# Stochastic Path Generation Algorithm for Route Choice Models[*]

E. Frejinger[†]      M. Bierlaire

May 16, 2007

## Abstract

Defining choice sets is necessary when modeling route choice behavior using random utility models. Since the number of paths between a given origin-destination pair may be intractable, path enumeration algorithms are used for this purpose.

In this paper, we present a new point of view on choice set generation. In contrast to existing approaches, we hypothesize that all paths connecting the origin to the destination belong to the "true" choice set. In this context, we view stochastic path enumeration algorithms as importance sampling of alternatives. For this type of sampling protocol it is necessary to correct the path utilities in order to obtain unbiased parameter estimates. We propose a stochastic path enumeration algorithm that makes the definition of such sampling correction possible. Some preliminary numerical results are presented.

# 1   Introduction

Path enumeration algorithms play an important role in route choice modeling with random utility models since choice sets are in general unobservable.

[†]Corresponding author.   Ecole Polytechnique Fédérale de Lausanne, Mobility and Transport Laboratory (TRANSP-OR), CH-1015 Lausanne, Switzerland.   E-mail: emma.frejinger@epfl.ch

Due to the often very large number of paths, estimating a model based on all elementary paths connecting a given origin-destination (OD) pair may not be possible. It is therefore necessary to enumerate a limited set of paths.

Recently, several researches have turned their attention to choice set generation and its effects on route choice model estimation results (e.g. Bekhor and Prato, 2006, Bekhor et al., 2006 and van Nes et al., 2006). Various heuristics have been proposed in the literature with the objective to generate the set of paths a traveler actually considers. This set should include all attractive paths but no unreasonable paths. The modeler defines attractiveness and reasonableness based on observed route choices and personal judgment.

In this paper we present a new point of view on path enumeration for route choice modeling. In contrast to existing literature, we hypothesize that the true choice set for a given OD pair is the universal one. That is, the set of all feasible paths. The objective of the choice set generation is to define choice sets such that the model estimation and prediction results are unbiased. For this purpose we consider stochastic path enumeration algorithms as importance sampling approaches. In order to obtain unbiased results, it is necessary to correct for this type of sampling protocol when estimating and applying route choice models. We propose a general stochastic choice set generation approach and a specific algorithm that allows the computation of sampling correction.

In the following section we present a review of existing choice set generation approaches and in Section 3 an overview of sampling of alternatives. In particular, we derive the correction for the sampling protocol corresponding to the proposed algorithm (described in Section 4). We give some preliminary numerical results in Section 5 before presenting conclusions and discussing topics for future work.

## 2 Choice Set Generation Approaches

For a given OD pair the number feasible paths (including paths with cycles) is unbounded. It is therefore always necessary to constrain route choice
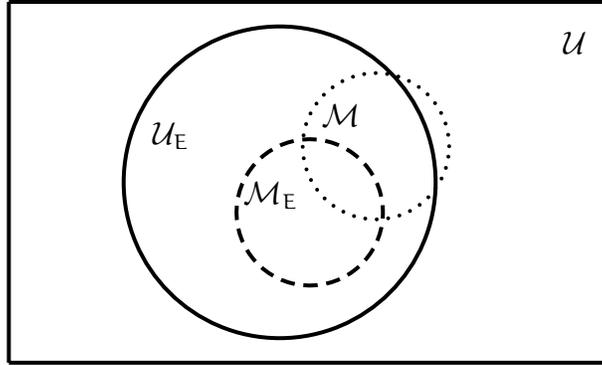
Figure 1: Illustration of path sets for a given OD pair

models to a limited number of alternatives. We illustrate in Figure 1 the different sets of paths for a given OD pair. In general only the set of elementary paths $\mathcal{U}_\text{E}$ is considered which is a subset of the unbounded universal set $\mathcal{U}$. The number of elementary paths is tractable but in real networks it is often too large to be enumerated. Existing path enumeration algorithms for route choice modeling generates a subset of elementary paths here denoted $\mathcal{M}_\text{E}$. The approach proposed in this paper produces a subset $\mathcal{M}$ that may contain paths with cycles.

Choice sets can be defined based on enumerated paths in two ways. Either, a deterministic way including all enumerated paths. Or, in a probabilistic way using the two-stage choice model proposed by Manski (1977) (see also Swait and Ben-Akiva, 1987, Ben-Akiva and Boccara, 1995, Morikawa, 1996 and Cascetta and Papola, 2001). This paper focuses on path enumeration and for the sake of simplicity we define choice sets deterministically.

Many heuristics for enumerating paths have been proposed in the literature. These can be divided into deterministic and stochastic approaches. The first category refers to algorithms always generating the same set of paths for a given OD pair. Examples of such approaches are link elimination (Azevedo et al., 1993), constrained k-shortest paths (e.g. van der Zijpp and Catalano, 2005), branch-and-bound (Friedrich et al., 2001, Hoogendoorn-Lanser, 2005 and Prato and Bekhor, 2006), labeled paths (Ben-Akiva et al.,

1984) and link penalty (de la Barra et al., 1993).

Two stochastic approaches have been proposed in the literature. Ramming (2001) used a simulation method that produces alternative paths by drawing link impedances from different probability distributions. The shortest path according to the randomly distributed impedance is calculated and introduced in the choice set. Recently, Bovy and Fiorenzo-Catalano (2006) proposed the so-called doubly stochastic choice set generation approach. Paths are enumerated by repeatedly computing shortest paths where the generalized cost function has both random parameters and random attributes. The algorithm has been applied to a multi-modal network.

# 3   Importance Sampling

The multinomial logit (MNL) model can be consistently estimated on a subset of alternatives. The probability that an individual $n$ chooses an alternative $i$ is then conditional on the choice set $\mathcal{C}_n$ defined by the modeler. This conditional probability is

$$P(i|\mathcal{C}_n) = \frac{e^{V_{in} + \ln q(\mathcal{C}_n|i)}}{\displaystyle\sum_{j \in \mathcal{C}_n} e^{V_{jn} + \ln q(\mathcal{C}_n|j)}} \tag{1}$$

and includes an alternative specific term, $\ln q(\mathcal{C}_n|j)$, correcting for sampling bias. This correction term is based on the probability of sampling $\mathcal{C}_n$ given that $j$ is the chosen alternative, $q(\mathcal{C}_n|j)$. See for example Ben-Akiva and Lerman (1985) or Train (2003) for detailed discussions on sampling of alternatives.

If all alternatives have equal selection probabilities, the estimation on the subset is done in the same way as the estimation on the full set of alternatives. Namely, $q(\mathcal{C}_n|i)$ is then equal to $q(\mathcal{C}_n|j)$ (uniform conditioning property, McFadden, 1978) and the correction for sampling bias cancels out in Equation (1). This simple random sampling protocol is however difficult to use in a path enumeration context. First of all, we are unaware of any

algorithm generating paths with equal probabilities without first enumerating the full set of paths. Second, due to the large (possibly intractable) number of paths, a simple random sample is likely to contain many alternatives that a traveler would never consider. Comparing the chosen path to a set of highly unattractive alternatives would not provide much information on the traveler's route choice.

Importance sampling is a more efficient scheme for path enumeration since it takes expected choice probabilities into account. Paths which are expected to have high choice probabilities have higher sampling probabilities than paths with lower expected choice probabilities. However, for this type of sampling protocol the correction terms in Equation (1) do not cancel out and $q(\mathcal{C}_n|j) \ \forall \ j \in \mathcal{C}_n$ must be defined. Note that if alternative specific constants are estimated, all parameter estimates except the constants would be unbiased even if the correction is not included in the utilities. In a route choice context it is in general not possible to estimate alternative specific constants and the correction for sampling is therefore essential.

We define a sampling protocol in the context of path enumeration as follows: a set $\widetilde{\mathcal{C}}_n$ is generated by drawing R paths with replacement from the universal set of paths $\mathcal{U}$ and adding the chosen path to it ($|\widetilde{\mathcal{C}}_n| = R + 1$). Each path $j \in \mathcal{U}$ has sampling probability $q(j)$ and $\sum_{j \in \mathcal{U}} q(j) \approx 1$. This approximation of the sum is based on the assumption that paths with cycles have very small probabilities.

The outcome of this protocol is $(\widetilde{k}_1, \widetilde{k}_2, \ldots, \widetilde{k}_J)$ where $\widetilde{k}_j$ is the number of times alternative j was drawn ($\sum_{j \in \mathcal{U}} \widetilde{k}_j = R$).

Following Ben-Akiva (1993) we derive the formulation of $q(\mathcal{C}_n|j)$ for this sampling protocol. The probability of an outcome is given by the multinomial distribution

$$P(\widetilde{k}_1, \widetilde{k}_2, \ldots, \widetilde{k}_J) = \frac{R!}{\prod_{j \in \mathcal{U}} \widetilde{k}_j!} \prod_{j \in \mathcal{U}} q(j)^{\widetilde{k}_j}. \tag{2}$$

The number of times alternative j appears in $\widetilde{\mathcal{C}}_n$ is $k_j = \widetilde{k}_j + \delta_{jc}$, where c denotes the index of the chosen alternative and $\delta_{jc}$ equals one if $j = c$ and zero otherwise. Let $\mathcal{C}_n$ be the set containing all alternatives corresponding

5

to the R draws ($\mathcal{C}_n = \{j \in \mathcal{U} \mid k_j > 0\}$). The size of $\mathcal{C}_n$ ranges from one to $R + 1$; $|\mathcal{C}_n| = 1$ if only duplicates of the chosen alternative were drawn and $|\mathcal{C}_n| = R+1$ if the chosen alternative was not drawn nor were any duplicates.

Using Equation (2), the probability of drawing $\widetilde{\mathcal{C}}_n$ given the chosen alternative i can be defined as

$$q(\widetilde{\mathcal{C}}_n|i) = \frac{R!}{(k_i - 1)! \prod\limits_{\substack{j \in \mathcal{C}_n \\ j \neq i}} k_j!} q(i)^{k_i - 1} \prod_{\substack{j \in \mathcal{C}_n \\ j \neq i}} q(j)^{k_j} = K_{\mathcal{C}_n} \frac{k_i}{q(i)} \tag{3}$$

where $K_{\mathcal{C}_n} = \frac{R!}{\prod_{j \in \mathcal{C}_n} k_j!} \prod_{j \in \mathcal{C}_n} q(j)^{k_j}$. We can now define the probability that an individual chooses alternative i given the set of draws $\widetilde{\mathcal{C}}_n$ as

$$P(i|\widetilde{\mathcal{C}}_n) = \frac{e^{V_{in} + \ln\left(\frac{k_i}{q(i)}\right)}}{\sum\limits_{j \in \mathcal{C}_n} e^{V_{jn} + \ln\left(\frac{k_j}{q(j)}\right)}}, \tag{4}$$

where $K_{\mathcal{C}_n}$ in Equation 3 cancels out since it is constant for all alternatives in $\mathcal{C}_n$.

In the following section we first present a general stochastic path enumeration approach that can be combined with various algorithms. Second we propose a biased random walk algorithm that allows for straightforward computation of path selection probabilities.

# 4   A Stochastic Path Enumeration Approach

This general stochastic approach for enumerating paths is based on the concept of subpaths where a subpath is a sequence of links. We define the probability of a subpath based on its distance to the shortest path. More precisely, its probability is defined by the double bounded Kumaraswamy distribution (Kumaraswamy, 1980) whose cumulative distribution function is $F(x_s|a, b) = 1 - (1 - x_s{}^a)^b$ for $x_s \in [0, 1]$. a and b are shape parameters and for a given subpath s with source node $v$ and sink node $w$, $x_s$ is defined as

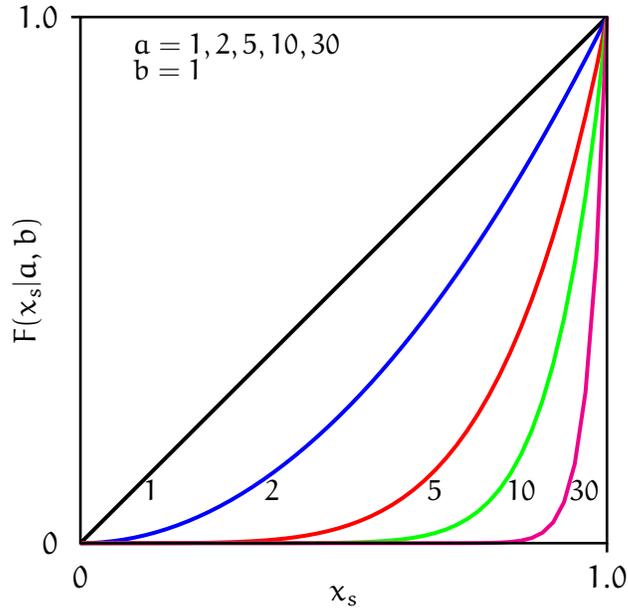$$x_s = \frac{SP(o, d)}{SP(o, w) + C(s) + SP(w, d)},$$

Figure 2: Kumaraswamy distribution - cumulative distribution function

where $C(s)$ is the cost of $s$, $o$ the origin, $d$ the destination and $SP(v_1, v_2)$ is the cost of the shortest path between nodes $v_1$ and $v_2$. Any generalized cost can be used in this context. Note that $x_s$ equals one if subpath $s$ is on the shortest path and $x_s \to 0$ as $C_s \to \infty$. In Figure 2 we show the cumulative distribution function for different values of $a$ when $b = 1$. The probabilities assigned to the subpaths can be controlled by the definition of the distribution parameters. High values of $a$ when $b = 1$ yield low probabilities for subpaths with high cost. Low values of $a$ have the opposite effect.

This is a flexible approach that can be used in various path enumeration algorithms including those presented in the literature. For example, in an algorithm similar to the link elimination approach but where the choice of subpaths (or links) to be eliminated is stochastic. Another example is a gateway algorithm, where a subpath is selected anywhere in the network, using the probability distribution described above. A generated path is composed of three segments: the shortest path from the origin to the source node of the subpath, the subpath itself, and the shortest path from the sink

7

node of the subpath to the destination. This gateway algorithm was used by Bierlaire et al. (2006) (see also Vrtic et al., 2006) for modeling long distance route choice behavior in Switzerland.

In this paper, we use a biased random walk algorithm that is described in the following section.

## 4.1 Biased Random Walk Algorithm

Starting from the origin, this algorithm selects a link using the probability distribution described previously. Another link starting at the sink node of the first one is then selected and this process is applied until the destination is reached and a complete path has been generated. The algorithm biases the random walk towards the shortest path in a way controlled by the parameters of the distribution. If a uniform distribution (special case of Kumaraswamy distribution with $a = 0$ and $b = 1$) is used then the algorithm corresponds to a simple random walk. Note however that a simple random walk does not generate a simple random sample of paths.

This algorithm has some nice properties that are important for an importance sampling approach. First, the path selection probabilities can be computed. The probability $q(j)$ of generating a path $j$ is the probability of selecting the ordered sequence of links $\Gamma_j$

$$q(j) = \prod_{\ell \in \Gamma_j} q(\ell | \mathcal{E}_v, a, b) \tag{5}$$

where $\ell$ denotes a link, $v$ its source node and $\mathcal{E}_v$ the set of outgoing links from $v$. In accordance with the approach presented previously $q(\ell | \mathcal{E}_v, a, b)$ is defined by the Kumaraswamy distribution using

$$x_\ell = \frac{SP(v, d)}{C(\ell) + SP(w, d)}.$$

A second property of this algorithm is that any path in $\mathcal{U}$ can potentially be generated, including paths with cycles.

8

# 5  Preliminary Numerical Results

In order to evaluate the effects of sampling correction we estimate models on synthetic data from two toy networks. Paths in the first network can have no overlap nor cycles. The second network is more complex and observations are generated using a probit model.

## 5.1  Network with Non-correlated Elementary Paths

Consider a network composed of 40 oriented links connecting two nodes (origin and destination). The universal choice set includes consequently 40 non correlated elementary paths. The links have different lengths (L) and some have a speed bump (SB). We associate a utility $U_j = \beta_L L_j + \beta_{SB} SB_j + \varepsilon_j$ with each path j, where $\beta_L = -0.6$, $\beta_{SB} = -0.3$ and $\varepsilon_j$ is distributed Gumbel (location parameter set to 0 and scale to 1). 500 observations have been generated by associating a choice with the highest utility for each set of draws of $\varepsilon_j \; \forall \; j \in \mathcal{U}$. The true model is hence MNL for this example.

We generate a set of paths for each observation using the biased random walk algorithm. The generalized cost function is the sum of length and number of speed bumps. Moreover, we make 40 draws using distribution parameters $a = 2$ and $b = 1$ which results in choice sets with 11.9 paths on average (maximum 18 and minimum 7).

The estimation results are reported in Table 1. We provide scaled coefficient estimates where the length coefficient has been normalized to its true value, $\widehat{\beta}_L = \beta_L = -0.6$. The scaled speed bump coefficient is significantly different from its true value $-0.3$ (t-test statistic 3.67) in the model without correction but this is not the case for the model with correction. This example confirms the theory on sampling of alternatives for path enumeration. Namely, a correction is necessary in order to obtain unbiased estimation results.

## 5.2  Network with Correlated Paths

The network is shown in Figure 3 where the origin and destination nodes are marked "O" and "D" respectively. All links have the length of one,

| Sampling correction | MNL without | MNL with |
|---|---|---|
| $\widehat{\beta}_L$ | -0.203 | -0.286 |
| Scaled estimate | -0.600 | -0.600 |
| Robust std. | 0.0193 | 0.019 |
| Robust t-test | -10.53 | -15.01 |
| $\widehat{\beta}_{SB}$ | -0.0194 | -0.143 |
| Scaled estimate | -0.0573 | -0.300 |
| Robust std. | 0.0662 | 0.0661 |
| Robust t-test | -0.29 | -2.17 |
| Null log-likelihood | -1069.453 | -1633.501 |
| Final log-likelihood | -788.42 | -759.848 |
| Adjusted $\bar{\rho}^2$ | 0.261 | 0.288 |
| BIOGEME (Bierlaire, 2005, Bierlaire, 2003) has been used for all model estimations. | | |

Table 1: Estimation Results for MNL Example

except the link in the upper left corner which has length three and the one in the lower right corner which has length two. Moreover, the links marked with SB have a speed bump. The network contains cycles, non elementary paths can therefore be enumerated with the biased random walk algorithm.

Path utilities are assumed to be link-additive and the utility for a link $\ell$ is $U_\ell = \beta_L L_\ell + \beta_{SB} SB_\ell + \sigma\sqrt{L_\ell}\nu_\ell$ with $\beta_L = -0.6$ and $\beta_{SB} = -0.4$. $\nu_\ell$ is distributed standard normal and the variance is assumed proportional to link length with a parameter $\sigma$ fixed to 0.8. In this case, observations can be generated according to a probit model (Burrell, 1968) by repeatedly computing the shortest path (minimizing $-U_a$) for each realization of the link utilities. Note that negative cycles are possible since $U_\ell$ can be positive. The shortest path algorithm cannot converge in the presence of negative cycles and these realizations of the link utilities are therefore ignored. 382 observations were generated using 500 realizations of the network.

We define a choice set for each observation in the same way as for the previous example but using 30 draws. The size of the choice sets ranges

Figure 3: Example Network

from 7 to 19 paths with an average of 13.5 paths.

We estimate MNL and path size logit (PSL) models (Ben-Akiva and Ramming, 1998, Ben-Akiva and Bierlaire, 2003, Frejinger and Bierlaire, 2007) with and without sampling correction. The results are reported in Table 2. The scaled $\widehat{\beta}_{\text{SB}}$ is significantly different from the true value $(-0.4)$ for both MNL models (t-test statistic of 6.18 and 6.28 respectively). A possible explanation is that the correlation is ignored in the MNL which biases the results. On the contrary, the scaled $\widehat{\beta}_{\text{SB}}$ is not significantly different from its true value for both PSL models. It seems that the path size term corrects for both sampling and correlation in this case. It is interesting to compare the standard deviation of the coefficient estimates between the two models. The estimates in the model with correction have smaller standard deviation. This supports the argument that the sampling bias is absorbed by the coefficients even thought this it does not significantly change the results for this example. Finally, note that $\widehat{\beta}_{\text{L}}$ is not significantly different from $-0.6$ for the models without correction. This is a coincidence since the scales of logit and probit models are different.

It is difficult to isolate the effects of sampling correction on the estimation results in the presence of correlation among alternatives. The reason is that we cannot estimate a model, such as probit, that is flexible enough to

11

capture the full correlation structure. The results are therefore necessarily biased since the PSL model approximates a nested logit model.

| Sampling correction | MNL without | MNL with | PSL without | PSL with |
|---|---|---|---|---|
| $\widehat{\beta}_L$ | -0.627 | -0.978 | -0.619 | -0.969 |
| Scaled estimate | -0.600 | -0.600 | -0.600 | -0.600 |
| Robust std. | 0.0397 | 0.032 | 0.0407 | 0.0358 |
| Robust t-test | -15.79 | -30.57 | -15.22 | -27.04 |
| $\widehat{\beta}_{SB}$ | -0.0822 | -0.0801 | -0.347 | -0.461 |
| Scaled estimate | -0.0787 | -0.0491 | -0.336 | -0.285 |
| Robust std. | 0.052 | 0.0559 | 0.182 | 0.158 |
| Robust t-test | -1.58 | -1.43 | -1.90 | -2.92 |
| $\widehat{\beta}_{PS}$ | | | 1.17 | 1.74 |
| Scaled estimate | | | 1.13 | 1.08 |
| Robust std. | | | 0.788 | 0.705 |
| Robust t-test | | | 1.49 | 2.47 |
| Null log-likelihood | -988.63 | -2769.959 | -988.63 | -2769.959 |
| Final log-likelihood | -676.111 | -653.396 | -674.481 | -649.268 |
| Adjusted $\bar{\rho}^2$ | 0.314 | 0.337 | 0.315 | 0.340 |
| BIOGEME (Bierlaire, 2005, Bierlaire, 2003) has been used for all model estimations. | | | | |

Table 2: Estimation Results for Example with Correlated Paths

# 6   Conclusions and Future Work

Defining choice sets is necessary for modeling route choice behavior with random utility models. In this paper we propose a new point of view on path enumeration. In contrast to existing literature, we hypothesize that all paths belong to the true choice set and view stochastic path generation as an importance sampling approach. In order to obtain unbiased parameter estimates it is necessary to correct path utilities for sampling bias.

We propose a stochastic path enumeration algorithm that allows the computation of path selection probabilities and sampling correction. Preliminary numerical results on two small networks are presented. This is ongoing research and several issues and questions remain to be investigated.

# 7    Acknowledgments

# References

Azevedo, J., Costa, M. S., Madeira, J. S. and Martins, E. V. (1993). An algorithm for the ranking of shortest paths, *European Journal of Operational Research* **69**: 97–106.

Bekhor, S., Ben-Akiva, M. E. and Ramming, S. (2006). Evaluation of choice set generation algorithms, *Annals of Operations Research* **144**(1).

Bekhor, S. and Prato, C. G. (2006). Effects of choice set composition in route choice modelling, *Paper presented at the 11th International Conference on Travel Behaviour Research*, Kyoto, Japan.

Ben-Akiva, M. (1993). Lecture notes on large set of alternatives. Massachusetts Institute of Technology.

Ben-Akiva, M., Bergman, M., Daly, A. and Ramaswamy, R. (1984). Modeling inter-urban route choice behaviour, *in* J. Vollmuller and R. Hamerslag (eds), *Proceedings of the 9th international symposium on transportation and traffic theory*, VNU Science Press, Utrecht, Netherlands, pp. 299–330.

Ben-Akiva, M. and Bierlaire, M. (2003). Discrete choice models with applications to departure time and route choice, *in* R. Hall (ed.), *Hand-*

*book of Transportation Science, 2nd edition*, Operations Research and Management Science, Kluwer, pp. 7–38. ISBN:1-4020-7246-5.

Ben-Akiva, M. and Boccara, B. (1995). Discrete choice models with latent choice sets, *International Journal of Research in Marketing* **12**: 9–24.

Ben-Akiva, M. and Lerman, S. (1985). *Discrete choice analysis: Theory and application to travel demand*, MIT Press, Cambridge, Massachusetts.

Ben-Akiva, M. and Ramming, S. (1998). Lecture notes: Discrete choice models of traveler behavior in networks. Prepared for Advanced Methods for Planning and Management of Transportation Networks. Capri, Italy.

Bierlaire, M. (2003). Biogeme: a free package for the estimation of discrete choice models, *Proceedings of the 3rd Swiss Transport Research Conference*, Ascona, Switzerland.

Bierlaire, M. (2005). An introduction to biogeme version 1.4. http://biogeme.epfl.ch.

Bierlaire, M., Frejinger, E. and Stojanovic, J. (2006). A latent route choice model in switzerland, *Proceedings of the European Transport Conference*, Strasbourg, France.

Bovy, P. H. and Fiorenzo-Catalano, S. (2006). Stochastic route choice set generation: behavioral and probabilistic foundations, *Paper presented at the 11th International Conference on Travel Behaviour Research*, Kyoto, Japan.

Burrell, J. (1968). Multipath route assignment and its application to capacity restraint, *Proceedings of the 4th International Symposium on the Theory of Road and Traffic Flow*.

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(4): 249–263.

de la Barra, T., Pérez, B. and Añez, J. (1993). Mutidimensional path search and assignment, *Proceedings of the 21st PTRC Summer Meeting*, pp. 307–319.

Frejinger, E. and Bierlaire, M. (2007). Capturing correlation with sub-networks in route choice models, *Transportation Research Part B* **41**(3): 363–378.

Friedrich, M., Hofsäss, I. and Wekeck, S. (2001). Timetable-based transit assignment using branch and bound, *Transportation Research Record* **1752**.

Hoogendoorn-Lanser, S. (2005). *Modelling Travel Behaviour in Multimodal Networks*, PhD thesis, Deft University of Technology.

Kumaraswamy, P. (1980). A generalized probability density function for double-bounded random processes, *Journal of Hydrology* **46**: 79–88.

Manski, C. (1977). The structure of random utility models, *Theory and decision* **8**: 229–254.

McFadden, D. (1978). *Modelling the choice of residential location*, Spatial interaction theory and residential location, North-Holland, pp. 75–96.

Morikawa, T. (1996). A hybrid probabilistic choice set model with compensatory and noncompensatory choice rules, *Proceedings of the 7th World Conference on Transport Research*, Vol. 1, pp. 317–325.

Prato, C. and Bekhor, S. (2006). Applying branch and bound technique to route choice set generation, *Presented at the 85th Annual Meeting of the Transportation Research Board*.

Ramming, M. (2001). *Network Knowledge and Route Choice*, PhD thesis, Massachusetts Institute of Technology.

Swait, J. and Ben-Akiva, M. (1987). Incorporating random constraints in discrete models of choice set generation, *Transportation Research Part B* **21**(2): 91–102.

Train, K. (2003). *Discrete Choice Methods with Simulation*, Cambridge University Press.

van der Zijpp, N. and Catalano, S. F. (2005). Path enumeration by finding the constrained k-shortest paths, *Transportation Research Part B* **39**(6): 545–563.

van Nes, R., Hoogendoorn-Lanser, S. and Koppelman, F. (2006). On the use of choice sets for estimation and prediction in route choice, *Paper presented at the 11th International Conference on Travel Behaviour Research*, Kyoto, Japan.

Vrtic, M., Schüssler, N., Erath, A., Axhausen, K., Frejinger, E., Bierlaire, M., Stojanovic, S., Rudel, R. and Maggi, R. (2006). Einbezug von reisekosten bei der modellierung des mobilitätsverhalten. Final report for SVI research program Mobility Pricing: Project B1, on behalf of the Swiss Federal Department of the Environment, Transport, Energy and Communications, IVT ETH Zurich, ROSO EPF Lausanne and USI Lugano.