# STOCHASTIC ROUTE CHOICE SET GENERATION: BEHAVIORAL AND PROBABILISTIC FOUNDATIONS

PIET H.L. BOVY[1] AND STELLA FIORENZO-CATALANO[2]

This article addresses the generation of choice sets for the purposes of route choice analysis and flow prediction in uni-modal and multi-modal networks. Ample attention is devoted to the implications for choice set generation of these various purposes. Based on a theory on choice behavior, a new model-based choice set generation approach will be elaborated, called the doubly stochastic choice set generation model, meant for establishing choice sets prior to the choice modeling step. Because of its stochastic principle, a typical property of the proposed generation approach is that the size and composition of the generated choice sets are stochastic variables. We will devote ample attention to these stochastic properties using theoretical derivations and experimental studies. The article reports on the calibration of this generation approach for multi-modal networks and illustrates the approach with a number of predictions in various networks.

KEYWORDS: Route choice sets, choice set generation, networks, multimodal

## 1. INTRODUCTION

In analyzing and predicting travel demand with discrete choice models, choice sets of the travel alternatives are needed. Choice sets are defined as the collection of travel options perceived available to an individual in satisfying his travel demand. From a variety of studies it is well-known that the size and composition of choice sets do matter in case of estimation and prediction (for the case of route choice see for example VanderWaerden et al., 2004). Correctness of choice parameter estimates and correctness of demand predictions depend on the quality of the adopted choice sets. From this it follows also that a clear distinction in purposes need to be made between establishing choice sets for estimation, that is estimating parameters of choice models, and for prediction, that is predicting flows in networks.

This article addresses the generation of choice sets for the purposes of route choice analysis and prediction in uni-modal and multi-modal networks. For an overview of this topic see among others Ramming (2002) and Hoogendoorn-Lanser (2005).

In this context, specific attention needs to be devoted to the various relevant choice set concepts (universal sets, consideration sets, etc.) and how these are related to the individual traveler's choice behavior on the one hand and the outside view of the modeler on the other hand.

Various approaches have been proposed in the literature for *explicit* generation of route choice sets prior to prediction. There are good reasons (see e.g. Bekhor et al., 2001) for a priori generation in the demand prediction instead of the widely applied so-called column generation during iterative network equilibrium modeling. These reasons ask for an in-depth assessment of choice set generation procedures. Based on a theory on choice behavior, we will elaborate below on the development of a new model-based choice set generation approach, called the doubly stochastic choice set generation model, meant for

---

[1] Transport & Planning Department, Faculty of Civil Engineering and Geosciences, Delft University of Technology, PO Box 5048, NL-2600 GA Delft, The Netherlands. Corresponding author (E-mail: P.H.L.Bovy@tudelft.nl).

[2] Transport & Planning Department, Faculty of Civil Engineering and Geosciences, Delft University of Technology, PO Box 5048, NL-2600 GA Delft, The Netherlands. Present address: TNO Bouw, P.O. Box 49, NL-2600 AA Delft, The Netherlands.

establishing choice sets prior to the choice modeling step. The basic hypothesis of this approach is that traveler's attribute preferences with respect to available alternatives highly determine the size and composition of the choice set considered in his choice decision. From this hypothesis we state that a traveler's trip utility function can be used as a basis for a generating function with which attractive alternatives can be extracted from the given network through optimal path search by stochastically varying trip attributes and attribute preferences.

Because of its stochastic principle, a typical property of the proposed generation approach is that the size and composition of the generated choice sets are stochastic variables. We will devote ample attention to these stochastic properties and will present some general probabilistic relationships derived for a few idealized cases.

The article will report on the calibration results achieved with observed choice sets from a survey of multi-modal train trips in the Dordrecht-Rotterdam Corridor in the Netherlands and will give a number of results on generated route choice sets in different types of networks such as the Dutch waterway network, the Dutch national road network, and a regional multi-modal public transport corridor (Rotterdam area).

## 2. NOTIONS OF CHOICE SETS

With the general notion of *choice set* we define the collection of travel options perceived available for making a trip by an individual or set of individuals. In this respect we need to distinguish between the traveler's perspective and the modeler's perspective. We define as the *subjective consideration set* the subset of alternatives available, feasible, and known to an individual traveler for a particular trip given his personal conditions. Such sets are usually very small and idiosyncratic given the particular conditions of individuals and specific trips, such as vehicle availability, time constraints etc. (For a more detailed elaboration on the notion of choice sets with particular reference to route choice see Bovy and Stern, 1990; Hoogendoorn-Lanser, 2005; Fiorenzo-Catalano, 2007).

Only in case of sufficient information about the traveler's specific conditions, the modeler can make probabilistic predictions of the consideration sets at hand, correspondingly called *objective* consideration sets, needed for choice modeling. This may hold for the case of estimation of choice functions based on individual observations of travel behavior. Usually, the objective consideration set will be larger that its subjective counterpart.

In the case of travel demand prediction with given choice models, usually no information is available on an individual level but rather on flows between areas. At most, some information is known at the level of the population of travelers, such as percentages of travelers with particular forms of vehicle availability or levels of income. Therefore, in the prediction case the modeler tries to establish a *collective* choice set for a particular interzonal travel demand that fits to that collection of travelers by trying to embrace the population of individual consideration sets into a single *collective* consideration set. Since preferences, information levels and other conditions usually are very different among individuals making a similar trip, the collective consideration set will be much larger and more varied than any individual choice set.

In the following we will elaborate on choice set generation for the establishment of route choice sets in networks for various purposes. In doing this we have to bear in mind a number of *specific characteristics of route choice sets*, such as among others, the following:

- the population of available routes (the universal set) in dense networks usually is very large and mostly not known; sometimes hundreds of routes can be identified for a single OD- pair;
- the set of feasible and attractive routes often is very large. In addition, this set is complex because of heterogeneous route composition and varying physical overlap among the routes;
- because of these aspects, it is difficult or even impossible for the analyst to enumerate the full set of routes attractive to the travelers and relevant in choice modeling;
- choice behavior among routes is ambiguous and may consist of various forms of choice process, such as sequential (from decision point to decision point), simultaneous (from origin to destination at once) or strategic (adaptive choice based on prevailing network conditions during the trip).

These issues require the adoption of special generation methods for establishing route choice sets. In addition, the sheer size of networks and number of OD-pairs ask for an efficient model-based computational procedure. Procedures developed for choices among modes or brands are not applicable to routes.

## 3. PURPOSES OF CHOICE SET GENERATION

Route choice sets may be generated mainly for the following three application purposes:
- Supply analysis of travel options in networks where the planner or researcher is interested to know the availability of travel alternatives, their number, their characteristics, their variety, their composition etc.;
- Demand model estimation (i.e. estimating behavioural parameters of utility functions of choice models at individual level);
- Prediction of choice probabilities in a demand analysis for determining route flow and link flow levels in networks using route choice models with known parameters derived from estimation, mostly at zonal level.

These different applications of choice sets appear to pose different requirements on size and composition of the choice sets to be used in estimation or prediction. In Bovy (2007), a more detailed elaboration can be found on these issues.

Whereas choice sets need not necessarily be exhaustive for estimation purposes, prediction choice sets must include at least all attractive routes. For estimation purposes, even if not all relevant alternatives are included in the choice set and small well-sampled choice sets are considered, it may nevertheless provide satisfactory estimation results (Ben-Akiva and Lerman, 1985). On the other hand for prediction purposes, the choice sets should include quite all realistic and reasonable alternatives in order to produce correct demand levels. In addition, the policy context of the travel demand analysis may require that particular alternatives regardless of their attractiveness need to be included in the prediction choice set. Whereas a lot of attention has been given to the estimation issue (see e.g. Ben-Akiva et al., 1985; Ben-Akiva and Bierlaire, 1999), the prediction case has been largely neglected, which is especially lacking in the route choice context.

In this article we will focus on the generation of choice sets for prediction of route and link flows, so we will devote our attention to the adequacy of choice sets for that purpose.

In the context of a prediction application, the analyst is interested to achieve satisfactory predictions of route and link flows, especially for those routes and links that have special policy relevance. Such a prediction involves calculating the choice probabilities of all non-zero OD-trips, maybe separately for user groups or trip purposes,

and then summing up the number of trips that will use each of the potentially feasible routes, and derived from this, through a network assignment, the use of links.

In most current network assignment studies, route flow predictions are not calculated explicitly but follow from application of simplified choice models and only implicitly generated shortest paths (so-called column generation). We are convinced however that following a two-step prediction approach distinguishing between an explicit choice set generation step followed by a choice modeling step using a probabilistic route choice model offers essential advantages.

This would require the specification of choice sets in which each route that may attract trips is included. In the case of predicting route flows the requirements on the quality of choice sets are very strict since in order to have correct route flows, predicted choice sets should include all relevant routes. Inclusion of some unattractive routes in the choice set if not too much overlapping with the attractive ones will not distort the demand predictions, nor will these have serious influence on computational efficiency.

Given this, prediction choice sets should be sufficiently large, maybe even including irrelevant unattractive routes, so that erroneous predictions of small route flows maybe compensate to a certain extent. Consequently, a prediction choice set should likely consist of all relevant routes with high probability of being chosen; inclusion of some unattractive routes is acceptable.

For a variety of reasons, *explicit generation* of route choice sets based on behavioral principles is advantageous both for estimation as well as prediction purposes. We stress that for prediction applications, explicit generation of route choice sets *prior* to route and link flow calculation instead of during the iterative flow calculations is strongly favoured for several reasons. It allows full control over desired properties of generated routes and of the size and composition of the choice set. A major advantage of a priori choice set generation for prediction purposes is that a priori given choice sets allow much more flexibility and realism in behavioural assumptions in the route choice models adopted. In that case, no restrictions exist on the type of choice model or utility function specification to be adopted. Applicability of advanced analytical approaches, easy consideration of route overlap, non-linear utility functions, and route-based attributes are only a few advantages of applying a priori generated choice sets. In addition, a priori enumeration in a network context not only offers a number of theoretical advantages, but also implementation and computational advantages in iterative network assignment approaches since repeated optimal route search no longer is necessary. This is especially true for dynamic and multi-user class network assignment types. Implicit choice set generation approaches (see e.g. Cascetta and Papola, 2001) typically do not specify the feasible alternatives and cannot offer these theoretical and practical benefits.

## 4. APPROACHES TO CHOICE SET GENERATION

In the past decades, several researchers have paid attention to the modeling and generation of route choice sets for estimation or prediction purposes. Well-known in this respect are for example the labeling approach of Ben-Akiva et al. (1984) and the implicit-availability approach of Cascetta and Papola (2001). In addition, from the operations research field a multitude of route generation algorithms were proposed to generate route sets useful for choice set composition. We stress that the establishment of useful route choice sets poses specific requirements to route generation algorithms. Some of the explicit route generation approaches have been compared and assessed by Ramming (2002) for the road network of Boston. Prato and Bekhor (2006) also present

an overview of deterministic and stochastic generation methods and establish a comparative performance analysis of a variety of generation approaches showing big differences in their performance. It appears from their analyses that simulation approaches, such as dealt with in this article, are very promising and show high levels of performance, even without a dedicated calibration of their parameters. All these methods are designed and applied to road networks. The works of Hoogendoorn-Lanser (2005) and Fiorenzo-Catalano (2007) specifically are dedicated to generation methods suitable for multi-modal transportation networks which are much more complex than uni-modal road networks. Whereas the first work develops a new deterministic branch and bound algorithm directed at choice sets for choice model estimation (see also Hoogendoorn-Lanser et al., 2006), the second one established a stochastic generation approach directed at explicit a priori generation of choice sets for demand prediction purposes. The latter method will be explained below.

## 5. DOUBLY STOCHASTIC CHOICE SET GENERATION APPROACH

We will now elaborate on the development of our model-based choice set generation approach, called the doubly stochastic choice set generation model, meant for establishing choice sets prior to the choice modeling step, be it in model estimation or demand prediction. The basic hypothesis of this approach is that traveler's attribute preferences with respect to available alternatives highly determine the size and composition of the choice set considered in his choice decision. From this hypothesis we state that trip utility functions known from literature can be used as the basis for a generating function with which attractive alternatives can be generated through optimal path search in the given network by stochastically varying network attributes and attribute preferences around their measured values. The adopted variances in network attribute values (travel time, waiting time, travel cost, etc.) reflect, among other matters, on the one hand the biased perception and cognition of the network by individual travelers and reflect on the other hand the variation in this perception of attributes among the travelers and their differences in knowledge about the network. The adopted variances in the related parameter values reflect the variation in the preferences for these attributes within the population of travelers acknowledging that there are time-sensitive, cost-sensitive, congestion-sensitive, etc types of travelers in the population having their own consideration sets.

The adoption of the stochastic optimal path approach does not mean that this is the prevalent choice mechanism, not at all! (the generation step is treated completely independent from the choice modeling step). This approach only is a means (an efficient one) to extract from the network attractive candidate routes for the choice set with sufficient variety in its composition so as to reflect the biased perceptions and cognition on the part of the individual traveler as well as the variation in perception, cognition and preferences among the group of travelers between the same OD-pair.

In practical application of this idea it means that on the one hand the link properties of the network at hand are randomized around their measured values $X$ according to certain rules, while on the other hand simultaneously the parameters of the generating function of travelers (the consideration utility function) are randomized around expected values $\beta$ according to certain rules. The expected parameter values and their corresponding random distributions of the generation model are subject to calibration based on route observations or other information. The choice set generation function then is given by the following general expression (1):

$$C_r^s = \sum_{a \in r} \sum_j \left( \beta_j^s + \varepsilon_j^s \right) \cdot \left( X_j^a + \varepsilon_j^a \right), \qquad (1)$$

where $C_r^s$ is the random cost of path r for user group $s$ (is shortest path search criterion), $\beta_j^s$ is the expected value of preference parameter for attribute $j$ for user class $s$, $\varepsilon_j^s$ is the random part of preference parameter $\beta_j^s$ drawn from distribution $P_j^s(0; \sigma_j^s)$, $X_j^a$ is the expected value of link attribute $j$ for link $a$, and $\varepsilon_j^a$ is the random part of link attribute $j$ of link $a$ drawn from distribution $P_j^a(0; \sigma_j^a)$.

Examples for attributes $X$ in road networks are distance or time, possibly differentiated by road type, traffic lights, lighting, etc. In public transport and multimodal networks many more attributes can be added related to waiting, transfers etc. We adopt specific generation functions for specific user classes $s$ if it is to be expected that their choice set formation will be different. This is for example relevant in modeling route choice in multi-modal networks where the feasibility of multi-modal travel options strongly depends on vehicle availability of the traveler. If car or bike are not available, then all travel options involving car use or bike use are not feasible.

The general expression (1) for the generation function meant for multi-modal networks may be simplified in special cases. For instance, in applications to the Dutch national main road network (see Bliemer and Taale, 2006, and Section 8) only the attribute 'time' was randomized while in an application to the Dutch waterway network only the attribute 'distance' was made stochastic. This can however easily be extended to for example the consideration of hierarchical choice patterns found in the population of travelers by differentiating the attributes by functional link type.

The choice set generation procedure then consists of repeated shortest path search under randomized conditions of the network and of the search criterion (parameters of generating function) governed by some termination criteria related to required size and composition of the choice set. The attributes and parameters used in the generating function are defined by probability distributions $P$ of which the expected values are related to estimates derived from exogenous choice model estimates or other exogenous information, while the distributional forms $P$ and variances $\sigma$ preferably follow from calibration of this choice set generation approach with observations.

The described doubly stochastic choice set generation algorithm may then be executed as follows (however, other sequences of calculations are also possible):
(1) repeat for given termination criteria (e.g. number of iterations);
(2) sample for each link attribute $X$ random values from a distribution specific for each attribute;
(3) for each distinguished traveler group $s$:
    a.  repeat for a given termination criterion (e.g. number of iterations);
    b.  sample for each parameter $\beta$ in the generation function random values from a distribution specific for each parameter;
    c.  compute for each link its random generalized cost $C$;
    d.  for each OD-pair, determine shortest path
    e.  insert path into path list if not yet found before, else ignore.

An important role is played by the termination criteria, since these determine the quality of the generated choice sets in terms of size, composition, variety, and the like. The termination criteria may relate to computational aspects (such as maximum number

of iterations for each separate randomization stage) or may relate to the generated choice sets (such as minimum or maximum size or to the changes of size or content of the choice sets in successive iteration steps).

A typical property of the proposed approach is that one and the same route maybe found several times as the best one, even with different cost values. These superfluous routes are ignored.

The shortest path search is the key operation in the procedure. Through this minimization principle the more attractive routes will most probably be generated by nature while disregarding the unattractive ones, while on the other hand it is a very efficient technique to enumerate routes. We repeat that this generation principle does not have any relationship to the actual choice mechanisms of travelers nor of that of modelers; it's only an efficient well-controlable means to extract candidate routes from the network.

The presented procedure offers ample opportunities to account for sufficient policy sensitiveness. First of all, changed network properties due to specific policies (e.g. new links, changed tolls, changed travel times) not only will influence the outcomes of the choice modeling but also will possibly lead to different (e.g. enlarged) choice sets, and thus influence the choice outcomes via this indirect way. Secondly, if policy analysis purposes require demand estimates for particular routes or links, these need to be present in the choice set. This can easily be accommodated by the proposed algorithm by making the nodes of such 'forced' links to be artificial origins and destinations.

## 5.1 Filtering process

Given the termination criteria, the described algorithm produces an exhaustive base route set (master set) that needs to be reduced in order to establish choice sets that efficiently suit their particular purpose (analysis, estimation, or prediction). To that end, a variety of selection constraints may be additionally applied. These constraints refer to requirements posed on individual routes (e.g. maximum detour, logical sequence of links by road type or mode, etc) and to requirements posed on comparative properties among selected routes (e.g. maximum overlap).

For estimation purposes, only a small choice set may suffice, however having a rich variety in terms of route attributes. Highly overlapping routes should be removed from the base set, while cases with dominant alternatives (better than others in every respect) should be removed as well. Also for prediction purposes, highly overlapping paths may be removed, while at the same time spatial variety is welcomed. In multi-modal network applications, a rich variety in multi-modal composition of the choice set is desirable, while at the same time non-feasible modal sequences (e.g. bus-car-bike-walk) should be removed from the base set. These filtering constraints are behavior-based following from travel observations. For a detailed account of selection criteria, see Fiorenzo-Catalano (2007).

As an illustration, example results of the route generation algorithm are shown in Figure 1. It illustrates that the number of generated routes for an OD pair is variable and depends on the OD pair and the network.

Figure 1(a) shows an OD pair with only two routes having a large mutual overlap. So, after filtering, only one route remains in the choice set. Figure 1(b) also shows an OD pair with two routes that however are clearly different and having a low overlap, hence both routes remain in the choice set. Finally, the OD pair in Figure 1(c) has many different routes in its choice set.
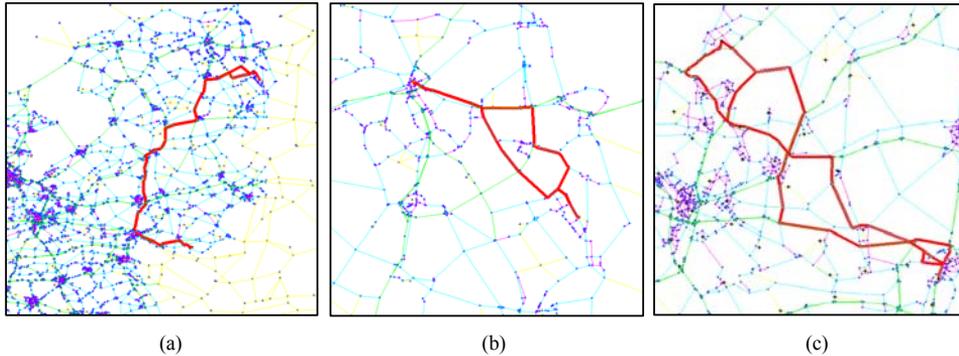
FIGURE 1: Example routes generated for different OD pairs (Dutch national road network)

Overlapping is a typical phenomenon of routes in networks and choice set generation methods always will produce overlapping paths to a certain extent. This is more a problem for the choice modeling step than for choice set generation (see Hoogendoorn-Lanser et al, 2005). In order to make choice model estimation and demand prediction efficient (with smallest acceptable choice set sizes) largely overlapping routes should be removed by the filtering step after generation (see Fiorenzo-Catalano, 2007; Hoogenddoorn-Lanser, 2005).

From earlier experiments in a multi-modal network (see Table 1) it is known that the doubly stochastic approach by far outperforms the single randomization approaches of either only network attributes or only preference parameters (see Fiorenzo-Catalano et al., 2004a,b). Given the same level of input variance, for the double randomization approach the prediction success rate of observed routes (full observed route is in generated route set) is about twice that (80% or more) of the single approaches (about 40% success rate).

TABLE 1: Some choice set generation results (base set) for sample of 37 observed OD-pairs (Rotterdam-Dordrecht corridor, The Netherlands) for 16 user groups combined (non-calibrated model)

| | Random attributes, 100 draws | Random Parameters, 100 draws | Double randomization 10×10 draws | Theoretical extreme |
|---|---|---|---|---|
| Total number of distinct multimodal routes generated | 286 | 283 | 701 | 1600×37 |
| Average choice set size | 8 | 8 | 19 | 1600 |
| Maximum choice set size | 18 | 16 | 36 | 1600 |
| Minimum choice set size | 3 | 4 | 10 | 1 |
| Prediction success rate of observed routes | 38% | 51% | 78% | 100% |

The results of the filtering process show that choice sets may dramatically be reduced (more than 50%) after adopting selection criteria such as a maximum overlapping constraint, modal sequence constraint, maximum detours, and the like. As an example, Table 2 shows choice set generation results from an application to the Dutch Inland Navigation Waterway Network for freight transport trips between 30 OD-pairs. After filtering of routes having a distance overlap of more than 90%, filtered choice sets on average were 73% smaller in size without significant loss in spatial variety.

TABLE 2: Impact of filtering on choice set size (Waterway network)

| | Maximum size | Minimum size | Average size | Total # of routes |
|---|---|---|---|---|
| Without filtering | 58 | 1 | 18.7 | 562 |
| With filtering | 12 | 1 | 5.1 | 154 |
| Reduction | | | | 408 = 73% |

## 6. PROBABILISTIC PROPERTIES

Because of its stochastic principle, a typical property of the presented generation approach is that its results are stochastic. This means that the size and composition of the generated choice sets are stochastic variables. The same holds for the order of route generation. It is a particular case of sampling with replacement. Whether a particular route alternative (e.g. the most attractive route) will have been generated (sampled) after a certain number of randomizations can only be said with a certain probability. Also whether the generated choice set will have a predefined minimum size after a certain number of randomizations can only be said with a certain probability. These probabilities on the one hand depend on the (mostly unknown) properties of the population of route alternatives in the network, while on the other hand they depend on the randomization properties (number of randomizations, seeds, variance levels, etc.).

Thus, in general, the relevant outcomes of the generation approach all are uncertain depending on its stochastic input. Therefore, in applying such generation approaches the question arises, for example, how many draws are needed to generate with a required probability a subset consisting of (at least) a certain number (say $k$) of different routes? Alternatively, the question might be what the probability of generating exactly $k$ different alternative routes ($k \leq n$) is given a predefined number of n draws from a network with $N$ alternative routes? Many more similar questions emerge in performing stochastic choice set generation approaches.

In the mathematical literature (Von Schelling, 1954), a specific instance of this type of generation problem is known as the Coupon Collectors Problem. It refers to the problem of a collector who tries to get a *complete* set of $N$ different desired objects (Collectors Items such as pictures of famous football teams). These objects are randomly contained in e.g. food products such as cans offered for sale in shops. The question then is how many of these products need to be purchased by a collector in order to have a complete set of the desired items. In this problem $k = N$ and required $n \geq N$, whereas in our choice set problem $k$ might be smaller than $N$.

Partly, such questions may be answered by monitoring the generation outcomes during the choice set generation process, such as for example the repetition of already selected alternatives. However, in planning the generation applications it is worthwhile to have some general rules at hand that can govern the experimental set up of the generation process such as the minimum number of random draws, or the best termination variable in the monitoring process.

On intuitive grounds, the following trivial probabilistic relationships can be stated:
a. with increasing number of randomizations:
   - choice set size will increase, albeit to a maximum;
   - choice set composition will be more stable;
   - # of selected attractive routes goes towards exhaustive;
   - # of selected unattractive routes is limited.
b. with increasing variance levels (of attributes and/or parameters):
   - choice set size will increase, not necessarily to a maximum size;

- choice set composition is less stable;
- # of selected attractive routes goes towards exhaustive;
- # of selected unattractive routes will increase.

These relationships indicate that variance levels should be kept modest.

We have made an in-depth theoretical analysis of the probabilistic properties of the presented choice set generation method (see Li et al., 2007).

A basic notion is the probability of having a particular alternative $r$ in the choice set after $n$ draws (of shortest routes) from a population of size $N$ in which route $r$ has a selection probability of $p_r$. This probability is:

$$P(S_r \leq n) = \sum_{m=1}^{n} (1-p_r)^{m-1} \cdot p_r = 1-(1-p_r)^n , \tag{2}$$

where $S_r$ is the number of random draws required to have for the first time alternative route $r$ in the choice set, $p_r$ is the selection probability of route r in population of routes of size $N$, and $n$ is the number of random draws (shortest path searches) taken from population of routes.

The expectation of the number of *different* routes ($K_n$) in the choice set for given $n$ and $N$ is given by:

$$E(K_n) = N - \sum_{r=1}^{N} (1-p_r)^n . \tag{3}$$

As long as route selection probabilities $p_r$ in the population show some regular patterns (for example are all equal to $1/N$, or regularly distributed on the probability axis) analytical expressions can be derived for most of the interesting questions. From these expressions a number of generally useful statements can be derived, such as the following sample may demonstrate (for details see Li et al., 2007):

- if (for limited universal sets) the ratio of sample size n and resulting choice set size $k$ is larger than 5, there is at least 90% confidence that this choice set size includes all alternatives with non-zero selection probability $p_r$.
- very large samples ($n \gg N$) are needed to guarantee sufficient stability in the composition of stochastically generated route choice sets;
- with large networks (many alternatives $N$) small choice sets (size $k$) can be generated already with small samples (size $n$).

The first finding implies for example that if during the generation process after more than 30 randomized draws 6 different routes have resulted, the process maybe stopped since the probability of finding additional different routes is very small (smaller than 0.1).

In addition to the theoretical findings, we show a number of results of a sensitivity analysis of the presented generation method to some of the randomization inputs by applying this method to a set of observations on interregional trips in the multi-modal network of the Rotterdam-Dordrecht area in The Netherlands (for details see Fiorenzo-Catalano, 2007). The generated choice set sizes presented below are generally very large because of the following reasons. All routes are unique at the finest level of spatial detail (links) and may overlap considerably because no filtering has been applied yet. In addition, in multi-modal networks, very many public transport alternatives exist because of the multiple modal combinations of the access, egress, and line haul trip parts.

The impact of randomization aspects on average generated choice set sizes (before filtering) is given in Table 3 showing how the number of randomized draws (from 10 to 30) from the link attribute distributions (rows) and from the preference parameter distributions (columns) influence the average generated size of the choice sets of 37

different trips. In each cell of the table three results achieved with three different seeds are given. The 3 choice set size values in the upper left cell refer to 100 random draws each (shortest path searches) whereas the rates in the lower right cell are based on 900 draws each. These choice set size predictions are stochastic outcomes.

TABLE 3: Average choice set size resulting from increasing numbers ($n$) of randomizations and 3 different random seeds S1 to S3 (same variances) ($N$=37 observed OD trips, link level, calibrated model)

| Parameters attributes | $n_\beta = 10$ S1  S2  S3 | $n_\beta = 15$ S1  S2  S3 | $n_\beta = 20$ S1  S2  S3 | $n_\beta = 25$ S1  S2  S3 | $n_\beta = 30$ S1  S2  S3 |
|---|---|---|---|---|---|
| $n_x$ =10 | 35.4–29.0–28.5 | 39.1–32.2–31.4 | 42.3–35.2–33.6 | 44.3–37.3–35.8 | 45.9–39.2–37.3 |
| $n_x$ =15 | 39.4–34.0–32.4 | 43.3–37.8–35.6 | 46.9–41.3–38.1 | 49.6–43.6–40.1 | 51.2–45.7–41.9 |
| $n_x$ =20 | 44.0–40.2–36.5 | 48.2–45.4–40.5 | 52.1–49.4–43.6 | 55.5–52.4–45.9 | 57.2–54.5–47.9 |
| $n_x$ =25 | 46.9– 42.1–39.0 | 51.1–47.5–42.9 | 54.8–51.8–46.2 | 58.1–54.8–48.8 | 60.1–56.9–50.9 |
| $n_x$ =30 | 50.1–44.8–42.2 | 54.6–50.8–46.5 | 58.2–55.5–50.2 | 61.7–58.6– 53.0 | 64.1–61.0–55.3 |

Expectedly, the generated choice set sizes increase with the number of randomizations though with diminishing growth rates. A doubling of average choice set size emerges when increasing the 10×10 randomization to a 30×30 level (9 fold). At the same time the random seed value appears to have a clear influence that happily diminishes with increasing numbers of randomizations. It appears that attribute randomization has a somewhat larger impact on changes in choice set sizes than parameter randomization.

Table 4 shows the impact of the number of randomized draws (from 10 to 30) from the link attribute distributions (rows) and from the preference parameter distributions (columns) on the prediction success rate (in percent correctly predicted modal sequences of chosen routes). In each cell of the table the results achieved with 3 different seeds are given. The 3 success rates in the upper left cell refer to 100 random draws each (shortest path searches) whereas the rates in the lower right cell are based on 900 draws each. In contrast to choice set sizes (Table 3), the spatial detail now is the so-called leg level of routes concerning only the correct prediction of the sequence of different modes used in the chosen multi-modal trips. A multi-modal route consists of 3 legs at least. These success rates are stochastic outcomes as well.

TABLE 4: Prediction success rates (in %) of chosen routes in the choice set for increasing number ($n$) of randomizations and 3 different seeds (same variances) ($N$=35 observed OD trips, leg level, calibrated model)

| Parameters attributes | $n_\beta = 10$ S1  S2  S3 | $n_\beta = 15$ S1  S2  S3 | $n_\beta = 20$ S1  S2  S3 | $n_\beta = 25$ S1  S2  S3 | $n_\beta = 30$ S1  S2  S3 |
|---|---|---|---|---|---|
| $n_x$ =10 | 74 – 77 – 86 | 77 – 80 – 86 | 80 – 80 – 86 | 80 – 83 – 86 | 80 – 86 – 86 |
| $n_x$ =15 | 80 – 83 – 86 | 83 – 83 – 86 | 86 – 83 – 86 | 86 – 86 – 89 | 86 – 89 – 89 |
| $n_x$ =20 | 83 – 86 – 89 | 86 – 86 – 89 | 89 – 86 – 89 | 89 – 89 – 91 | 89 – 89 – 91 |
| $n_x$ =25 | 83 – 86 – 89 | 86 – 86 – 89 | 89 – 86 – 89 | 89 – 89 – 91 | 89 – 89 – 91 |
| $n_x$ =30 | 83 – 86 – 89 | 86 – 86 – 91 | 89 – 86 – 91 | 89 – 89 – 94 | 89 – 89 – 94 |

In contrast to choice set size, prediction success rates appear less sensitive to the level of randomization. Already after 100 randomizations a fairly high level of success rate has been achieved. This is to be expected because the chosen route naturally is an attractive one that will be predicted early in the random generation process. Also in this case the impact of attribute randomization appears larger than parameter randomization. There is a clear difference in results between the 3 seeds that diminishes though with increasing level of randomization.

If instead of chosen routes observed consideration sets are used to measure the prediction success rates these appear expectedly to be somewhat lower while the sensitivity to randomizations appear to be similar. The outcomes presented in Tables 3 and 4 suggest that convergence is not yet achieved and that increasing the randomization further might slightly improve the results. In conformity with theoretical results, about 1000 randomizations per user group may give optimal results (convergence).

Finally let us look at the impact of randomization at the level of individual OD trips. For 35 of these trips (same data set from Rotterdam-Dordrecht region), we apply a calibrated (see Section 7) choice set generation procedure inclusive a few filtering steps such as on vehicle availability constraints. The resulting sets might e.g. be adopted for estimation of route choice models. We show the resulting choice set sizes of three randomizations (different seeds) with of 20 randomizations for attributes $X_j$ and parameters $\beta_j$ each.

Table 5 shows a few statistics on average outcomes demonstrating that averaged over the sample of OD pairs the outcomes are very close for the different randomizations. After applying the filter more than half of the routes have were removed (compare Table 3).

TABLE 5: Comparative statistics of choice set generation with three different randomizations (with calibrated model and after filtering)

| $N = 35$ OD trips | Total unique routes generated (leg level) | Average choice set size | Minimum choice set size | Maximum choice set size | Prediction success rate (at leg level) |
|---|---|---|---|---|---|
| Seed 1 | 816 | 23.3 | 7 | 47 | 89 % |
| Seed 2 | 723 | 20.7 | 7 | 47 | 86 % |
| Seed 3 | 782 | 22.3 | 7 | 50 | 89 % |

However, more interesting are the individual outcomes given in Figure 2 showing that the generated choice set sizes for individual OD trips are very close, sometimes even identical, for the different randomizations.
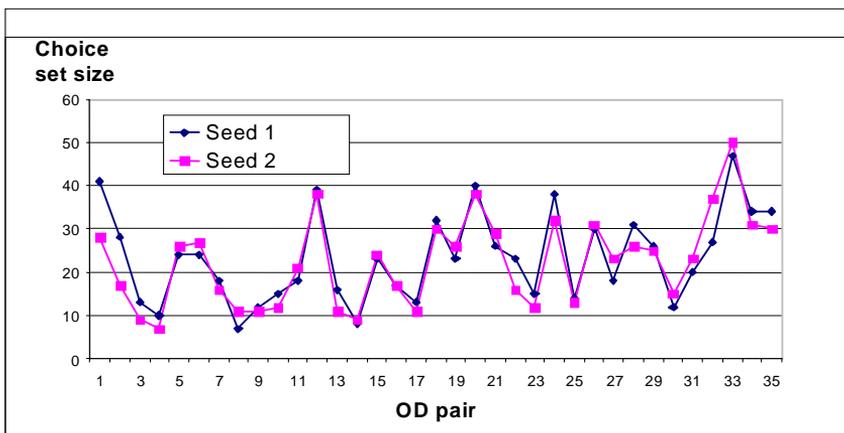


FIGURE 2: Choice set sizes for individual trips 1 to 35, compared for two random seeds

From these outcomes one may hypothesize that the compositions of the choice sets in terms of unique routes most probably are also very alike.

The various outcomes shown above indicate that the performance quality of the doubly stochastic choice set generation approach successfully competes with the best approaches known so far, as shown in Prato and Bekhor (2006).

## 7. CALIBRATION OF THE DOUBLY STOCHASTIC CHOICE SET GENERATION METHOD

In general, the attributes $X$ and parameters $\beta$ used in the generating function (1) are defined by probability distributions characterized by a distribution type, an expectation, and a variance.

The expected values of the attributes (link costs, link times, link distances, etc) in most cases may stem from readily available data and need not be calibrated. For the distributions of the random attribute values of links, the adoption of positive uniform distributions at the level of the links seems justified since the summation of random link values to corresponding route values then results in a smooth positive distribution of approximately normal form.

The variances of the uniform attribute distributions in principle are subject to calibration. In line with our behavioral hypothesis (see Section 5) that the attribute randomization reflects perception errors and the like, we confined the variance calibration to those time attributes known to be most important in travelers trip choices (e.g. in-vehicle times of car, PT, and bicycle, PT waiting time, etc.).

The expected values of the $\beta$ parameters as well as their distribution types and variances are generally subject to calibration. Because of the multitude of attributes not all parameters are calibrated. Many of them are determined by engineering judgment supported by the travel choice behavior literature. Especially the parameters of the time attributes are considered relevant for calibration. Also for parameter distributions positive uniform distributions are assumed after some experimenting with other distributional forms such as truncated normal. In a first step prior to calibration, default values are set for expectations and variances of which the default expectations were derived from exogenous choice model estimates published in literature.

The calibration then involves mainly a set of variance values to be determined by maximizing some measure of resemblance with observations of individual consideration sets or chosen routes. Various calibration performance measures are used for that purpose, such as:
- the number of chosen routes or routes in the consideration sets reproduced exactly by the generation algorithm (strong criterion);
- spatial coverage measures expressing to what extent the observed routes fully or partly are covered by the generated set of routes (weak criterion).

It is tempting to perform the calibration as a constrained optimization problem according to sound statistical rules. However, in-depth analysis shows that the objective function is a complex non-continuous function of the unknown decision variables (parameters and variances). This is the more true for multi-modal networks. Therefore, by a trial-and-error procedure the default values have been adapted to maximize the first criterion (leaving a more sophisticated method for future research).

Table 6 shows some of the calibration results achieved with observed consideration sets from a survey of multi-modal train trips in the Dordrecht-Rotterdam Corridor in the Netherlands.

TABLE 6: Subset of calibrated parameters of time attributes (all uniformly distributed)

| β parameters | Mode | E(β) | Max. dev. σ | minimum (E(β)-σ) | maximum (E(β)+σ) | Coeff. of variation |
|---|---|---|---|---|---|---|
| PT in-vehicle time | Bus/Tram | 1.0 | 0.4 | 0.6 | 1.4 | 23 % |
| | Metro | 0.8 | 0 | - | - | - |
| | IC Train | 1.0 | 0 | - | - | - |
| | Expr/Stop Train | 1.0 | 0.2 | 0.8 | 1.2 | 12 % |
| PT board/alight time | All modes | 1.0 | 0 | - | - | - |
| PT waiting time | Bus/Tram | 1.5 | 0.7 | 0.8 | 2.2 | 27 % |
| | Metro | 1.2 | 0.2 | 1.0 | 1.4 | 10 % |
| | IC Train | 1.2 | 0.2 | 1.0 | 1.4 | 10 % |
| | Expr/Stop Train | 0.9 | 0.3 | 0.6 | 1.2 | 19 % |
| Walk time | | 0.6 | 0.2 | 0.4 | 0.8 | 19 % |
| Car in-vehicle time | Secondary | 1.0 | 0.4 | 0.6 | 1.4 | 23 % |
| | Motorway | 1.0 | 0 | - | - | - |
| Car access/egress time | Secondary | 1.0 | 0 | - | - | - |
| Bicycle time | | 1.0 | 0.2 | 0.8 | 1.2 | 12 % |
| Bicycle access/egress time | | 1.0 | 0.2 | 0.8 | 1.2 | 12 % |

The parameter variances (last column) look modest. The decisive variance level in the application however is the product of the variances of the parameters and the attribute values, which is much larger.

The calibration exercise succeeded in improving the generation performance significantly from a prediction success rate of the chosen routes (leg level) of 78% before calibration to more than 90% after calibration.

## 8. APPLICATIONS AND ASSESSMENT

The developed doubly stochastic choice set generation approach has been successfully applied in a number of different networks (road, waterway, multimodal); for an overview see Fiorenzo-Catalano (2007).

Recent applications are done on the Dutch national main road network in the course of traffic flow predictions with a dynamic assignment model (see Bliemer and Taale, 2006). The stochastic choice set generation is a separate explicit step prior to the route choice modeling and route flow assignment. This assignment problem is characterized by the following demand and network figures: 345 zones (trip entry and exit points), 109.000 non-zero OD-flows, 25,341 one-directional links, and 17,963 nodes.

The route generation algorithm adopted the following parameters: maximum allowed route overlap = 80%, number of (shortest path) iterations = 25, path search criterion = travel time, and coefficient of variation in link travel time = 66%.

The generation algorithm produced the following outcomes: 469,000 unique routes, average choice set size = 4.3, large variation of choice set sizes among OD-pairs, and only few (13) minutes computation time.

This application demonstrated the following. This way of explicit prior route generation is feasible, even for very large networks. The attractive paths have to be calculated only once (at the beginning) in stead of repeatedly in classical assignment

approaches, which strongly reduces computation time. Starting the iterative assignment process with a set of routes instead of a single route (which is common for most assignment approaches) enables the distribution of OD flows over multiple routes already from the first iteration. This appears to speed up the convergence of the dynamic assignment significantly. Sets of (loaded) routes are available for further subsequent analyses.

The doubly stochastic choice set generation approach has also been applied in multi-modal networks, namely specifically in the Rotterdam-Dordrecht region in The Netherlands. Sections 6 and 7 show some relevant outcomes (choice set sizes, generation performance) related to a set of observations in that region. However, also regionwide applications have been performed. These show the following. Even with strict constraints (concerning overlap, modal sequences, detours and the like) mostly large choice sets result (about 25 routes). The generated choice sets show a large variety of unimodal and multimodal route alternatives (see Figure 3 for an example). Especially the public transport alternatives appear to be frequent and are manifold. They constitute the majority of generated alternatives. This is partly due to the good public transportation provision in the region and the multitude of different modes available to most travelers.

| Route Nr | Modes used and the related distances traveled | | | | | | | | | Total | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Access | | | **Main** | | | Egress | | | | |
| | Mode | Dist | Mode | Dist | Mode | Dist | Mode | Dist | Mode | Dist | Dist | Time |
| 1 | | | | | **Car** | 26.0 | | | | | 26.0 | 35.3 |
| 2 | | | | | **Car** | 24.9 | | | | | 24.9 | 37.8 |
| 3 | | | | | **Car** | 26.2 | | | | | 26.2 | 35.3 |
| 4 | | | | | **Car** | 27.1 | Metro | 6.5 | Walk | 0.2 | 27.1 | 42.3 |
| 5 | | | Bicycle | 3.7 | **Train** | 19.8 | Bicycle | 3.3 | | | 26.8 | 58.8 |
| 6 | | | Bicycle | 3.5 | **Train** | 19.8 | Bicycle | 3.3 | | | 26.6 | 58.1 |
| 7 | | | Bicycle | 3.5 | **Train** | 19.8 | Bicycle | 3.3 | | | 26.6 | 58.0 |
| 8 | | | Bicycle | 3.5 | **Train** | 19.8 | Bicycle | 3.3 | | | 26.7 | 58.3 |
| 9 | | | Bicycle | 3.5 | **Train** | 19.8 | Bicycle | 3.3 | | | 26.7 | 58.2 |
| 10 | | | Bicycle | 3.5 | **Train** | 19.8 | Bicycle | 3.3 | | | 26.7 | 58.3 |
| 11 | | | Bicycle | 3.5 | **Train** | 19.8 | Bicycle | 3.3 | | | 26.7 | 58.2 |
| 12 | | | Bicycle | 0.7 | **Train** | 22.2 | Bicycle | 3.3 | | | 26.3 | 58.6 |
| 13 | | | Bicycle | 0.7 | **Train** | 22.2 | Bicycle | 3.3 | | | 26.3 | 58.7 |
| 14 | | | Bicycle | 0.7 | **Train** | 22.2 | Tram | 3.1 | Walk | 0.2 | 26.3 | 61.1 |
| 15 | | | Bicycle | 3.7 | **Train** | 19.8 | Tram | 3.1 | Walk | 0.2 | 26.8 | 61.3 |
| 16 | | | Bicycle | 3.5 | **Train** | 19.8 | Tram | 3.1 | Walk | 0.2 | 26.7 | 60.6 |
| 17 | | | Bicycle | 3.5 | **Train** | 18.1 | Metro | 3.3 | Walk | 0.2 | 25.3 | 62.0 |
| 18 | | | Bicycle | 5.0 | **Train** | 16.1 | Metro | 3.3 | Walk | 0.2 | 24.7 | 66.4 |
| 19 | | | Walk | 0.7 | **Train** | 20.6 | Metro | 3.3 | Walk | 0.2 | 24.9 | 63.6 |
| 20 | | | Walk | 0.7 | **Train** | 22.2 | Tram | 3.1 | Walk | 3.1 | 26.3 | 66.2 |
| 21 | Walk | 0.7 | Bus | 2.9 | **Train** | 19.8 | Tram | 3.1 | Walk | 0.2 | 26.7 | 59.4 |
| 22 | Walk | 0.7 | Bus | 2.9 | **Train** | 18.1 | Metro | 3.3 | Walk | 0.2 | 25.3 | 60.8 |

FIGURE 3: Example multimodal choice set generated for an OD-trip (all modes available) from the City of Dordrecht to the City of Rotterdam

## 9. CONCLUSIONS

In establishing and assessing choice set generation methods, consideration should be given to the purpose of the choice sets, that is analysis of supply conditions, choice model estimation, or demand prediction. These purposes pose different requirements to the choice sets in terms of size, composition and variety.

Choice model estimation results are strongly dependent on the quality of the applied choice sets. For demand prediction, explicit generation of route choice sets has a number of advantages. A priori given choice sets allow nearly unlimited flexibility and freedom in choice models to be adopted. In addition to theoretical advantages, they offer significant computational advantages in iterative demand calculation procedures since repeated optimal path search no longer is necessary.

The proposed doubly stochastic choice set generation approach is goal-oriented and based on an easy understandable behavioral hypothesis. Despite the randomizations, its minimization principle (repeated shortest path search) guarantees that attractive alternatives will be in the choice set with very high probability while unattractive ones will have negligible probabilities precluding problems in choice modeling. Theoretical and experimental results show that despite the stochastic principle of the method, its outcomes in terms of size and composition of generated choice sets are fairly stable already at modest numbers of randomization iterations. The repeated shortest path search principle makes the method computationally very efficient, which has been proven by various applications in very large networks.

The various exercises performed with the doubly stochastic choice set generation approach have shown the following results:
- The resulting choice sets show the desired properties in terms of choice set size (minimum and maximum size) and composition (variety in trip properties); both the face validity (plausibility of generated alternatives) and empirical validity (conformity with observations) are very high; in nearly all cases the reported chosen alternative is member of the generated choice set;
- The achieved coverage at the level of complete individual trips is very high (about 90% in terms of numbers of modal legs covered) and is even higher  for particular link types such as access and egress legs of multi-modal train trips;
- The doubly stochastic approach by far outperforms the singly stochastic approaches (only randomizing attributes or parameters).

## ACKNOWLEDGEMENT

## REFERENCES

Bekhor,S., Ben-Akiva, M.E. and Ramming, S. (2001) Route choice: choice set generation and probabilistic choice models. Proceedings 4th TRISTAN Conference, Azores, Portugal.

Ben-Akiva, M.E., Bergman, M.J., Daly, A.J. and Ramaswamy, R. (1984) Modeling inter-urban route choice behavior. In J. Volmuller and R. Hamerslag (eds.) Proceedings 9th International Symposium on Transportation and Traffic Theory. VNU Science Press, pp. 299-330.

Ben-Akiva, M.E. and Bierlaire, M. (1999) Discrete choice methods and their applications to short term travel decisions. In R.W. Hall (ed.) Handbook of Transportation Science, Kluwer Academic Publishers, 1999, pp. 5-61.

Ben-Akiva, M.E. and Lerman, S. (1985) Discrete Choice Analysis: Theory and Application to Travel Demand, The MIT Press.

Bliemer, M.C.J. and Taale, H. (2006) Route generation and dynamic traffic assignment for large networks. Paper DTA Conference Leeds, June 2006.

Bovy, P.H.L. (2007) Modeling route choice sets in transportation networks: a preliminary synthesis. Proceedings VI TRISTAN Conference, Phuket, June 2007.

Bovy, P.H.L. and Stern, E. (1990) Route Choice: Wayfinding in Transport Networks, Kluwer Academic Publishers.

Cascetta E. and Papola, A. (2001) Random utility models with implicit availability/ perception of choice alternatives for the simulation of travel demand. Transportation Research Part C, 9, 249-263.

Fiorenzo-Catalano, S. (2007) Choice set generation in multi-modal transport networks. PhD Thesis, Delft University of Technology, TRAIL, The Netherlands.

Fiorenzo-Catalano, S., van Nes, R. and Bovy, P.H.L. (2004a) Choice set generation for multi-modal travel analysis. Proceedings of TRB Annual Meeting 2004, CD-ROM.

Fiorenzo-Catalano, S., van Nes, R. and Bovy, P.H.L. (2004b) Choice set generation for multi-modal travel analysis. European Journal of Transport Infrastructure Research, 4, 195-209.

Hoogendoorn-Lanser S. (2005) Modeling travel behavior for multi-modal transport networks. TRAIL Thesis Series T2005/4, TRAIL, The Netherlands.

Hoogendoorn-Lanser, S., van Nes, R. and Bovy, P.H.L. (2005) Path size and overlap in multi-modal transport networks. In H. Mahmassani (ed.) Transportation and Traffic Theory: Flows, Dynamics and Human Interaction, Elsevier, pp. 63-84.

Hoogendoorn-Lanser, S., van Nes, R. and Bovy, P.H.L. (2006) A rule-based approach to multi-modal choice set generation. Paper IATBR, Kyoto 2006.

Li, H., Bovy, P.H.L. and Hooghiemstra, G. (2007) Probabilistic properties of stochastic choice set generation. Delft University of Technology, TRAIL Studies in Transportation Science.

Prato, C.G. and Bekhor, S. (2006) Applying branch and bound technique to route choice set generation. Paper TRB Annual Meeting, CD-ROM.

Ramming, M.S. (2002) Network knowledge and route choice. PhD thesis, MIT, Cambridge.

VanderWaerden, P., Borgers, A. and Timmermans, H. (2004) Choice set composition in the context of pedestrian's route choice modeling. Proceedings TRB 2004 Annual Meeting CD-ROM, Paper # 04-2601.

Von Schelling, H. (1954) Coupon collecting for unequal probabilities. American Mathematical Monthly, 61, 304-311.