 Introduction

The Transport Survey of Residents of Canada (TSRC) is a monthly, cross-sectional survey collected by Statistic Canada to measure the volume, characteristics and economic impact of domestic travel. The survey provides a large quarterly sample of performed trips within Canada, along with socio-economic data and the activities and expenditures performed on each trip. Results are released yearly, with the data available at a monthly temporal resolution.

The TSRC was designed to measure the size and economic impacts of Canada's domestic tourism industry. It was first performed in 2005, and replaces the Canadian Travel Survey. In 2011, the survey was redesigned to bring the questionnaire more in line with the World Tourism Organisation guidelines, and align the recorded activities with the International Travel Survey (ITS).

The TSRC acts as the main data source for the estimation and calibration of the destination choice model presented by this Thesis. Hence, this section provides an overview of the aspects of the TSRC and its design that relate to the development of a destination choice model. In particular, the methodology behind the survey is discussed, and the relevant variables and weightings available in the resultant microdata are highlighted.

2.1.2 Method

The survey is performed as a voluntary supplement to the compulsory Labour Force Survey (LFS). The LFS is a mandatory household survey of around 54,000 households to measure employment, and has a 90% response rate. The LFS sample consists of the entire civilian, non-institutionalized population over 15 years of age. A sub-sample of these households is selected to answer the TSRC, excluding residents of the Yukon, the Northwest Territories and Nunavut and people living on Native Reserves. A respondent is randomly selected from the household and asked to complete the travel survey. The survey is a computer-assisted telephone interview (CATI) available in both of Canada's official languages, English and French. 15 minutes are allocated for each respondent, with as many trips being collected as possible in that time.

2.1.3 Data

2.1.3.1 Spatial Resolution

All spatial data points, namely those for home location, trip origins and destinations and stopovers are provided in the microdata at three resolutions, Province or Territory, Census Division, and Census Metropolitan Agglomeration (CMA). Canada is made up of ten provinces and three territories, the largest of which is Ontario, the focus of this thesis.

Census Divisions are the next largest geographical area in Canada. Census Divisions represent groups of neighboring municipalities combined to aid regional planning and the provision of common services. After the provinces and territories, they are the most stable spatial unit. They were last modified for the 2011 census, and therefore are consistent between each TSRC dataset since the revised version was introduced. In most provinces and territories, these census divisions are defined in legislation, however in Newfoundland and Labrador, Manitoba, Saskatchewan, Alberta, Yukon, Northwest Territories and Nunavut, provincial or territorial law does not provide for these administrative geographic areas. In these cases, the census divisions are allocated by Statistics Canada.

Census sub-divisions are the next smallest geographical area, representing individual municipalities. These are recorded as part of the survey, however are not available in the TSRC microdata. The finest level of aggregation available is that of the Census Metropolitan Areas (CMA) and Census Agglomerations (CA). CMAs and CAs represent certain clustered areas of population around an urban core. More specifically, to be defined as a CMA, an area must have a total population of at least 100,000, with half

of those living in the core urban area. CAs, which related to CMAs but require a core population of only 10,000, are not recorded in the TSRC data.

Since CMAs do not topographically cover the whole Canadian study area, but only identify particular dense urban areas, census divisions are the most detailed resolution available for consistent use when working with the TSRC data. CMAs are only recorded for 51.5% of trip origins, and 48.3% of trip destinations. However, CMA's can be used to better allocate trips to transport zones in urban areas, as discussed in section (CITE).

2.1.3.2 Error Detection and Imputation

The computer-assisted nature of the survey allows for real-time error detection and consistency checking during the interview process. One example is that the program will inform the interviewer if the number of nights recorded for a trip does not match the number of nights recorded in various types of accommodations. Dont Know and Refused are also valid options for many questions, to prevent false answers been recorded. Sanity checks against extreme values are also performed, and the coding of geographical areas is mostly performed automatically.

Two forms of imputation are performed for the survey, for trip details and expenditure amounts respectively. Since the survey only allows 15 minutes for the recording of trip details, the details of non-selected trips are imputed from other trips recorded for that resident. This imputation process is multi-staged, and is performed per respondent. A donor pool of trips is selected that are similar to the non-selected trip. A distance function is then used to select the closest donor-trip to the recipient, and the detailed variables (activities, expenditures, etc) are copied over to the recipient trip.

2.1.3.3 Weighting

The weighting of records is particularly important when working with survey data. They allow the researcher to scale up the results for a sample to build an accurate representation of population, taking into account under- and over- represented groups within the survey. In total, four weightings are provided for the TSRC, with two relating to trip avriables: full-sample person weights first-month person weights person-trip weights trip weights

As the TSRC sample is based on the LFS survey, person weights are applied from the LFS and recalibrated to reflect subsampling, non-response, and known control groups.

After the 2011 redesign, respondents are asked about same-day trips that ended in the previous month, but overnight trips that ended in the previous two months. This means that effectively only half the sample is asked about same-day trips. To account for this, two weights are provided for each person record. A first month weight, that can be used for any person variable, and a second "full sample" weight that can be applied to person characteristics and overnight travel variables.

The person-trip weight, used to estimate trip volume, is then calculated by accommodating for identical trips, declared and reported trips, missing data and non-response. These weights are treated for outliers and recall bias. In calculating the person-trip weight, the person weight is also multiplied by the number of identical trips that this trip represents. The person-trip weight (WTEP) can be used against all socio-economic characteristics, as well all trip and visit variables, excluding expenditures. Trip weight (WTTP) is then calculated by dividing the number of household members that went on the trip, is only used to calculate expenditures, and as such is not relevant to the model design.

2.1.4 Microfile Format

The results from the TSRC are provided as yearly collections, separated into individual files for persons, trips and visits. The survey results are provided as fixed width delimited .dat files. A code book and data dictionaries are provided to decode the values stored in each line. The schema for encoded variables such as province are consistent across files and years (i.e. Ontario is always coded as 35), meaning that once read from the correct position on a line, values don't need to be decoded before being compared with each other.

Each person record is associated with one or more trips. Not all persons recorded in the person microdata necessarily have a trip recorded for a particular time period, as the survey records the travel behavior of both travelers and non-travelers.

Each recorded trip record has at least two associated visit records, or more if intermediate overnight stops were recorded. Visits are classified into two types, origins or destination/airport. Each Trip has one origin visit record, and at least one destination record. Where the main mode of travel for the trip is "Air", two or more airports are specified as visit records, along with the 3-digit airport code for the respective Canadian airport. The survey codebook notes that these airport records may be adjusted to protect respondent privacy.

2.1.4.1 Trip Datafile

Trips covered by the TSRC include same-day trips of more than 40km and overnight trips with at least one night in Canada. Domestic same-day and overnight trips are recorded in full. International trips with no nights in Canada are not recorded in the TSRC. For trips with an overnight destination, but some nights in Canada, only the domestic portion of the trip is recorded, with the point of departure from Canada recorded in the MDxxx variables for trip destination. The TR_D11 variable records the number of times this trip was performed in the reference month, and must be taken into account when estimating trip frequencies.

Socio-economic variables for the traveller are recorded for each trip record; namely age, gender, education level, employment status and income. The number of household members who participated on the trip is also recorded.

Trip purpose is recorded at two categorical levels. In the first, purposes are split into four options:

- Holidays, leisure or recreation (TODO frequency tables 2011-2014)
- Visit friends or relatives
- Business - All business and work related trips, except routine travel which is a regular part of the job
- Other - All trips for other reasons except regular household chores

A second variable is more specific, with seven categories:

- Holidays, leisure or recreation (TODO frequency tables 2011-2014)
- Visit friends or relatives
- Shopping, non-routine
- Personal conference, convention or trade show
- Other personal reasons
- Business conference, convention or trade
- Other Business

The point to point trip distance from origin to destination, excluding intermediate stops is also recorded. The distance variable is discussed further in section (CITE)

The main mode of travel is also provided, In the following categories:

- Car or Truck
- Commercial Aircraft
- Camper or RV
- Bus
- Train
- Ship, ferry
- Boat
- Not stated
- Other

2.1.4.2 Visit Datafile

The visit data file provides a stops performed on each trip, which can be linked to the relevant trip by the Public Use Microdata File Number (PUMFID) and the Trip Identification Number (TRIPID). Each trip has at least two visits associated with it, an origin and a destination visit, differentiated by the VISRECFL variable. The AIRFLAG variable is used to identify visit records that refer to an airport entry or exit.

If a location is visited twice during a single trip, only one visit is recorded for that location. The visits are not guaranteed to be recorded in the chronological order of visitation, even though the visits are collected in chronological order during the survey process. This lack of order prevents the visit records from being used model trip chaining, as discussed in section (CITE).

2.2 Filtering of Trip Records

For the model input, the TSRC trip records from 2011 to 2014 were collated together, giving 220,512 trip records. Not all these trips were relevant to the estimation of the destination choice model. Firstly, records were removed where:

- Either an origin or destination is not stated
- The trip purpose is not leisure, visit or business
- A distance is not recorded
- The mode is recorded as air and the destination and origin airports are identical

The TSRC trip files provide trip records not just for Ontario, but for all of Canada. However, as a model for Ontario, we are only concerned with the following categories of trips that influence travel in Ontario:

- Internal trips within Ontario - Internal (II)
- Trips entering Ontario - Incoming (EI)
- Trips leaving Ontario - Outgoing (IE)
- Trips that cross Ontario - External (EE)

Any trips that didn't fit one of these categories is also excluded from the trip dataset used to estimate the destination choice model. Internal, incoming and outgoing trips don't need to be filtered, and all trips in these categories are retained in the trip set. External trips, on the other hand, are filtered to remove trips that don't cross Ontario. Excluding such external trips is important to make sure that the estimated model reflects the behaviour of travel in Ontario, which could be different to the behaviours in other provinces.

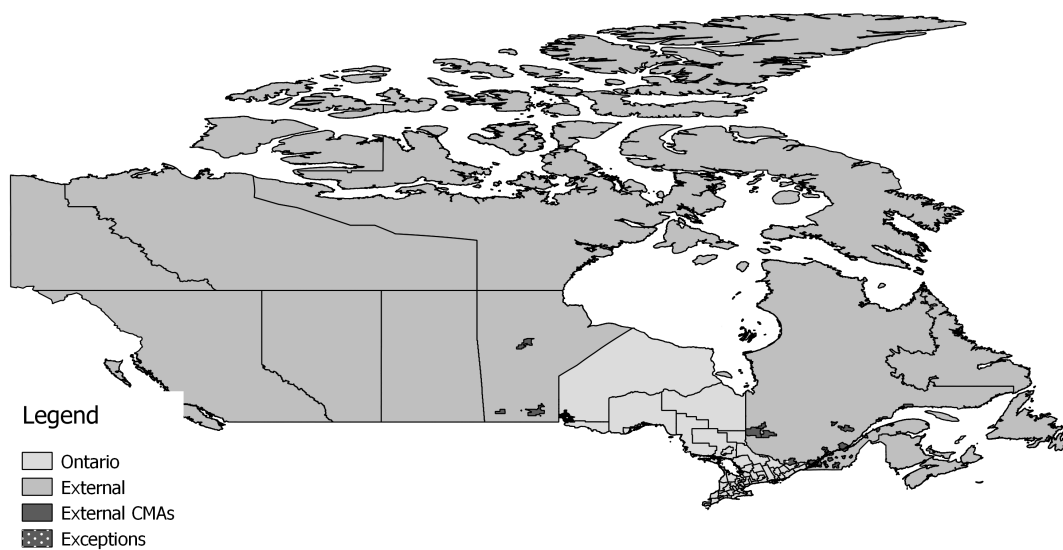


FIGURE 2.1: Dividing external zones into east and west

The unique geography of the Canadian provinces greatly restricts the number of external origin-destination pairs that need to be considered when excluding unwanted external trips. Ontario acts as land bridge between the eastern and western parts of Canada, see figure 2.1, dividing the external zones into two groups, east and west. Trips originating in one group and arriving in another have to pass through Ontario. And the converse is true for trips within a group. Hence all trips that don't go between east and west can be removed. There are two zones which are the exception to this, zones 85 and 117 in western Quebec. Journeys between these zones and other zones in Quebec may pass through Ontario. For example, figure 2.2 illustrates a journey from Gatiéau, Quebec to Montreal Airport takes around 2 hours when passing through Ontario, and 2 hours and 30 minutes otherwise. Trips between these two exception zones and all other zones in Quebec are therefore retained.

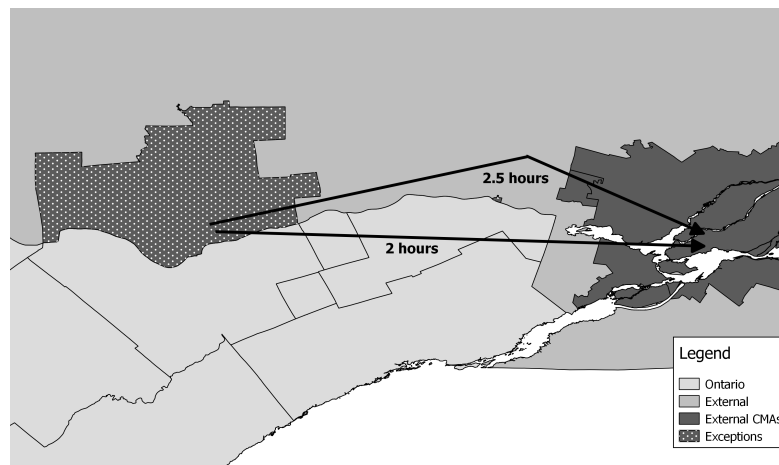


FIGURE 2.2: An example of an external origin-destination pair that passes through Ontario

Figure 2.3 illustrates how our the trip distribution of the estimated trips fits the observed distribution much better after undesired external trips are removed. In total 69,328 individual trip records remain from the TSRC dataset for model estimation, representing 40,177,841 weighted trips.

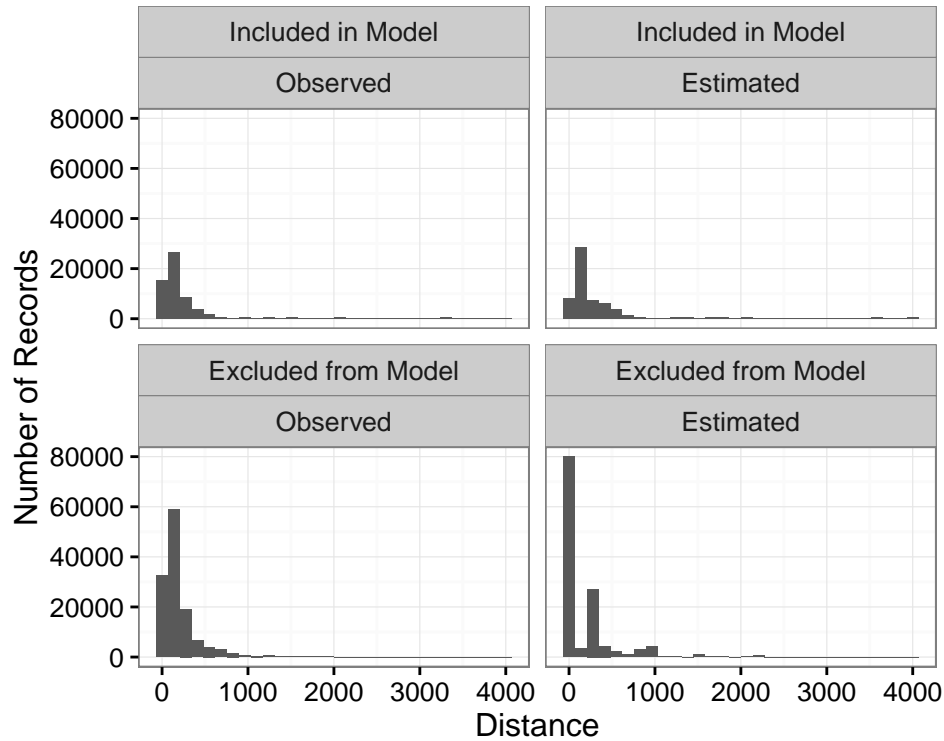


FIGURE 2.3: Estimated vs. Observed Trip Distance

2.3 Zone System

This chapter discusses the definition of the choice set of alternatives for the Destination Choice model. Numerous factors need to be considered when designing the choice set. Firstly, The sample size of the data available to estimate the model coefficients is an important restraint. With a small sample set relative to the size of the destination choicet, not enough records are available to calculate the parameter coefficients with high confidence. Additionally, if many destinations have a low number of trip records, or even none at all, then the impact of the attractiveness of such destinations on the model is not considered. When working with a large destination choice set, a second element must be considered. When deciding on a destination, a traveller cannot reasonably compare the utility of thousands of destinations. Some models (CITE) have used nested destination choice models, where a region is first selected, from which a destination is selected. Large destinations choice sets also lead to very long computation times when estimating the model coefficients. A balance needs to be found between the detail represented in the choice set, and the computability and validity of the coefficients.

For this particular destination choice model, a zone system was already provided, consisting of 6671 Traffic Analysis Zones (TAZ). The TAZs can be grouped into 4 categories;

6495 internal zones for Ontario, 48 and 121 external zones representing the rest of Canada and North America respectively, and 7 zones for remaining world-wide destinations. As this thesis is only concerned with domestic travel within Canada, only the Internal zones for Ontario, and 48 external zones within Canada are considered. The external zones are not modified as TSRC origins and destinations are directly translatable to the external zones.

For the Internal zones (TAZs) within Ontario are allocated using a gradual raster based zone approach, based on the method presented by Moeckel and Donnelly (2015). The 6495 generated TAZs vary in size from $0.879km^2$ to $3600km^2$, with smaller cells defined for more populous areas, and larger cells for regional areas. The gradual zone system is designed on the premise that it is desirable to have larger zones in rural areas where there is less population, and hence, less activity. This method reduces the number of TAZs, and hence, the complexity of the model, while only removing detail where it is least required.

However, the TSRC trip origins and destinations don't match this custom zone system, being only available at a much broader resolution. Origins and Destinations in the TSRC microfiles can be identified using three variables of varying spatial resolution; Province, Census Division (CD) and Census Metropolitan Agglomeration (CMA). As discussed in 2.1.3, Provinces and Census Divisions cover the national study area completely, while CMAs are only defined for areas of urban population. Hence, either internal zones need to be aggregated to the level available in the TSRC data, or the trip origins and destinations need to be allocated to TAZs using a weighted random selection. For this thesis, the first approach is chosen. The allocation of destinations will be performed at a later stage in the transport model, and the proposed method is discussed further in section resection:implementation.

Defining a zone system for Ontario based on the TSRC data As a first step, the zones are defined by the census divisions comprising Ontario, of which there are 49, as this provides the highest resolution available that fully covers the study area. However, even in rural areas, the TAZs are much smaller than the size of a Census Division. When the zone system is defined purely using the Census Divisions within Ontario, over 50% of Census Divisions have more than 75 TAZs, with a large spread (see figure 2.4.

Although CMAs are defined only for selected urban areas of Canada, they can be considered alongside the CDs when allocating zones to improve the spatial resolution of the zone system. The concept of a CMA aligns closely with the objective behind the gradual raster cell size of the provided zone system for Ontario. CMAs identify areas of denser population around an urban core that may be of particular significance to geographers and modellers. By simply including CMAs as zones in the aggregated zoning model, the

number of zones is increased to 57, allowing trip origins and destinations in urban areas to be more accurately assigned during disaggregation to TAZs. As seen in figure 2.4, this significantly lowers the mean number of TAZs per zone, and also reduces the spread of values, indicating a much improved zone system.

This approach has one drawback, as a large outlier, representing the Toronto CMA is now observable, consisting of over 2000 TAZs. This outlier corresponds to the CMA of Toronto, the most populous in Ontario, both a very large generator and attractor of trips. In 2014, Toronto represented 13.4% and 10% of trip origins and destinations respectively. It should be clear that this is a very undesirable occurrence. While it is hard to say whether this outlier would significantly affect the destination choice model, it would regardless provide significant challenges when allocating trips to individual TAZs, affecting the overall quality of the model. It is worth noting briefly that the choice of zone system can have cascading effects throughout the whole transport model, which need to be considered during its creation.

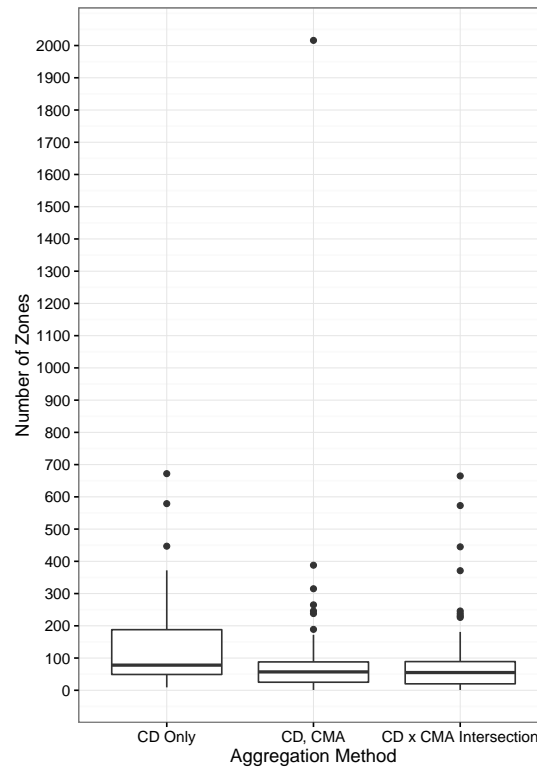


FIGURE 2.4: Different methods of aggregating internal zones to match the TSRC spatial resolution.

This lone outlier was not present when only the CDs were considered as destination alternatives. Since CMAs often overlap multiple CDs, rather than simply including CMAs and CDs independently, we can overlap the CDs and CMAs to fully reflect the number of destination choices available in the TSRC data. This is done as followed; the

area of each CD that does not intersect with a CMA is selected as a zone. Then, each unique combination of CMA and CD is recorded as a new zone. An example of this process is shown in figure (CITE). This process gives us 69 zones for Ontario, at 41% increase over the simplest approach that only considers CDs.

Figure 2.4 illustrates the difference between these methods. When only Census Divisions are used, a significant number of CDs have a large number of assigned TAZs. When the CMA's are considered, the results are clearly better. A lower average number of TAZs per aggregate zone will give better results when trips origins and destinations are disaggregated. Taking the intersection of CDs and CMAs has little effect on most zones, but still improves the overall result in a very significant way. The CMA of Toronto overlaps with 7 separate CDs, and can with this method be divided into seven smaller zones.

This third method has another advantage, in that the a distinction between urban and rural areas is now encoded into the zone system. This will be important in the estimation process as 51.5% of trips in the filtered TSRC survey originated in a CMA, and 48.3% had destination recorded as a CMA. Not only is it clear that urban areas are important drivers of long distance travel, but that, interestingly, CMAs are more likely to be origins than destinations.

2.4 Aggregating of Zonal Data

All the data on distances, population and employment was provided at the TAZ level. This section describes how they were allocated to the zone system. The TAZs themselves were assigned to the zone which intersected their centroid. Where the centroid of the TAZ didn't intersect any zones, the first intersecting zone was used. When the TAZ did not intersect with any part of the Canadian census boundaries at all, it was assigned manually to the nearest zone.

Socio-economic variables, namely population and employment were aggregated from the TAZ level to the zone system using a summation.

A distance matrix was provided for the original nset of TAZs. It was calculated without congestion using the Canadian road network and intra-zonal travel times were not included. This skim matrix needed to be aggregated to the zone system used in the destination choice model. This was done by taking a distance d between each pair of sub zones weighted by the multiplied populations p of the origin and destination.

$$d_{ij} = \frac{\sum_{k \in i, l \in j} d_{kl} \cdot p_k \cdot p_l}{\sum_{k \in i, l \in j} p_k \cdot p_l}$$

2.5 Foursquare

Chapter 3

Model Design

Chapter 4

Gravity Model

4.1 Design

4.2 Calibration

4.3 Results

Chapter 5

Destination Choice Model

5.1 Estimation

5.1.1 Socio-economic Variables

5.1.2 Origin-Destination Interactions

5.1.3 Incorporating LBSN Data

5.1.4 Model Subsets

Seasons

Income Strata

5.1.5 Final Model

Chapter 6

Implementation

6.1 MTO Model Integration

6.2 Algorithm

6.3 Disaggregation to TAZs

Chapter 7

Validation

7.1 Case study of a new Ski Resort

Chapter 8

Epilogue

8.1 Conclusions

8.2 Discussion

8.3 Future Work

Appendix A

An Appendix

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Vivamus at pulvinar nisi. Phasellus hendrerit, diam placerat interdum iaculis, mauris justo cursus risus, in viverra purus eros at ligula. Ut metus justo, consequat a tristique posuere, laoreet nec nibh. Etiam et scelerisque mauris. Phasellus vel massa magna. Ut non neque id tortor pharetra bibendum vitae sit amet nisi. Duis nec quam quam, sed euismod justo. Pellentesque eu tellus vitae ante tempus malesuada. Nunc accumsan, quam in congue consequat, lectus lectus dapibus erat, id aliquet urna neque at massa. Nulla facilisi. Morbi ullamcorper eleifend posuere. Donec libero leo, faucibus nec bibendum at, mattis et urna. Proin consectetur, nunc ut imperdiet lobortis, magna neque tincidunt lectus, id iaculis nisi justo id nibh. Pellentesque vel sem in erat vulputate faucibus molestie ut lorem.

Quisque tristique urna in lorem laoreet at laoreet quam congue. Donec dolor turpis, blandit non imperdiet aliquet, blandit et felis. In lorem nisi, pretium sit amet vestibulum sed, tempus et sem. Proin non ante turpis. Nulla imperdiet fringilla convallis. Vivamus vel bibendum nisl. Pellentesque justo lectus, molestie vel luctus sed, lobortis in libero. Nulla facilisi. Aliquam erat volutpat. Suspendisse vitae nunc nunc. Sed aliquet est suscipit sapien rhoncus non adipiscing nibh consequat. Aliquam metus urna, faucibus eu vulputate non, luctus eu justo.

Donec urna leo, vulputate vitae porta eu, vehicula blandit libero. Phasellus eget massa et leo condimentum mollis. Nullam molestie, justo at pellentesque vulputate, sapien velit ornare diam, nec gravida lacus augue non diam. Integer mattis lacus id libero ultrices sit amet mollis neque molestie. Integer ut leo eget mi volutpat congue. Vivamus sodales, turpis id venenatis placerat, tellus purus adipiscing magna, eu aliquam nibh dolor id nibh. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Sed cursus convallis quam nec vehicula. Sed vulputate neque eget odio fringilla ac sodales urna feugiat.

Phasellus nisi quam, volutpat non ullamcorper eget, congue fringilla leo. Cras et erat et nibh placerat commodo id ornare est. Nulla facilisi. Aenean pulvinar scelerisque eros eget interdum. Nunc pulvinar magna ut felis varius in hendrerit dolor accumsan. Nunc pellentesque magna quis magna bibendum non laoreet erat tincidunt. Nulla facilisi.

Duis eget massa sem, gravida interdum ipsum. Nulla nunc nisl, hendrerit sit amet commodo vel, varius id tellus. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Nunc ac dolor est. Suspendisse ultrices tincidunt metus eget accumsan. Nullam facilisis, justo vitae convallis sollicitudin, eros augue malesuada metus, nec sagittis diam nibh ut sapien. Duis blandit lectus vitae lorem aliquam nec euismod nisi volutpat. Vestibulum ornare dictum tortor, at faucibus justo tempor non. Nulla facilisi. Cras non massa nunc, eget euismod purus. Nunc metus ipsum, euismod a consectetur vel, hendrerit nec nunc.

Bibliography

- Abdulazim, Tamer et al. (2015). “Framework for Automating Travel Activity Inference Using Land Use Data: The Case of Foursquare in the Greater Toronto and Hamilton Area, Ontario, Canada”. In: *Transportation Research Record: Journal of the Transportation Research Board* 2526, pp. 136–142.
- Adler, Thomas and Moshe Ben-Akiva (1979). “A theoretical and empirical model of trip chaining behavior”. In: *Transportation Research Part B: Methodological* 13.3, pp. 243–257. ISSN: 0191-2615.
- Anas, Alex (1983). “Discrete choice theory, information theory and the multinomial logit and gravity models”. In: *Transportation Research Part B: Methodological* 17.1, pp. 13–23. ISSN: 0191-2615.
- Ben-Akiva, Moshe E (1974). *Structure of passenger travel demand models*. ISBN: 0309023734.
- Canadian Transport Commission (1971). “Intercity passenger transport study”. In: Casey, HJ (1955). “Applications to traffic engineering of the law of retail gravitation”. In: *Traffic Quarterly* 9.1, pp. 23–35.
- Cheng, Zhiyuan et al. (2011). “Exploring Millions of Footprints in Location Sharing Services.” In: *ICWSM* 2011, pp. 81–88.
- Daly, Andrew (1982). “Estimating choice models containing attraction variables”. In: *Transportation Research Part B: Methodological* 16.1, pp. 5–15. ISSN: 0191-2615.
- Jin, Peter et al. (2014). “Location-Based Social Networking Data: Exploration into Use of Doubly Constrained Gravity Model for Origin-Destination Estimation”. In: *Transportation Research Record: Journal of the Transportation Research Board* 2430, pp. 72–82.
- Kitamura, Ryuichi (1984). “Incorporating trip chaining into analysis of destination choice”. In: *Transportation Research Part B: Methodological* 18.1, pp. 67–81. ISSN: 0191-2615.
- Lindqvist, Janne et al. (2011). “I’m the mayor of my house: examining why people use foursquare-a social-driven location sharing application”. In: *Proceedings of the SIGCHI conference on human factors in computing systems*. ACM, pp. 2409–2418.

- McFadden, Daniel (1973). "Conditional logit analysis of qualitative choice behavior". In:
- Miller, Eric (2004). "The trouble with intercity travel demand models". In: *Transportation Research Record: Journal of the Transportation Research Board* 1895, pp. 94–101. ISSN: 0361-1981.
- Mishra, Sabyasachee et al. (2013). "Comparison between gravity and destination choice models for trip distribution in Maryland". In: *Transportation Research Board 92nd Annual Meeting*.
- Moeckel, Rolf and Rick Donnelly (2015). "Gradual rasterization: redefining spatial resolution in transport modelling". In: *Environment and Planning B: Planning and Design* 42.5, pp. 888–903.
- Moeckel, Rolf, Rhett Fussell, and Rick Donnelly (2015). "Mode choice modeling for long-distance travel". In: *Transportation Letters* 7.1, pp. 35–46. ISSN: 1942-7867.
- Noulas, Anastasios et al. (2012). "A tale of many cities: universal patterns in human urban mobility". In: *PloS one* 7.5, e37027.
- Oi Walter; Schuldiner, Paul (1962). *An Analysis of Urban Transportation Demand*. Evanston, IL: Northwestern University Press.
- Outwater, Maren L et al. (2015). "Tour-Based National Model System to Forecast Long-Distance Passenger Travel in the United States". In: *Transportation Research Board 94th Annual Meeting*.
- Outwater, Maren et al. (2010). "California statewide model for high-speed rail". In: *Journal of Choice Modelling* 3.1, pp. 58–83. ISSN: 1755-5345.
- SA, Rokib et al. (2015). "Origin-Destination Trip Estimation from Anonymous Cell Phone and Foursquare Data". In: *Transportation Research Board 94th Annual Meeting*. 15-2379.
- Schiffer, RG (2012). "NCHRP Report 735: Long-Distance and Rural Travel Transferable Parameters for Statewide Travel Forecasting Models". In: *Transportation Research Board of the National Academies, Washington, DC*.
- Spear, Bruce D (1977). *Applications of new travel demand forecasting techniques to transportation planning. A study of individual choice models*. Report. URL: <http://ntl.bts.gov/DOCS/SICM.html>.
- Stephanedes, Yorgos J, Vijaya Kumar, and Balakrishnan Padmanabhan (1984). "A fully disaggregate mode-choice model for business intercity travel". In: *Transportation Planning and Technology* 9.1, pp. 13–23. ISSN: 0308-1060.
- Train, Kenneth E (2009). *Discrete choice methods with simulation*. Cambridge university press. ISBN: 1139480375.
- Warner, Stanley Leon (1962). "Stochastic choice of mode in urban travel: A study in binary choice". In: *METROPOLITAN TRANSPORTATION SERIES*.

Wilson, Alan G (1967). “A statistical theory of spatial distribution models”. In: *Transportation research* 1.3, pp. 253–269. ISSN: 0041-1647.