# Use of Generalized Maximum Entropy Estimation for Freight Flows Modelling and an Application

Esra Satici[1]* and Haydar Demirhan[2]

[1]*General Directorate of Turkish Highways* and [2]*Hacettepe University*

*Abstract*: In this paper, freight transportation is taken into account. One of the models used for modelling "Origin-Destination" freight flows is log-regression model obtained by applying a log-transformation to the traditional gravity model. Freight flows between ten provinces of Turkey is analyzed by using generalized maximum entropy estimator of the log-regression model for freight flow. The data set is gathered together from the axle load survey performed by Turkish Directorate of Highways and other socioeconomic and demographic variables related with provinces of interest. Relations between considered socioeconomic and demographic variables and freight flows are figured out and results are discussed.

*Key words*: Bootstrap method, freight demand modelling, generalized maximum entropy, gravity model.

## 1. Introduction

Freight transportation is one of the most important economic activities in a country. Therefore, freight demand modelling takes an important place in the transportation engineering and social sciences. Analysis of freight flows is also important for policy making. Models used for freight demand analysis can be categorized as: macro economic models, spatial interaction models, and micro economic models.

Spatial interaction models are used in the literature for modelling of flows between origins and destinations. An objective of this type of modelling is to explain variation in the level of flows between origin-destination (OD) pairs. One of the commonly used spatial interaction models in practice is the gravity model. The gravity model is originally founded on Newton's Law of Gravity, which states that two bodies attract each other in proportion to their masses and inversely by

---
*Corresponding author.

the square of the distance between them. The simplest formulation of the gravity model has the following functional form (Ortuzar and Willumsen, 2001):

$$T_{ij} = \frac{\alpha P_i P_j}{d_{ij}^2},\qquad(1.1)$$

where $P_i$ and $P_j$ are the population of the region of origin and destination, respectively, $d_{ij}$ is the distance between these regions, and $\alpha$ is a proportionality factor. The first rigorous use of the gravity model is shown by Casey (1955) for shopping trips between towns in a region (Ortuzar and Willumsen, 2001). Later, various forms of gravity model obtained by including different variables depending on research topic.

In freight generation modelling, relations between socioeconomic characteristics of regions, use of land, and freight generation amounts of regions are considered. Freight generation is an example for cross-sectional data since it depends on gross domestic product, vehicle ownership, population and difference in geographic location. Regarding this idea, there are many studies about spatial econometrics. Spatial econometrics is a subfield of econometrics that deals with the treatment of spatial interaction (spatial autocorrelation) and spatial structure (spatial heterogeneity) in regression models for cross-sectional and panel data (Anselin, 1988).

The estimation and testing of spatial econometric models are studied by Whittle (1954). Anselin (1988) introduces the spatial econometric regression models and specification testing methods. LeSage and Pace (2008) propose the family of models that rely on a spatial auto-regression filtering for flows between regions and introduce maximum likelihood estimation of this family of models. To use this family of models when the sample size is small and parametric assumptions are violated, LeSage and Pace (2008) propose the generalized maximum entropy (GME) estimator of this family of models and illustrate this approach by using commodity flows between Italy and European countries of the Balkanise area.

After the entropy concept is developed as a measure of uncertainty by Shannon in 1948, the maximum entropy principle is formulated by Jaynes (1957) as a method for estimation and inference particularly for ill-posed and/or ill-conditioned problems. More recently, Golan, Judge and Miller (1996) develop the GME estimator in the context of non-normal disturbances. Eruygur (2005), compares the GME estimator of unknown parameters of the general linear models with the ordinary least squares (OLS) estimator by Monte Carlo simulations and concludes that the performance of the GME estimator is remarkably good, especially for small sample sizes. Ciavolino and Al-Nasser (2009) compare the GME estimator with the partial least squares estimator in the presence of outliers, missing data and multicollinearity by Monte Carlo simulations and show that the results of the GME outperform the partial least squares in terms of mean

squared error.  As seen from the literature, the GME estimation method has several advantages over the conventional maximum likelihood and least squares formulations.  The main advantages are that it is more efficient, avoids strong parametric assumptions, works well when the sample size is small, and uses prior information.

In this study, we propose to use the GME estimation method for modelling the freight flow data over log-regression models.  It is difficult to fulfil strong parametric assumptions required by conventional maximum likelihood and least squares formulations for freight flow data, and sample size is small in general.  It is very probable that the problem of freight flow modelling will be ill-posed for most of the cases.  In fact, there are a small number of countries that have great number of provinces.  Thus, we need an approach such as the GME estimation that is successful in freight flow modelling when the sample size is small.  We utilize from the advantages of the GME estimation method to model freight flow data of Turkey.  In addition, it is possible to use this modelling approach to make predictions for OD pairs, for which we have not observed any data yet.

In the second section, a description of used data set is given.  In the third section, the GME estimation method is introduced and solutions for parameters of a general linear model are given.  In the fourth section, using the axle load survey data performed by Turkish Directorate of Highways, GME estimators of a log-regression model of freight flows between top ten provinces of Turkey from the respect of the amount of generated and attracted freights are obtained, and a discussion is given in the last section.

## 2. Data Description

The data set of interest is collected by the Republic of Turkey General Directorate of Highways (GDH).  Each year, the GDH applies axle weight inspections on 45 points of highways of Turkey.  In these inspections, heavy vehicles are stopped at the roadside and weighed.  Also, surveys including freight and travel information are applied to drivers. Detailed statistics can be found in the website of GDH[1]. Our data is collected during these inspections between 2007 and 2009.  The data for these three years are aggregated by using growth coefficients obtained over 2009 annual average daily traffic density (AADT). AADT is a measure used primarily in transportation planning and transportation engineering. It is the total volume of vehicle traffic of a highway or road for a year divided by 365 days. Used growth coefficient is the ratio of 2009 AADT of the relevant road segment to sample size. If an OD pair is seen in more than one inspection, the maximum freight amount is taken for this OD pair. Included provinces are Adana,

---

[1]http://www.kgm.gov.tr/Sayfalar/KGM/SiteEng/Root/MainPageEnglish.aspx

Ankara, Bursa, Hatay, İçel, İstanbul, İzmir, Kayseri, Kocaeli, and Konya. Our data set includes, freight amount (tons/day), population, gross domestic product (GDP), employment, registered number of trucks (RNT), and distance between provinces in each OD pair. These socioeconomic variables are included because they are found to be highly correlated with total amount of transported freight (Ünal, 2009).

## 3. Generalized Maximum Entropy Estimation

Let $Y$ be the $N \times 1$ dependent variable vector, $X$ be the $N \times K$ known matrix of explanatory variables. The linear regression model is defined as below:

$$Y = X\beta + \varepsilon, \tag{3.1}$$

where $\beta$ is the $K \times 1$ vector of unknown parameters and $\varepsilon$ is the $N \times 1$ vector of unknown errors. The standard least squares estimation of $\beta$ vector of parameters is the solution of the following optimization problem:

$$\min_{\beta} \left\{ \sum_{i=1}^{N} \varepsilon_i^2, \;\; \varepsilon_i = Y_i - X_i\beta, \;\; \forall i \right\}.$$

The objective is to minimize the quadratic sum of squares function for $\beta$. The maximum entropy approach is based on the entropy objective function $H(p)$ instead of the quadratic sum of squares objective function. In order to be able to use entropy principle, the unknown parameter vector should be written in terms of probabilities. Each unknown parameter $\beta_k$ is reparameterized for $M \geq 2$ as follows (Golan, Judge and Miller, 1996):

$$\beta_k = \sum_{m=1}^{M} z_{km} p_{km}, \quad k = 1, 2, \cdots, K, \quad M \geq 2. \tag{3.2}$$

Define $Z$ as a $K \times K \cdot M$ diagonal matrix of support points. Then $\beta$ is rewritten as follows:

$$\beta = Zp^{\beta} = \begin{bmatrix} z_1' & 0 & 0 & \cdots & 0 \\ 0 & z_2' & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & z_k' \end{bmatrix} \begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_k \end{bmatrix}, \tag{3.3}$$

where $p^{\beta}$ is a $KM \times 1$ vector of probabilities or weights on support points. As can be seen, the implementation of the maximum entropy formalism allowing for unconstrained parameters starts by choosing a set of discrete points by the

researcher based on his a priori information about the value of parameters to be estimated (Eruygur, 2005). If the researcher has no prior information about the sign and magnitude of the unknown $\beta_k$, support space should be defined uniformly symmetric around zero with end points of large magnitude. For instance, for $M = 5$ and for a scalar $C$, $z_k' = [-C, -C/2, 0, C/2, C]$.

Similarly, the unknown error vector $\varepsilon$ is reparameterized as follows (Golan, Judge and Miller, 1996):

$$\varepsilon_t = \sum_{j=1}^{J} v_{tj} p_{tj}, \quad t = 1, 2, \cdots, T, \tag{3.4}$$

where, $v_t' = [v_{t1}, v_{t2}, \cdots, v_{tj}]$ is support space and $p_t' = [p_{t1}, p_{t2}, \cdots, p_{tj}]$ is a vector of unknown probabilities. $V$ is defined as a $T \times T \cdot J$ diagonal matrix of support points $v_{ij}$. Then, the unknown error vector $\varepsilon$ is rewritten as follows:

$$\varepsilon = V p^\varepsilon = \begin{bmatrix} v_1' & 0 & 0 & \cdots & 0 \\ 0 & v_2' & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & v_j' \end{bmatrix} \begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_j \end{bmatrix}. \tag{3.5}$$

In practice, discrete support spaces for both parameters and errors, supplied by the researcher, are based on economic or other prior information. The support space of errors is defined according to the Chebyshev's inequality $(-v\sigma, +v\sigma)$. Golan, Judge and Miller (1996) recommend using the "three-sigma rule" to establish bounds on the error components: the lower bound is taken as $-3\sigma_y$ and the upper bound is taken as $3\sigma_y$, where $\sigma_y$ is the standard deviation of the dependent variable. For example, if $J = 5$, then $v_t' = [-3\sigma_y, -1.5\sigma_y, 0, 1.5\sigma_y, 3\sigma_y]$ is used. With the assumption that unknown weights on the parameters and the error support for the linear regression model are independent, the unknown parameters and errors are obtained by solving the constrained optimization problem of $\max H(p^\beta, p^\varepsilon) = -p^{\beta'} \ln p^\beta - p^{\varepsilon'} \ln p^\varepsilon$.

The data constrained GME estimator of the linear regression model is defined by the following constrained maximum entropy problem (Golan, Judge and Miller, 1996):

$$\max_{p^\beta, \, p^\varepsilon} H(p^\beta, p^\varepsilon) = -\sum_{k=1}^{K} \sum_{m=1}^{M} p_{km}^\beta \ln p_{km}^\beta - \sum_{t=1}^{T} \sum_{j=1}^{J} p_{tj}^\varepsilon \ln p_{tj}^\varepsilon, \tag{3.6}$$

subject to the constraints:

$$\sum_{k=1}^{K} \sum_{m=1}^{M} x_{tk} z_{km} p_{km}^\beta + \sum_{j=1}^{J} v_{tj} p_{tj}^\varepsilon = y_t, \quad t = 1, 2, \cdots, T, \tag{3.7}$$

$$\sum_{m=1}^{M} p_{km}^{\beta} = 1, \quad k = 1, 2, \cdots, K, \tag{3.8}$$

$$\sum_{j=1}^{J} p_{tj}^{\varepsilon} = 1, \quad t = 1, 2, \cdots, T. \tag{3.9}$$

The solutions for $\hat{p}_{km}^{\beta}$ and $\hat{p}_{ij}^{\varepsilon}$ are obtained by the method of Lagrange multipliers as follows (Eruygur, 2005):

$$\hat{p}_{km}^{\beta(\text{GME})} = \frac{e^{-\sum_{t=1}^{T} \hat{\lambda}_t z_{km} x_{tk}}}{\sum_{m=1}^{M} e^{-\sum_{t=1}^{T} \hat{\lambda}_t z_{km} x_{tk}}}, \quad \Omega_k^{p^{\beta}}(\hat{\lambda}_t) = \sum_{m=1}^{M} e^{-\sum_{t=1}^{T} \hat{\lambda}_t z_{km} x_{tk}}, \tag{3.10}$$

$$\hat{p}_{tj}^{\varepsilon(\text{GME})} = \frac{e^{-\hat{\lambda}_t v_{tj}}}{\sum_{j=1}^{J} e^{-\hat{\lambda}_t v_{tj}}}, \quad \Omega_t^{p^{\varepsilon}}(\hat{\lambda}_t) = \sum_{j=1}^{J} e^{-\hat{\lambda}_t v_{tj}}. \tag{3.11}$$

Substituting (3.10) and (3.11) in (3.2) and (3.4), respectively, GME estimators of $\beta_k$ and $\varepsilon_t$ are found as:

$$\hat{\beta}_k^{\text{GME}} = \sum_{m=1}^{M} \hat{p}_{km}^{\beta(\text{GME})} z_{km}, \quad k = 1, 2, \cdots, K, \tag{3.12}$$

and

$$\hat{\varepsilon}_t^{\text{GME}} = \sum_{j=1}^{J} \hat{p}_{tj}^{\varepsilon(\text{GME})} v_{tj}, \quad t = 1, 2, \cdots, T. \tag{3.13}$$

Under the conditions that error support is symmetric around zero and the errors are independent, the GME estimator is consistent and asymptotically normal (Joshi, Hanrahan, Murphy and Kelley, 2010). In order to obtain estimates of $\beta_k^{\text{GME}}$ and $\varepsilon_t^{\text{GME}}$ over the equations (3.12) and (3.13), numerical optimization techniques should be employed to obtain a solution to the system defined by the equations (3.10) and (3.11) considering the constraints given in (3.7)-(3.9).

## 4. Analysis of the Transportation Data

In this section, we focus on the application of the GME estimation to intercity freight demand modelling of 10 provinces of Turkey. Explanatory variable $(Y)$ of this data set contains the amount of freight transported from each of $N$ origin provinces to $N$ destination provinces. Therefore, there is a strong dependence structure within $Y$, it is unsuitable to analyze this data set with the approaches based on MLE. We apply GME estimates of econometric model given by Anselin (1988) to intercity freight movement of 10 provinces of Turkey. Contrary to the

other econometric models, the model given by Anselin (1988) includes $N^2$ origin and destination pairs. The model is as formulated follows:

$$Y = \beta_{\mathrm{O}} X_{\mathrm{O}} + \beta_{\mathrm{D}} X_{\mathrm{D}} + \gamma D + \varepsilon,$$

where $Y$ is a $100 \times 1$ vector of the amount of freight transported from each of N origin provinces to $N$ destination provinces. $Y$ can be generated as origin or destination based (LeSage and Pace, 2008). We follow the origin based manner. $X_{\mathrm{O}}$ is the explanatory variable matrix and corresponding to origin provinces, $X_{\mathrm{D}}$ is the explanatory variable matrix corresponding to destination provinces, $D$ is a $100 \times 1$ explanatory variable vector including distances between OD pairs, and $\varepsilon$ includes error terms. $\beta_{\mathrm{O}}$, $\beta_{\mathrm{D}}$ and $\gamma$ include regression coefficients corresponding to origin and destination provinces, and distances between OD pairs, respectively.

Elements of vector of distances are obtained by using gravity model. Distance is taken as zero for intra-provincial transportations and distances between centres of origin and destination provinces are squared for inter-provincial transportations. To get a better linear relationship, we apply the logarithmic transformation to dependent and independent variables.

Values of $M$ and $J$ are taken as 5 following Golan, Judge, and Miller (1996), who states that the largest progress is obtained for $M = J = 5$ in sensitivity of estimates for limited data. Therefore, matrices $Z$ and $V$ are defined for $M = J = 5$. Elements of $Z$ are determined by referring the multiple regression coefficients for the amount of produced freight of provinces given by Ünal (2009). The $3 - \sigma$ rule is followed to determine elements of the matrix $V$. Thus, (3.1) is reparameterized as follows:

$$Y = XZp^{\beta} + Vp^{\varepsilon}.$$

(3.6) is maximized subject to the constraints given in (3.7)-(3.9) by using a computer program prepared in Matlab 7.

In order to test the model and determine significant variables on the amount of transported freight, we employ the bootstrap method (see, Efron and Tibshirani (1993) for details of the bootstrap method). The bootstrap estimates obtained for 1000 and 1250 bootstrap samples are very close. Thus, the number of bootstrap samples is taken as 1000. So as to test the overall model, achieved significance level (ASL) is obtained. We conclude that our overall model is significant because ASL for the overall model is less than 0.001. Bootstrap estimates of Akaike information criterion (AIC) and Schwarz's Bayes criterion (SBC) are obtained as $-349.48$ and $134.49$, respectively. Bootstrap estimates of model parameters, corresponding standard errors (SE) and obtained ASL values for each parameter are given in Table 1.

Table 1: Bootstrap estimates of model parameters and corresponding achieved significance levels

| Variable | $\beta_k^{\mathrm{GME}}$ | SE | ASL |
|---|---|---|---|
| Log(Population$_\mathrm{O}$) | 0.0546 | 0.0249 | < 0.001 |
| Log (GDP$_\mathrm{O}$) | 0.0565 | 0.0263 | < 0.001 |
| Log (RNT$_\mathrm{O}$) | 0.0378 | 0.0179 | 0.026 |
| Log (Employment$_\mathrm{O}$) | 0.0486 | 0.0224 | < 0.001 |
| Log (Population$_\mathrm{D}$) | 0.0544 | 0.0248 | < 0.001 |
| Log(GDP$_\mathrm{D}$) | 0.0563 | 0.0259 | < 0.001 |
| Log(RNT$_\mathrm{D}$) | 0.0376 | 0.0179 | < 0.001 |
| Log(Employment$_\mathrm{D}$) | 0.0484 | 0.0223 | < 0.001 |
| Log(Distance$^2$) | 0.0302 | 0.1279 | 0.999 |

According to ASL values presented in Table 1, effects of all explanatory variables are significant on the amount of transported freight but natural logarithm of square of the distance. We study the prominent provinces with respect to the generation and the attraction of freight, and as well these provinces are far from each other within the economic distance context in a way of international highway transportation. That is why; it is found that the contribution of the distance is found insignificant. We drop the natural logarithm of square of the distance from the model and obtain new bootstrap estimates.

The overall model without the effect of distance is significant according to the obtained ASL value (< 0.001). Bootstrap estimates of model parameters, corresponding standard errors and obtained ASL values for each parameter of the second model are given in Table 2.

Table 2: Bootstrap estimates of parameters of the second model and corresponding achieved significance levels

| Variable | $\beta_k^{\mathrm{GME}}$ | SE | ASL |
|---|---|---|---|
| Log(Population$_\mathrm{O}$) | 0.0610 | 0.0096 | < 0.001 |
| Log (GDP$_\mathrm{O}$) | 0.0630 | 0.0104 | < 0.001 |
| Log (RNT$_\mathrm{O}$) | 0.0421 | 0.0078 | < 0.001 |
| Log (Employment$_\mathrm{O}$) | 0.0541 | 0.0089 | < 0.001 |
| Log (Population$_\mathrm{D}$) | 0.0607 | 0.0101 | < 0.001 |
| Log(GDP$_\mathrm{D}$) | 0.0628 | 0.0112 | < 0.001 |
| Log(RNT$_\mathrm{D}$) | 0.0421 | 0.0080 | < 0.001 |
| Log(Employment$_\mathrm{D}$) | 0.0540 | 0.0094 | < 0.001 |

According to ASL values presented in Table 2, effects of all explanatory variables are significant on the amount of transported freight. We obtain smaller SE estimates than those obtained for the first model. The cause of this situation is possibly that we have multicollinearity between some explanatory variables and the distance. Exclusion of the effect of distance makes the results more precise. Bootstrap estimates of the AIC and SBC are obtained as $-423.69$ and $57.67$, respectively. Thus, the model not including the effect of distance is better in the estimation of the amount of transported freight than the first model.

Standardized estimates of regression coefficients for origin and destination provinces are given in Table 3 to compare impacts of explanatory variable on the amount of transported freight.

Table 3: Standardized estimates of regression coefficients

| Variable | Standardized $\beta_k^{\mathrm{GME}}$ |
|---|---|
| Log(Population$_O$) | 6.363149 |
| Log (GDP$_O$) | 6.048194 |
| Log (RNT$_O$) | 5.368396 |
| Log (Employment$_O$) | 6.106079 |
| Log (Population$_D$) | 6.030926 |
| Log(GDP$_D$) | 5.608030 |
| Log(RNT$_D$) | 5.238004 |
| Log(Employment$_D$) | 5.776144 |

It is seen from Table 3 that the effects of explanatory variables on amount of freight are almost the same level. However, population$_O$, GDP$_O$, employment$_O$ and population$_D$ are seen to be come forward. When the amount of freight is examined in terms of characteristics of the provinces, it is seen that the characteristic features of the province, in which the freight is produced, are more effectual than those of the destination province. The most effectual explanatory variable is population of the origin province while the least effectual explanatory variable is RNT of the destination province. According to this, population$_O$ is 1.2 times more effective than RNT$_D$. Under the results of our analysis, GDP of the origin province and population of the destination province are almost same effective on the amount of freight.

Lastly, we obtain OLS estimates of the regression parameters to see what the results are if we employ OLS method instead of GME estimation. OLS estimates of model parameters, standard errors of parameters, $p$-values corresponding to the $t$-statistics are given in Table 4.

Table 4: OLS estimates standard errors of model parameters and corresponding $p$-values

| Variable | $\beta_k^{\text{OLS}}$ | SE | $p$-value |
|----------|------|------|-----------|
| Constant | $-15.888$ | 3.857 | $< 0.001$ |
| Log(Population$_O$) | $-1.042$ | 1.085 | 0.339 |
| Log (GDP$_O$) | 0.902 | 0.309 | 0.004 |
| Log (RNT$_O$) | $-0.038$ | 0.556 | 0.946 |
| Log (Employment$_O$) | 0.576 | 1.249 | 0.646 |
| Log (Population$_D$) | 4.170 | 1.085 | $< 0.001$ |
| Log(GDP$_D$) | $-0.441$ | 0.309 | 0.158 |
| Log(RNT$_D$) | $-0.083$ | 0.556 | 0.882 |
| Log(Employment$_D$) | $-2.792$ | 1.249 | 0.028 |
| Log(Distance$^2$) | $-0.096$ | 0.013 | $< 0.001$ |

It draws attention in Table 4 that according to the $p$-values, effects of population$_D$, GDP$_O$, employment$_D$ and distance are significant and those of the rest are insignificant at the 5% significance level. We suspect from a multicollinearity and obtain tolerance and variance inflation factors (VIF), and they are presented in Table 5.

Table 5: Tolerance and VIF values fort he model parameters

| Variable | Tolerance | VIF |
|----------|-----------|-----|
| Log(Population$_O$) | 0.015 | 68.11 |
| Log (GDP$_O$) | 0.125 | 8.01 |
| Log (RNT$_O$) | 0.059 | 16.90 |
| Log (Employment$_O$) | 0.011 | 92.21 |
| Log (Population$_D$) | 0.015 | 68.12 |
| Log(GDP$_D$) | 0.125 | 8.01 |
| Log(RNT$_D$) | 0.059 | 16.90 |
| Log(Employment$_D$) | 0.011 | 92.21 |
| Log(Distance$^2$) | 1 | 1 |

As expected, all the variables but GDP$_O$, GDP$_D$ and distance are affected by certain multicollinearity patterns. The results of hypothesis tests conducted over OLS estimates are unreliable. Therefore, one should use another method for estimation.

An alternative to OLS method would be the well-known ridge regression technique. We also apply the ridge regression. For the ridge regression model, values of AIC and SBC are obtained as $-209.3$ and 274.7, respectively. These values

are even greater than those obtained for the GME estimation of our first model. Therefore, we do not present the parameter estimates obtained from the ridge regression.

As seen here, the GME estimation with bootstrapping produces more reliable and precise results than those obtained by the OLS or the ridge regression techniques.

## 5. Conclusion

In this study, we use the GME estimation to model and analyse the origin − destination matrices of freight flow. Modelling origin − destination matrices has an important role in transportation planning. Models constructed by the method of GME are used effectively when the number unknown parameters of a model is small or there are violations of assumptions.

We model the freight flow between top ten provinces of Turkey from the respect of load production and capture, including various socioeconomic variables by using the method of GME. Effects of included socioeconomic variables are found as positive. In general, characteristic features of the origin province are more effectual than those of the destination province.

In our later studies, we intend to widen our analyses by adding a spatial lag term to the model and including new road inspection data.

## Acknowledgements

## References

Anselin, L. (1988). *Spatial Econometrics*: *Methods and Models*. Kluwer Academic, Boston.

Casey, H. J. (1955). Applications to traffic engineering of the law of retail gravitation. *Traffic Quarterly* **9**, 23-35.

Ciavolino, E. and Al-Nasser, A. D. (2009). Comparing generalized maximum entropy and partial least squares methods for structural equation models. *Journal of Nonparametric Statistics* **21**, 1017-1036.

Efron, B. and Tibshirani, R. J. (1993). *An Introduction to the Bootstrap.* Chapman & Hall, New York.

Eruygur H. O. (2005). Generalized maximum entropy (GME) estimator: formulation and a Monte Carlo study. Paper presented at the VII National Symposium on Econometrics and Statistics, Istanbul, Turkey, May 26-27.

Golan, A., Judge, G. and Miller, D. (1996). *Maximum Entropy Econometrics*: *Robust Estimation With Limited Data*. Wiley, New York.

Jaynes, E. T. (1957). Information theory and statistical mechanics. *Physics Review* **106**, 620-630.

Joshi, S. R., Hanrahan, K., Murphy, E. and Kelley, H. (2010). Estimating an import demand systems using the generalized maximum entropy method. Paper presented at the Twelfth Annual Conference, Lausanne, September 9-11.

LeSage, J. P. and Pace, R. K. (2008). Spatial econometric modelling of origin-destination flows. *Journal of Regional Science* **48**, 941-967.

Ortúzar, J. de D. and Willumsen, L. G. (2001). *Modelling Transport*, 3rd edition. Wiley, New York.

Ünal, L. (2009). *Modeling of Freight Transportation on Turkish Highways*. Ph.D. Dissertation, Middle East Technical University, Ankara.

Whittle, P. (1954). On stationary process in the plane. *Biometrika* **41**, 434-449.

Esra Satici
Directorate of Strategy Development
General Directorate of Turkish Highways
İnönü Blv. No. 14, 06100, Ankara, Turkey
esra.satici@gmail.com

Haydar Demirhan
Department of Statistics
Hacettepe University
Beytepe, 06800, Ankara, Turkey.
haydarde@hacettepe.edu.tr