

Prediction of AADT on a nationwide network based on an accessibility-weighted centrality measure

Working Paper**Author(s):**

Sarlas, Georgios; Axhausen, Kay W. 

Publication date:

2015

Permanent link:

<https://doi.org/10.3929/ethz-b-000102909>

Rights / license:

[In Copyright - Non-Commercial Use Permitted](#)

Originally published in:

Arbeitsberichte Verkehrs- und Raumplanung 1094

1 **Prediction of AADT on a nationwide network based on an accessibility-**
2 **weighted centrality measure**

3 Georgios Sarlas*

4 PhD Student

5 Institute for Transport Planning and Systems

6 ETH Zurich

7 HIL F 51.3, Stefano-Francini-Platz 5

8 8093, Zurich, Switzerland

9 Phone: +41 44 633 37 93

10 E-mail: georgios.sarlas@ivt.baug.ethz.ch

11

12 Kay W. Axhausen

13 Professor

14 Institute for Transport Planning and Systems

15 ETH Zurich

16 HIL F 31.3, Stefano-Francini-Platz 5

17 8093, Zurich, Switzerland

18 Phone: +41 44 633 39 43

19 E-mail: axhausen@ivt.baug.ethz.ch

20

21

22 * Corresponding author

23

24

25

26

27

28

29

30

31

32

33 Total Word Count: **5348 (Text) + 7 (Figures/Tables) * 250 + References (500) = 7598 words**

34

35

36

37

38 Submitted for Publication and Presentation at the Transportation Research Board 95th Annual

39 Meeting January 10-14, 2015, Washington, D.C.

40 **ABSTRACT**

41 In the present paper a direct demand modelling approach for AADT prediction on a nationwide
42 network is presented. A particular focus is given on the construction of a variable that can capture
43 the interregional demand patterns by taking into account the direction of potential interactions
44 over space, called accessibility-weighted centrality, by applying different modifications on the
45 stress centrality measure, tailored for the task of AADT prediction. It is exhibited that it can lead
46 to a significant enhancement on the accuracy of the models. In addition to the already tested
47 models in the literature, two SAR models are estimated and it is shown that GWR and kriging are
48 more appropriate for interpolation purposes, while spatial error and OLS models might give
49 slightly worse results but they have the potential to be applied both for interpolation and
50 forecasting since their estimated parameters are unbiased and consistent. A comparison of models
51 predictive accuracy to the output of a traditional four-step model is conducted to show that direct
52 demand models on nationwide scale can constitute a trustworthy alternative to more advanced,
53 but definitely more data demanding and computationally burdensome models.

54 **INTRODUCTION**

55 Many studies in the field of transport modelling have dealt with the issue of annual average daily
56 traffic (AADT) prediction, developing different methodologies to tackle the problem. In general,
57 two main streams of literature can be found. One that exploits different modelling techniques
58 aiming at resolving the issues of spatial dependence and heterogeneity, while in the second
59 stream the construction and the inclusion of more variables describing the demand patterns in
60 models is investigated. The employed methodologies vary from the aspatial regression techniques
61 to the statistical techniques accounting for the spatial effects. In particular, the later encompass
62 two different approaches. The first one is utilizing a data-driven approach of spatial statistics
63 called kriging, while the second one utilizes the geographically weighted regression (GWR) of
64 the class of spatial econometric models. Nevertheless, the majority of the studies developed
65 various methodologies tailored for small, or medium, scale level of analysis in terms of network
66 size, having mainly the purpose to interpolate AADT from known to unmeasured locations.

67 **Literature review**

68 Xia et al.(1) developed a multiple regression model for estimating AADT on non-state roads of
69 Florida and found that the most important contributing predictors are the roadway characteristics
70 along with the area type, while socioeconomic variables were found to have an insignificant
71 impact on AADT. Similarly, Mohamad et al.(2) developed a multiple regression model for
72 AADT prediction for county roads in Indiana, incorporating various demographic variables
73 which were found to be significant. In a similar context, Desylas et al. (3) developed a multiple
74 regression analysis model for pedestrian flows.

75 The plausibility of applying the GWR model for estimating AADT was demonstrated in another
76 study (4) and it was shown that it can lead to the enhancement of the prediction accuracy,
77 compared to the aspatial ordinary linear regression. Eom et al. (5) exploited ordinary kriging for
78 interpolating AADT for non-freeway facilities in Wake County, North Carolina, and concluded
79 that its predictive capability is much better than the ordinary regression models. Along the same
80 line of thought, Wang and Kockelman (6) applied kriging-based methods for AADT prediction at
81 unmeasured locations, making use of Texas highway count data, and highlighted further the
82 capability of applying kriging for prediction purposes on a statewide network. Selby and
83 Kockelman (7) explored the application of two spatial methods for prediction of AADT on the
84 same statewide network (universal kriging and GWR), and they concluded that both methods
85 reduce prediction errors over aspatial regression techniques whereas the predictive capabilities
86 of kriging exceed those of GWR. Interestingly, the estimation of the kriging parameters taking
87 into account network distances, instead of Euclidean, showed no enhanced performance.

88 Furthermore, Pulugurtha and Kusam (8) developed Generalized Estimating Equations models to
89 estimate AADT using integrated spatial data from multiple network buffer bandwidths. Spatial
90 data included off-network characteristics such as demographic, socio-economic and land use
91 characteristics, captured over multiple network buffer bandwidths around a link and integrated by
92 the employment of distance decreasing weights. The methodology was applied on a city level
93 (Charlotte, North Carolina). As a continuation of the previous study (9), the authors exploited the
94 application of the principle of demographic gravitation to estimate AADT based on land-use
95 characteristics on the same network. A negative binomial model was estimated along with neural
96 network models. Interestingly, the results obtained showed that the developed models gave
97 significantly lower errors in comparison to outputs from traditional four-step method used by
98 regional modelers.

99 In a recent study by Lowry (10), a new method for interpolating AADT was presented, tailored
100 for communities where attributes such as roadway characteristics, land-use etc., are uniform over
101 space, and thus their inclusion in the model bears no explanatory power. The new method used
102 novel explanatory variables that are derived through a modified form of stress centrality, a
103 network analysis metric that quantifies the topological importance of a link in a network. The
104 case study showed high quality results. The same methodology found application as well for
105 estimating directional bicycle volumes (11).

106 **Description of the framework of the paper**

107 The objective of the current paper is to develop a direct demand modelling approach for
108 prediction of AADT on a nationwide network, which has not been addressed in the existing
109 literature. The particularity of the nationwide network level case stems from the incapability of
110 the spatial densities of different socioeconomic data to capture adequately the demand patterns
111 that occur on the links, since they fail to bear explanatory power with respect to high volume of
112 interregional through traffic. Naturally, the construction of a variable that can account for

113 interregional flows necessitates, taking into account the direction of potential interactions,
114 allowing us to capture the demand capacity interaction at the core of transport modelling. More
115 specifically, a variable called accessibility-weighted centrality measure is constructed, building
116 upon the work of Lowry (10) who showed that the use of stress centrality constitutes a huge
117 improvement over traditional direct-demand models, and modifying accordingly the stress
118 centrality measure for the particular problem at hand.

119 In addition to the already tested models in the literature, the family of spatial simultaneous
120 autoregressive (SAR) models is exploited with their capability to be applied for AADT prediction
121 purposes. The advantage of such models is that they can resolve spatial dependence issues,
122 offering a structural explanation of the AADT and since their estimated coefficients are unbiased
123 and consistent, they can fulfill both interpolation and forecasting purposes which is important for
124 policy evaluation and project appraisal purposes. In summary, a set of different models is
125 estimated and evaluated in order to draw sound conclusions on the newly constructed variable
126 and also on models' capabilities to be employed for AADT prediction purposes and thus
127 highlight in a quantifiable way their strengths and weaknesses. At last, a comparison of models
128 predictive accuracy to the output of a traditional four-step model is conducted to show to what
129 extent such models can constitute a trustworthy alternative to more advanced, but definitely more
130 data demanding and computationally burdensome, models.

131 **METHODOLOGY**

132 **Accessibility-weighted centrality measure**

133 The construction of a new variable capturing the interregional demand patterns, taking into
134 account the direction of potential interactions over space, is of high importance for the estimation
135 of AADT models on a nationwide network. Making use of the graph theory, centrality is an index
136 that aims to identify the most influential persons in the context of a social network. Different
137 centrality indices have been introduced over the years, aiming at the identification and the
138 quantification of the importance of a particular person in a social network. In general, centrality
139 indices take into account the number of shortest paths that pass by a given link/node, either for
140 given pairs of nodes, or for all pair of nodes within the network. In the case where a capacity
141 constraint exists in the form of a particular weight/cost associated with each link/node, then this
142 weight should be taken into account in the routing algorithm for the identification of the shortest
143 paths.

144 Departing from the social sciences questions, centrality indices are meaningful for all networks'
145 analysis. From this viewpoint, centrality indices are meaningful for the analysis of transport
146 networks as well and can provide a quantifiable measure of the importance of links, taking into
147 account the network structure and the cost of traversing each link (distance or time). In the case

148 of transportation, networks correspond to directed networks, given the allowed and prohibited
 149 turning movements on its vertices (nodes), and are modelled as higher level networks in order to
 150 account for them. Stress centrality index was introduced by Shimbel (12) and is defined as the
 151 number of shortest paths connecting all pairs of nodes of the network that pass from a link.

$$152 \quad \text{Stress centrality}_e = \sum_{i,j \in V} \sigma_{ij}(e) \quad (1)$$

153 Where e is any link of the network, V the set of all nodes, σ_{ij} the shortest path from node i to node
 154 j , and $\sigma_{ij}(e)$ is equal to one if the link e is part of the shortest path connecting i and j nodes.

155 By definition, higher hierarchical links have high centrality values, while that might be the case
 156 as well for lower hierarchical links given the network structure. In the case of transport networks,
 157 the hierarchy is given by the functional class of the roads where their importance is normally
 158 matched by the number of trips using the given link. Naturally, two issues with respect to the
 159 application of the stress centrality index for transport networks come to the surface. First, the
 160 issue of travel demand since not all nodes are attracting or producing the same number of trips
 161 and thus this should be taken into account in the centrality formulation. Second, interaction
 162 between nodes tends to diminish and becomes very small as the distance between them increases,
 163 which should be accounted for in a modified stress centrality formulation.

164 Addressing the aforementioned issues takes place in three steps. At first, the issue of trip
 165 production and attraction is addressed by making the assumption that production is related to the
 166 economically active population in the vicinity of the origin node, and attraction at the
 167 employment positions at the destination node. Second, the interaction intensity between the nodes
 168 should be associated with a function that diminishes by network distance. The distance decay
 169 function embedded in the measure of travel accessibility is employed for that reason, since
 170 accessibility is a measure of how far people are willing, or able, to travel on the course of their
 171 daily life and quantifies how interaction opportunities decrease over the distance (13). Two
 172 variations of distance decay function are checked to identify the one that fits the data better (14).
 173 Last, a restriction has to be imposed with respect to the direction of potential interactions by
 174 standardizing the accessible opportunities from each node to each node, by the total number of
 175 opportunities accessible by the origin node in total. The incorporation of these changes in the
 176 stress centrality index and the derivation of the constructed index, called *accessibility-weighted*
 177 *centrality*, is presented below. It should be noted that the constructed variable mirrors to a great
 178 extent the first two steps of the traditional four-step model, however this is inevitable due to the
 179 nature of the relationships that we need to capture in the variable.

$$180 \quad \text{Accessibility – weighted centrality}_e = \sum_{i,j \in V} \sigma_{ij}(e) \quad (2)$$

$$181 \quad \sigma_{ij}(e) = \sum_{i,j \in V} \text{Popul}_i \frac{\text{Employment}_j * f(\text{cost}_{ij})}{\text{Travel Accessibility}_i} \quad (3)$$

182
$$\text{Travel Accessibility}_i = \sum_i^j \text{Employ}_j * f(\text{cost}_{ij}) \quad (4)$$

183
$$f(\text{cost}_{ij}) = \begin{cases} e^{\beta * \text{cost}_{ij}} \\ e^{\beta * \text{cost}_{ij}^a} \end{cases} \quad (5)$$

184 The parameters of the distance-decay function can be either estimated if data availability allows
185 it, or taken from another study, if required.

186 **Modelling approaches**

187 In order to test the predictive accuracy of models for AADT prediction, the application of
188 different models is examined. In particular, the classical ordinary least square (OLS) model
189 constitutes the starting point due to its simplicity, where the dependent variable Y is described by
190 a linear function of independent variables X with the parameters β being the least squares
191 estimates. One of the main assumptions of the model requires that the error should be spherical,
192 meaning that they should be homoscedastic and not auto-correlated.

$$Y = \beta X + \varepsilon \quad (6)$$

193 where Y is a vector with N values of the dependent variable, β is a vector with the regression
194 coefficients, X is a matrix with the independent variables and ε a vector of error terms.

195 However, the application of the OLS estimator for the statistical analysis of spatial data results to
196 residuals that are not independent, but spatially correlated, leading to the violation of the
197 assumptions of the OLS estimator.

198 Spatial econometrics was popularized by Anselin (15) and are defined as the use of regression
199 models by accounting for the impact of spatial effects (spatial dependence and heterogeneity) in
200 their specification and estimation, avoiding the statistical problems such as unreliable statistical
201 tests and biased and inconsistent estimated parameters. This is facilitated by the inclusion of a
202 spatial weight matrix (W) in the model specification that incorporates information about the
203 extent of the neighborhood, the type of the adjacency, and the relative weight that should be
204 assigned on the neighboring locations. In the transport network case, it specifies the expected
205 direction and mechanism of influence.

206 In the case of the spatial dependence, SAR models can account for it by the inclusion of relevant
207 spatial autoregressive components (16). In particular, the spatial error model assumes that the
208 spatial dependence exists in the error term of the model, and thus the spatial autoregressive
209 process is applied to it.

$$Y = \beta X + u \quad (7)$$

210
$$\text{with } u = \lambda W u + \varepsilon \quad (8)$$

211 where u the error term, λ the spatial autoregressive coefficient, W a matrix with the contiguity
 212 structure having dimensions $N \times N$, and ε a vector of independent and identically distributed (iid)
 213 error terms.

214 The spatial lag model assumes that the spatial dependence exists in the response variable and
 215 applies the spatial autoregressive process to the response variable, treating it as a lagged variable.
 216 The formulation of the model is:

$$Y = \rho WY + \beta X + \varepsilon \quad (9)$$

217 where ρ is the spatial autocorrelation parameter, and WY is the term for the lagged variable.

218 On the front of spatial heterogeneity, geographically weighted regression constitutes a technique
 219 which allows different relationships to exist in space, instead of a global relationship, and
 220 provides localized estimates of the coefficients (17).

$$Y(z) = \beta_i(z)X + u \quad (10)$$

221 Where the notation $\beta_i(z)$ indicates that the parameter describes a relationship around location u
 222 and is specific to that location (17).

223 Kriging is a geostatistical technique used for interpolation purposes. In the case of ordinary
 224 kriging, the assumption is that the unobserved value is decomposed into two terms, the local
 225 trend βX , and the error terms which are spatially correlated and their variance is assumed to
 226 follow a semivariogram relation $\gamma(h_{ij})$, as a function of the distance h between the points. In
 227 previous studies of AADT (e.g. (7)), three semivariogram functions are evaluated.

228 Exponential: $\gamma(h_{ij}; c_0, c_e, a_s) = c_0 + c_e \left(1 - e^{-\frac{h_{ij}}{a_s}} \right) \quad (11)$

229 Gaussian: $\gamma(h_{ij}; c_0, c_e, a_s) = c_0 + c_e \left(\frac{1.5h_{ij}}{a_s} - 0.5 \left(\frac{h_{ij}}{a_s} \right)^3 \right) \quad (12)$

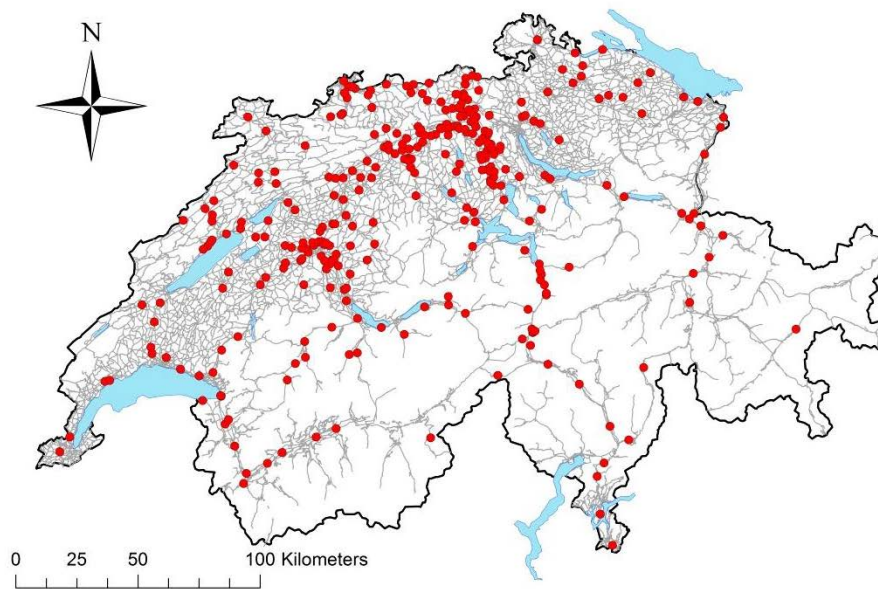
230 Spherical: $\gamma(h_{ij}; c_0, c_e, a_s) = c_0 + c_e \left(1 - e^{-\frac{h_{ij}}{a_s}} \right) \quad (13)$

231 Last, the negative binomial regression is widely used along with the Poisson regression, for the
 232 modelling of count data, accounting properly for their non-negative nature.

233 CASE STUDY

234 In order to assess the plausibility of applying a direct demand modelling approach for prediction
 235 of AADT on a nationwide network, and evaluate the capability of the accessibility-weighted

236 centrality measure to enhance the predictive accuracy of such models, a case study is designed
237 and conducted. More specifically, the network of Switzerland is employed as the study network
238 (ARE; National Transport Model, 2010), where the Federal Roads Office collects count data at
239 various locations of the network and calculates AADT values. As the basis year, the year 2010 is
240 chosen in order to be comparable with the output of the latest version of the National Transport
241 Model. In particular, for the basis year AADT data on 314 links exist which are used for the
242 model estimation as dependent values. A map of the study network along with the spatial
243 distribution of the count locations can be seen in Figure 1.



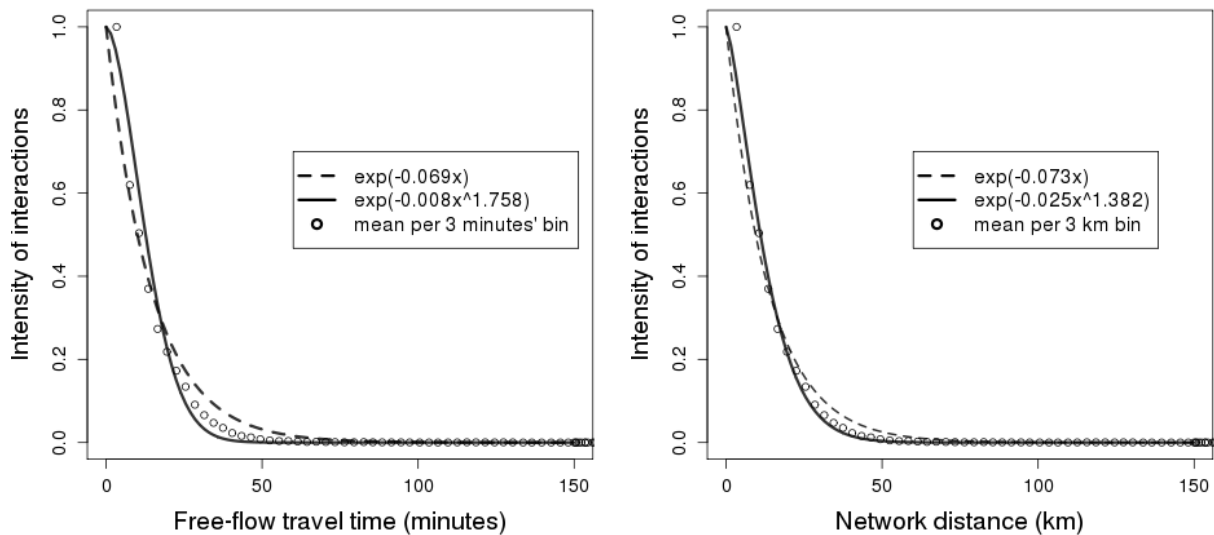
244
245 **FIGURE 1 Case study network and count locations (source: ARE, National Transport Model,**
246 **2010)**

247 **Accessibility-weighted centrality measure**

248 The first step is to proceed to the construction of the accessibility-weighted centrality measure for
249 the study network, according to the defined methodology. As mentioned before, the new measure
250 includes a distance decay function which serves the purpose of capturing the diminishing
251 intensity interactions over distance and two variations of distance decay function are checked to
252 identify the one that fits better the data, in line with a previous study (14). Obviously, different
253 parameters are associated with different trip purposes; e.g. people are willing to travel shorter
254 distances for shopping activities than for commuting to work. In our case, the interregional
255 commuting to work trips are the ones mainly contributing to the available AADT values and thus
256 the estimated parameters should correspond to this trip purpose.

257 In order to facilitate the estimation of the parameters of these two functions, we make use of an
 258 existing Origin-Destination (O-D) matrix, that corresponds to the demand for trips among all the
 259 municipalities of Switzerland in year 2010 (ARE; National Transport Model, 2010). The travel
 260 cost among all municipalities is calculated by identifying their shortest paths on the employed
 261 weighted directed network, both in terms of distance and travel time (free-flow travel time).

262 Subsequently, the portion of total daily commuters/trips that lie within each bin of given length is
 263 calculated. The length of the bin is chosen to be three minutes and three kilometers respectively.
 264 These portions take values between 0 and 1 and they are referred to as the interaction intensity. A
 265 normalization of the aforementioned portions by the percentage of the first bin (maximum)
 266 follows to ensure that we have values covering the whole range of potential values. The next step
 267 is to quantify how interaction intensity decreases over space, which actually corresponds to the
 268 parameters of the distance decay functions. The nonlinear least-squares estimates of the
 269 parameters are calculated by following the Gauss-Newton algorithm. The estimated parameters
 270 and the shape of the distance decay functions are presented in Figure 2, where the function with
 271 the two parameters is found to fit better to the data, for both cases of distance and travel time, and
 272 thus is the chosen one.



273
 274 **FIGURE 2** Estimated parameters of the distance decay functions

275 Alternatively, these parameters could be taken from previous studies as long as the employed cost
 276 metric is consistent with the one of the case study.

277 The next step is to define the origin and the destination nodes of the network that their shortest
 278 paths are accounted in the calculation of the centrality measure. Given the interregional character
 279 of the trips, a convenient choice is to employ a zonal level according to the administrative level

280 of municipalities. In this case, a node close to the centroid of each zone serves as the origin and
281 destination node for the trips of each zone, associating on it the population and the employment
282 positions of each zone. The advantage of that choice is the availability of socioeconomic data
283 aggregated on this level while the methodology can be easily applied if more disaggregated data
284 (e.g. on a hectare level) exist along with the identification of different population and employment
285 clusters, which can then replace the employed zonal analysis level.

286 Finally, the calculation of the accessibility-weighted centrality value takes place for all the links
287 where count data exists for both metric costs of network distance and travel time. For
288 computational reasons, given the finding that zones with distances more than 60 kilometers or
289 minutes between them (Figure 2), have interaction intensity close to zero, we restrict the
290 time/distance window around each link to these values. Essentially that means that only the
291 shortest paths among the origins and destinations within a radius of 60 kilometers or 60 minutes
292 around each link are found and taken into account.

293 **Independent variables**

294 In essence, the regression yields two components; one that captures the impact of supply on
295 AADT, and one that captures the impact of demand allowing to model their interaction. On the
296 supply side, variables describing the road capacity are put to use. More specifically, the
297 functional class of the road, the number of lanes, and the speed limit are the chosen explanatory
298 variables. On the demand side, the constructed accessibility-weighted centrality measure is
299 introduced for incorporating information about the magnitude and the direction of the spatial
300 interactions, serving as an approximation of spatial flows. Additional spatial variation is added on
301 the demand side by the inclusion of the public transport network density in the vicinity of each
302 road (density of public transport stops within 5 km radius), as indicative of the intensity of local
303 activities, and thus of local demand. The summary statistics of the included variables are
304 presented in Table 1. As it can be seen, in conjunction with the box-plot in Figure 3, the new
305 variable has a similar magnitude as the AADT while their correlation is close to 0.75, taking it as
306 evidence that the new variable has the capability of reproducing satisfyingly the variation of
307 demand over space.

308 **TABLE 1 Summary statistics of variables**

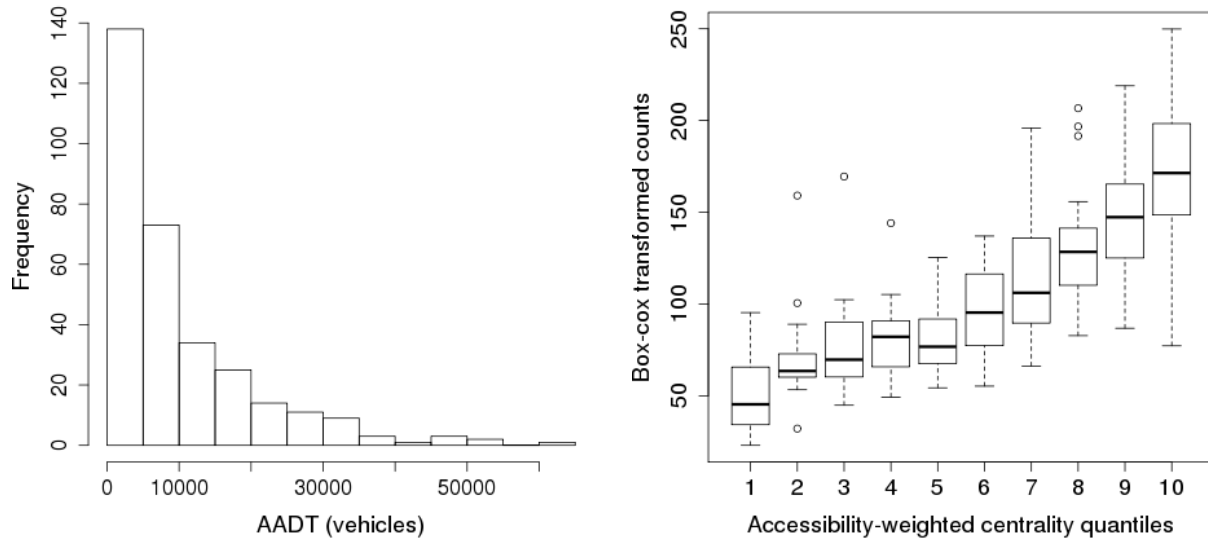
Variables	Unit	Mean	Median	St. Dev.
AADT (before transformation)	Vehicles	9834	5851	10399
Collector road	Dummy	34	-	-
Alpine road	Dummy	21	-	-
Rural major road	Dummy	28	-	-
Major road	Dummy	112	-	-
Freeway-Highway	Dummy	119	-	-
Two-lane road	Dummy	92	-	-
Three-lane road	Dummy	13	-	-
Free-flow speed	km/hr	85.75	80	22.56
Accessibility-weighted centrality	Empl. opportunities	9846	5400	11737
Public transp. density: 5km radius	stops/ sq. km	1.1	0.82	0.86

309 **AADT transformation**

310 The particularity of using count data as the dependent variable in the context of linear regression
 311 models, stems from their non-negative character which can lead to a number of shortcomings
 312 (18). In this case, either models accounting for it should be employed such as Poisson or negative
 313 binomial regression models, or the dependent variable should be transformed to conform to the
 314 assumptions of normality and/ or homoscedasticity of variance (19). Based on that, the Box-Cox
 315 transformation(20) is applied on the AADT data in order to allow the estimation of linear
 316 regression models. The transformation form is presented below while the identified ξ value for
 317 the AADT data is found to be equal to 0.414, indicating a transformation somewhere between the
 318 square and the third root.

$$Y_{tr} = \begin{cases} \frac{Y^\xi - 1}{\lambda}, & \xi \neq 0 \\ \ln Y, & \xi = 0 \end{cases} \quad (14)$$

319 Given the high correlation of the centrality variable with the AADT, we choose to apply an
 320 identical Box-Cox transformation to it, to maintain their strong linear relation in the model. The
 321 histogram of the AADT values before the transformation is presented in Figure 3 (left side),
 322 while on the right side the box-plot of the transformed centrality quantiles are plotted to show
 323 their clear linear correlation with the transformed AADT values.



324
 325 **FIGURE 3 Histogram of the AADT data and box-plot of the centrality with respect to**
 326 **AADT**

327 It should be noted that the involved data processing, models estimation, and network processing
 328 are undertaken with the statistical programming language R (21), making use of different
 329 available packages (22–24).

330 **MODEL ESTIMATION - RESULTS**

331 In this section, a set of different models is estimated and evaluated in order to draw safe
 332 conclusions on both the newly constructed variable and also on models’ capabilities. In addition
 333 to the already tested models in the literature, the family of spatial simultaneous autoregressive
 334 (SAR) models is tested as well. An assessment of models predictive accuracy and comparison to
 335 the output of a traditional four-step model is conducted to show to what extent such models can
 336 constitute a trustworthy alternative.

337 At first, an OLS model is estimated to serve as the basis for the comparison and also for
 338 examining the existence of spatial autocorrelation in the residuals and thus justify if the need for
 339 the estimation of spatial regression models arises. The spatial autocorrelation is calculated in
 340 terms of the Moran’s I measure which shows that there is statistically significant autocorrelation
 341 of 0.21. The implication of this, as mentioned before, is that the estimates are biased and
 342 inconsistent and more (or less) explanatory power is attributed to them than it should.

343 Therefore, the estimation of spatial error and lag models necessitates in order to account for the
 344 autocorrelation issues. Driven by this, three spatial weight matrices are constructed based on
 345 Euclidean distance, and network cost, both in terms of time and distance, in order to evaluate the

346 direction that correlation occurs. The identification of the the spatial extent of autocorrelation in
 347 the OLS residuals is used as an indicator to define the extent of the neighborhood. In particular,
 348 for the Euclidean and the network distance, the Moran's I measure exhibits that the
 349 autocorrelation exists up to a radius of 20 and 30 kilometers respectively. In the case of network
 350 time, the autocorrelation remains significant up to a radius of 25 minutes of free-flow travel time.
 351 The last part of the construction of the spatial weight matrices is to determine the weight that
 352 should be assigned to each neighboring location. Making use again of the Moran's I measure, we
 353 conclude that the inverse distance metric along with a normalization of the sum of the weights of
 354 the neighboring locations to one, is the more appropriate to capture the spatial structure. The
 355 estimated coefficients for the OLS and the spatial regression models are presented in Table 2.

356 **TABLE 2 Estimated coefficients for the different models**

Indep. Variables	OLS		Sp. error model		Sp. lag model	
	Estimate	Sign.	Estimate	Sign.	Estimate	Sign.
Collector road	23.50	***	22.17	***	9.42	
Alpine road	26.59	***	29.92	***	15.88	.
Rural major road	39.59	***	37.79	***	24.98	***
Major road	41.05	***	39.91	***	25.43	***
Freeway-Highway	53.43	***	53.36	***	39.84	***
Two-lane road	28.35	***	26.49	***	24.93	***
Three-lane road	76.20	***	73.65	***	74.20	***
Free-flow speed	0.16	.	0.18	*	0.14	.
Acc.-weighted centrality	0.26	***	0.25	***	0.24	***
Public transp. density: 5km	7.16	***	6.89	***	5.69	***
lamda	-		0.55	***	-	
rho	-		-		0.20	***
<i>Adjusted R</i>	0.978		-	-		
<i>Akaike Inf. criterion</i>	2663		2625		2641	
<i>Moran's I measure</i>	0.21***		-0.006		0.13	***
<i>No. of observations</i>						314

*Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1*

357

358 In summary, the OLS coefficients of the functional class variables have the expected order of
 359 magnitude, while the impact of the number of lanes and the free-flow speed is in line with
 360 expectations. The demand relevant variables, public transport density and accessibility-weighted
 361 centrality, have positive impact and they are statistically significant. It should be mentioned that
 362 they centrality value with the distance decay function as a relationship of the travel time distance
 363 is found to be slightly more statistically significant, and thus the one employed. The same pattern
 364 can be observed in the estimated coefficients of the spatial models, with the spatial autoregressive
 365 and autocorrelation parameters found to be statistically significant. In terms of goodness-of-fit

366 measures, the Akaike information criterion shows that the spatial error model is the best one
367 among the three.

368 The next estimated model is GWR, which aims to resolve spatial heterogeneity issues and it is
369 calculated by taking into account an adaptive bandwidth. The results are reported in Table 3.

371 Interestingly, the statistics of the centrality variable's coefficient show that it has relatively low
372 variation over space, providing further evidence on its ability to approximate interregional
373 demand patterns.

374 **TABLE 2 Estimated coefficients for GWR model**

Variable	Min.	1 st Quant.	Median	3 rd Quant.	Max.
Collector road	2.94	21.28	27.46	27.45	35.26
Alpine road	0.59	23.92	29.73	29.59	37.03
Rural major road	16.43	33.75	42.43	41.64	50.13
Major road	14.54	37.16	44.27	44.76	54.77
Freeway-Highway	23.10	45.31	58.42	58.48	69.13
Two-lane road	18.57	24.55	30.60	30.45	34.86
Three-lane road	53.07	67.87	76.61	76.16	84.54
Free-flow speed	-0.31	-0.02	0.12	0.12	0.28
Acc.-weighted centrality	0.08	0.23	0.26	0.27	0.31
Public transp. density: 5km	-3.25	4.67	6.60	6.34	9.31
Local R square	0.86	0.90	0.92	0.92	0.94

375
376 The negative binomial regression results are not reported, but the estimates exhibit the same
377 patterns as in the OLS model. In this particular case, the untransformed AADT and centrality
378 variables are employed.

379 Evaluation of predictive accuracy of models

380 The developed models are evaluated in terms of their predictive accuracy, both for in-sample and
381 out-of-sample. For the out-of-sample, an 80% share of the count locations are randomly chosen
382 and used for the estimation of the model while the remaining 20% is used for the validation part.
383 Given the relatively low number of observations, the out-of sample predictive accuracy of the
384 model exhibits variation and in order to account for it, a number of 100 replications is performed
385 in order to draw safer conclusions and the corresponding mean values are reported.

386 The following five accuracy measures are calculated in order to allow the evaluation to take
387 place. Mean percentage error (MPE) and mean absolute percentage error (MAPE) are easily
388 interpretable measures, having the main disadvantage though that they are influenced by outliers.
389 Symmetric mean absolute percentage error (SMAPE) is a similar measure which has the

390 advantage that it corrects for outlier's influence. Median absolute percentage error (MdAPE) has
 391 the advantage that it is not influenced by outliers and can provide an overview of the distribution
 392 of the errors in conjunction with MPE. Mean squared error (MSE) because of the quadratic term
 393 is influenced heavily by the outliers. An overview of the employed accuracy measures is given by
 394 Makridakis and Hibon (25), where they conclude that for forecasting purposes MSE and SMAPE
 395 are found to be the more preferable measures. It should be noted that AADT predicted values are
 396 reversely transformed before the calculation of the measures. The formulas of the accuracy
 397 measures are given below with \hat{Y}_i the predicted value, while the results are reported in Table 4.

$$MPE = \frac{1}{n} \sum_{i=1}^n \frac{\hat{Y}_i - Y_i}{Y_i} * 100 \quad (15)$$

398 $MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{\hat{Y}_i - Y_i}{Y_i} \right| * 100 \quad (16)$

$$SMAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{Y_i - \hat{Y}_i}{\frac{Y_i + \hat{Y}_i}{2}} \right| * 100 \quad (17)$$

$$MdAPE = median \left(\left| \frac{\hat{Y}_i - Y_i}{Y_i} \right| * 100 \right) \quad (18)$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (19)$$

399 A comparison of the accuracy measures reveals similar patterns for both in-sample and out-of-
 400 sample. In particular, among the variations of SAR models, the one that employs a spatial matrix
 401 based on the free-flow time distance gives slightly better results. Among the kriging models, it
 402 can be concluded that the one with the exponential semivariogram has better accuracy.

403 The negative binomial model yields the results with the lower predictive accuracy, providing
 404 support to the argument of the necessity of transforming the dependent variable that does not
 405 conform to the assumptions of normality.

TABLE 3 In-sample and out-of-sample predictive accuracy of estimated models

	Model	MdAPE	MPE	MAPE	MSE	SMAPE
In-sample predictive accuracy	OLS	25.27	12.98	35.38	1.38E+07	7.59
	Negative binomial	23.25	24.61	45.01	2.54E+07	8.35
	Sp. error: Eucl. distance	22.10	12.17	33.96	1.18E+07	7.35
	Sp. error: Netw. distance	21.41	11.95	33.56	1.14E+07	7.26
	Sp. error: Netw. fftt	21.67	11.91	33.44	1.13E+07	7.24
	Sp. lag: Eucl. distance	23.87	12.05	33.95	1.33E+07	7.36
	Sp. lag: Netw. distance	24.18	12.05	34.09	1.31E+07	7.36
	Sp. lag: Netw. fftt	24.23	11.95	34.09	1.29E+07	7.37
	GWR	19.83	8.32	27.86	7.69E+06	6.22
	National model (4-step)	17.45	18.12	19.79	1.29E+07	4.36
out-of-sample predictive accuracy	OLS	26.13	14.58	37.68	1.56E+07	7.97
	Negative binomial	26.04	27.32	48.40	4.31E+07	8.83
	Sp. error: Eucl. distance	26.12	15.14	38.35	1.55E+07	8.03
	Sp. error: Netw. distance	26.06	15.20	38.38	1.56E+07	8.04
	Sp. error: Netw. fftt	26.13	14.36	38.23	1.57E+07	8.06
	Sp. lag: Eucl. distance	26.23	14.32	37.05	1.56E+07	7.87
	Sp. lag: Netw. distance	26.26	13.93	37.11	1.59E+07	7.90
	Sp. lag: Netw. fftt	26.27	13.52	37.00	1.58E+07	7.92
	Kriging: Spherical	24.58	12.66	35.29	1.44E+07	7.66
	Kriging: Gaussian	25.19	12.79	35.68	1.53E+07	7.75
	Kriging: Exponential	24.54	12.48	35.11	1.39E+07	7.63
	GWR	24.80	13.10	36.50	1.45E+07	7.84
	National model (4-step)	17.74	18.06	19.86	1.27E+07	4.37

408
409 Among the estimated models, GWR has the highest in-sample accuracy while kriging has the
410 highest out-of-sample. However, in terms of MdAPE the out-of-sample difference between
411 kriging and SAR models is 1.5%, while in terms of MAPE is almost 3 percent. In terms of
412 SMAPE, all models besides negative binomial regression yield similar out-of-sample results.
413 Moreover, taking into account the fact that GWR and kriging models are aimed for interpolation
414 purposes, it can be concluded that the spatial error model gives similar results, while having the
415 advantage that it can be applied for forecasting purposes since its parameters are unbiased and
416 consistent. Interestingly, OLS out-of-sample accuracy is slightly better than spatial error model,
417 which is not the case in-sample.

418 Attempting a comparison with the Swiss national model's accuracy, which corresponds to the
419 state-of-practice four-step model used for AADT estimation, the national model outperforms the
420 estimated models. However, it has significantly higher MPE than the other models, which is of
421 similar magnitude as MAPE, revealing that it systematically overestimates the AADT. In general,
422 national transport model has higher accuracy than the other models but at the same it has to be
423 pointed out that its higher MPE value raises some concerns, given the much more data and

424 complicated models it employs. In addition, a potential source of introduced might have resulted
425 from not accounting for international commuters which can lead to underestimation of AADT
426 close to the borders.

427 Attempting a comparison with the results of a study of a similar scale (7) where kriging models
428 were estimated and the MAPE was calculated to be close to 60%, the difference in the magnitude
429 of the accuracy can be attributed to a great extent to the inclusion of the centrality measure. In the
430 case of the study conducted by Lowry for a community network, the reported MdAPE values of
431 28%, are slightly larger but of similar magnitude with our results.

432 **CONCLUSIONS**

433 In the present paper a direct demand modelling approach for AADT prediction on a nationwide
434 network is presented. It is exhibited that the construction of a variable that can account for
435 interregional flows, such as the accessibility-weighted centrality measure, can lead to a
436 significant enhancement on the accuracy of the models. In addition to the already tested models
437 in the literature, the spatial error model is estimated and it is shown that GWR and kriging
438 models are more appropriate for interpolation purposes while spatial error and OLS models have
439 the potential to be applied for forecasting purposes as well since they are estimated parameters
440 are unbiased and consistent. Under this consideration, spatial error model and OLS can be used
441 within a structural equation framework to make statements about the speed and the AADT on a
442 link level, accounting for both their well-known interdependencies and the spatial autocorrelation
443 (26). These two constitute the minimum requirements for the transport project appraisal.

444 At last, a comparison of models predictive accuracy to the output of a traditional four-step model
445 is conducted to show that direct demand models can constitute a trustworthy alternative to more
446 advanced, but definitely more data demanding and computationally burdensome models.
447 Conceptually, it is arguable that a simplified approach cannot exhibit the predictive accuracy and
448 the sensitivity of the existing approaches (four-step or agent-based models). However, the higher
449 sensitivity might allow to address more issues, but then raises the issue if the forecast is better, as
450 there are more independent variables to forecast/fix. Furthermore, it cannot be overseen that
451 when it comes to the appraisal of public transport projects, as Flyvbjerg et al. (27) argue, the
452 quality of the demand forecasts has not been improved over the years even though more complex
453 and advanced models have been employed.

454 The developed methodology can be easily applied to different scales of network, where a finer
455 zonal analysis level and the identification of clusters of trip production and attraction can be used.
456 Moreover, it requires only publicly available socioeconomic data and can utilize different
457 available networks (e.g. Open street map).

458 **ACKNOWLEDGEMENTS**

459 This paper is based on an ongoing research project funded by the Swiss National Science
460 Foundation entitled “Models without (personal) data?”.

461 **REFERENCES**

- 462 1. Xia, Q., F. Zhao, Z. Chen, L. Shen, and D. Ospina. Estimation of Annual Average Daily
463 Traffic for Nonstate Roads in a Florida County. *Transportation Research Record: Journal*
464 *of the Transportation Research Board*, Vol. 1660, Jan. 1999, pp. 32–40.
- 465 2. Mohamad, D., K. Sinha, T. Kuczek, and C. Scholer. Annual Average Daily Traffic
466 Prediction Model for County Roads. *Transportation Research Record: Journal of the*
467 *Transportation Research Board*, Vol. 1617, Jan. 1998, pp. 69–77.
- 468 3. Desyllas, J., E. Duxbury, J. Ward, and A. Smith. Pedestrian demand modelling of large
469 cities: an applied example from London. 2003.
- 470 4. Zhao, F., and N. Park. Using geographically weighted regression models to estimate
471 annual average daily traffic. *Transportation Research Record: Journal of the*
472 *Transportation Research Board*, Vol. 1879, 2004, pp. 99–107.
- 473 5. Eom, J., M. Park, T.-Y. Heo, and L. Huntsinger. Improving the Prediction of Annual
474 Average Daily Traffic for Nonfreeway Facilities by Applying a Spatial Statistical Method.
475 *Transportation Research Record*, Vol. 1968, No. 1968, 2006, pp. 20–29.
- 476 6. Wang, X., and K. M. Kockelman. Forecasting Network Data. *Transportation Research*
477 *Record: Journal of the Transportation Research Board*, Vol. 2105, No. -1, Dec. 2009, pp.
478 100–108.
- 479 7. Selby, B., and K. M. Kockelman. Spatial prediction of traffic levels in unmeasured
480 locations: applications of universal kriging and geographically weighted regression.
481 *Journal of Transport Geography*, Vol. 29, May 2013, pp. 24–32.
- 482 8. Pulugurtha, S., and P. Kusam. Modeling Annual Average Daily Traffic with Integrated
483 Spatial Data from Multiple Network Buffer Bandwidths. *Transportation Research Record:*
484 *Journal of the Transportation Research Board*, Vol. 2291, Dec. 2012, pp. 53–60.
- 485 9. Duddu, V., and S. Pulugurtha. Principle of demographic gravitation to estimate annual
486 average daily traffic: Comparison of statistical and neural network models. *Journal of*
487 *Transportation Engineering*, No. June, 2013, pp. 585–595.
- 488 10. Lowry, M. Spatial interpolation of traffic counts based on origin–destination centrality.
489 *Journal of Transport Geography*, Vol. 36, Apr. 2014, pp. 98–105.

- 490 11. McDaniel, S. Using Origin-Destination Centrality to Estimate Directional Bicycle
491 Volumes. *assets.conferencespot.org*, 2014, pp. 1–16.
- 492 12. Shimbel, A. Structural parameters of communication networks. *Bull. Math. Biophys*, Vol.
493 15, No. 4, 1953, pp. 501–507.
- 494 13. Hansen, W. G. How Accessibility Shapes Land Use. *Journal of the American Institute of*
495 *Planners*, Vol. 25, No. 2, May 1959, pp. 73–76.
- 496 14. Halás, M., P. Klapka, and P. Kladivo. Distance-decay functions for daily travel-to-work
497 flows. *Journal of Transport Geography*, Vol. 35, Feb. 2014, pp. 107–119.
- 498 15. Anselin, L. *Spatial Econometrics: Methods and Models*. Springer Netherlands, Dordrecht,
499 1988.
- 500 16. Kissling, W. D., and G. Carl. Spatial autocorrelation and the selection of simultaneous
501 autoregressive models. *Global Ecology and Biogeography*, Vol. 17, No. 1, Jun. 2007, pp.
502 59–71.
- 503 17. Charlton, M., and S. Fotheringham. Geographically Weighted Regression, White paper.
504 *National Centre for Geocomputation. National University of Ireland Maynooth*, Sep. 2009.
- 505 18. Winkelmann, R. *Econometric analysis of count data*. 2008.
- 506 19. Osborne, J. W. Improving your data transformations: Applying the Box-Cox
507 transformation. *Practical Assessment, Research & Evaluation*, Vol. 15, No. 12, 2010, pp.
508 1–9.
- 509 20. Box, G. E. P., and D. R. Cox. An analysis of transformations. *Journal of the Royal*
510 *Statistical Society. Series B (Methodological)*, Vol. 26, No. 2, 1964, pp. 211–252.
- 511 21. R Development Core Team. R: A Language and Environment for Statistical Computing, *R*
512 *Foundation for Statistical Computing*, Vienna, Austria, 2011.
- 513 22. Csárdi, G., and T. Nepusz. The igraph software package for complex network research.
514 *InterJournal Complex Systems*, Vol. 1695, 2006, p. 1695.
- 515 23. Bivand, R., L. Anselin, O. Berke, A. Bernat, M. Carvalho, Y. Chun, C. F. Dormann et al.
516 *spdep: Spatial dependence: weighting schemes, statistics and models*, 2011.
- 517 24. Pebesma, E. J. Multivariable geostatistics in S: The gstat package. *Computers and*
518 *Geosciences*, Vol. 30, No. 7, 2004, pp. 683–691.
- 519 25. Makridakis, S., and M. Hibon. Evaluating accuracy (or error) measures. *Insead*, 1995, pp.
520 1–41.

- 521 26. Sarlas, G., and K. W. Axhausen. Localized speed prediction with the use of spatial
522 simultaneous autoregressive models. *Paper presented at the 94th Annual Transportation*
523 *Research Board Meeting, Washington D.C, 2015.*
- 524 27. Flyvbjerg, B., M. K. Skamris Holm, and S. L. Buhl. How (In) accurate are demand
525 forecasts in public works projects. *Journal of the American planning association*, Vol. 71,
526 2005.
- 527