

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

**ASSESSING THE QUALITY OF ORIGIN-DESTINATION MATRICES
DERIVED FROM ACTIVITY/TRAVEL SURVEYS:
RESULTS FROM A MONTE CARLO EXPERIMENT**

Mario Cools, Elke Moons, Geert Wets *

Transportation Research Institute
Hasselt University
Wetenschapspark 5, bus 6
BE-3590 Diepenbeek
Belgium
Fax.:+32(0)11 26 91 99
Tel.:+32(0)11 26 91 {31, 26, 58}
Email: {mario.cools, elke.moons, geert.wets}@uhasselt.be

* Corresponding author

Number of words = 4,342
Number of Figures = 8
Number of Tables = 4
Words counted: $4,342 + 12 \cdot 250 = 7,342$ words

Revised paper submitted: November 13, 2009

ABSTRACT

To support policy makers combating travel-related externalities, quality data is required for the design and management of transportation systems and policies. To this end, large amounts of money have been spent on collecting household and person-based data. The main objective of this paper is to assess the quality of origin-destination matrices derived from household activity/travel surveys. To this purpose, a Monte Carlo experiment is set up to estimate the precision of OD-matrices given different sampling rates. The Belgian 2001 census data, containing work/school-related travel information for all 10,296,350 residents, are used for the experiment. For different sampling rates, 2000 random stratified samples are drawn. For each sample, three origin-destination-matrices are composed: one at municipality level, one at district level, and one at provincial level. The correspondence between the samples and the population is assessed by using the Mean Absolute Percentage Error (MAPE) and a censored version of the MAPE (MCAPE). The results show that no accurate OD-matrices can be directly derived from these surveys. Only when half of the population is queried, an acceptable OD-matrix is obtained at provincial level. Therefore, it is recommended to use additional information to better grasp the behavioral realism underlying destination choices and to collect information about particular origin-destination pairs by means of vehicle intercept surveys. In addition, the results suggest using the MCAPE next to traditional criteria to examine dissimilarities between different OD-matrices. An important avenue for further research is the investigation of the effect of sampling proportions on travel demand model outcomes.

1 BACKGROUND

In the modern cosmopolite society, travel is a cornerstone for human development, both for personal and commercial reasons: travel is not only regarded as one of the boosting forces behind economic growth, but is also seen as a social need providing people the opportunity for self-fulfillment and relaxation. As a result of the continuous evolution of modern society (e.g. urban sprawl, increasing female participation in labor, decline in traditional household structures), transportation challenges have accrued and have become more complex (1). Consequently, combating environmental (e.g. greenhouse-emissions such as CO₂, methane, NO_x; noise, odor annoyance and acid precipitation), economic (e.g. use of nonrenewable energy sources; and time lost due to congestion) and societal (e.g. health problems such as cardiovascular and respiratory diseases; traffic casualties; community severance and loss of community space) repercussions is a tremendous task (2).

To support policy makers in addressing these externalities, quality data are required for the design and management of transportation systems and policies (3). To this end, during the last four decades, large amounts of money have been spent on collecting household and person-based data. For most metropolitan areas, the largest part of planning budgets (an estimated \$7.4 million per year) was devoted to the conduct of household and person travel surveys (4). The data collected by these surveys are used for a wide variety of applications, including traffic forecasting, transportation planning and policy, and system monitoring (3).

The main objective of this paper is to assess the quality of origin-destination matrices derived from travel surveys. Mark that origin-destination matrices are core components in both traditional four-step and modern activity-based travel demand models. A sample size experiment is set up to estimate the precision of the OD-matrices given different sampling rates. Thus, an assessment of the appropriateness of travel surveys for deriving origin-destination relations can be made. Note that different types of travel surveys exist: Cambridge Systematics (5) distinguished seven different commonly used types of surveys (household activity/travel surveys; vehicle intercept and external surveys; transit onboard surveys; commercial vehicle surveys; workplace and establishment surveys; hotel/visitor surveys; and parking surveys). Each of these survey types provides a unique perspective for input into travel demand models. In this paper, the term ‘travel survey’ is confined to the first category, namely the household activity/travel surveys.

In a household activity/travel survey, respondents are queried about their household characteristics, the personal characteristics of the members of the household, and about recent activity/travel experiences of some or all household members. For most regions, household activity/travel surveys remain the best source of trip generation and distribution data, and therefore, are an important building block for travel demand models. In addition to model building purposes, these surveys are also used to poll specific target populations (such as transit users and non-users), to assess the potential demand and level of public support for major infrastructural projects, and to create a deeper understanding of travel behavior in the region (5). For a more elaborate discussion concerning travel surveys the reader is referred to (3,5,6). Recent trends in household travel surveys are discussed by Stopher and Greaves (7).

The remainder of this paper is organized as follows. Section 2 provides an extended discussion on the set-up of the sample size experiment. The relationship between sampling rates and the precision of a general statistic (i.e. the proportion of the commuting population) is highlighted in the first part of Section 3. The second part of Section 3 provides the results and

1 corresponding discussion of the statistical analysis of the main sample size experiment. Finally,
 2 some general conclusions will be formulated and avenues for further research indicated.

3 4 **2 SET-UP OF THE SAMPLE SIZE EXPERIMENT**

5
 6 As mentioned in the introduction, the main goal of this paper is the assessment of the quality of
 7 origin-destination matrices derived from household activity/travel surveys and, consequently,
 8 providing an answer to the question of how large a sample size should be to provide accurate
 9 OD-information in a region. To this end a Monte Carlo experiment is set up to estimate the
 10 precision of the OD-matrices given different sampling rates. A Monte Carlo experiment involves
 11 the use of random sampling techniques and computer simulation to obtain approximate solutions
 12 to mathematical problems. It involves repeating a simulation process, using in each simulation a
 13 particular set of values of random variables generated in accordance with their corresponding
 14 probability distribution functions (8). A Monte Carlo experiment is a viable approach for
 15 obtaining information about the sampling distribution of a statistic (in this study the precision of
 16 an origin-destination matrix) of which a theoretical sampling distribution may not be available
 17 due to the complexity. Monte Carlo simulation is generally suitable for addressing questions
 18 related to sampling distribution, especially when a) the theoretical assumptions of the statistical
 19 theory are violated; b) the theory about the statistic of interest is weak; or c) no theory exists
 20 about the statistic of interest (9). The latter is the case in this study (i.e. the precision of OD-
 21 matrices given different sampling rates).

22 The Monte Carlo experiment reported in this paper focuses on commuting (i.e. work and
 23 school related) trips made in Belgium. The 2001 census data will be used for the experiment. In
 24 particular, the census queried information about the departure and arrival times and locations of
 25 work/school trips (when applicable) for all 10,296,350 residents. For different sampling rates,
 26 ranging from one (the full population) to a millionth, 2,000 random stratified samples were
 27 drawn (2,000 for each sampling rate). Note that this number of samples is common in
 28 transportation oriented simulation experiments (e.g. 10,11). To ensure that the persons in the
 29 samples were geographically distributed in the country, the sample was stratified by
 30 geographical area: three nested stratification levels were taken into account, namely province,
 31 district and municipality. The sample was proportionately allocated to the strata. In other words,
 32 the sample in each stratum was selected with the same probabilities of selection (12).

33 For each sample, the proportion of persons making commuting trips was calculated, and
 34 three corresponding (morning commute) origin-destination-matrices (OD-matrices) were
 35 composed: one OD-matrix on municipality level (589 by 589 matrix), one OD-matrix on district
 36 level (43 by 43 matrix), and one on provincial level (11 by 11 matrix). A side-note has to be
 37 made for the latter OD-matrix: actually there are only 10 provinces in Belgium, but the Brussels
 38 metropolitan capital area (accounting for about 1/10 of the entire population) was treated as a
 39 separate province. The correspondence of the sample proportion and sample OD-matrices with
 40 the population (census) proportion and OD-matrices was then tabulated.

41 The correspondence between the sample and the population is assessed by using the
 42 Mean Absolute Percentage Error (MAPE) and an accommodated version of the MAPE. The
 43 MAPE is the mean of the Absolute Percentage Errors (APE) and is calculated by:

44
$$MAPE_{ij} = \sum_i \sum_j APE_{ij} / N, \text{ with } APE_{ij} = \left| \frac{A_{ij} - E_{ij}}{A_{ij}} \right| \times 100,$$

1 where A_{ij} is the population count for the morning commute from origin i to destination j , E_{ij} the
 2 sample count (scaled up to population level) for this morning commute, and N the total number
 3 of origin-destination cells. Despite its widespread use, the MAPE has several disadvantages.
 4 Armstrong and Collopy (13) for instance, argued that the MAPE is bounded on the low side by
 5 an error of 100% (origin-destination counts are all positive integers), but there is no bound on
 6 the high side. In response to this comment, Makridakis (14) proposed a modified MAPE
 7 (MDAPE), which is often referred to as SAPE (smoothed absolute percentage error) or SMAPE
 8 (symmetric mean absolute percentage error). This modified MAPE (MDMAPE) is given by:

$$9 \quad MDMAPE_{ij} = \sum_i \sum_j MDMAPE_{ij} / N, \text{ where } MDMAPE_{ij} = \left| \frac{A_{ij} - E_{ij}}{(A_{ij} + E_{ij})/2} \right| \times 100.$$

10 Although this modification accommodates the above described problem, it treats large positive
 11 and negative errors very differently (15). Therefore, in this paper, a new modification of the
 12 MAPE is proposed, named the Mean Censored Absolute Percentage Error (MCAPE). This new
 13 statistic takes into account the above described comments by limiting the positive values to a
 14 maximum of 100. Mathematically, the MCAPE is given by the following formula:

$$15 \quad MCAPE_{ij} = \sum_i \sum_j MCAPE_{ij} / N, \text{ where } MCAPE_{ij} = \min \left\{ 100, \left| \frac{A_{ij} - E_{ij}}{A_{ij}} \right| \times 100 \right\}.$$

16 When A_{ij} in the above formulae would be equal to zero, the different criteria would be
 17 undefined. This has been remedied by equalizing the APE_{ij} , $MDAPE_{ij}$ and $MCAPE_{ij}$ to zero in
 18 these occasions. After all, when the true population count equals zero (no person in the full
 19 population corresponds to the considered origin-destination pair) the up-scaled sample count
 20 also equals zero, and thus the true zero is correctly estimated.

21 The correspondence between the sample proportion (p) of persons making commuting
 22 trips and population proportion (π) is calculated by simply calculating the Absolute Percentage
 23 Error (APE):

$$24 \quad APE = \left| \frac{\pi - p}{\pi} \right| \times 100.$$

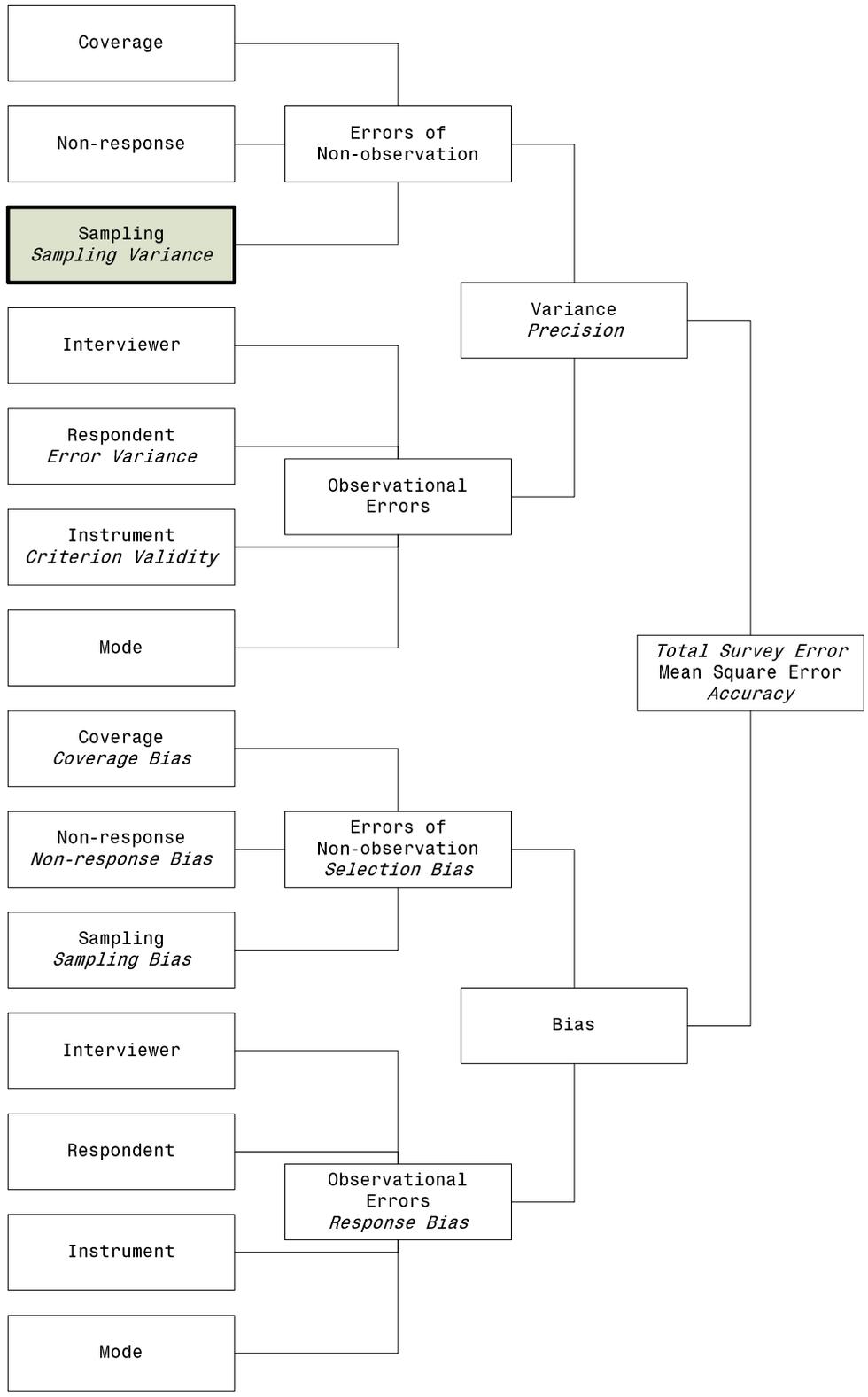
25 No accommodation of this APE was required, as the population proportion (π) was equal to
 26 62.59 percent, and consequentially the APE could not exceed 100.

27 To recapitulate, for each sampling rate, 2000 MAPE values and MCAPE values are
 28 calculated for the OD-matrix on municipality level, for the OD-matrix on district level, and for
 29 the OD-matrix on province level. In addition 2000 APE values are computed for the commuting
 30 proportion. For each of these sets of 2000 values, the 2.5th percentile, the 5th percentile, the 95th
 31 percentile and the 97.5th percentile was calculated. The k^{th} percentile is that value x such that the
 32 probability that an observation drawn at random from the population is smaller than x , equals k
 33 percent (16). The 2.5th percentile and 97.5th percentile are used to construct the 95% percentile
 34 interval which will be illustrated graphically as lower and upper bounds for the median. The 5th
 35 percentile and 95th percentile will be displayed in the corresponding tables because one is most
 36 often only interested in the one-sided alternative. In addition, the median (the 50th percentile)
 37 and the arithmetic mean are also computed.

38 To guarantee that the Monte Carlo experiment is really estimating the precision of the
 39 OD-matrices in function of different sample rates, rather than in function of other (unobserved)
 40 effects, one could take a look at the different sources of errors and biases in surveys. Groves (17)
 41 distinguished different sources of inaccuracy in surveys, of which an overview is given in Figure

1 1. Since in this experiment the true population values are known, and samples are drawn under
2 ideal circumstances (no response bias, no selection bias, no observation errors, no non-response
3 and perfect coverage), the resulting variations in the experiment are only a consequence of the
4 sampling variance (indicated with a gray box, framed with a tick black line in Figure 1). Thus, as
5 intended, the relationship between different sample sizes and the precision/accuracy (sample
6 variance) of the quantities under study are investigated.

7



1
2

FIGURE 1 Potential sources of the total survey error.

1 **3 RESULTS**

2

3 **3.1 Proportion of the Commuting Population**

4

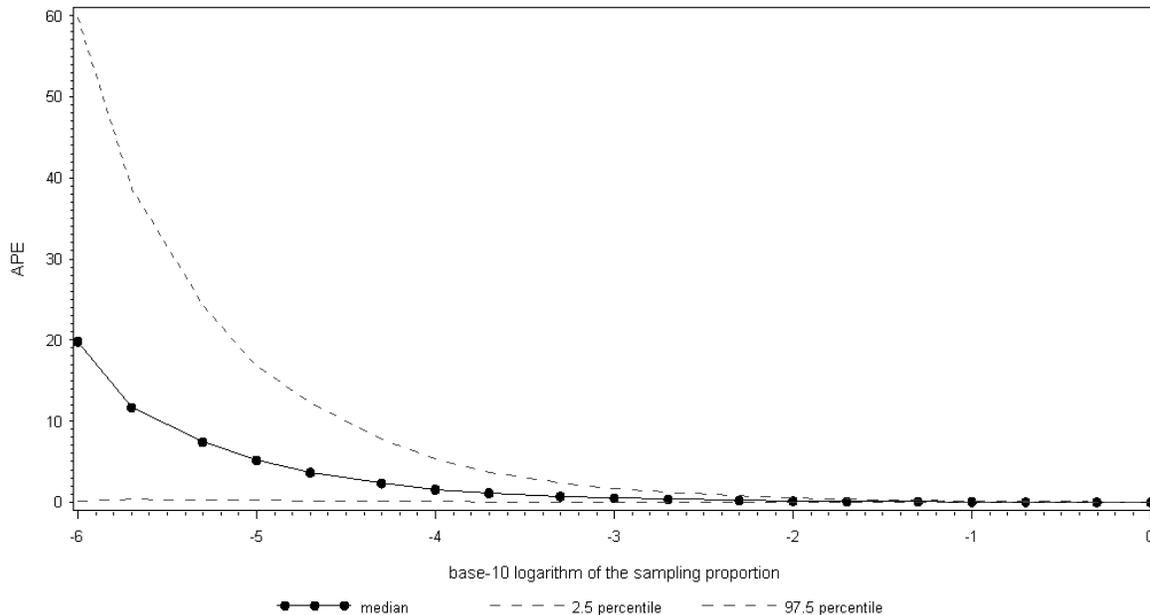
5 Before elaborating on the quality of OD-matrices in the second part of this Section, in this first
 6 part, an assessment of the appropriateness of travel surveys for deriving traditional indices, such
 7 as the mean number of trips made or the mean number of activities performed by
 8 individuals/households, or the proportion of the population making work/school-related trips, the
 9 latter being subject of the Monte Carlo experiment, is made. For traditional indices such as the
 10 mean number of trips made / activities performed by individuals/households, classical sample
 11 size calculations can be used to determine optimal sample sizes. Cools et al. (18), for instance,
 12 calculated the required number of households for a household activity survey using the
 13 following formula:

14
$$n \geq \frac{z^2 p(1-p)}{md^2},$$

15 where n equals the sample size, p the sample (survey) proportion, md the maximal deviation and
 16 z the z -value of the desired confidence interval. For the ‘safest’ case (i.e. $p = 0.5$), a maximal
 17 deviation of 2% and a confidence level of 95% would require a minimum of at least 2,401
 18 households. This example illustrates that for aggregate indices, such as the proportion of the
 19 commuting population, a clear theory exists and Monte Carlo simulation is not per se required.
 20 Notwithstanding, an investigation of the relationship between sampling rates and precision
 21 (sample variance) is still valuable, and especially contributes to the literature when the focus is
 22 turned to the different percentiles that are examined.

23 Results from the Monte Carlo experiment for the proportion of commuters in the
 24 population are graphically displayed in Figure 2 and numerically represented in Table 1.

25



26

27 **FIGURE 2 Relationship between absolute percentage error and sampling rate for**
 28 **commuting proportion.**

Figure 2 shows a clear relationship between the Absolute Percentage Error (APE) and the sampling proportion. As expected, the additional improvement in precision decreases as the sampling rate increases: for instance the increase in precision (decrease in APE) from a sampling rate of one millionth (base-10 logarithm of the sampling proportion equals minus 6) to one hundred-thousandth (base-10 logarithm equals minus 5) is considerably larger than the increase in precision from a sampling rate of one thousand to one hundred. This is especially so for the upper bound of the 95% percentile interval (97.5th percentile).

TABLE 1 APE-Statistics for the Commuting Proportion given Different Sampling Rates¹

Sampling Rate (SR)	Log ₁₀ SR	Mean	P5	Median	P95
0.000001	-6.00	20.270	1.670	19.826	46.744
0.000002	-5.70	13.915	1.096	11.684	34.377
0.000005	-5.30	8.642	0.659	7.412	21.509
0.000010	-5.00	6.083	0.474	5.192	14.849
0.000020	-4.70	4.293	0.343	3.671	10.618
0.000050	-4.30	2.741	0.213	2.317	6.793
0.000100	-4.00	1.879	0.137	1.574	4.586
0.000200	-3.70	1.330	0.097	1.140	3.219
0.000500	-3.30	0.853	0.072	0.723	2.097
0.001000	-3.00	0.602	0.052	0.508	1.485
0.002000	-2.70	0.430	0.040	0.362	1.054
0.005000	-2.30	0.269	0.022	0.227	0.659
0.010000	-2.00	0.190	0.015	0.160	0.462
0.020000	-1.70	0.129	0.010	0.109	0.313
0.050000	-1.30	0.080	0.007	0.067	0.197
0.100000	-1.00	0.053	0.004	0.045	0.132
0.200000	-0.70	0.034	0.003	0.029	0.082
0.500000	-0.30	0.016	0.001	0.013	0.038
1.000000	0.00	0.000	0.000	0.000	0.000

¹ 'P' stands for the percentile, e.g. P5 stands for the 5th percentile.

The results also show that when the full population is sampled, an absolute precision is obtained (absence of all variation). By definition this result should be obtained. When an average deviation of 5 percent is considered acceptable, a sample rate between 1 and 2 hundred-thousandth is required (5 percent lies between the mean values 4.293 and 6.083). On the other hand, from the median value one could conclude that in 50% of the cases the maximal deviation (APE) is smaller than 5.192 percent. A more cautious approach entails the use of the 95th percentiles. Suppose that only 5% of the cases the APE was allowed to exceed 2, then a sampling rate of about 5 ten-thousandth would be required, which roughly corresponds to sampling 5000 persons.

3.2 Precision of Origin-Destination Matrices

In this part of the result section, an assessment of the appropriateness of household activity/travel surveys for deriving OD-matrices is made. Recall that a Monte Carlo simulation is particularly suitable for addressing the questions concerning the distribution of the precision of these OD-matrices, as no real theoretical background of this distribution exists. First, attention will be paid to OD-matrices at municipality level. Afterwards, the focus is laid on OD-matrices at district and provincial level.

3.2.1 OD-Matrices at Municipality Level

Before expanding on the results of the Monte Carlo experiment, it is important to mention that the true OD-matrix (OD-matrix composed from the full population) is a very large and sparse matrix: of the 346,921 origin-destination pairs (589 times 589), 77.8% are zero-cells. As zero-cells in the full population are by definition correctly predicted by taking a sample from this population, the actual overall precision is significantly boosted by the sparseness of the true OD-matrix. Therefore, the decision was made to present the results based on the 76,882 non-zero-cells. To derive the values that include the zero-cells, one only needs to divide the MAPE and MCAPE values by 4.512 (= all cells / (all cells – zero-cells)).

TABLE 2 MAPE and MCAPE for OD-Matrices Derived at Municipality Level

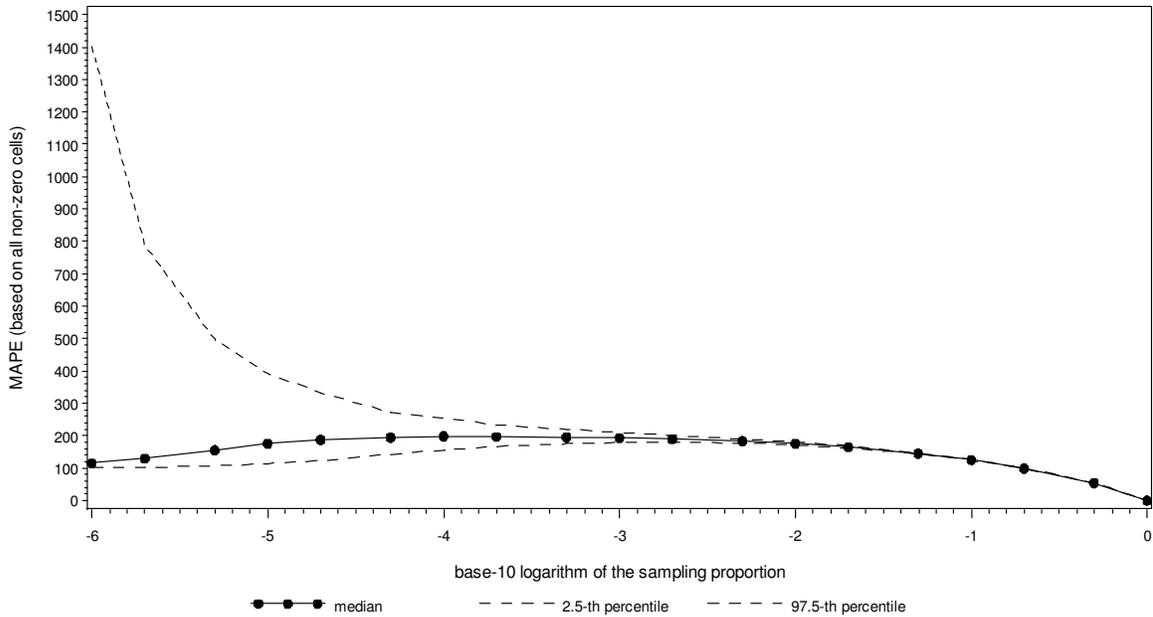
Sampling Rate (SR)	Log ₁₀ SR	MAPE				MCAPE			
		Mean	P5	Median	P95	Mean	P5	Median	P95
0.000001	-6.00	205.996	100.902	115.559	753.493	100.000	100.000	100.000	100.000
0.000002	-5.70	200.992	102.528	129.575	754.266	100.000	100.000	100.000	100.000
0.000050	-5.30	199.861	109.279	155.540	431.119	99.999	99.998	99.999	100.000
0.000010	-5.00	199.089	116.768	175.588	352.896	99.997	99.995	99.997	99.999
0.000020	-4.70	198.647	127.682	187.707	304.125	99.994	99.992	99.994	99.996
0.000050	-4.30	198.338	148.261	194.723	259.440	99.977	99.972	99.977	99.981
0.000100	-4.00	199.452	160.661	198.325	243.313	99.927	99.919	99.927	99.936
0.000200	-3.70	197.459	170.778	196.746	226.279	99.784	99.769	99.784	99.798
0.000500	-3.30	195.930	178.611	195.344	215.156	99.414	99.391	99.414	99.437
0.001000	-3.00	193.624	181.284	193.511	206.508	98.999	98.973	98.999	99.026
0.002000	-2.70	190.182	181.664	189.982	199.411	98.393	98.360	98.392	98.427
0.005000	-2.30	183.425	177.916	183.465	188.907	97.033	96.990	97.033	97.077
0.010000	-2.00	175.654	171.907	175.659	179.551	95.352	95.298	95.353	95.407
0.020000	-1.70	164.993	162.373	165.044	167.739	92.797	92.730	92.798	92.865
0.050000	-1.30	145.263	143.745	145.271	146.836	87.514	87.424	87.515	87.598
0.100000	-1.00	124.970	124.078	124.960	125.866	81.369	81.273	81.369	81.469
0.200000	-0.70	99.172	98.724	99.169	99.631	72.293	72.193	72.294	72.392
0.500000	-0.30	54.108	54.089	54.108	54.128	54.108	54.089	54.108	54.128
1.000000	0.00	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000

1 Inspection of Table 2 immediately reveals that no accurate OD-matrices are obtained at
2 municipality level, even if zero-cells are taken into account: a survey that would query half of
3 the population still would have an average absolute percentage error of 11.99 percent when zero-
4 cells are included and correspondingly of 54.11 % when only the actual predictions (non-zero-
5 cells) are taken into account. This clearly indicates that the direct derivation of origin-destination
6 matrices from household activity/travel surveys should be avoided. Notwithstanding, origin-
7 destination matrices derived from household activity/travel surveys are very valuable: even a
8 simple gravity model with the inverse squared distance as deterrence function, taking into
9 account the productions and attractions derived from the surveys, already results in a clear
10 improvement of the OD-matrices. This is certainly a plea for travel demand models that
11 incorporate the behavioral underpinnings of destination choices (activity location choices) given
12 a certain origin, like for instance models that make use of space-time prisms, e.g. (19), and
13 models that combine data from different sources, such as data integration tools, e.g. (20). In
14 addition, OD-matrices derived from travel surveys form a good basis for OD-matrices derived
15 from traffic counts: as multiple OD-matrices can be derived from the same set of traffic counts,
16 OD-matrices derived from travel surveys provide a good basis for constraining the matrices
17 derived from traffic counts (21). A thorough look at Table 2 also reveals that when half the
18 population is sampled, the values for the MAPE and MCAPE are the same. This can be
19 explained by the fact that when using half of the population none of the 2000 samples has a
20 MAPE higher than 1.

21 When the general tendency of the precision of the OD-matrices derived from travel
22 surveys is discussed, Figures 3 and 4 provide a clear insight in the relationship between the
23 precision and the sampling rate. From Figure 3 one can clearly see that the median MAPE first
24 increases when samples are becoming larger, and then starts to decrease. The increase in median
25 MAPE for the smallest sampling rates can be accounted for by the fact that on average more
26 cells are seriously overestimated, whereas the maximum underestimations are bounded by
27 100%. This effect is filtered out by using the MCAPE, as can be seen from Figure 4; a clear
28 decreasing relationship is visible here. Next to the difference in relationships between the MAPE
29 and MCAPE, one could also observe a clear difference between the percentile interval for the
30 MAPE and the percentile interval for the MCAPE. By condensing the APE to a maximum of
31 one (i.e. the CAPE), almost all variability around the median value is filtered out: the 2.5th and
32 97.5th percentiles almost coincide with the median values in case of the MCAPE.

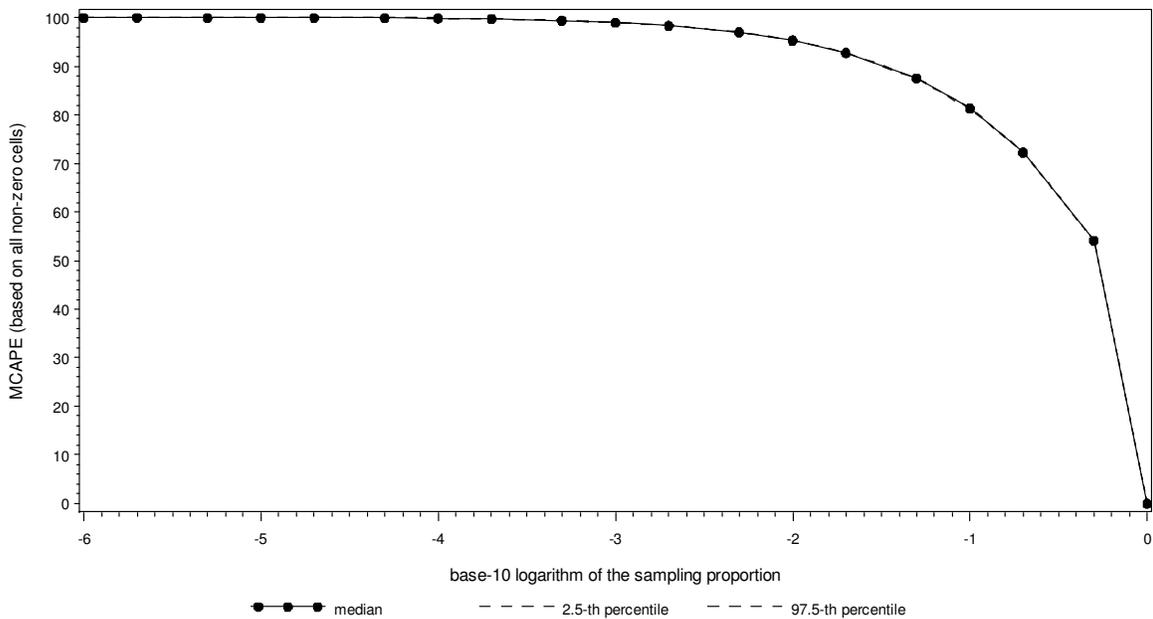
33 When this decreasing pattern of the MCAPE (Figure 4) is compared to the one of the
34 proportions (Figure 2), a clear contrast in the tendency can be seen: while the pattern for
35 proportion is a convex decreasing function, for the OD-matrices this is a concave decreasing
36 function. This difference in pattern, as well as the difference in precision, can be explained by
37 the fact that proportions are aggregate indices, and that surveys are extremely suitable for
38 capturing these aggregate figures, while in OD-matrices all individual information is used.

39
40



1
2
3

FIGURE 3 Relationship between MAPE and base-10 logarithm of the sampling rate.



4
5
6

FIGURE 4 Relationship between MCAPE and base-10 logarithm of the sampling rate.

7 *3.2.2 OD-Matrices at District Level*

8
9
10
11
12

Similar to the true OD-matrix at municipality level, the true OD-matrix at district level (43 by 43) comprises a non-negligible amount of zero-cells. Nonetheless, the number of non-zero-cells is considerably smaller: 10.4% of the 1,849 origin-destination pairs are zero-cells. Recall that zero-cells in the full population are by definition correctly predicted by taking a sample from this

1 population. Therefore, similar to the previous paragraph, the results are based on the 1,657 non-
 2 zero-cells. The values that include the zero-cells can be calculated by dividing the MAPE and
 3 MCAPE values by 1.116.

4 A thorough look at Table 3 shows that also at district level no accurate OD-matrices can
 5 be derived. Even if zero-cells are included in the calculations, surveying half of the population
 6 would result in an average absolute percentage error of 22.25 percent (24.832 divided by 1.116).
 7 When compared to the results of OD-matrices derived at municipality level, the results including
 8 the zero-cells are worse at district level than at municipality level (an average MAPE of 22.25
 9 percent versus one of 11.99 percent). This is due to the fact that at municipality level a much
 10 larger share (77.8 percent versus 10.4 percent) of zero-cells is automatically correctly predicted.
 11 In contrast, when the results of only the non-zero-cells are compared, the precision of the OD-
 12 matrices derived at district level is higher than the precision of the OD-matrices derived at
 13 municipality level. This result confirms that predictions on a more aggregate level are more
 14 precise.

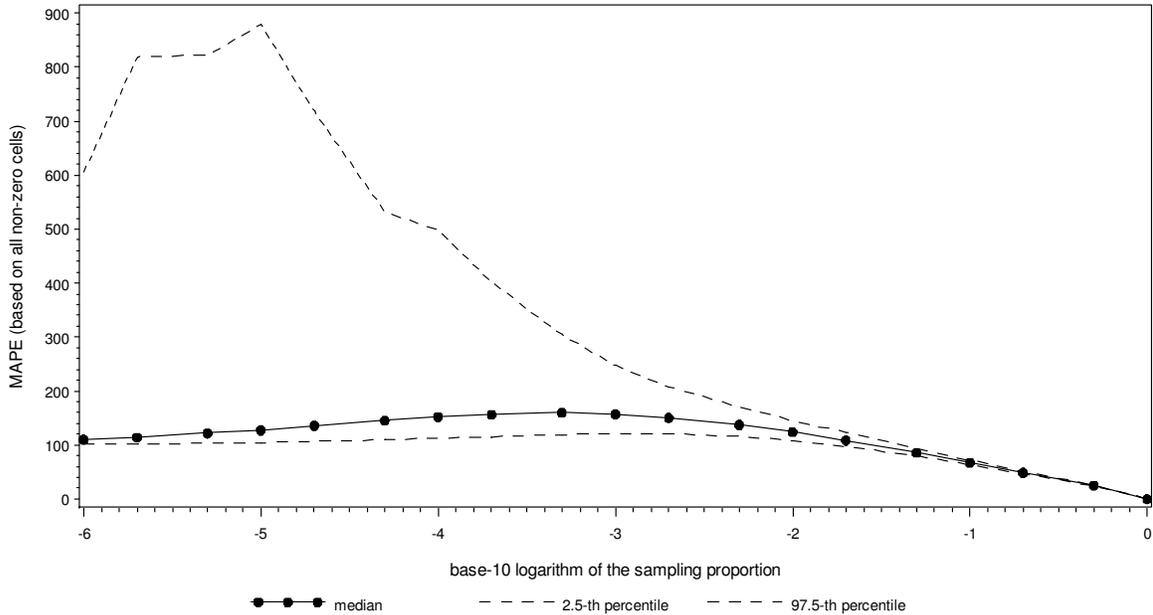
15 **TABLE 3 MAPE and MCAPE for OD-Matrices Derived at District Level**

16

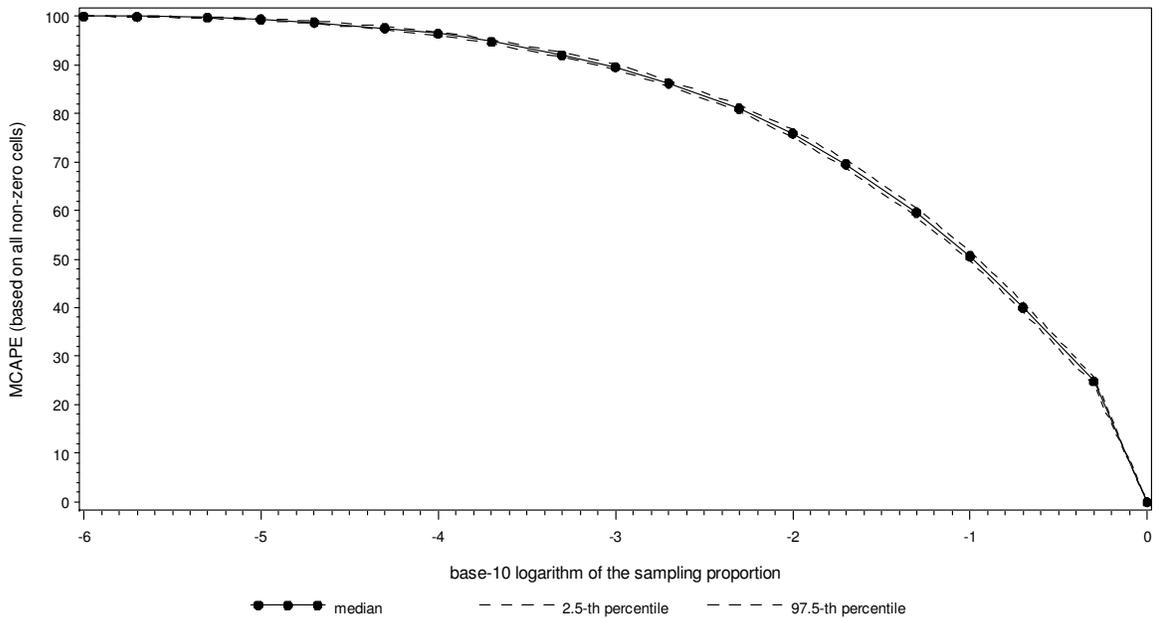
Sampling Rate (SR)	Log ₁₀ SR	MAPE				MCAPE			
		Mean	P5	Median	P95	Mean	P5	Median	P95
0.000001	-6.00	171.196	100.978	110.178	299.050	99.996	99.987	99.993	100.000
0.000002	-5.70	206.097	101.875	113.822	385.419	99.937	99.869	99.937	100.000
0.000050	-5.30	197.675	104.054	121.673	442.676	99.716	99.589	99.717	99.841
0.000010	-5.00	200.483	105.752	127.559	481.848	99.352	99.174	99.351	99.535
0.000020	-4.70	191.212	108.363	136.036	417.340	98.750	98.524	98.749	98.970
0.000050	-4.30	186.436	111.940	145.249	383.240	97.578	97.331	97.572	97.860
0.000100	-4.00	187.614	114.861	152.048	359.850	96.444	96.145	96.449	96.738
0.000200	-3.70	177.715	118.461	156.391	307.348	94.788	94.416	94.791	95.131
0.000500	-3.30	172.565	122.683	160.022	274.422	92.023	91.596	92.027	92.443
0.001000	-3.00	164.271	124.944	157.032	230.024	89.557	89.113	89.557	90.004
0.002000	-2.70	154.307	123.783	150.441	196.883	86.232	85.711	86.228	86.753
0.005000	-2.30	138.952	118.396	137.611	164.985	81.028	80.476	81.014	81.605
0.010000	-2.00	124.538	110.153	123.828	141.608	75.874	75.224	75.874	76.529
0.020000	-1.70	109.048	99.198	108.610	120.450	69.540	68.823	69.548	70.242
0.050000	-1.30	85.890	80.198	85.743	92.146	59.623	58.825	59.617	60.407
0.100000	-1.00	67.276	63.761	67.238	70.863	50.649	49.818	50.645	51.482
0.200000	-0.70	48.653	46.789	48.640	50.602	40.040	39.253	40.042	40.890
0.500000	-0.30	24.832	24.110	24.826	25.529	24.832	24.110	24.826	25.529
1.000000	0.00	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000

17
 18 A visual representation of the relationship between the precision of the OD-matrices
 19 derived at district level and the sampling rate is provided in Figures 5 and 6. Inspection of Figure
 20 5 reveals a pattern very similar to the one observed in Figure 3: the MAPE first increases when
 21 samples are becoming larger, and then starts to decrease. Recall that the increase in median
 22 MAPE for the smallest sampling rates can be accounted for by the fact that more cells are

1 seriously overestimated on average, while the maximum underestimations are bounded by
2 100%. By analogy with the results at municipality level, this effect is filtered out by using the
3 MCAPE, as could be noticed from Figure 6. Moreover, the relationship between the MCAPE
4 and sampling proportion is a concave decreasing function, similar to the relationship between the
5 MCAPE and sampling rate at municipality level.
6



7
8 **FIGURE 5 Relationship between MAPE and base-10 logarithm of the sampling rate.**



10
11 **FIGURE 6 Relationship between MCAPE and base-10 logarithm of the sampling rate.**

12

3.2.3 OD-Matrices at Provincial Level

In contrast to the true OD-matrices at municipality and district level, the true OD-matrix at provincial level (11 by 11) only comprises non-zero-cells. Examination of Table 4 reveals that at provincial level, barely any accurate OD-matrices can be derived. Nonetheless, in contrast to the results at municipality and district level, for the largest sample sizes acceptable results are obtained: sampling half of the population would result in an average absolute percentage error of 3.4 percent, and surveying one fifth of the population results in an average absolute percentage error of 7.4%. Results from Table 4, also confirm that predictions related with a more aggregate level are more precise. Notwithstanding, results at provincial level confirm the finding unraveled at the lower levels (municipality and district) that the direct derivation of origin-destination matrices from household activity/travel surveys should be avoided.

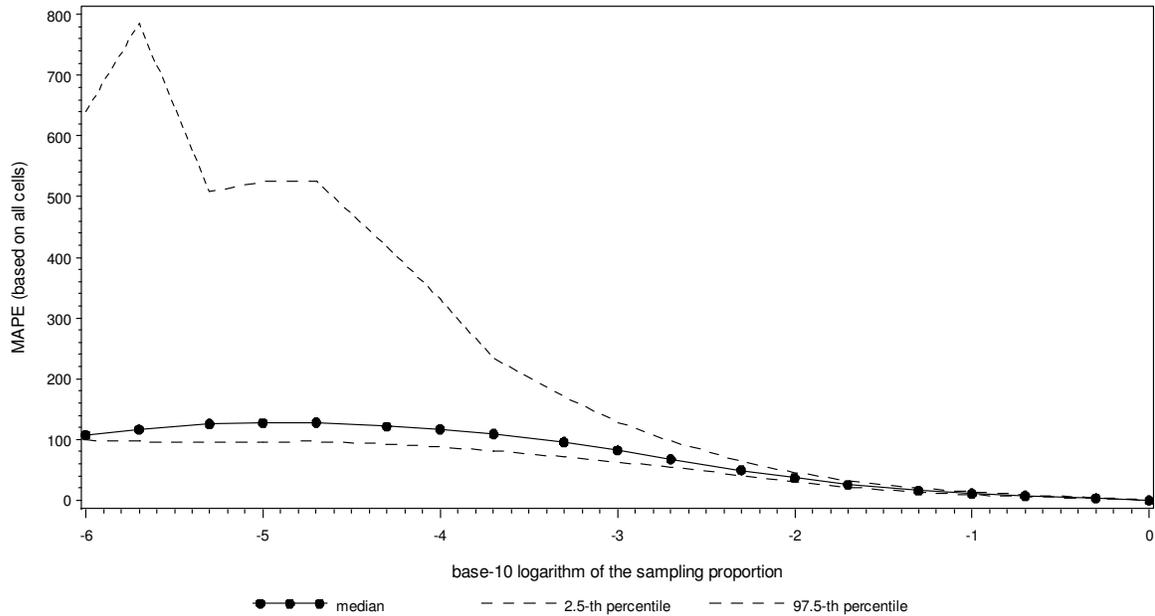
TABLE 4 MAPE and MCAPE for OD-Matrices Derived at Provincial Level

Sampling Rate (SR)	Log ₁₀ SR	MAPE				MCAPE			
		Mean	P5	Median	P95	Mean	P5	Median	P95
0.000001	-6.00	176.769	99.263	107.557	394.608	99.063	98.121	99.098	100.000
0.000002	-5.70	203.155	97.373	116.674	433.429	97.457	96.026	97.482	98.892
0.000050	-5.30	181.113	96.627	126.086	350.822	95.265	93.815	95.268	96.704
0.000010	-5.00	168.445	98.339	127.604	348.163	93.433	92.038	93.444	94.797
0.000020	-4.70	168.536	99.464	128.362	360.363	91.485	90.201	91.460	92.835
0.000050	-4.30	151.833	95.531	121.858	311.175	86.839	84.890	86.849	88.801
0.000100	-4.00	139.293	89.910	117.180	260.735	81.456	78.907	81.492	83.867
0.000200	-3.70	122.260	83.171	109.350	199.592	75.496	72.706	75.491	78.216
0.000500	-3.30	102.579	73.950	95.857	153.082	66.131	63.291	66.100	68.976
0.001000	-3.00	86.550	65.717	82.586	121.220	58.695	55.863	58.703	61.425
0.002000	-2.70	70.179	55.819	67.913	91.566	50.581	47.480	50.567	53.506
0.005000	-2.30	50.419	42.091	49.683	61.219	40.069	37.270	40.045	42.948
0.010000	-2.00	37.062	31.187	36.759	44.062	31.583	28.460	31.522	34.604
0.020000	-1.70	26.160	22.031	25.987	30.902	23.833	20.970	23.816	26.754
0.050000	-1.30	16.138	13.530	16.044	19.084	15.566	13.331	15.566	17.840
0.100000	-1.00	11.165	9.297	11.150	13.148	11.018	9.271	11.021	12.841
0.200000	-0.70	7.417	6.252	7.374	8.784	7.404	6.243	7.370	8.721
0.500000	-0.30	3.449	2.903	3.440	4.050	3.449	2.903	3.440	4.050
1.000000	0.00	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000

The visualization of the relationship between the precision of the OD-matrices derived at provincial level and the sampling proportion is shown in Figures 7 and 8. Analogous to the relationships between the sample rate and the precision of the OD-matrices at municipality and district level, the MAPE first increases when samples are becoming larger, and then starts to decrease (Figure 7). Again, the use of the MCAPE filters out this effect. In contrast to the results at municipality and district level, the relationship between the precision and the MCAPE reveals an s-shaped decreasing function: for the smallest sampling rates the relationship is concavely

1 decreasing, similar to the OD-matrices at municipality and district level; but for the larger
 2 sampling rates the relationship is a convex decreasing function. Moreover, the 95 percentile
 3 interval is much wider than for the OD-matrices at less aggregated levels. The most important
 4 reasons for this are the degree of sparseness and size of the matrix: for the less aggregate levels
 5 (municipality and district level), a lot of the variability of the precision is taken away by the
 6 large amounts of (zero-)cells.

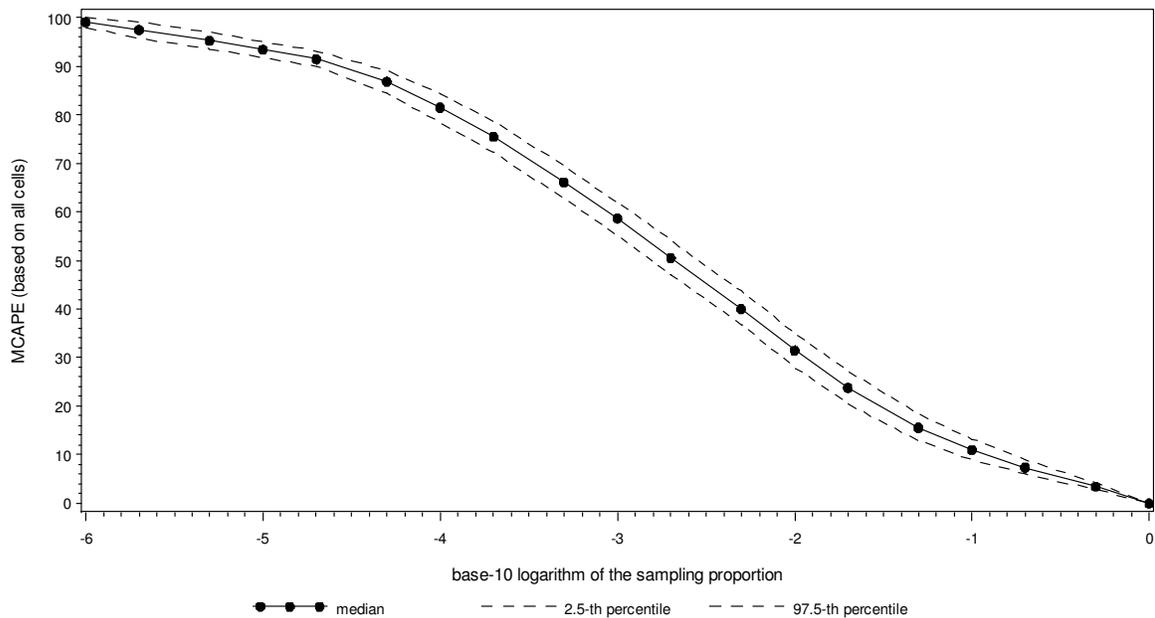
7



8

9 **FIGURE 7 Relationship between MAPE and base-10 logarithm of the sampling rate.**

10



11

12 **FIGURE 8 Relationship between MCAPE and base-10 logarithm of the sampling rate.**

4 CONCLUSIONS

In this paper, an assessment of the quality of origin-destination matrices derived from household activity/travel surveys was made. The results showed that no accurate OD-matrices can be directly derived from these surveys. Only when half of the population is queried, an acceptable OD-matrix is obtained at provincial level. Therefore, it is recommended to use additional information to better grasp the behavioral realism underlying destination choices. This is certainly a plea for travel demand models that incorporate the behavioral underpinnings of destination choices (activity location choices) given a certain origin. Moreover, matrix calibration techniques could seriously improve the quality of the matrices derived from these household activity/travel surveys. In addition, it is recommended to collect information about particular origin-destination pairs by means of vehicle intercept surveys rather than household activity/travel surveys, as these vehicle intercept surveys are tailored for collecting specific origin-destination data. Mark that the results presented in this paper do not negate the value of travel surveys as was shown in the example of deriving the commuting population, but indicate that sophistication is needed in the manner in which the data are employed.

A second important finding in this paper is that traditional methods to assess the comparability of two origin-destination matrices could be enhanced: the MCAPE index that was proposed has clear advantages over the traditional indices. The most important advantage being the fact that the MCAPE filters out the noise created by the asymmetry of the traditional criteria. Therefore, when dissimilarities between different OD-matrices are investigated, the use of the MCAPE index next to traditional criteria is highly recommended.

An important avenue for further research is the investigation of the relationship between the variability in the outcomes of travel demand models and underlying survey data. Triangulation of both travel demand modeling and small area estimation models could prove to a pathway for success. An empirical investigation of the effect of sampling proportions in household activity/travel surveys on final model outcomes would further illuminate the quest for optimal sample sizes. A thorough examination of the minimum required sampling rate of a household travel survey such that trip distribution models (e.g. a gravity model) could help fill in the full trip table certainly is an important step in further analyses. Model complexity and computability will certainly be key challenges in this pursuit.

5 ACKNOWLEDGEMENTS

The authors would like to thank Katrien Declerq for her advice on the implementation of the experiment.

6 REFERENCES

- (1) Haustein, S., and M. Hunecke. Reduced use of environmentally friendly modes of transportation caused by perceived mobility necessities: an extension of the theory of planned behavior. *Journal of Applied Social Psychology*, Vol. 37, No. 8, 2007, pp. 1856-1883.
- (2) Steg, L. Can public transport compete with the private car. *IATSS Research*, Vol. 27, No. 2, 2003, pp. 27-35.

- 1
2 (3) TRB Committee on Travel Survey Methods. *The On-Line Travel Survey Manual: A*
3 *Dynamic Document for Transportation Professionals*. Provided by the Members and
4 Friends of the Transportation Research Board's Travel Survey Methods Committee
5 (ABJ40), Washington, D.C., 2009. <http://trbtsm.wiki.zoho.com>. Accessed July 22, 2009.
6
- 7 (4) Stopher, P., R. Alsnih, C. Wilmot, C. Stecher, J. Pratt, J. Zmud, W. Mix, M. Freedman, K.
8 Axhausen, M. Lee-Gosselin, A. Pisarski, and W. Brög. *Standardized Procedures for*
9 *Personal Travel Surveys*. National Cooperative Highway Research Program Report 571.
10 Transportation Research Board, Washington, 2008.
11
- 12 (5) Cambridge Systematics. *Travel Survey Manual*. Prepared for U.S. Department of
13 Transportation and the U.S. Environmental Protection Agency. Travel Model Improvement
14 Program (TMIP), Washington, D.C., 1996.
15
- 16 (6) Tourangeau, R., M. Zimowski, and R. Ghadialy. An Introduction to Panel Surveys in
17 Transportation Studies. Report DOT-T-84, Prepared for the Federal Highway
18 Administration. National Opinion Research Center, Chicago, I.L., 1997.
19
- 20 (7) Stopher, P.R., and S.P. Greaves. Household travel surveys: Where are we going?
21 *Transportation Research Part A: Policy and Practice*, Vol. 41, No. 5, 2007, pp. 33-40.
22
- 23 (8) Rubinstein, R.Y. *Simulation and the Monte Carlo Method*. John Wiley and Sons, Inc., New
24 York, 1981.
25
- 26 (9) Fan, X., A. Felsövályi, S.A. Sivo, and S.C. Keenan. *SAS[®] for Monte Carlo Studies: A*
27 *Guide for Quantitative Researchers*. SAS Institute. Cary, N.C., 2000.
28
- 29 (10) Patel, A., and M. Thompson. Consideration and Characterization of Pavement
30 Construction Variability. In *Transportation Research Record: Journal of the*
31 *Transportation Research Board*, No. 1632, Transportation Research Board of the National
32 Academies, Washington, D.C., 1998, pp. 40-50.
33
- 34 (11) Awasthi, A., S.S. Chauhan, S.K. Goyal, and J.-M. Proth. Supplier selection problem for a
35 single manufacturing unit under stochastic demand. *International Journal of Production*
36 *Economics*, Vol. 117, No. 1, 2009, pp. 229-233.
37
- 38 (12) Groves, R.M., F.J. Fowler, M. Couper, J.M. Lepkowski, E. Singer, and R. Tourangeau.
39 *Survey Methodology*. John Wiley and Sons, Inc., Hoboken, N.J., 2004.
40
- 41 (13) Armstrong, J.S., and F. Collopy. Error measures for generalizing about forecasting
42 methods: empirical comparisons. *International Journal of Forecasting*, Vol. 8, No. 1,
43 1992, pp. 69-80.
44
- 45 (14) Makridakis, S. Accuracy measures: theoretical and practical concerns. *International*
46 *Journal of Forecasting*, Vol. 9, No. 4, 1993, pp. 527-529.

- 1
2 (15) Goodwin, P., and R. Lawton. On the asymmetry of the symmetric MAPE. *International*
3 *Journal of Forecasting*, Vol. 15, No. 4., 1999, pp. 405-408.
4
5 (16) Good, P.I. *Resampling Methods: A Practical Guide to Data Analysis, Third Edition*.
6 Birkhäuser, Boston, 2006.
7
8 (17) Groves, R.M. *Survey Errors and Survey Costs*. John Wiley and Sons, Inc., Hoboken, N.J.,
9 1989.
10
11 (18) Cools, M., E. Moons, T. Bellemans, D. Janssens, and G. Wets. Surveying activity-travel
12 behavior in Flanders: assessing the impact of the survey design. In Macharis, C., and L.
13 Turcksin (eds.), *Proceedings of the BIVEC-GIBET Transport Research Day 2009, Part II*.
14 VUBPRESS, Brussels, 2009, pp. 727-741.
15
16 (19) Pendyala, R.M., T. Yamamoto, and R. Kitamura. On the formulation of time-space prisms
17 to model constraints on personal activity-travel engagement. *Transportation*, Vol. 29, No.
18 1, 2002, pp. 73-94.
19
20 (20) Nakamya, J., E. Moons, S. Koelet, and G. Wets. Impact of Data Integration on Some
21 Important Travel Behavior Indicators. In *Transportation Research Record: Journal of the*
22 *Transportation Research Board*, No. 1993, Transportation Research Board of the National
23 Academies, Washington, D.C., 2007, pp. 89-94.
24
25 (21) Abrahamsson, T. *Estimation of Origin-Destination Matrices Using Traffic Counts: A*
26 *Literature Survey*. IIASA Interim Report IR-98-021/May, 1998.
27
28
29
30
31
32
33
34