

1 **TRAVEL PATTERN AND VARIABILITY PREDICTION**

2
3 Ruohan Li

4 Department of Civil, Architectural and Environmental Engineering
5 The University of Texas at Austin
6 lruohan_322@yahoo.com
7

8 Kara M. Kockelman

9 (Corresponding Author)

10 E.P. Schoch Professor of Engineering
11 Department of Civil, Architectural and Environmental Engineering
12 The University of Texas at Austin
13 kcockelm@mail.utexas.edu
14 Phone: 512-471-0210 & FAX: 512-475-8744
15

16 Under review for publication in the *Journal of the Transportation Research Forum*
17 January 2017
18

19 **ABSTRACT**

20 It is often important to know the travel demand of an area, which can be represented by the
21 annual vehicle-miles travelled (VMT) of the population in vehicles belonging to this area.
22 However, VMT can be very costly and difficult to track. Surveys are one of the most effective
23 tools in gathering travel demand data, by interviewing survey respondents about their travel
24 mileage across one or two consecutive days which can ultimately be used to model annual VMT.
25 In this study, using the data from Puget Sound Regional Council (PSRC), a regression of annual
26 VMT over both daily and 2-day VMTs ran for 20 times with randomly selected dates, resulting
27 in an average R-squared value of 0.1928 for annual vs. daily and 0.2233 for annual vs. 2-day.
28 Demographic variables include household income, age of the head of the household, number of
29 children, number of drivers per vehicle, as well as the month and week day of the selected date.
30 In order to keep track of the variation in travel among individual days across the year, the Gini
31 coefficient of each vehicle's travel pattern is also determined. The 214 vehicles have a mean Gini
32 coefficient of 0.2465. However, the adjusted R-square value for this regression turns out to be
33 0.1242, indicating that it's not an easy task to predict the Gini coefficient of a vehicle from
34 variables such as annual VMT, household income, age of head of household, number of children,
35 and number of drivers per vehicle.
36

37 **INTRODUCTION AND MOTIVATION**

38 Vehicle-miles traveled (VMT) is a metric that can be used to represent travel demand (Cervero et
39 al., 2002). In order to find the VMT for a population of vehicles, single day surveys are
40 performed in which each household participating completes a trip diary capturing all trips
41 undertaken during the 24-hour survey period. Surveys, are still the most prevalent among the
42 possible methods used to calculate VMT, due to the high drop-off rate and respondent fatigue in
43 multi-day surveys (Stopher et al., 2008). People's travel patterns, however, can vary
44 considerably over time (Pendyala and Pas, 2000). There can be days of extremely heavy travel as
45 well as days on which no travel takes place at all. Compared to obtaining only one day of trip
46 data, surveying for two days has the advantage of better capturing this variation. As a result,

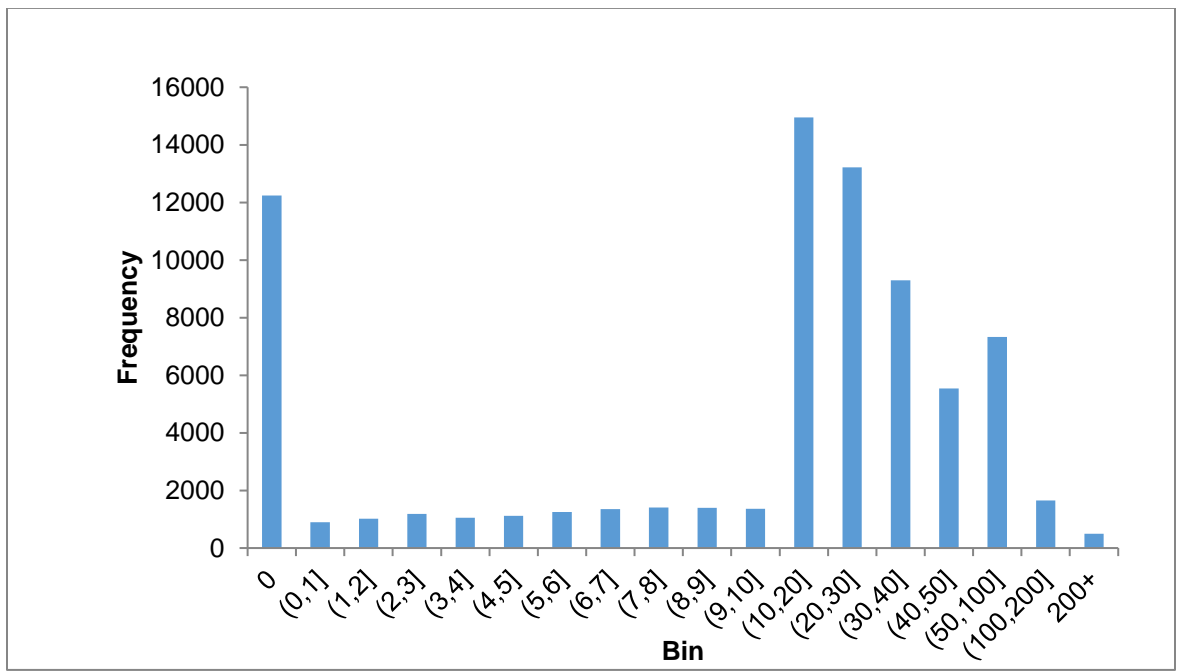
1 there is an increasing use of multi-day surveys of travel behavior usually two to three days in
2 length, aiming at capturing more variance than one day can provide (Axhausen et al., 2000).
3 Thus, when evaluating how descriptive of the annual VMT the short term travel records are in
4 this paper, the primary focus are on single day and 2-day values.

5
6 The variability, then, can be expressed by another term called Gini coefficient, initially derived
7 from economics, typically for income distribution. It is the area between the line of equality and
8 the Lorenz curve over the area of the triangle formed by the x-axis and the line of equality. The
9 Lorenz curve is a curve of the cumulative percentage of total travel plotted over the domain of 0
10 to 1. The line of equality is the Lorenz curve when the distribution is completely even across the
11 population. The value of the Gini coefficient varies between 0 and 1, and the smaller its value,
12 the less variation, and thus the more stable and predictable the travel pattern. A regression of the
13 Gini coefficient over annual VMT and a series of demographic variables is then run, to see how
14 they are correlated.

15
16 **DATA SET**

17 The data came from the Puget Sound Regional Council (PSRC) when it conducted the Traffic
18 Choices Study by placing GPS tolling meters in the vehicles of volunteer households. The final
19 data set contains 329 unique households and 484 vehicles. In order to get rid of the correlation of
20 travel among different vehicles in the same household, one vehicle per household was used, and
21 after removing certain households due to a low tracking period or missing demographic
22 information, the study was carried out with 214 vehicles, each from a different household.

23
24 **Figure 1. Histogram for Daily VMT**



26
27 Figure 1 shows the histogram for daily VMTs of all the vehicles. It can be seen from the graph
28 that 0 takes up a heavy proportion. Another peak in the histogram takes place between 10 miles
29 and 20 miles per day. This indicates that although no travel happening on a day at all is a very

1 common phenomenon, if a car does travel, a very probable amount it travels on one day falls
 2 between 10 miles and 20 miles. However, this falls short when compared to the average annual
 3 VMT of 19,850 miles found in the 2009 National Household Travel Survey (NHTS), and the
 4 daily VMT of 28.97 (Santos, 2009). One explanation can be that the sample being surveyed in
 5 the data set used in this paper might not be the most representing group of the overall travel
 6 pattern nationally, and might on average travel less.

7
 8 A regression is run of annual VMT (used in the study is the sum of 360 days of travel) over the
 9 travel mileage of one day or two consecutive days of each vehicle along with demographic
 10 information including household income, number of children, age of the head of the household,
 11 and number of drivers per vehicle, as well as date and day of the week the travel happened. After
 12 running the regression with randomly picked dates, the t-statistic and P-value of each variable is
 13 examined.

14
 15 **DATA ANALYSIS**

16 **Table 1. OLS of Annual VMT over Daily and 2-Day VMT** ($n_{obs} = 214$ vehicles)
 17

	<i>Coefficients</i>	<i>t Stat</i>	<i>P-value</i>		<i>Coefficients</i>	<i>t Stat</i>	<i>P-value</i>
Intercept	7128.3	16.54	0.000	Intercept	6757.4	15.81	0.000
Daily VMT	56.5	7.44	0.000	Two-Day VMT	34.4	8.60	0.000
Number of Children	485.4	1.98	0.049	Number of Children	351.3	1.48	0.139
Thursday	-1343.8	-1.97	0.050	Thursday	-1455.7	-2.21	0.028
Saturday	-1508.0	-2.20	0.029	Saturday	-833.9	-1.28	0.201
February	1430.5	1.73	0.085	February	1573.9	1.97	0.050
March	2354.8	2.49	0.013	March	2601.0	2.86	0.005
April	2468.4	3.18	0.002	April	2316.6	3.09	0.002
June	1706.1	1.61	0.108	June	1993.0	1.96	0.051
July	891.5	1.06	0.291	July	813.0	1.00	0.319
August	1697.7	2.31	0.022	August	1061.2	1.48	0.140

18 Adjusted $R^2 = 0.273$

Adjusted $R^2 = 0.322$

19
 20 Table 1 above is the regression result for one set of random dates. Because of the instability of
 21 the correlation between annual VMT and the selected parameters, the regression is repeated 20
 22 times. Table 2 contains all 20 runs' R^2 values, suggesting that Table 1's example is a relatively
 23 high fit result..

24
 25 One notable finding is that the intercepts are very high, and the coefficient for the daily travel
 26 amount is significantly less than 360. This might be due to the zeros in the sample. Since there
 27 are days that the vehicles sit still without any travel, and these zero days can be sampled into the
 28 daily VMT. When being plotted, they occupy the positive y axis, and thus bringing the intercept
 29 up and the slope down.

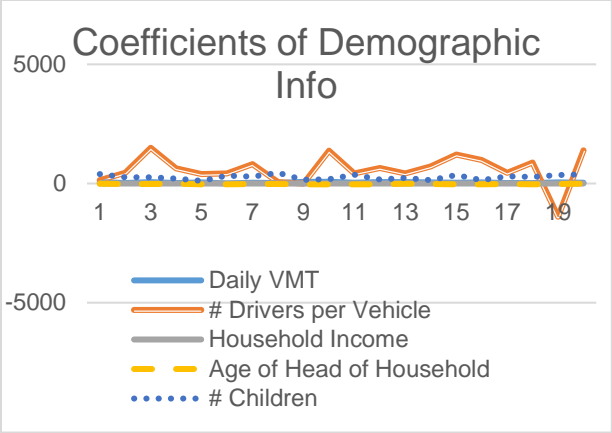
30
 31 **Table 2. R-Squared Value for 20 Additional Regressions** ($n = 214$ vehicles)

Single day	2-day
------------	-------

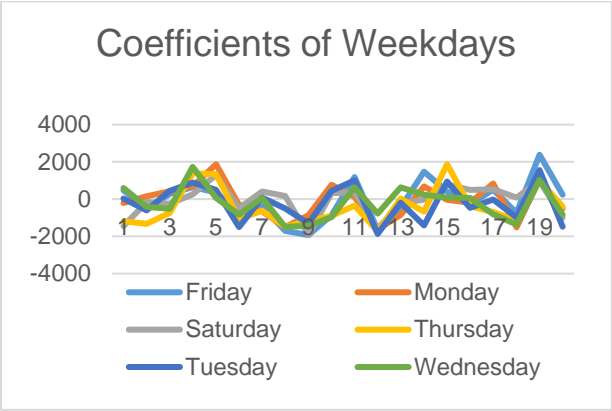
0.319	0.362
0.100	0.070
0.224	0.308
0.116	0.215
0.221	0.103
0.086	0.100
0.134	0.254
0.185	0.225
0.287	0.356
0.282	0.334
0.144	0.197
0.247	0.284
0.215	0.255
0.167	0.237
0.170	0.205
0.287	0.312
0.121	0.089
0.178	0.258
0.278	0.090
0.210	0.213

1
2 In the 20 pairs of regression models, the R-squared values vary significantly, but, as expected,
3 most (80 percent) of the pairs delivered a higher R-squared value for the 2-day sample. This
4 indicates that the prediction result is unstable, and also varies with the date being selected.
5 However, within the same sample, 2 consecutive days of travel tends to be more descriptive of
6 how much the vehicle travels throughout the year than only one day.
7 But the R-squared value for 2-day regressions can drop very low, one of the reasons being that
8 when two days are added together, variability can fall, but can also increase when both days
9 sampled are unusual (e.g., a summertime driving vacation or two consecutive days of zero
10 driving). There 12,595 zero-VMT days in the data set (16.3% of the total), and half (6,263) of
11 these zero-VMT days are directly followed by a zero-VMT day.
12 Plotting the coefficients can make evident which parameters are of more significance. The
13 following plots are divided into three categories: demographic information, months, and
14 weekdays.

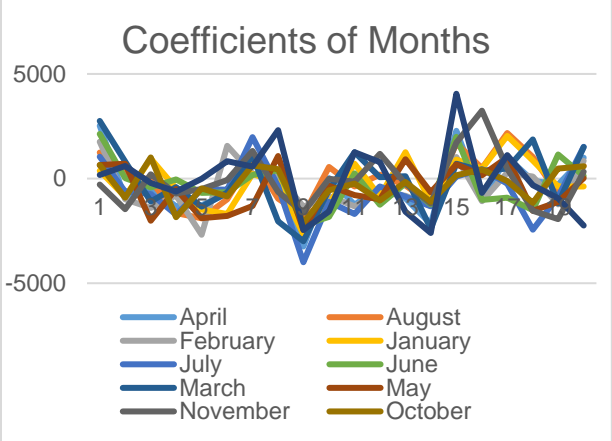
15
16 **Figure 2(a). Coefficients of each Variable in 3 Categories for Annual VMT vs. Daily VMT**



1



2



3

4

5

6

7

8

9

10

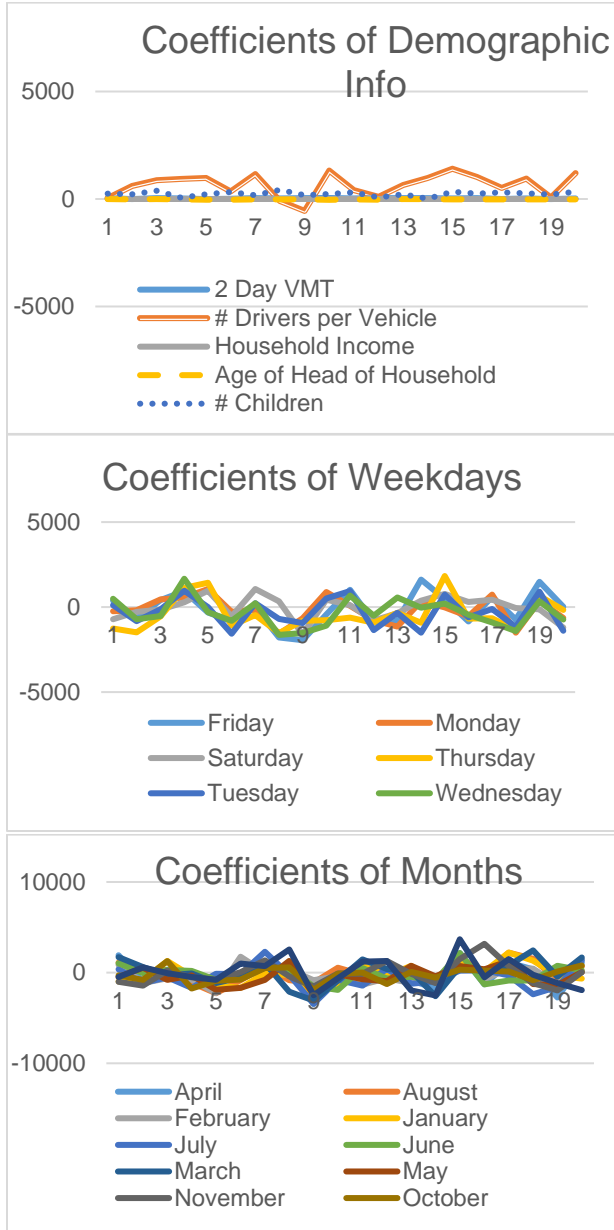
11

12

13

14

1 **Figure 2(b). Coefficients of each Variable in 3 Categories for Annual VMT vs. 2-Day VMT**
 2



3
4
5
6
7 When looking at the results of the regression, not many consistent patterns can actually be
 8 drawn, other than the fact that the age coefficient of the head of the household tends to be
 9 negative and the number of children coefficient is usually positive. This is due to the variability
 10 of people's travel patterns. Even if the miles a vehicle travels on a day and all household
 11 information are known, it is difficult to determine whether that day is a day of heavy or light
 12 travel for the household. This is analogous to the fact that randomly selecting one sample from
 13 the population is often not relevant enough to predict the entire normal distribution. Surveying
 14 for more than one day can eliminate this type of error to some extent, but the extended survey
 15 period reduces reporting accuracy as well. It is possible, however, with the technology of GPS
 16 devices, to keep track of the surveyed vehicles' travel pattern over a longer period of time, like

1 an entire week, without sacrificing accuracy. As shown in Table 3 and as expected, R2 values
 2 for a week's worth of VMT data out-perform a day or two of such data, in predicting the
 3 household's annual VMT with that vehicle. It should be noted though, that the month to which a
 4 week is considered to belong is determined by which month the first day of the week falls in (so
 5 the month may change, and there is a greater chance of overlap in some of the sampled days
 6 across the 6 regression results, due to a week-long sampling period).

7
 8
 9 **Table 3. R-Squared Value for 6 Annual VMT vs. Weekly VMT Regressions**

	Adj R ²
Run 1	0.389
Run 2	0.454
Run 3	0.385
Run 4	0.445
Run 5	0.389
Run 6	0.499
Average	0.427

11
 12
 13 The next task, then, would be to sort out the demographic patterns of households that might have
 14 more regularly traveling vehicles, and an indication for this would be the Gini coefficient.

15
 16 **GINI COEFFICIENT**

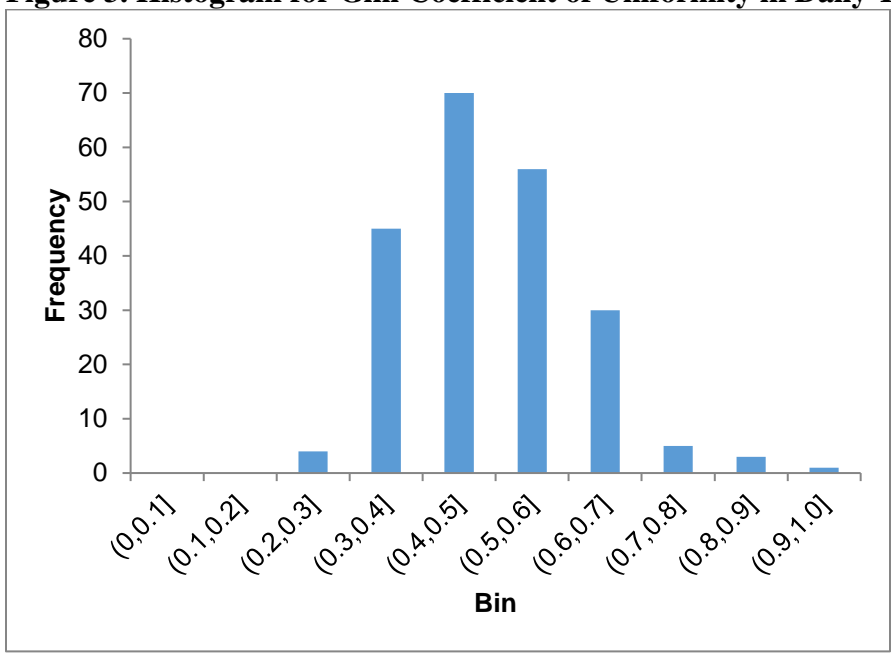
17 The Gini coefficient is a term originated from economics and has been used as a conventional
 18 measure of income inequality (Dorfman, 1979). However, it can be used in other areas as well.
 19 In this study, for example, it is possible to look at the Gini coefficient of the daily VMT of a
 20 certain vehicle across a year. It is defined as the area enclosed by the line of equality and the
 21 Lorenz curve (Turrell and Mathers, 2001). In economics, the Lorenz curve represents the
 22 cumulative distribution of income while income units such as individuals or households are
 23 arranged in an ascending order from left to right along the x-axis (Kakwani, 1977). The line of
 24 equality is a line in the first quadrant passing through the origin that forms a 45 degree angle
 25 with the positive x-axis, and it is the Lorenz curve when income is evenly distributed among all
 26 individuals. The income inequality Gini index of the United States was 0.480 in 2014, increasing
 27 by 5.9 percent since 1993 (DeNavas-Walt and Proctor, 2015). In this study, the line of equality
 28 shows the cumulative distribution of a car's daily VMT over a year's period if it travels the same
 29 amount every day throughout the year, while the Lorenz curve is the actual cumulative
 30 distribution across days of the week arranged from left to right in the order of ascending daily
 31 VMT. The area between the 2 curves then, is the Gini coefficient.

32
 33 In this study, since it is picked 360 days as a year, there are a sufficiently large number of points
 34 on the x-axis to use the rectangular approximation method. The rectangular approximation
 35 method is used to approximate the area under the curve as the sum of the areas of n rectangles
 36 each having a width of $\frac{1}{n}$, while n is the number of days being looked at. In this case, the distance
 37 between the points on the line of equality and the Lorenz curve on each day is multiplied by the

1 width of $\frac{1}{360}$ to get the area of 360 rectangles, and summing up the areas of them would be the
 2 approximated Gini coefficient.

3

4 **Figure 3. Histogram for Gini Coefficient of Uniformity in Daily Travel Distances**



5

6 Figure 3 shows the histogram of the vehicles' Gini coefficient for daily VMT. Among the 214
 7 Gini coefficients calculated, more than a third are between 0.4 and 0.5, indicating that a large
 8 percentage of Lorenz curves form an area with the line of equality approximately half the size of
 9 the triangle formed by the line of equality with the x and y axes. It is very rare that this value
 10 drops below 0.3 or goes beyond 0.8.

11

12 **Figure 4. Lorenz Curves for n = 194 Vehicles over 360 Days versus Line of Equality**

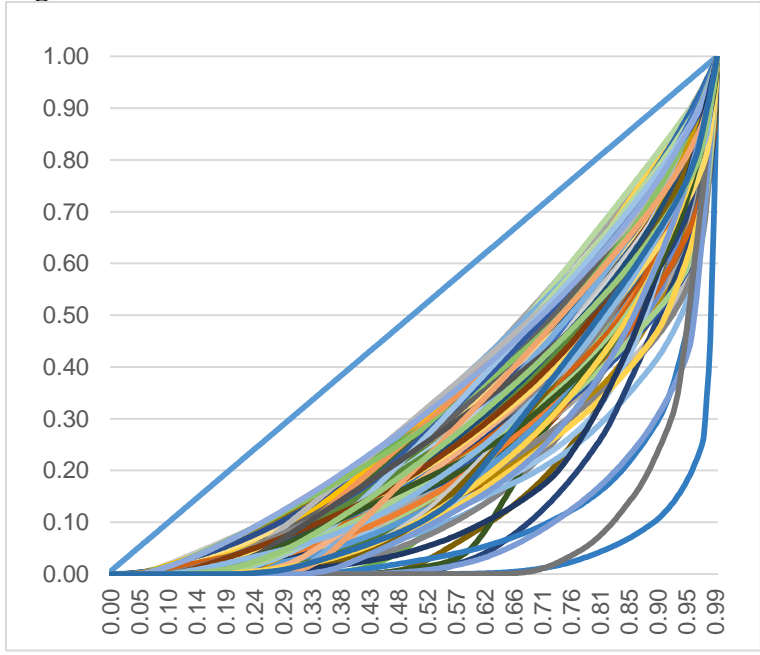


Figure 4 displays the plot for the Lorenz curve of 194 different vehicles, and the area between each Lorenz curve and the line of equality is the Gini coefficient of that specific vehicle's daily VMT throughout the year. This area reaches its maximum when all travel occurs on one day. As the total number of days increases, the maximum Gini coefficient approaches 1. When the travel pattern is homogeneous every day throughout the year, the Gini coefficient reaches its minimum of 0, and the Lorenz curve coincides with the equality line. Since there are vehicles whose travel mileage are not kept track of for a full 360-day period, they are removed from the list when graphing, and that's why the number of vehicles on the graph are less than that included in the regressions.

When surveying households, the lower the Gini coefficient of travel, the easier it is to predict the annual VMT from the daily value, since it indicates a more homogeneous travel. If the Gini coefficient is high because of how far the vehicle travels per day varies widely throughout the year a one day survey doesn't provide enough information on where this day falls on the spectrum of travel distance for this vehicle, and might be an indicator that suggests coming back to the household for another time, or asking them to provide travel details over a longer period of time.

What demographic characteristics may deliver more day-to-day VMT variability and a higher Gini coefficient for the vehicle? Table 2 provides the results of such a regression, for $y = \text{Gini coefficient}$ versus household demographic information, the following result is shown in Table 2.

Table 2. OLS of Gini Coefficient over Annual VMT and Demographic Variables (n = 214 vehicles)

	Coefficients	t Stat	P-value
Intercept	0.2679	10.45	0.000
Annual VMT	0.0000	-5.09	0.000
Annual Household Income	0.0000	1.07	0.287
Number of Children	-0.0019	-0.43	0.666
Age of Household Head	0.0003	0.87	0.384
#Drivers per Vehicle	0.0058	0.30	0.761

From the regression, however, it can be concluded that the Gini coefficient heavily depends on annual VMT. This discovery makes sense, since generally the regular travelers that go to certain places (for example, school or work) every single day accumulate a higher annual VMT. Due to the low R-squared value of annual VMT over both single day and 2-day VMTs, using a predicted annual VMT to find out the Gini coefficient may cause even more significant error. But it can be possible to read from the odometer of the vehicle to extract the value of the annual VMT, and use this value to estimate the Gini coefficient.

Household income is another important variable in predicting the Gini coefficient. There is a positive coefficient linked to this variable, indicating that the Gini coefficient tends to be higher when the surveyed household has a high income. This pattern makes sense as well, since usually

1 wealthier people have more control over their time, including what to do and where to go. They
2 are more likely to take vacations and go on long journeys. Long, driven journeys would appear in
3 the data as very large daily VMTs, while non-driven (or rented-car) journeys can result in a
4 series of zeros, since the owners are away and their personal vehicles are left unattended. Either
5 situation delivers a higher Gini coefficient.

6
7 The Gini coefficient tends to increase with age of the household head, which may be due to
8 retirees not having a daily work or school commitment, and thus less regular travel needs.

9
10 The number of children in a household and the drivers-to-vehicles ratio seems are not
11 statistically significant predictors of the Gini coefficient, but a larger data set will tend to result
12 in smaller p-values on any coefficient estimate. Other variables that may be quite helpful to have
13 include location of the household within its region (e.g., local population and jobs densities, land
14 use balance metrics, and distance to worker workplaces), distances to nearby major cities,
15 household income, number of students in the household, and ages of all household members.

16 17 **CONCLUSION**

18 People's travel patterns vary from day to day, and after knowing how much one travels on a
19 specific day, it is may be difficult to predict how much a person travels due to the discrepancy in
20 how one day of surveying might compare to other days in the year. Longer or more repetitive
21 surveys, then, can be carried out in order to make a closer estimation, but the respondent's own
22 accuracy in reporting trips decrease as the length of survey period increases. There are certain
23 households that have vehicles that follow a more regular travel pattern, of which we can take the
24 advantage to survey for only one day to reduce the cost of longer surveys and the loss of
25 accuracy due to respondent fatigue. Higher annual VMT, lower household income, and lower
26 age of the head of the household are all such indicators. If not all variables appear to indicate a
27 low Gini coefficient, the researchers can consider surveying another day. But if it occurs that all
28 variables suggest a high Gini coefficient, it might be a better idea to turn towards another
29 household since the extended survey period to make up for the variability might result in even
30 lower accuracy due to the survey length. One major difference between the daily VMT used in
31 the study and that obtained in traffic surveys is that in the study, the daily VMT is collected by
32 the GPS device, while in actual traffic surveys, this value is often self-reported. However, it was
33 found that self-reported car trips compare well with the actual distance, especially when the trip
34 is short (Salon, 2013). Also, because of the use of 360 days instead of 365 days when selecting
35 data, the generated value can be scaled up by $\frac{365}{360}$ in order to get the travel mileage of a full 365-
36 day year.

37 38 **REFERENCES**

- 39 Axhausen, K. W., Zimmermann, A., Schönfelder, S., Rindsfuser, G., & Haupt, T. (2002).
40 Observing the Rhythms of Daily Life: A Six-Week Travel Diary. *Transportation*, 29(2), 95-124.
41 Cervero, R., & Hansen, M. (2002). Induced Travel Demand and Induced Road Investment: A
42 Simultaneous Equation Analysis. *Journal of Transport Economics and Policy (JTEP)*, 36(3),
43 469-490.
44 DeNavas-Walt, C., and B. D. Proctor. "US Census Bureau. Income and poverty in the United
45 States: 2014." *Current Population Reports* (2015): 60-252.

- 1 Dorfman, R. (1979). A Formula for the Gini Coefficient. *The Review of Economics and*
2 *Statistics*, 61(1), 146-149.
- 3 Kakwani, N. (1977). Applications of Lorenz Curves in Economic Analysis. *Econometrica*, 45(3),
4 719-727.
- 5 Pendyala, R. M., & Pas, E. I. (2000). *Multi-Day And Multi-Period Data for Travel Demand*
6 *Analysis and Modeling* (No. E-C008,).
- 7 Salon, D. (2014). Comparison of Self-Reported to Network-Calculated Trip Distances for the
8 California Add-on to the 2009 National Household Travel Survey. In Transportation Research
9 Board 93rd Annual Meeting (No. 14-5389).
- 10 Santos, A., McGuckin, N., Nakamoto, H. Y., Gray, D., & Liss, S. (2011). *Summary of Travel*
11 *Trends: 2009 National Household Travel Survey* (No. FHWA-PL-11-022).
- 12 Stopher, P., Kockelman, K., Greaves, S., & Clifford, E. (2008). Reducing Burden and Sample
13 Sizes in Multiday Household Travel Surveys. *Transportation Research Record: Journal of the*
14 *Transportation Research Board*, (2064), 12-18.
- 15 Turrell, Gavin, and Colin Mathers. "Socioeconomic Inequalities in All-Cause and Specific-Cause
16 Mortality in Australia: 1985–1987 and 1995–1997." *International Journal of Epidemiology* 30,
17 no. 2 (2001): 231-239.