

---

# Los Alamos National Laboratory

## TRANSIMS REPORT SERIES

Creating Synthetic  
Baseline Populations

Richard J. Beckman  
Keith A. Baggerly  
Michael D. McKay

The logo for the Travel Model Improvement Program (TMIP) features the letters 'TMIP' in a large, bold, italicized, outlined font. The letters are contained within a horizontal rectangular frame that has a slight 3D effect with a top and bottom bar.

---

**Travel  
Model  
Improvement  
Program**

Department of Transportation  
Federal Highway Administration  
Bureau of Transportation Statistics  
Federal Transit Administration  
Assistant Secretary for Policy Analysis

Environmental Protection Agency

---



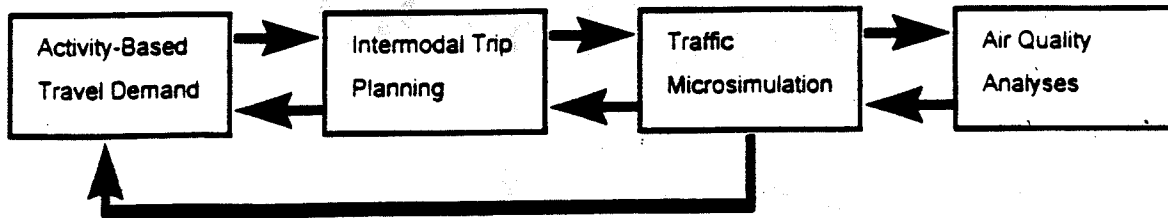
U.S. Department of  
Transportation



U.S. Environmental  
Protection Agency

# WHAT IS TRANSIMS?

TRANSIMS is a new system of travel forecasting models being developed by the Los Alamos National Laboratory for the Travel Model Improvement Program. The TRANSIMS models are a wholly new approach to travel forecasting, specifically designed to meet the needs of today's transportation decision makers for more accurate information on traffic congestion, differential impacts of transportation and motor vehicle emissions. TRANSIMS is composed of four basic modules.



*Activity-Based Travel Demand* estimates the number, characteristics and locations of activities in which individuals will participate during the forecast period. Activities are work, shopping, recreation, etc. These activity estimates are based on characteristics of individuals, their households and vehicles determined by a *synthetic population generator*.

*Intermodal Trip Planning* computes combined route and mode trip plans to accomplish individuals' desired activities. Intermediate activities such as shopping may occur during the routing of a principal trip such as work. TRANSIMS maintains the identities and characteristics of individual drivers, vehicles and other travelers throughout their trips. Trips are identified by specific geographic points of origin and destination.

*Traffic Microsimulation* computes the movement of persons, goods and vehicles on the simulated transportation network second by second during the forecast period. The microsimulation continuously computes the operating status of all vehicles and engines throughout the trips, including locations, speeds, acceleration or deceleration. Every motor vehicle in the study area is monitored in this way, thereby indicating areas and times of traffic congestion and emission concentrations.

*Air Quality Analyses* identify the kinds of emissions and calculate the effects of emissions on the atmosphere in the study area. The air quality module estimates the nature, amount and conditions of emissions by each motor vehicle. The output of the emissions modules are consistent with airshed models.

TRANSIMS is a considerable departure from traditional, four-step travel forecasting procedures. The new technical approaches in TRANSIMS permit equity analyses of transportation alternatives, service reliability and forecast uncertainty. Moreover, all of the traditional analyses conducted by the best current four-step models can be conducted with TRANSIMS.

# CREATING SYNTHETIC BASELINE POPULATIONS

by

Richard J. Beckman, Keith A. Baggerly, Michael D. McKay  
Statistics Group,

Los Alamos National Laboratory, Los Alamos, NM 87545, U.S.A.

## Abstract

To develop activity-based travel models individual travelers and households must be considered. Methods for creating baseline synthetic populations of households and persons using 1990 census data are given. Households are generated on either a census block group or census tract basis by judicious selection of households from the associated Public Use Microdata Sample (PUMS) of the Census Bureau. The procedures are validated by creating pseudo census tracts from PUMS samples and considering the joint distribution of the size of households and the number of vehicles in the households. It is shown that the joint distributions created by these methods do not differ substantially from the true values.

## 1. Introduction

Activity-based transportation models, such as those outlined and/or reviewed in Recker *et al.* (1986a, 1986b), Kitamura (1988), Axhausen and Garling (1991), Bhat and Koppelman, (1993), Garling, Kwan and Colledge (1994), Recker (1994) and Smith *et al.* (1995), require that individual travelers and households rather than aggregates be considered. The purpose of this paper is to outline a methodology for the creation of a synthetic baseline population of individuals and households which can be used in such models. While the populations developed here are necessary for activity-based models, aggregated demographic characteristics of these populations can also be used in the traditional four step process to estimate travel demand.

The techniques given here for the construction of synthetic populations rely on 1990 census data. They can, however, be modified as new census data becomes available. The census data includes the Census Standard Tape File 3A (STF-3A) (see Census, 1992a) and the Public Use Microdata Sample (PUMS) (see Census, 1992b). STF-3A is a collection of summary tables of demographics, such as the number of persons per household, for census tract or census block group sized areas and is often used in transportation studies. Most tables in STF-3A summarize one demographic characteristic, but a few cross-classified summary tables of two or three demographics are also given.

The PUMS is beginning to be utilized in transportation studies (for example see Purvis, 1994). The PUMS file consists of a 5% representative sample of almost complete census records (addresses and other unique identifiers are missing) from those contained in a collection of census tracts or other small geographic census areas, which collectively is called a Public Use Micro Area (PUMA). A PUMA is constructed so that it contains approximately 100,000 individuals. Since essentially complete demographics are given for each individual and household in the PUMS, entire multiway summary tables which are not available in STF-3A can be created for the PUMA area.

The basic algorithm for the construction of a synthetic population is as follows. Each census tract or part of a census tract which contributes to a given PUMA is considered. (Block groups can and possibly should be used, but without loss of generality here we talk exclusively about census tracts.) Summary tables for a selected set of demographics from STF-3A are assembled for each of the census tracts. The multiway demographic table of these demographics is constructed from the PUMS of the corresponding PUMA. Then, a multiway table is estimated for each census tract where the marginal totals in the constructed tables match the marginal totals given by STF-3A and the correlation structure in the multiway table constructed from the PUMS is maintained. In the last step of the algorithm households to makeup the synthetic population

are drawn from the PUMS in proportion to the estimated entries in the multiway tables for the census tracts.

The form of the data in STF-3A and the exact tables used for this procedure are given in the next section. Section 3 illustrates methodologies for the creation of the estimated multiway table for each census tract in a PUMA, while techniques for the creation of the synthetic population are given in the next Section. Methodologies for the validation and verification of the technique are given in Section 5. Discussion and Summary Sections follow.

## **2. Census Data**

Census tract summary tables in STF-3A used in the creation of synthetic households are divided into three categories: family households, nonfamily households, and group quarters. Family households are those households with two or more related persons. Persons living alone or unrelated persons living together are nonfamily households. Group quarters are dwellings such as prisons or college dormitories. Since the summary tables in STF-3A are different for each of the three types, the corresponding synthetic populations are generated separately. Family households are considered first. The summary tables in STF-3A which concern family households are :

1. P24: Age of the Householder,
2. P107: Family Income,
3. P112: Number of Workers in the Family,
4. P124A&B: Poverty Status (which is not used here) x Race x Family Type x Presence and Age of Children.

Not all categories that are given for the above STF-3A tables are used in the procedure. For example, the summary table P107 of family income has 25 categories. These are collapsed to 7 categories here, as the properties of the resulting populations using all 25 categories are almost identical with those using 7 categories. This does not preclude the use of all 25 categories in practice.

Examples of the four summary tables for census tract 1216.04 in Tarrant County, Texas are given in Table 1. The Family Class demographic given in Table 1 is derived from the Family Type and Age of Children categories given in summary table P124A&B. The 12 Family Classes are:

1. Married Couple: Children under age 5 only,
2. Married Couple: All children between 5 and 17,
3. Married Couple: Children under 5 and 5 to 17,
4. Married Couple : No children under 18,
5. Male Householder-No Wife Present: Children under age 5 only,
6. Male Householder-No Wife Present: All children between 5 and 17,
7. Male Householder-No Wife Present: Children under 5 and 5 to 17,
8. Male Householder-No Wife Present: No children under 18,
9. Female Householder-No Husband Present: Children under age 5 only,
10. Female Householder-No Husband Present: All children between 5 and 17,
11. Female Householder-No Husband Present: Children under 5 and 5 to 17,
12. Female Householder-No Husband Present: No children under 18.

The poverty level in summary table P124A&B is not considered. Data in the categories of "below the poverty level" and "above the poverty level" were summed, yielding the resulting race by family class summary table given in Table 1.

Summary tables in STF-3A for nonfamily households are:

1. P17: Household Type and Gender,
2. P20: Race x Household Type x Presence and Age of Children. (The race of nonfamily householders is derived from this table and is the only demographic used from this table.),
3. P24: Age of Nonfamily Householder,
4. P110: Nonfamily Household Income,
5. P127: Poverty Status (not used here) x Age of Householder x Household Type.

These tables for census tract 1216.04 in Tarrant County, Texas are shown in Table 2. Once again some of the categories shown in Table 2 are collapsed from the full set of categories given in STF-3A. However, the three categories of age shown for the P127 summary table are exactly those given in STF-3A.

There are only two summary tables for group quarters in STF-3A. These are:

1. P40: Group Quarters,
2. P41: Group Quarters x Age.

These two tables for census tract 1216.04 of Tarrant County, Texas are given in Table 3.

The second major source of census data used here is the PUMS. The samples given in the PUMS are a representative 5% sample of households and group quarters from the PUMA. Weights are assigned to each household and person in the sample so that weighted summary statistics can be formed.

In the technique presented here, weighted multiway summary tables corresponding to the demographic variables and categories in Tables 1, 2 and 3 are formed from the PUMS for family households, nonfamily households and group quarters for each PUMA in the metropolitan area under study. The tables are constructed from the PUMS by adding the household weights for households in each category of the multiway table. For group quarters, the person weights are summed for each category.

### 3. Estimating the Cross-Classified Table

In this section a method for the construction of cross-classified tables of demographics for each census tract in the area of study is discussed. These multiway tables are constructed to satisfy the marginal summaries of STF-3A. Additionally, the estimate of the correlation structure of the census tract multiway table as given by the PUMS is maintained. For all of the procedures developed here, it is assumed that all of the census tracts and census places that contribute to a PUMA have the same correlation structure.

Correlation in a multiway table is measured by odds ratios. For a 2 by 2 table

P <sub>1,1</sub>	P <sub>1,2</sub>
P <sub>2,1</sub>	P <sub>2,2</sub>

the odds ratio is

$$\phi = \frac{P_{1,1} P_{2,2}}{P_{1,2} P_{2,1}}$$

Odds ratios for multiway tables or  $n \times m$  two-way tables are given by a simple extension of the formula given above for the  $2 \times 2$  table. For an  $n_1 \times n_2 \times \dots \times n_m$  table the odds ratios have the general form

$$\phi = \frac{(P_{i_1, i_2, \dots, i_m})(P_{i_1, \dots, i_j + c_1, \dots, i_k + c_2, \dots, i_m})}{(P_{i_1, \dots, i_j + c_1, \dots, i_k, \dots, i_m})(P_{i_1, \dots, i_j, \dots, i_k + c_2, \dots, i_m})}$$

The primary tool used to complete the multiway table for each census tract is iterative proportional fitting (IPF) (Deming and Stephan, 1940). It can be shown that in situations where the marginal totals of a multiway table are known and a sample from the population which generated these totals is provided IPF gives a constrained maximum entropy estimate of the true proportions in the population multiway table (Ireland and Kullback, 1968). Additionally, IPF estimates maintain the same odds ratios as those in the sample table in the absence of any marginal information to the contrary (see for example Little and Wu, 1991).

To describe IPF the following notation is used. Let the proportion of observations in a sample from a  $m$ -way marginal table from the PUMS be denoted by

$$P_{i_1, i_2, \dots, i_m} = \frac{n_{i_1, i_2, \dots, i_m}}{n}$$

where  $i_j = 1, 2, \dots, n_j$  represents the observed value of the  $j^{\text{th}}$  demographic with  $n_j$  categories (for example age of the householder has  $n_j=7$  categories),  $n$  is the total number of observations in the table and  $n_{i_1, i_2, \dots, i_m}$  is the number of counts in cell  $(i_1, i_2, \dots, i_m)$ . Also, let  $T_k^{(j)}$  be the known marginal totals for the  $k^{\text{th}}$  category of the  $j^{\text{th}}$  demographic. The total number  $n$  is

$$n = \sum_{k=1}^{n_j} T_k^{(j)} \text{ for all } j.$$

Let  $p_{i_1, i_2, \dots, i_m}^{(t)}$  denote the estimated proportions in cell  $(i_1, i_2, \dots, i_m)$  at iteration  $t$  of the IPF procedure and let

$$p_{i_1, \dots, i_j = k, \dots, i_m}^{(t)} = \sum_{i_1=1}^{n_1} \dots \sum_{i_m=1}^{n_m} p_{i_1, i_2, \dots, i_j = k, \dots, i_m}^{(t)}$$



where the above sums are not over the index corresponding to the fixed category  $k$ .

Iterative proportional fitting begins by letting

$$p_{i_1, i_2, \dots, i_m}^{(0)} = p_{i_1, i_2, \dots, i_m}$$

At iteration  $t$  the estimated proportions  $p_{i_1, i_2, \dots, i_m}^{(t)}$  are derived using the following procedure. For each margin in turn, update the estimated proportions (for all values of  $i_1, i_2$  etc., and the  $k^{\text{th}}$  category of the  $j^{\text{th}}$  marginal) by

$$p_{i_1, i_2, \dots, i_j=k, \dots, i_m}^{(\text{new})} = p_{i_1, i_2, \dots, i_j=k, \dots, i_m}^{(\text{old})} \frac{T_k^{(j)}/n}{P_{\dots, \dots, i_j=k, \dots}^{(\text{old})}}$$

where for the first marginal for  $p^{(\text{old})}$  corresponds to  $p^{(t-1)}$ , the resulting estimates from the last iteration. For the second and later marginals  $p^{(\text{old})}$  is set equal to the  $p^{(\text{new})}$  estimated for the previous marginal. Finally,  $p^{(t)}$  is set equal to  $p^{(\text{new})}$  resulting from the last marginal. The iterations continue until the relative change between iterations in each estimated  $p_{i_1, i_2, \dots, i_m}^{(t)}$  is small. In practice we find this procedure converges in 10 to 20 iterations.

Minor adjustments must be made to the IPF routine in order to handle marginal tables in the forms given in Tables 1, 2 and 3. Two way marginals such as the race  $\times$  family class table in Table 1 present no problem to the IPF routine. Such marginal tables are converted to a single demographic whose categories are all the combinations of two demographics involved and IPF proceeds as usual. For example the  $5 \times 12$  table for race  $\times$  family class in Table 1 is considered as a table of one demographic variable with 60 categories. If two marginal tables contain a common demographic variable (e.g. the alone-not alone demographic in tables P17 and P127 in Table 2) the procedure is not altered. The fitting proceeds as above treating each marginal separately. For the case where one demographic variable is in two summary tables and has fewer categories in one than the other (e.g. summary tables P24 and P127 shown in Table 2) an additional step is required. The procedure uses only one marginal table at a time.

When the table with the “collapsed” marginal is considered, the procedure updates the cells as above where all of the cells that contribute to the individual “collapsed” categories are updated by the same proportion.

As stated before it is assumed that each of the census tracts contained in a PUMA has the same set of odds ratios. Without this or similar assumptions the estimation of odds ratios for individual tables is intractable. Even with this assumption, however, the PUMS sample does not reflect the correlation structure of the individual tables. This is easily seen by considering two  $2 \times 2$  tables with sample sizes  $n_1$  and  $n_2$  and the same odds ratio  $\phi$ . Let the proportions in the cells of the two tables be

Table-1	
$P_{1,1}^{(1)}$	$P_{1,2}^{(1)}$
$P_{2,1}^{(1)}$	$P_{2,2}^{(1)}$

Table-2	
$P_{1,1}^{(2)}$	$P_{1,2}^{(2)}$
$P_{2,1}^{(2)}$	$P_{2,2}^{(2)}$

With the assumption of equal odds ratios in the tables we have

$$\phi = \frac{P_{1,1}^{(1)} P_{2,2}^{(1)}}{P_{1,2}^{(1)} P_{2,1}^{(1)}} = \frac{P_{1,1}^{(2)} P_{2,2}^{(2)}}{P_{1,2}^{(2)} P_{2,1}^{(2)}}$$

Now, the odds ratio for the combined tables and thus for the sample is

$$\phi' = \frac{P_{1,1} P_{2,2}}{P_{1,2} P_{2,1}}$$

where

$$P_{i,j} = n_1 P_{i,j}^{(1)} + n_2 P_{i,j}^{(2)}$$

From the above equations it is easy to see that in general  $\phi' \neq \phi$ .

The fact that  $\phi' \neq \phi$  implies that correct statistical techniques for the estimation of  $\phi$  must simultaneously take into account the data from all of the census tracts or parts of census tracts that contribute to the PUMS.

The following two step IPF procedure considers all census tracts and parts of tracts that contribute to the PUMS. It results in estimated multiway tables for which the odds ratios in each table are identical, and when the tables are combined the odds ratios are the odds ratios of the PUMS. First, the marginal tables for all  $n$  census tracts in the PUMS are added. Then an  $m$ -dimensional multiway table denoted by  $T_0$  is obtained by IPF of the PUMS against the summed marginals. The second step may be viewed as the construction of an  $(m+1)$ -dimensional table. The first  $m$  dimensions of the table are the  $m$  marginals. The last dimension is created by "stacking" the  $n$  tables for the census tracts. IPF is used a second time where the known marginal tables are the combined marginals for the  $n$  tables and  $T_0$  is taken as an additional marginal table for the "stacked" tables. The final estimated tables are obtained by IPF of a  $(m+1)$ -dimensional table with entries of 1 against these marginal tables. The following example of two  $2 \times 2$  tables illustrates this procedure.

Suppose that there are two  $2 \times 2$  tables with the following marginals.

TABLE 1			
	v2=1	v2=2	Total
v1=1	?	?	1700
v1=2	?	?	1050
Total	1505	1245	2750

TABLE 2			
	v2=1	v2=2	Total
v1=1	?	?	1405
v1=2	?	?	905
Total	700	1610	2310

The corresponding total marginal tables and the "PUMS" are :

TABLE 1 + TABLE 2			
	v2=1	v2=2	Total
v1=1	?	?	3105

"PUMS"			
	v2=1	v2=2	Total
v1=1	45	108	153

v1=2	?	?	1955
Total	2205	2855	5060

v1=2	63	37	100
Total	108	145	253

IPF of the "PUMS" against the TABLE 1 + TABLE 2 marginal totals gives table T0:

T0			
	v2=1	v2=2	Total
v1=1	949	2156	3105
v1=2	1256	699	1955
Total	2205	2855	5060

T0 is used as one of the marginal tables for a 2x2x2 table constructed by "stacking" the two original tables on top of one another. Two additional marginal tables are those given by the existing marginals of the two tables. These are:

Marginal: v1 x Tables		
	Table 1	Table 2
v1=1	1700	1405
v1=2	1050	905

Marginal: v2 x Tables		
	Table 1	Table 2
v2=1	1505	700
v2=2	1245	1610

To obtain the final estimated tables, the three 2x2 marginal tables are used as marginals in the IPF of a 2x2x2 table with cell entries equal to 1. The resulting IPF estimates are:

TABLE 1: Final Estimate			
	v2=1	v2=2	Total
v1=1	701	999	1700
v1=2	804	246	1050
Total	1505	1245	2750

TABLE 2: Final Estimate			
	v2=1	v2=2	Total
v1=1	248	1157	1405
v1=2	452	453	905
Total	700	1610	2310

It is interesting to note that each of the above estimated tables has an odds ratio of .21; the odds ratio of these two tables combined is .24, which is also the odds ratio of the "PUMS". If the two tables are fitted individually to the PUMS by IPF, the odds ratio of each table is the

same as the "PUMS" odds ratio, .24. The odds ratio of the combined tables is then .28 which is not the odds ratio of the "PUMS".

#### 4. Creating the Synthetic Population

The synthetic population of households is constructed by selecting entire households from the PUMS in proportion to the estimated probabilities given in the multiway table obtained by the technique in Section 3. It is not required that the selection procedure of Section 3 be used to estimate the multiway table; any estimation scheme may be used. However, it is the only routine considered here.

The number of households to be generated of each demographic type (having a specific set of demographics) is determined for each census tract. These numbers can be obtained either by multiplying the total number of households by the probabilities in the estimated multiway table, or by drawing the numbers at random according to these probabilities. The first case, which we call "in expectation", requires the addition or deletion of a few households due to rounding.

Once the number of households of each demographic type to be selected is determined, households with different demographics are considered separately. For a combination of demographic characteristics a set of probabilities is assigned to each household in the PUMS (after it has been split into family and nonfamily households and group quarters) where PUMS samples "close" to the combination of desired demographic characteristics are assigned higher probabilities. These probabilities are computed by considering the "distance" between a PUMS household, indicated by  $p$ , and the households characterized by a cell,  $c$ , in the multiway table for census tract. The following function is used to calculate the probabilities.

$$D(p, c) = w_p \prod_{i \in J} \left( 1 - |(d_i^p - d_i^c)/r_i|^k \right) \cdot \prod_{i \notin J} (1 - \Delta(d_i^p, d_i^c))$$

where for family households

1.  $J$  is the set of ordinal variables such as {income, age, workers} for family households and  $\sim J$  is the set of categorical variables such as {family class, race} for the family households.
2.  $d_i^p$  is the value of the  $i^{\text{th}}$  demographic for household  $p$  from the PUMS,
3.  $d_i^c$  is the value of the  $i^{\text{th}}$  demographic of the household of cell type  $c$  from the census tract,
4.  $r_i$  is the range of demographic  $i$  in the PUMS,
5.  $w_p$  is the weight associated with household  $p$  from the PUMS,
6.  $\Delta(d_i^p, d_i^c) = \begin{cases} \alpha & d_i^p = d_i^c \\ 1 - \alpha & d_i^p \neq d_i^c \end{cases}$

When  $\alpha = 0$  and  $k \rightarrow 0$ ,  $D(p,c)$  is a 0-1 loss function. In this case, samples from the PUMS are considered only if there is an exact match in the demographics  $d_i^p$  and  $d_i^c$ .

To acquire households for the synthetic population, the probability of selecting household  $p$  for the synthetic population for a household with demographics  $c$  is given by

$$\Pr\{\text{Selecting Household } p\} = D(p, c) / \sum_j D(j, c).$$

Using the procedure given in Section 3, there is always at least one sample in the PUMS which exactly matches the table's demographics. Therefore, the 0-1 loss function of step 7 above may be used to select PUMS households for the synthetic population. Other procedures for estimating the multiway table, such as assuming independence between the marginals and multiplying marginal probabilities to create the table, may not have this property. We show in the next section that for the method of Section 3, a 0-1 loss function is desired.

Since there is randomness in the selection of the households by the above method, multiple populations are constructed for each study. These multiple populations may then be used to investigate the uncertainty in the results of the study which is due to the construction of the synthetic populations.

## 5. Validation Studies

One method of validating the scheme for creating a synthetic population compares demographic characteristics of the synthetic populations with those of the true population using

variables not involved in the generation of the population. For example, to judge the bias and the variance of the generation technique, the distribution of the number of people per household as summarized in STF-3A can be compared with its distribution in the synthetic population. This variable is available in STF-3A for all households and is not used in the generation of the synthetic family and nonfamily households.

Utilizing a 0-1 loss function, 100 synthetic populations of family and nonfamily households for census tract 1216.04 of Tarrant County, Texas were created. Figure 1 shows the differences in the true distribution of household size for this census tract and those distributions in the synthetic populations. Each line in the figure represents the difference between one of the synthetic population and the known values for the census tract. The box in the lower right hand of the figure shows the actual distribution of households and the average number of synthetic households for each household size.

The largest discrepancy between the synthetic populations and the actual population is in the number of households with 5 or more persons where on the average 14 (with a base of 157) too many households of this size are generated. In each population of synthetic households exactly 464 households with exactly one person were produced. The number of generated one person households is always exactly correct with a 0-1 loss function due to the fact that the marginal totals in Table 2 for P17 and P127 are exactly satisfied. The average total population for the synthetic population is 5628 persons, while the true number of persons residing in the census tract is 5592 persons. The difference of only 36 persons is remarkable considering that the total number of persons is not controlled by the procedure to create the synthetic population.

In a second validation scheme a validation population is constructed from the existing PUMS data. The samples in 20 to 30 PUMS data sets are considered as the full set of demographics of a "census tract". A "super PUMS" is constructed as a 5% random sample from this collection. Synthetic populations are generated for each of the validation "census tracts". Since the entire

set of demographics is known for the validation population of "census tracts", comparisons of the synthetic population with these demographics can be made.

In the study here 22 PUMSs were selected from the San Francisco Bay Area to serve as pseudo census tracts. A "super PUMS" was created by taking a 5% sample from the collection of 22 PUMS data sets. The average number of persons and vehicles per household was computed for each synthetic population using a 0-1 loss function. The results are given in Figure 2. For each of the 22 pseudo census tracts, the absolute differences between the proportion of synthetic households of sizes 1 to 5+ and with 0 to 3+ vehicles and the true proportions in the population are highlighted. The darker areas on the figures show the most discordant regions. It is seen in this figure that the procedure does well in the prediction of this joint distribution. This most discrepant results are in the first three "census tracts" where the number of vehicles for small household sizes are overestimated. In particular, the proportion of single person households with no vehicles is greatly underestimated. These three "census tracts" were known to be different from the rest as they are the PUMAs from the City of San Francisco, where due to parking difficulties and the availability of a good mass transit system the number of vehicles is lower than in the general population. They were added to the validation set to see if a few heterogeneous census tracts in a PUMA would invalidate the procedure. The results shown in Figure 1 demonstrate the robustness of the procedure for the majority (and most homogeneous) "census tracts". In practice, however, the true census tracts from a PUMA are more homogeneous than those constructed in the validation study by combining PUMAs, and the resulting estimated populations will be closer to the true values for all census tracts.

One potential methodology for the construction of a synthetic population is to forgo the fitting of the multiway table with IPF and draw the population directly from the PUMS according to the number of households given in STF-3A for block groups or census tracts. The improvement by using IPF for both family and nonfamily households is shown in Figure 3. In this figure the



mean absolute deviation between the proportions of household sizes by the number of vehicles estimated using IPF for the 22 "census tracts" is plotted against the mean absolute deviation of the proportions in the "census tracts" and the true proportions in the "super PUMA". The latter is equivalent to the selection of households directly from the PUMS without IPF. Most points on the two plots in Figure 3 are above the line. These show the "census tracts" where the IPF routine does a better job in the prediction of the joint distribution of household size and the number of vehicles.

## 6. Discussion

We now briefly discuss two additional facets of the procedure. First, the multiway table generated from the PUMS is sparse. There are 11,760 cells in the multiway table for families (obtained by multiplication of the number of categories for each of the marginals) while there are only 609 cells with nonzero entries in the multiway table constructed from the PUMS. The IPF algorithm estimates a zero proportion for all cells that are zero in the sample. Since the PUMS is a sample, many of the approximately 11,000 empty cells might not be empty in the population. Therefore, one may wish to "tweak" each of the empty cells with a partial count of .1 or .01 before using the IPF routine.

Second, given the estimated multiway table from the IPF routine, the potential exists for the imputation of an additional demographic from the PUMS. For example, one demographic that could be imputed is the number of people in the household which is not one of the demographic characteristics used to construct the multiway tables for either family or nonfamily households.

Both imputation and tweaking were investigated and neither improved on the results obtained by using the basic algorithms given above. Synthetic family populations for census tract 1.07 of Bernalillo County, New Mexico were generated with and without both tweaking and imputing the persons in the household. Zero cells in the PUMS were replaced by .01 and imputation of the

household size was accomplished using Classification And Regression Trees (CART) (Breiman, et. al., 1984). Both of these investigations required the use of a non 0-1 loss function since in these cases there is no guarantee of PUMS samples matching exactly the nonzero cells of the multiway table generated by IPF. The values of  $k$  and  $\alpha$  in  $D(p,c)$  of Section 4 were set to .1 and .05 respectively.

Tweaking badly biases the statistics for the marginal tables and is not recommended. Imputation does little more than add variability to the synthetic populations. It does not improve the results from the basic algorithm. These results are not surprising. Tweaking places an entry in every cell of the estimated multiway table. Some of these cells may be much different from the closest samples in the PUMS. Therefore, selection of even the closest PUMS samples using  $D(p,c)$  will tend to bias the results. On the other hand the household size was imputed from the PUMS and the closest cells are not too distant from the nonzero cells in the multiway table. Hence, in this case no appreciable bias is added to the sample.

The use of a non 0-1 loss function was also investigated for populations generated with no imputation or tweaking. The results are similar to the imputation results where increased variability but no bias was noticed.

To maintain the proper odds ratio structure, the statistically correct procedure for IPF as given in Section 3 is to create an  $(m+1)$  dimensional table, with the table number as the additional variable is the table number and to fit the  $(m+1)$  tables simultaneously. Investigations show that this procedure only marginally outperforms the simpler procedure of fitting the multiway tables individually. Consequently, the simpler procedure can probably be used without much harm.

One could assume that the variables that make up the multiway table are independent. The marginal proportions for these variables could be multiplied and the synthetic households could be drawn from the PUMS using  $D(p,c)$ . However, this procedure has the same problem as "tweaking" in that it badly biases the marginal distributions of the individual variables, not to

mention the joint distributions of variables such as household size and the number of vehicles. Assuming independence between the marginal variables is not recommended.

## 7. Summary

A method has been given for the generation of synthetic populations on a census tract or block group basis. The technique uses only census data and it reproduces the existing population in a reasonable way. There are two steps in the methodology. First a multiway demographic table of proportions is estimated. Here, it is estimated using iterative proportional fitting. However, any reasonable statistical method for this estimation would be acceptable. Secondly, a synthetic population of households is drawn from the PUMS so as to match the proportions in the estimated table.

We have shown by validation that synthetic populations generated by this procedure have desirable characteristics. In the synthetic population the marginal distribution of variables used in the construction of the multiway table match the truth exactly (within rounding). The distribution of variables not used in the construction of the multiway table, such as household size, are reasonably estimated as evidenced in Figure 1. Also, the joint distributions of two or more variables not used in the multiway table construction, such as household size and the number of vehicles, are estimated well by this procedure.

## 8. References

- Axhausen, K. W., and Garling, T. (1991) *Activity-based approaches to travel analysis: conceptual frameworks, models and research problems*. U.S. Department of Commerce, National Technical Information Service, TSU Ref:628.
- Bhat, C. R., and Koppelman, F. S. (1993) A conceptual framework of individual activity program generation. *Transportation Research A*, 27A, 433–446.
- Breiman, L., Friedman, J. H., Olshen, R. A. and Stone, C. J. (1984) *Classification and Regression Trees*, Wadsworth and Brooks: Monterey, CA.
- Census (1992a) Census of Population and Housing, 1990; Summary Tape File 3 on CD-ROM Technical Documentation / prepared by the Bureau of the Census. —Washington: The Bureau, 1992.
- Census (1992b) Census of Population and Housing, 1990; Public Use Microdata Sample U.S. Technical Documentation / prepared by the Bureau of the Census. —Washington: The Bureau, 1992.
- Deming, W. E. and Stephan, F. F. (1940) On a least squares adjustment of a sampled frequency table when the expected marginal tables are known. *The Annals of Mathematical Statistics*, 11, 427–444.
- Garling, T., Kwan, M. and Colledge, R. G. (1994) Computational-Process modelling of household activity scheduling. *Transportation Research B*, 28B, 355–364.
- Ireland, C. T. and Kullback, S. (1968) Contingency tables with given marginals. *Biometrika*, 55, 179–188.
- Kitamura, R. (1988) An evaluation of activity-based travel analysis. *Transportation*, 15, 9–34.
- Little, R. J. A. and Wu, M. M. (1991) Models for contingency tables with known marginals when target and sampled populations differ. *Journal of the American Statistical Association*, 86, 87–95.
- Purvis, C. L. (1994) Using 1990 Census Public Use Microdata Sample to estimate demographic and automobile ownership models. *Transportation Research Record*, 1443, 21–29.
- Recker, W. W. (1994) A household activity pattern problem: general formulation and solution. *Transportation Research B*, 29B, 61–77.
- Recker, W. W., McNally, M. G. and Root, G. S. (1986a) A model of complex travel behavior: Part I—theoretical development. *Transportation Research A*, 20A, 307–318.
- Recker, W. W., McNally, M. G. and Root, G. S. (1986b) A model of complex travel behavior: Part II—an operational model. *Transportation Research A*, 20A, 319–330.

Smith, L., Beckman, R., Anson, D., Nagel, K. and Williams, M. (1995) *TRANSIMS: Transportation ANalysis and SIMulation System*. Los Alamos National Laboratory Unclassified Report, LA-UR-95-1664, Los Alamos, NM 87544.

**TABLE 1.**

**Family Summary Statistics For Census Tract 1216.04 Tarrant County, TX**

<b>P24: AGE OF HOUSEHOLDER</b>							
Age	15-24	25-34	35-44	45-54	55-64	65-74	>74
n	100	445	382	283	164	78	39

<b>P107: FAMILY INCOME</b>							
Income	<\$10K	\$10-15K	\$15-25K	\$25-35K	\$35-50K	\$50-100K	>\$100K
n	147	117	216	324	371	267	49

<b>P112: WORKERS IN FAMILY</b>				
Workers	0	1	2	>2
n	89	489	792	121

<b>P124A&amp;B FAMILY CLASS x RACE</b>					
	White	Black	A. Indian	Asian	Other
(1) Couple Child<5	73	7	0	0	14
(2) Couple Child 5-17	276	23	0	6	18
(3) Couple Child'n <5 and 5-17	150	0	0	0	0
(4) Couple No Child'n <18	533	0	0	0	0
(5) Male HH Child<5	26	0	0	0	0
(6) Male HH Child 5-17	0	15	0	0	0
(7) M. HH Child'n <5 and 5-17	0	0	0	0	0
(8) Male HH No Child'n <18	19	8	0	11	0
(9) Fem. HH Child<5	28	13	0	0	0
(10) Fem. HH Child 5-17	119	45	0	0	11
(11) F. HH Child'n <5 and 5-17	11	0	0	0	0
(12) Fem. HH No Child'n <18	85	0	0	0	0

TABLE 2.

Nonfamily Summary Statistics for Census Tract 1216.04 Tarrant County, TX

<b>P17: HOUSEHOLD TYPE x GENDER</b>		
Type/Gender	Male	Female
Living Alone	243	403
Not Living Alone	120	61

<b>P20: RACE OF HOUSEHOLDER</b>					
Race	White	Black	Am. Indian	Asian	Other
n	793	34	0	0	0

<b>P24: AGE OF HOUSEHOLDER</b>							
Age	15-24	25-34	35-44	45-54	55-64	65-74	>74
n	139	146	105	97	123	74	143

<b>P110: HOUSEHOLD INCOME</b>							
Income	<\$10K	\$10-15K	\$15-25K	\$25-35K	\$35-50K	\$50-100K	>\$100K
n	250	69	184	140	75	101	8

<b>P127: AGE x HOUSEHOLD TYPE</b>			
Type/Age	15-64	65-74	>74
Living Alone	439	64	143
Not Living Alone	171	10	0

TABLE 3

Group Quarters Summary Statistics Census Tract 1216.04 Tarrant County, TX

P40: PERSONS IN GROUP QUARTERS		
Institutionalized:		
	Correctional Institutions	0
	Nursing Homes	211
	Mental Hospitals	0
	Juvenile Institutions	0
	Other	0
Other:		
	College Dormitories	0
	Military Quarters	0
	Emergency Shelters	0
	Visible in Street	0
	Other	0

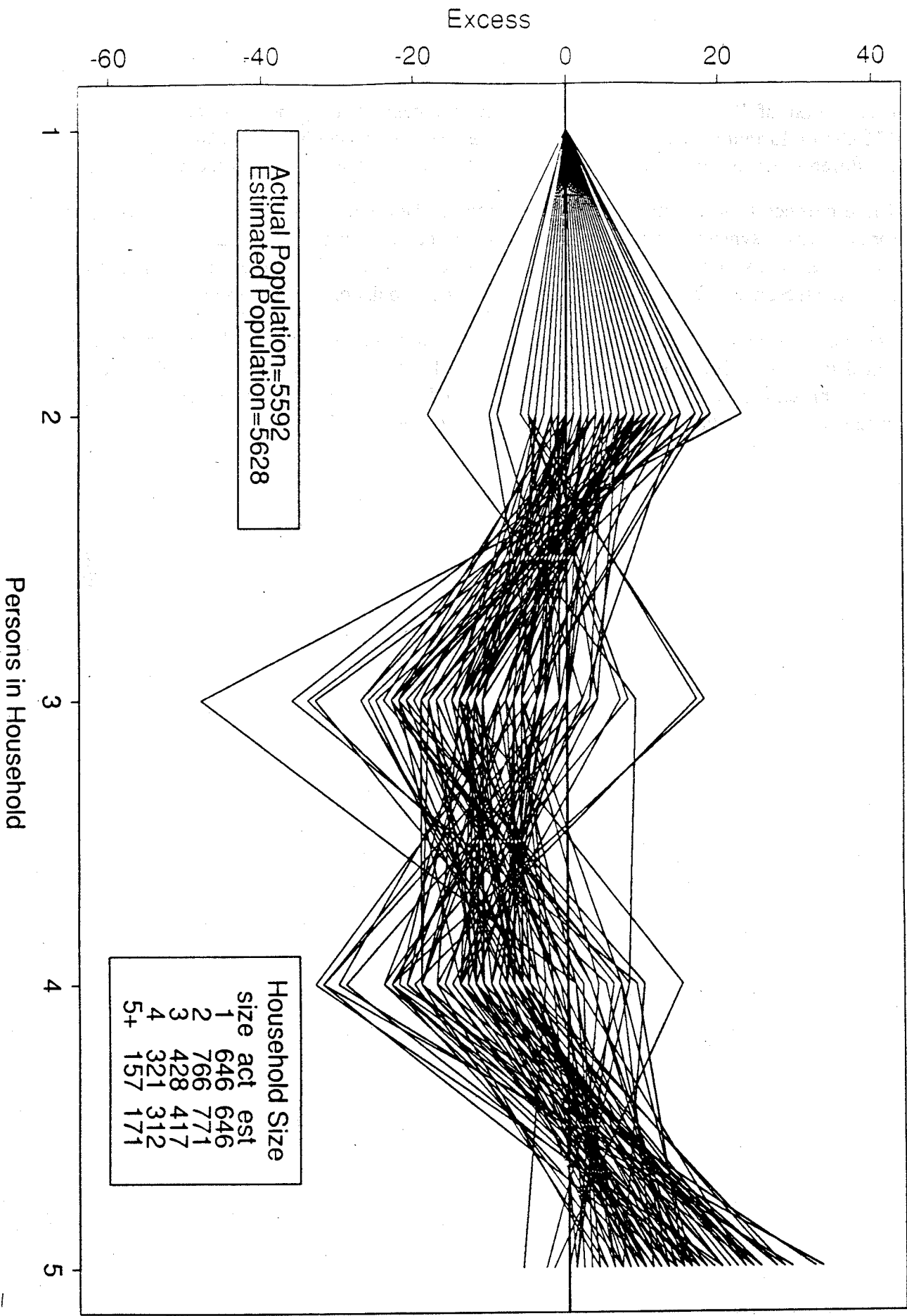
P41: GROUP QUARTERS x AGE			
Type/Age	<18	18-64	>64
Institutionalized:	0	0	211
Other:	0	0	0



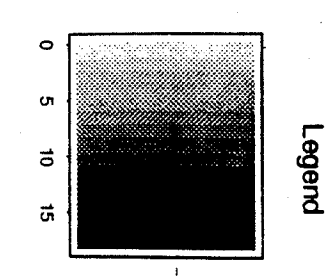
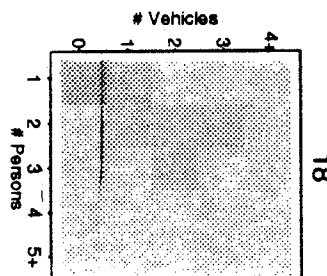
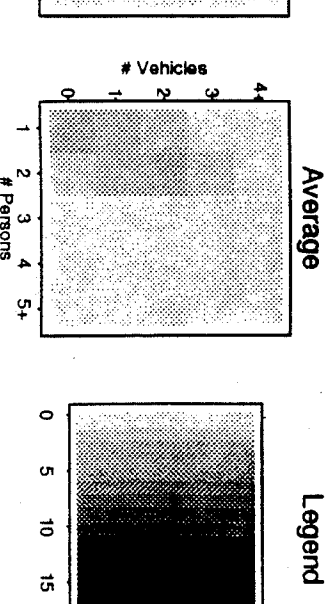
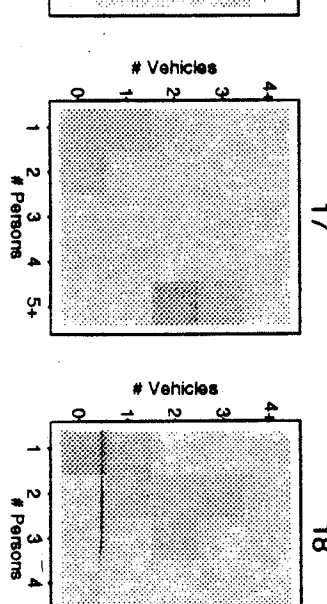
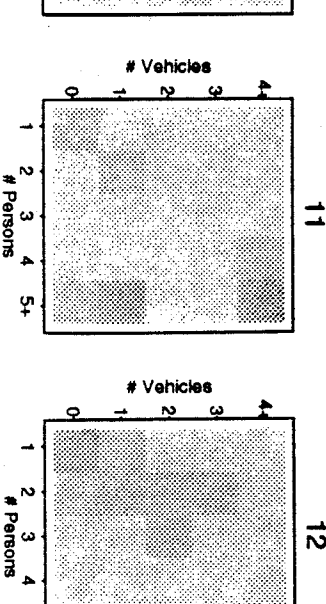
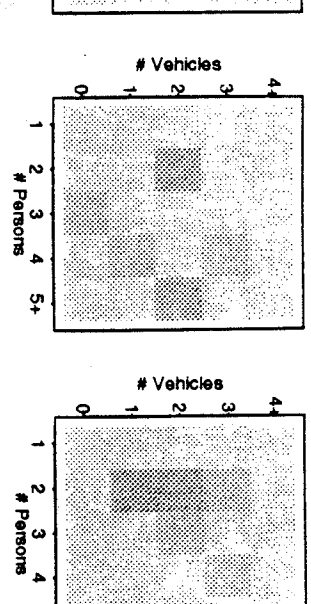
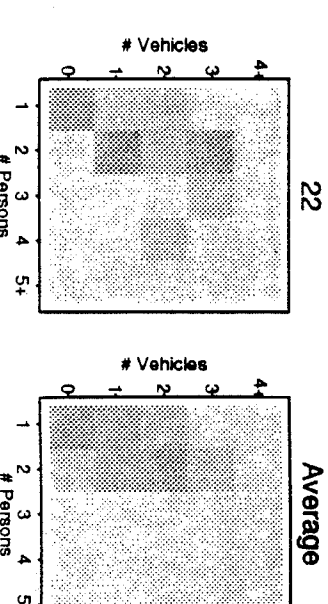
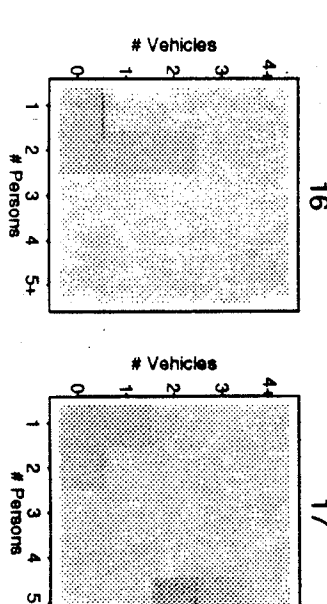
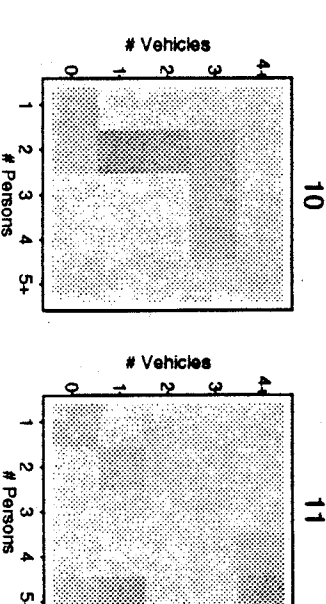
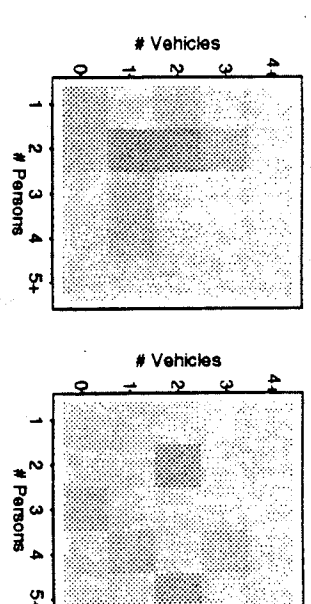
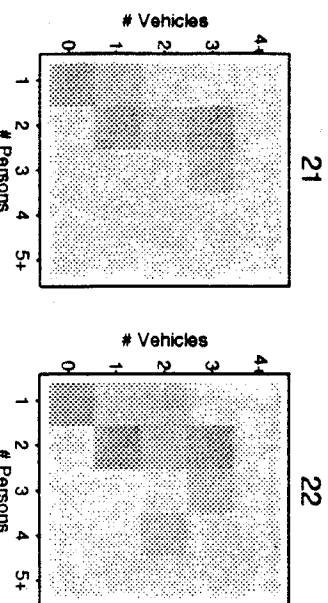
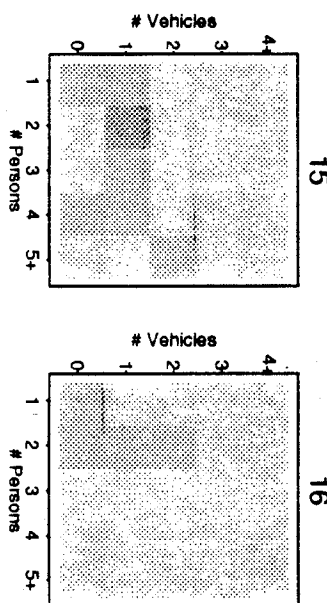
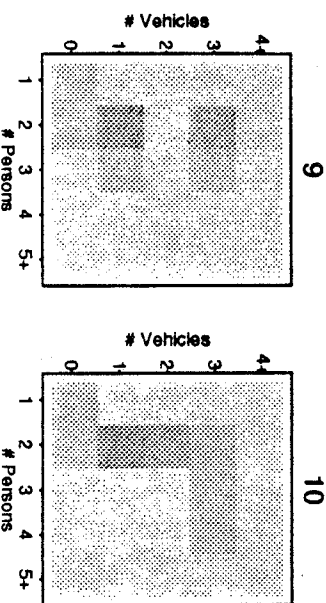
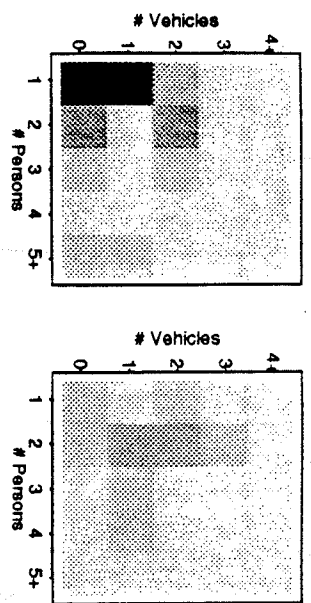
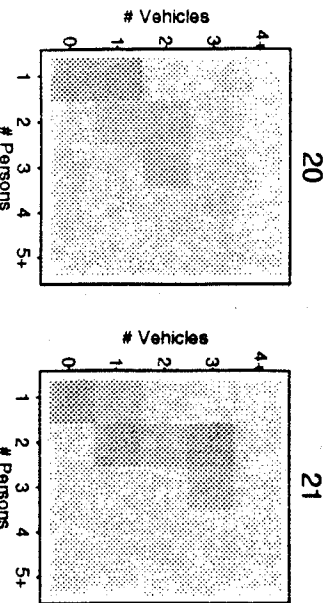
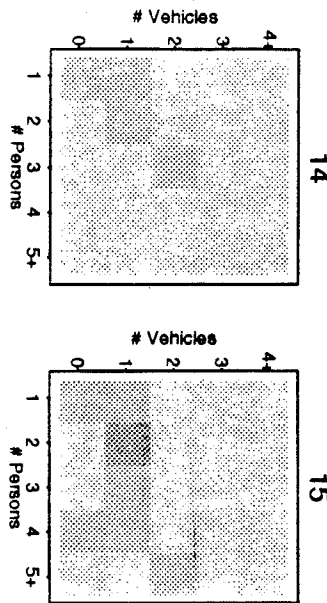
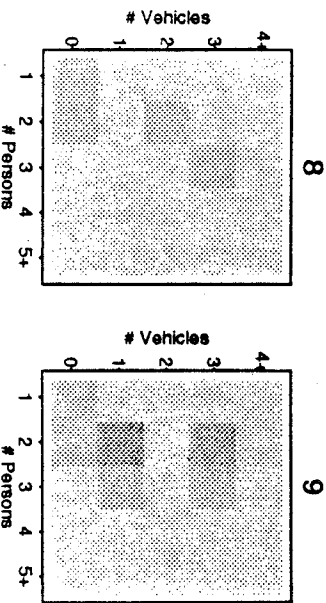
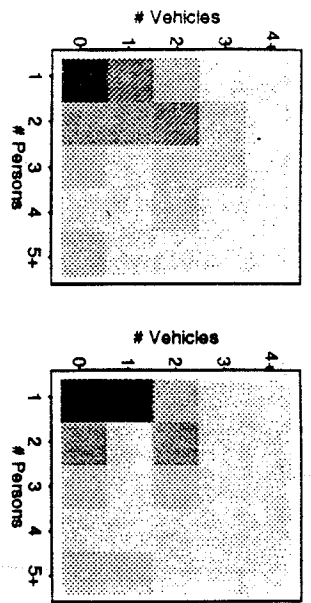
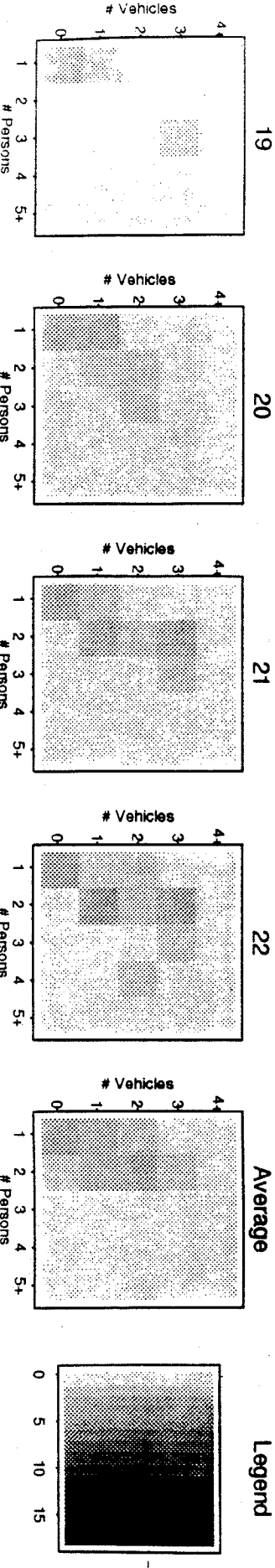
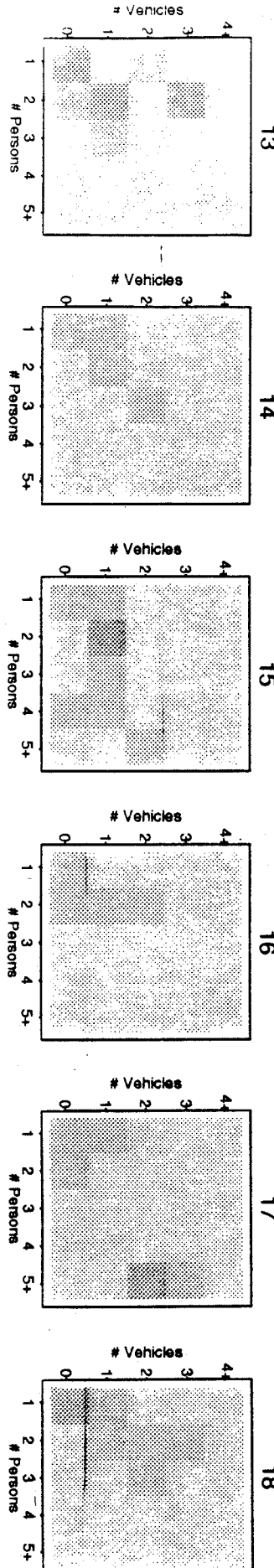
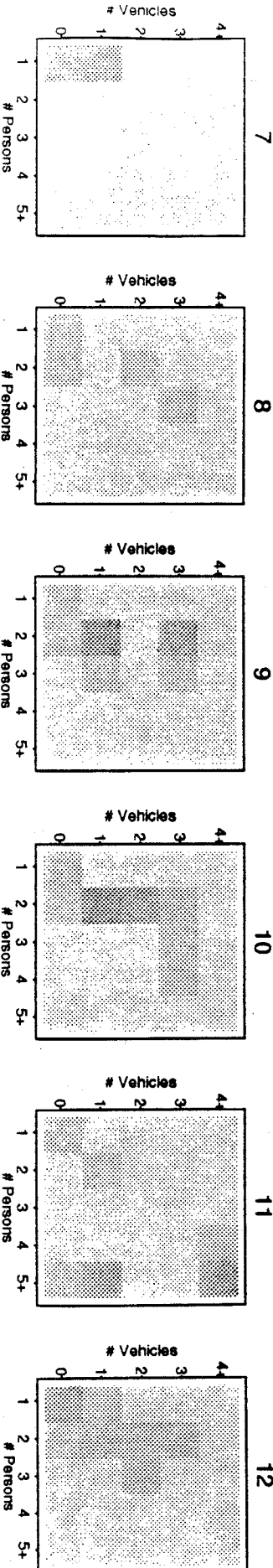
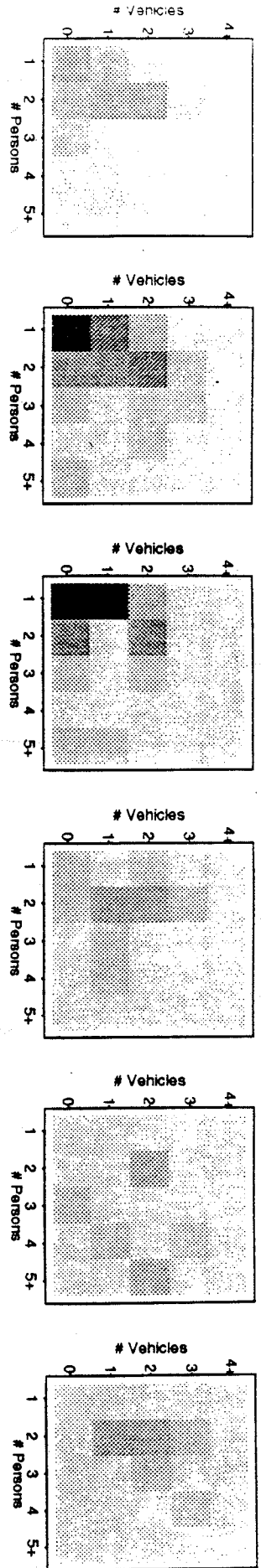
### Figure Captions

1. Distribution of household size for 100 synthetic populations generated for census tract 1216.05 of Tarrant County, Texas. Each line represents one synthetic population and shows the difference between the number of households of a particular size and the true number.
2. The difference from the truth in the proportions of household size and number of vehicles for an average synthetic population. Each box on the plot represents 1 of 22 "census tracts" which were constructed from San Francisco Bay area PUMS. The darker the shading in the box the larger the difference between the synthetic population and the truth.
3. The mean absolute deviation of the proportions of households with household sizes 1 to 5+ and 0 to 4+ vehicles in the synthetic population from the true proportions. Points on the plots represent these differences using only the PUMS and those differences using IPF. The differences are smaller for those populations generated with IPF.

# Distribution of Household Size

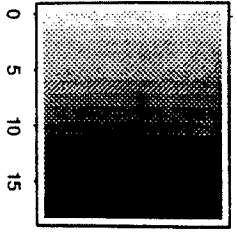


# Combined Fam/Nonfam households - Cell (%) differences from Estimate

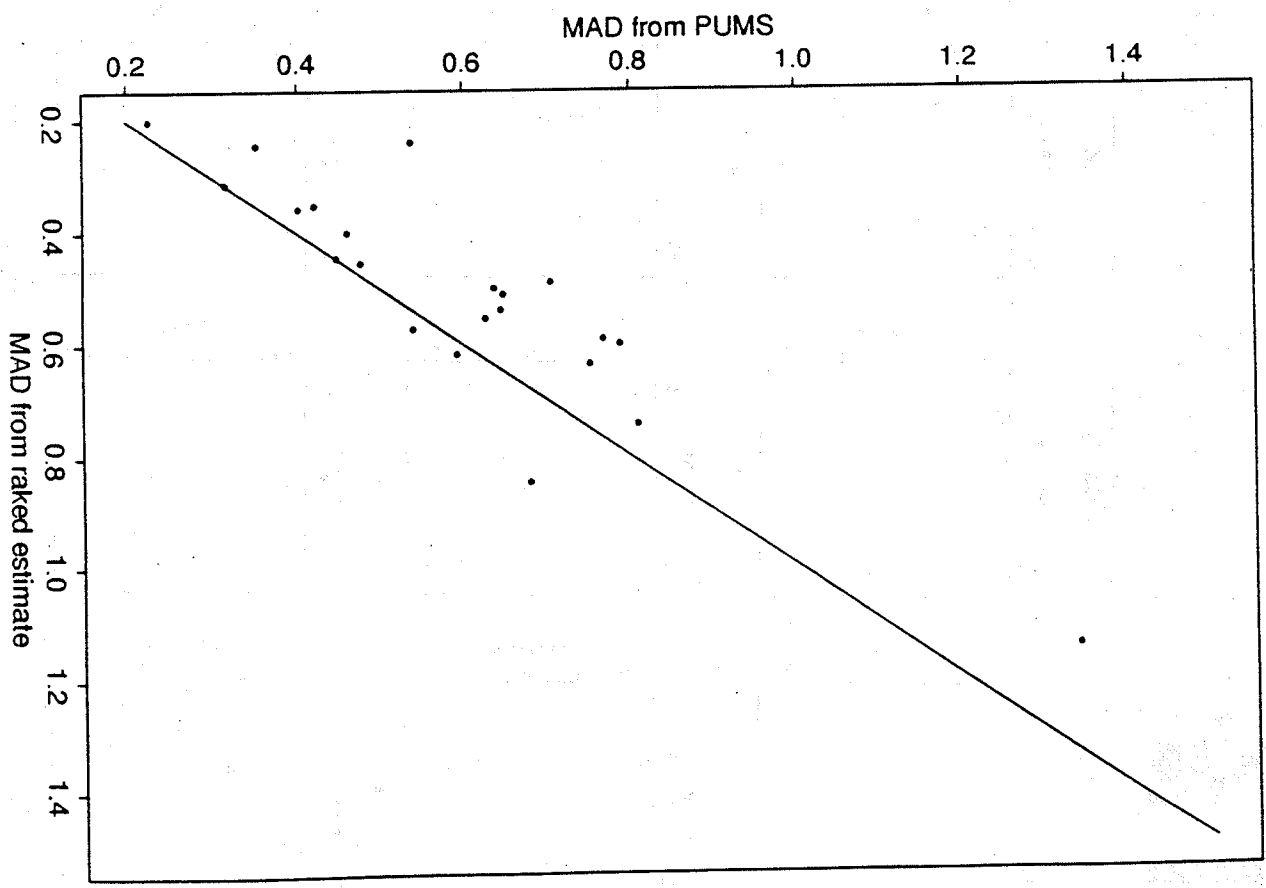
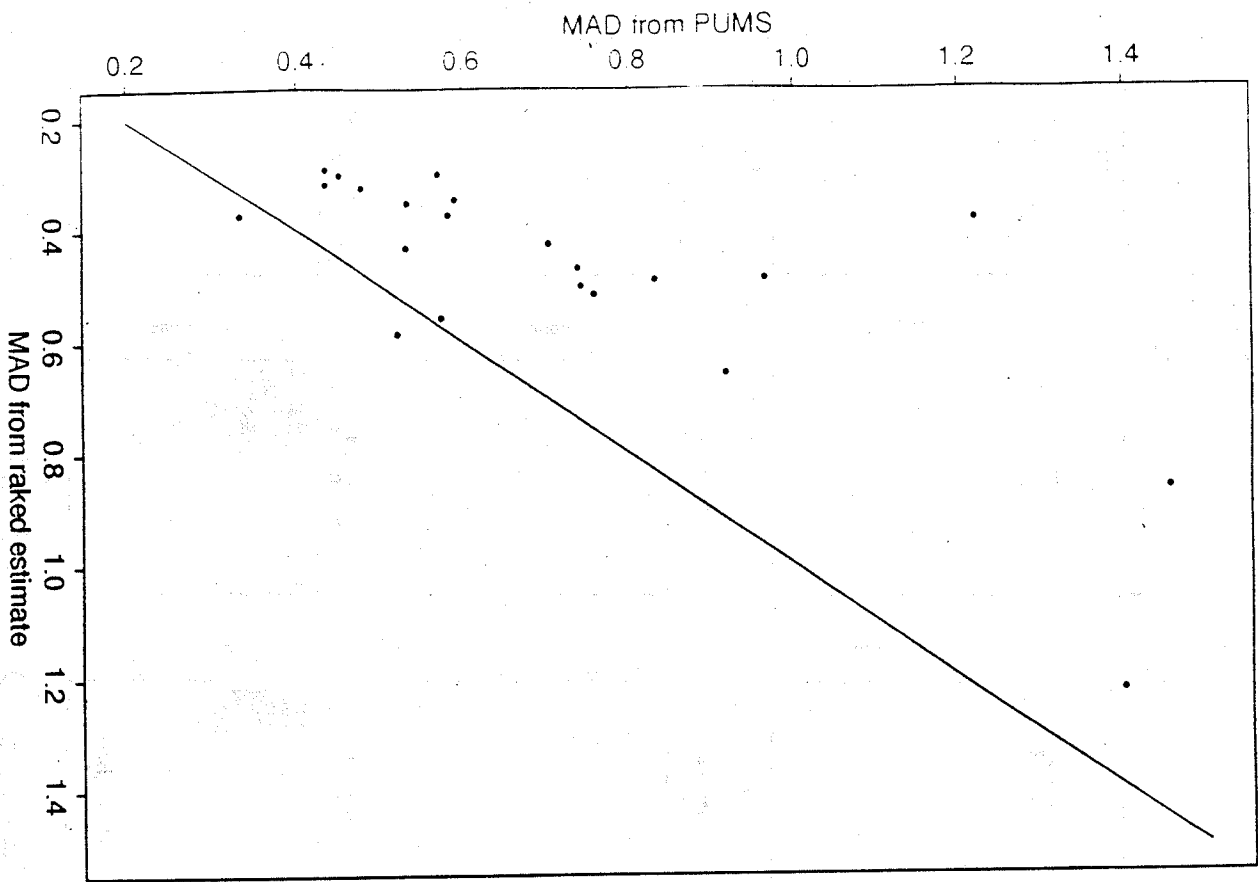


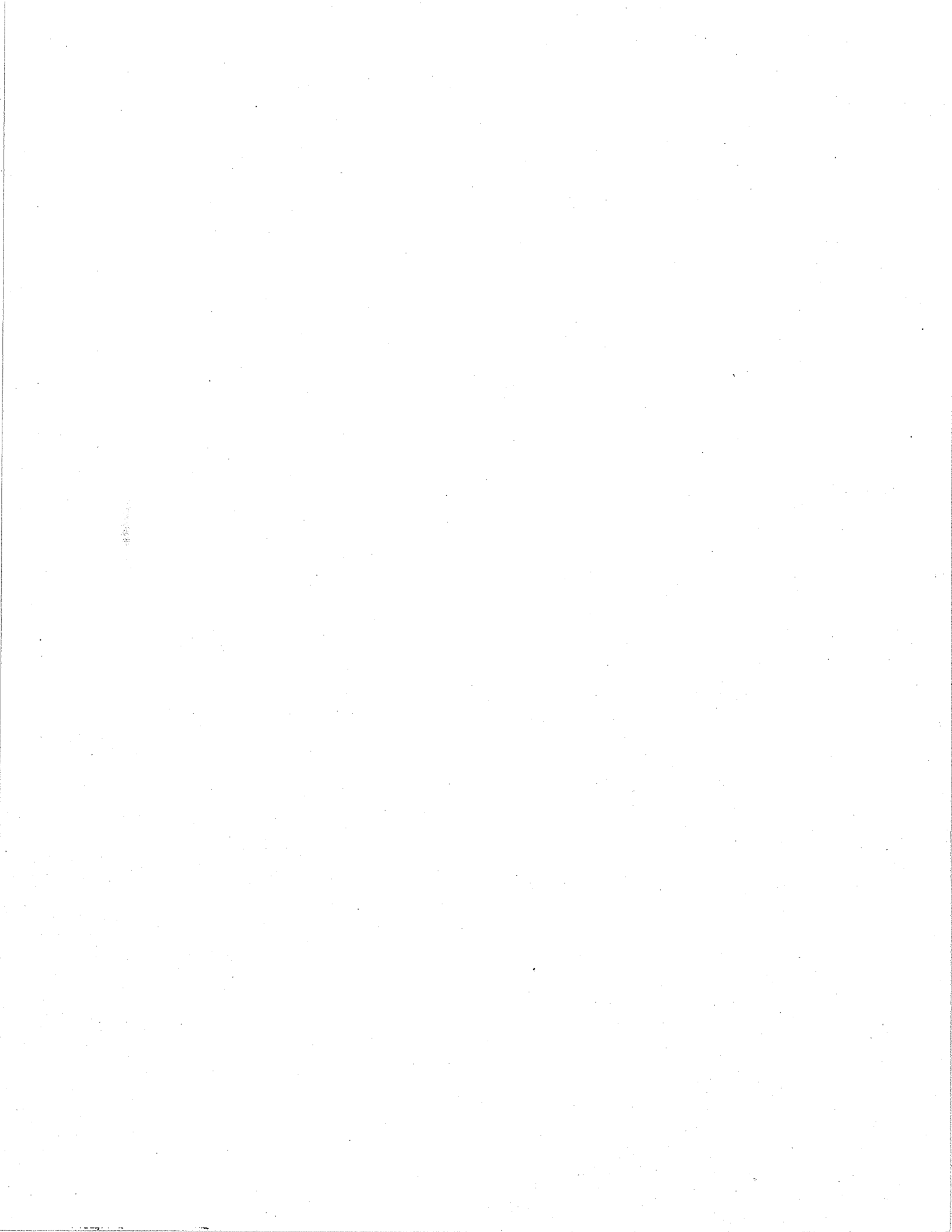
Average

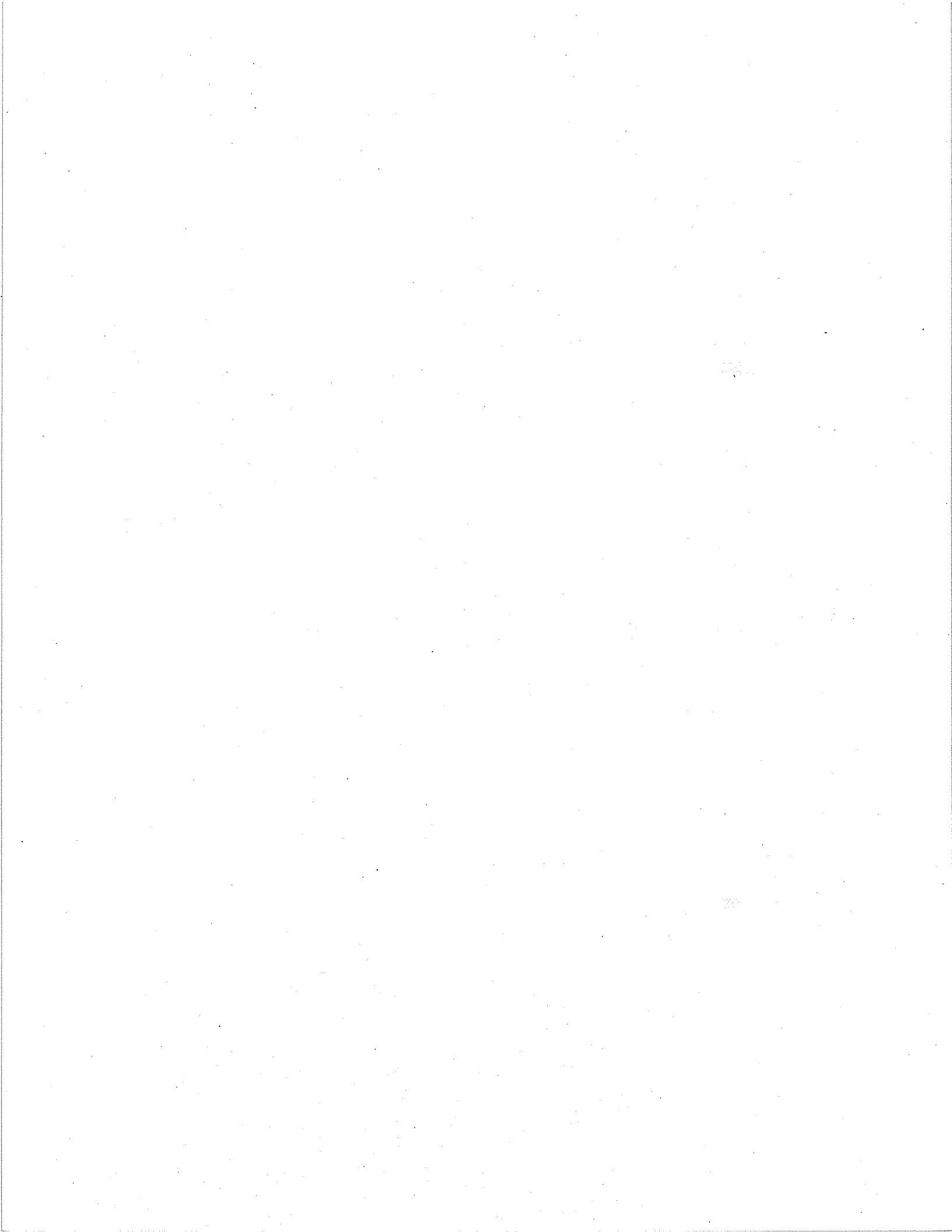
Legend



Mean Absolute Deviations for fitted Person X Auto %  
Family  
Nonfamily







## **Travel Model Improvement Program**

**The Department of Transportation, in cooperation with the Environmental Protection Agency, has embarked on a research program to respond to the requirements of the Clean Air Act Amendments of 1990 and the Intermodal Surface Transportation Efficiency Act of 1991. This program addresses the linkage of transportation to air quality, energy, economic growth, land use and the overall quality of life. The program addresses both analytic tools and the integration of these tools into the planning process to better support decision makers. The program has the following objectives:**

- 1. To increase the ability of existing travel forecasting procedures to respond to emerging issues including: environmental concerns, growth managements, and lifestyles along with traditional transportation issues,**
- 2. To redesign the travel forecasting process to reflect changes in behavior, to respond to greater information needs placed on the forecasting process and to take advantage of changes in data collection technology, and**
- 3. To integrate the forecasting techniques into the decision making process, providing better understanding of the effects of transportation improvements and allowing decisionmakers in state governments, local governments, transit agencies, metropolitan planning organizations and environmental agencies the capability of making improved transportation decisions.**

**This research was funded through the Travel Model Improvement Program.**

**Further information about the Travel Model Improvement Program may be obtained by writing to:**

**TMIP Information Request  
Metropolitan Planning (HEP-20)  
Federal Highway Administration  
U.S. Department of Transportation  
400 Seventh Street, SW  
Washington, D.C. 20590**

1917

Dear Sir,

I have the honor to acknowledge the receipt of your letter of the 14th inst.

and in reply to inform you that the same has been forwarded to the proper authorities for their consideration.

I am, Sir, very respectfully,  
Yours faithfully,  
[Signature]

[Name]  
[Title]

[Address]

[Additional information]

[Closing remarks]