

# **Spatial Interaction Models: Formulations and Applications**

by

**A.S. Fotheringham**

*State University of New York at Buffalo, USA*

and

**M.E. O'Kelly**

*The Ohio State University, USA*



**KLUWER ACADEMIC PUBLISHERS**

DORDRECHT / BOSTON / LONDON

Library of Congress Cataloging in Publication Data

Fotheringham, A. Stewart.

Spatial interaction models : formulations and applications / A.  
Stewart Fotheringham, Morton E. O'Kelly.

p. cm. -- (Studies in operational regional science)

Bibliography: p.

Includes index.

ISBN 0-7923-0021-1 (U.S.)

1. Space and economics--Mathematical models. 2. Geography,  
Economic--Mathematical models. 3. Industry--Location--Mathematical  
models. I. O'Kelly, Morton E., 1955- . II. Title. III. Series.  
HB199.F65 1988

338.6'042'0724--dc19

88-29416

CIP

ISBN 0-7923-0021-1

---

Published by Kluwer Academic Publishers,  
P.O. Box 17, 3300 AA Dordrecht, The Netherlands.

Kluwer Academic Publishers incorporates  
the publishing programmes of  
D. Reidel, Martinus Nijhoff, Dr W. Junk and MTP Press.

Sold and distributed in the U.S.A. and Canada  
by Kluwer Academic Publishers,  
101 Philip Drive, Norwell, MA 02061, U.S.A.

In all other countries, sold and distributed  
by Kluwer Academic Publishers Group,  
P.O. Box 322, 3300 AH Dordrecht, The Netherlands.

All Rights Reserved

© 1989 by Kluwer Academic Publishers

No part of the material protected by this copyright notice may be reproduced or  
utilized in any form or by any means, electronic or mechanical  
including photocopying, recording or by any information storage and  
retrieval system, without written permission from the copyright owner

Printed in The Netherlands

# CHAPTER 3

---

## THE CALIBRATION OF SPATIAL INTERACTION MODELS

### 3.1 Introduction

In the previous chapter, four specific forms of spatial interaction model were derived from the general formulation in equation (1.7). These were the unconstrained, production-constrained, attraction-constrained and doubly constrained models. Common to each is the need to obtain accurate estimates of the model's parameters. This is necessary in order to forecast interactions in different systems and it is also useful in providing information on the system under investigation (see Chapter 1). The process of obtaining estimates of a model's parameters is known as model calibration.

To calibrate an interaction model it is necessary to have a known interaction matrix. The question that might be asked is "if we already have the interaction matrix, why do we need to calibrate an interaction model?" There are three major reasons why we might want to undertake such an operation:

- (i) To predict interactions at some future time period or in a different spatial system.
- (ii) To forecast the effects on interaction patterns of planned changes in spatial structure. For example, if the interaction data consisted of journeys-to-work in an urban area, it would be possible to forecast the changes in traffic patterns that would result if a major industrial development took place in some part of the city or if a new road were built.
- (ii) Model calibration yields parameter estimates that can provide information on the system under investigation and it is possible to draw conclusions about interaction behaviour from a comparison of parameter estimates. For example, if a model were calibrated in the same spatial system with migration data from two time periods, it might be possible to conclude, from a comparison of the two distance-decay parameters, that migrants are becoming less or more constrained by distance over time. The research task would then focus on explaining why this change is occurring.

This chapter discusses techniques that can be employed to obtain estimates of the parameters of the four interaction models described in Chapter 2. Because of the similarity of the production-constrained and attraction-constrained models, however, only the former model is discussed. Two different methods of calibrating each model are outlined—ordinary least squares regression (OLS) and maximum likelihood estimation (MLE). It is not the task of this book to discuss these techniques in detail (Lewis-Beck (1980) and Achen (1982) provide introductions to regression while Hanushek and Jackson (1977) provide a useful description of maximum likelihood estimation); rather, we concentrate on the application of these techniques to the calibration of spatial interaction models.

### 3.2 The Calibration of Spatial Interaction Models by Regression

In order to be calibrated by regression, a model must be in a linear format, that is, linear in terms of its parameters. Hence, the models described in equations (2.20), (2.21), (2.26) and (2.33) must be transformed into the following format:

$$y = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n \quad (3.1)$$

where, in standard regression terminology,  $y$  is the dependent variable (a function of interaction in this case), the  $x$ 's represent independent variables (functions of the  $v$ 's,  $w$ 's and  $c_{ij}$ 's),  $\alpha_0$  is a scale parameter and the remaining  $\alpha$ s represent the relationship between particular independent variables and the dependent variable. The necessary linear transformation is now described for each of the spatial interaction models.

### 3.2.1 The Unconstrained Model

The transformation of the unconstrained model,

$$T_{ij} = kv_i^{\mu} w_j^{\alpha} c_{ij}^{\beta} \quad (3.2)$$

into a linear format is achieved very easily by simply taking logarithms of both sides of the equation so that equation (3.2) becomes:

$$\ln T_{ij} = \ln k + \mu \ln v_i + \alpha \ln w_j - \beta \ln c_{ij} \quad (3.3)$$

where  $\ln$  denotes a natural logarithm. Note that purely for pedagogic purposes, the simplest form of the model is represented. The addition of other variables would not alter the basic transformations that are the focus of this discussion. Where the functional forms of the variables in each model are made explicit, as above, only the transformations using power functions are reported in the interests of brevity. For instance, if equation (3.2) contained an exponential cost function, the linear transformation would clearly be

$$\ln T_{ij} = \ln k + \mu \ln v_i + \alpha \ln w_j - \beta c_{ij} \quad (3.4)$$

If the assumptions of OLS regression are met,  $\mu'$ ,  $\alpha'$  and  $\beta'$  are unbiased and consistent estimates of  $\mu$ ,  $\alpha$  and  $\beta$ , respectively, but  $\exp(\ln k)$  is a biased estimate of  $k$  [see Heien (1968) and Haworth and Vincent (1979)]. In fact,  $k$  will always be underestimated when obtained by OLS unless the model fit is perfect. The underestimation of  $k$  results in

$$\sum_i \sum_j T_{ij}' \leq \sum_i \sum_j T_{ij} \quad (3.5)$$

so that a more accurate estimate of  $k$ ,  $k'(\text{new})$ , should be obtained after the regression in the following manner:

$$k'(\text{new}) = k'(\text{old}) \cdot \sum_i \sum_j T_{ij} / \sum_i \sum_j T_{ij}' \quad (3.6)$$

A simple example of the underestimation of total flows that can occur is given below in section 3.3 where this important, but often ignored, feature of logarithmic regression is discussed further along with two other problems that can arise in the regression-based calibration of spatial interaction models.

Potential heteroscedacity problems can occur in the calibration of equation (3.3) by regression due to the logarithmic transformation (Stronge and Schultz, 1978) although this problem can usually be solved by calibrating the model by weighted or generalised least squares regression, with a weight equal to  $(\ln T_{ij})^{1/2}$ , rather than by OLS regression. The OLS estimators are unbiased and consistent but generally have larger variances (that is, are less efficient) than GLS estimators so that use of the former is less likely to detect significant relationships. Empirical research reported by van Est and van Setten (1979) and by Nakanishi and Cooper (1975), however, questions whether there is any practical difference between the two techniques in terms of parameter estimation. Useful descriptions of the heteroscedacity problem and weighted least squares regressions are given by

Hanushek and Jackson (1977), Gujarati (1978) and Cooper and Weekes (1983).

### 3.2.2 The Production-Constrained (Attraction-Constrained) Model

Consider the following general form of the production-constrained interaction model:

$$T_{ij} = O_i \exp[\sum_n \alpha_n f_n(x_{ijn})] / \sum_k \exp[\sum_n \alpha_n f_n(x_{ikn})] \quad (3.7)$$

where  $x_{ijn}$  is the  $n$ th explanatory variable and  $f_n$  is the functional form of that variable in the model. The transformation of this general structure (equivalent to that found in the attraction-constrained model) into a form linear in parameters was first described for a specific model by Nakanishi and Cooper (1974) and has been used by Stetzer (1976) and van Est and van Setten (1979) and others. To understand the transformation, multiply together the set of flows emanating from each origin to the  $n$  destinations so that

$$\prod_j T_{ij} = O_i^n \prod_j \exp[\sum_n \alpha_n f_n(x_{ijn})] / (\sum_k \exp[\sum_n \alpha_n f_n(x_{ikn})])^n. \quad (3.8)$$

Take the  $n$ th root of both sides of the equation,

$$\frac{(\prod_j T_{ij})^{1/n} = O_i (\prod_j \exp[\sum_n \alpha_n f_n(x_{ijn})])^{1/n}}{\sum_k \exp[\sum_n \alpha_n f_n(x_{ikn})]}. \quad (3.9)$$

Divide both sides of the equation into  $T_{ij}$  and substitute for  $T_{ij}$  in the right-hand side:

$$T_{ij} / (\prod_j T_{ij})^{1/n} = \exp[\sum_n \alpha_n f_n(x_{ijn})] / (\prod_j \exp[\sum_n \alpha_n f_n(x_{ijn})])^{1/n} \quad (3.10)$$

and then by taking logarithms of both sides and rearranging, this equation is made linear in terms of its parameters

$$\ln T_{ij} - (1/n) \sum_j \ln T_{ij} = \sum_n \alpha_n [f_n(x_{ijn}) - (1/n) \sum_j f_n(x_{ijn})]. \quad (3.11)$$

Note that if  $f_n$  is a unitary function, so producing an exponential function in the original model, then the expression  $(1/n) \sum_j f_n(x_{ijn})$  is an arithmetic mean. If  $f_n$  is a logarithmic function, so producing a power function in the original model, the expression is a geometric mean.

While equation (3.11) is the form of the production-constrained model that is most often calibrated by regression, it can be rearranged to the following:

$$\ln T_{ij} = \sum_n \alpha_n [f_n(x_{ijn})] + (1/n) \sum_j \ln T_{ij} - (1/n) \sum_n \alpha_n \sum_j f_n(x_{ijn}) \quad (3.12)$$

which can be simplified to:

$$\ln T_{ij} = k_i + \sum_n \alpha_n [f_n(x_{ijn})] \quad (3.13)$$

which is merely an unconstrained model with an origin-specific constant term. Cesario (1975b) was amongst the first to note this relationship.

Apart from allowing the production-constrained model (and hence the attraction-constrained model) to be calibrated by regression, the linearising technique described above has another useful property. It can be used to obtain estimates of either origin propulsiveness or destination attractiveness. For example, suppose the patronage of state parks is of interest and the number of people visiting the parks from various places is known, but an accurate

measure of the attractiveness of each park is not available (a recurring problem in recreation and tourism studies). Simply using the number of visitors to each park as a measure of attractiveness would be misleading since some parks will obviously be closer to larger centers of population than others. A measure of attractiveness is needed that is independent of location and this can be provided as follows. Let the unknown attractiveness of a destination be denoted by  $a_j$  and let spatial separation be represented as a power function of distance. The production-constrained model can then be written as:

$$T_{ij} = O_i a_j d_{ij}^\beta / \sum_k a_k d_{ik}^\beta \quad (3.14)$$

which, in linearised form, becomes,

$$\begin{aligned} \ln T_{ij} - (1/n) \sum_j \ln T_{ij} = \\ (\ln a_j - (1/n) \sum_j \ln a_j) - \beta (\ln d_{ij} - (1/n) \sum_j \ln d_{ij}) . \end{aligned} \quad (3.15)$$

When equation 3.15 is calibrated by OLS, the constant term in the regression is an estimate of  $(\ln a_j - (1/n) \sum_j \ln a_j)$  which can be used to yield values of the unknown  $a_j$ s. Dummy variables must be used in order to obtain estimates of all the  $a_j$ s and, to avoid perfect multicollinearity, one of the destinations must be excluded from the regression. The attractiveness of that destination is then obtained from the constant term alone. For the other destinations, the attractiveness is derived from the constant term plus the relevant dummy variable parameter. If the attractiveness of the excluded destination is defined as some constant, say 1.0, then this defines  $(1/n) \sum_j \ln a_j$  and all of the other attractiveness values can then be determined relative to the one set at 1.0. In essence, this technique is akin to estimating a regression model with an origin-specific constant term [Cesario (1975b)]. Baxter and Ewing (1981) provide empirical examples of the estimation of attractiveness in this way and an example is provided below in section 3.6.

### 3.2.3 The Doubly Constrained Model

At first glance it might appear impossible to linearise the doubly constrained model from section 2.3.4 written here with a power function of distance:

$$T_{ij} = A_i O_i B_j D_j d_{ij}^\beta \quad (3.16)$$

where,

$$A_i = 1 / \sum_j B_j D_j d_{ij}^\beta \quad (3.17)$$

and

$$B_j = 1 / \sum_i A_i O_i d_{ij}^\beta . \quad (3.18)$$

However, two techniques have recently been described that achieve this task. Sen and Soot (1981) and Gray and Sen (1983) have provided a technique that separates the estimation of the distance (cost) parameter from the calculation of balancing factors. The procedure is termed the odds ratio method and involves taking ratios of interactions so that the  $A_i O_i$  and  $B_j D_j$  terms in the model cancel out. From equation (3.16),

$$T_{ij} = A_i O_i B_j D_j d_{ij}^\beta \quad (3.19)$$

$$T_{ji} = A_j O_j B_i D_i d_{ji}^\beta \quad (3.20)$$

$$T_{iu} = A_i O_i B_u D_u d_{iu}^\beta \quad (3.21)$$

$$T_{uj} = A_u O_u B_j D_j d_{uj}^\beta \quad (3.22)$$

so that

$$(T_{ij}/T_{ii}).(T_{jj}/T_{jj}) = [(d_{ij}/d_{ii}).(d_{jj}/d_{jj})]^\beta \quad (3.23)$$

The equation is then made linear in its parameters by taking logarithms,

$$\begin{aligned} \ln T_{ij} + \ln T_{jj} - \ln T_{ii} - \ln T_{jj} = \\ -\beta(\ln d_{ij} + \ln d_{jj} - \ln d_{ii} - \ln d_{jj}) \end{aligned} \quad (3.24)$$

Sen and Soot suggest using weighted least squares to counteract the heteroscedastic error terms caused by the logarithmic transformation with the weight being  $(T_{ij}^{-1} + T_{jj}^{-1} + T_{ii}^{-1} + T_{jj}^{-1})^{0.5}$ .

Once  $\beta$  is estimated, the balancing factors of the model can be obtained by iterating equations (3.17) and (3.18).

Two potential problems can arise with the above linear transform of the doubly-constrained model. The interaction matrix has to be square and it is necessary for the intrazonal interactions ( $T_{ii}$  and  $T_{jj}$ ) to be non-zero. An alternative method of linearising the doubly-constrained gravity model that does not suffer from these problems is given by Sen and Soot (1981) and uses the relationship that,

$$\begin{aligned} \ln T_{ij} - (1/n)\sum_j \ln T_{ij} - (1/m)\sum_i \ln T_{ij} + (1/mn)\sum_i \sum_j \ln T_{ij} \\ = -\beta [\ln d_{ij} - (1/n)\sum_j \ln d_{ij} - (1/m)\sum_i \ln d_{ij} + (1/mn)\sum_i \sum_j \ln d_{ij}] \end{aligned} \quad (3.25)$$

where  $(1/n)\sum_j \ln T_{ij}$  is the row mean of the  $\ln T_{ij}$ s;  $(1/m)\sum_i \ln T_{ij}$  is the column mean; and  $(1/mn)\sum_i \sum_j \ln T_{ij}$  is the grand mean.

Again, weighted least squares may be preferable to ordinary least squares. In this instance, the weight is simpler, being  $T_{ij}^{0.5}$ . This linearisation technique can be used for rectangular matrices. However, it needs to be modified when intrazonal flows are zero (Sen and Pruthi, 1983).

### 3.3 Cautionary Notes on the Calibration of Interaction Models by Regression

The transformations described above provide a useful mechanism for calibrating interaction models by regression which has the advantage of being a readily available calibration technique with well-known properties. There are, however, three potential problems that can arise with the technique. The first concerns the estimation of the constant term in an unconstrained model, which as already discussed, tends to be biased downwards; the second concerns goodness-of-fit calculations; and the third concerns the replacement of zero flows. Each is now examined in more detail.

#### 3.3.1 Bias in the Constant Term

The constraint operating in regression is that,

$$\sum_i \sum_j e_{ij} = 0 \quad (3.26)$$

where  $e_{ij}$  represents the error term in the regression. This implies that

$$\sum_i \sum_j y_{ij}' = \sum_i \sum_j y_{ij} \quad (3.27)$$

where  $y_{ij}$  represents the transformed interaction variable in one of the transformations described above and  $y_{ij}'$  is the predicted value of  $y_{ij}$ . In the case of the unconstrained model, for example,

$$\sum_i \sum_j \ln T_{ij}' = \sum_i \sum_j \ln T_{ij} \quad (3.28)$$

which does not imply that the desired constraint,

$$\sum_i \sum_j T_{ij}' = \sum_i \sum_j T_{ij} \quad (3.29)$$

is met. In fact, the constraint in (3.28) ensures that, unless the model is perfectly accurate,

$$\sum_i \sum_j T_{ij}' < \sum_i \sum_j T_{ij} \quad (3.30)$$

that is, the total flow volume predicted will be less than the actual total flow volume. If the model is imperfect, the variance of the predicted values will be less than the variance of the actual values. Consequently, there is a tendency for small flows to be overpredicted and large flows to be underpredicted. While these over- and underpredictions cancel each other out in logarithms, in terms of real flows, the underpredictions of large flows will be greater than the overprediction of small flows. Hence, the discrepancy described in (3.30) arises.

This can be demonstrated in an alternative way. The average predicted value from the model,  $\exp\{[1/(mn)]\sum_i \sum_j \ln T_{ij}'\}$ , is actually the geometric mean of the  $T_{ij}'$  values. By a proof given by Beckenbach and Bellman (1961), the geometric mean of a set of values is always less than or equal to the arithmetic mean. Hence,

$$\sum_i \sum_j \exp(\ln T_{ij}') \leq \sum_i \sum_j T_{ij} \quad (3.31)$$

To remove this discrepancy, it is necessary to revise the estimate of the constant term in the model through the use of equation (3.6) or, as Heien (1968) and Haworth and Vincent (1979) suggest, by performing the following transformation on the exponential of the estimated constant from a logarithmic regression:

$$k'(\text{new}) = k'(\text{old}).\exp(s^2/2) \quad (3.32)$$

where  $s^2$  is the sample variance of the logarithmic error terms from the regression. Since  $s^2 \geq 0$ ,  $\exp(s^2/2) \geq 1$  and  $k'(\text{new}) \geq k'(\text{old})$ .

### 3.3.2 Goodness-of-fit Calculations and Logarithmic Regressions

Generally, the calibration of models by regression is undertaken with statistical computer software such as SPSS-X, SAS, and the various micro-based packages. It is important to realise that the goodness-of-fit statistic(s) reported in such packages, which are generally based on an R-squared measure, relate the accuracy of the model in replicating the dependent variable; the latter not being interaction but one of the transformations of interaction described above. Thus, the goodness-of-fit statistics reported in such packages can be misleading: we are interested in predicting interaction, not some transformation of it. It is therefore necessary to transform the predicted dependent variable into a predicted interaction and then calculate a goodness-of-fit statistic with those predicted interactions. A discussion on goodness-of-fit statistics suitable for assessing the performance of interaction models is given below in section 3.7.



### 3.3.3 The Problem of Zero Interactions

Because all of the transformations described above involve taking a logarithm of interaction, a problem arises whenever the sampled interaction between any two points is zero since the logarithm is then undefined. Several solutions to this problem, of varying degrees of satisfaction, are possible. Perhaps the most obvious solution is simply to remove all zero interactions from the analysis. However, the resulting parameter estimates would not reflect the low volumes of interaction that occur between certain origins and destinations and so would be misleading. A second solution is to remove from the analysis all origins and destinations associated with zero interactions. However, a great deal of useful information can be lost in this way, and in particularly sparse matrices, there may be no origin that has a non-zero interaction to every destination.

The third solution is by far the most commonly used in dealing with zero interactions and it involves adding a constant to elements of the interaction matrix. Two possibilities exist here: one is to add the constant to every flow in the matrix; the other is to add the constant only to the zero flows. In practice, there seems little difference between the two methods in terms of the resulting parameter estimates. In both cases, some uncertainty exists over the value of the constant to be added. Probably the most frequently encountered method of dealing with zero interactions is to add one to every zero flow and this can be justified on the grounds that the flows recorded are generally integer and one is the closest approximation to zero. It also ensures that the logarithmic interactions have a minimum of zero. However, Sen and Soot (1981) provide a theoretical justification for the addition of 0.5 rather than one although their justification does not take into account possible effects on parameter estimates.

None of the above problems is encountered in the second method of calibration, maximum likelihood estimation (MLE), which is now discussed. There is a trade-off, however, in that MLE calibration routines tend to be less accessible than their regression counterparts but in section 3.8 a computer programme developed especially for the calibration of spatial interaction models by maximum likelihood is discussed.

### 3.4 The Calibration of Models by Maximum Likelihood Estimation

In essence, the technique of MLE is to find parameter estimates that maximise the likelihood of observing a sample set of interactions from a theoretical distribution. The steps involved in the calibration include identifying a theoretical distribution for the interactions, maximising the likelihood function of this distribution with respect to the parameters of the interaction model, and then deriving equations that ensure the maximisation of the likelihood function. For convenience, the logarithm of the likelihood function is usually used since this is at a maximum whenever the likelihood function is at a maximum (see Pickles, 1986, for a demonstration of this). Parameter estimates that maximise the likelihood function are termed maximum likelihood estimates. Maximum likelihood estimators have several desirable properties: they are consistent, asymptotically efficient and are asymptotically normally distributed. This latter property is particularly useful in significance testing and further comment on this subject is made below.

The method of obtaining parameter estimates by MLE is described for each of the four interaction models although again the method is identical for the production-constrained and attraction-constrained models so is only discussed in terms of the former.

#### 3.4.1 The Unconstrained Model

Flowerdew and Aitkin (1982) state that interactions can be considered to be the outcome of a Poisson process if it is assumed that there is a constant probability of any individual

in  $i$  moving to  $j$ , that the population of  $i$  is large, and that the number of individuals interacting is an independent process. Consequently, the probability that  $T_{ij}$  is the number of people recorded as moving between  $i$  and  $j$  is given by,

$$p(T_{ij}) = \frac{\exp(-T_{ij}') \cdot T_{ij}'^{(T_{ij})}}{T_{ij}!} \quad (3.33)$$

where  $T_{ij}'$  is the expected outcome of the Poisson process. Note that this is distinguished from the observed value  $T_{ij}$ , the latter being subject to sampling and measurement errors and therefore fluctuates around the expected value,  $T_{ij}'$ . Since  $T_{ij}'$  is unknown and unobservable, it has to be estimated from some theoretical model such as the unconstrained model in equation (3.2).

Consider the log-likelihood of a set of observed flows  $\{T_{ij}\}$  where each flow is the outcome of a particular Poisson process. This log-likelihood,  $L^*$ , can be represented as:

$$L^* = \sum_i \sum_j \ln[\exp(-T_{ij}') \cdot T_{ij}'^{(T_{ij})} / T_{ij}!] \quad (3.34)$$

which is equivalent to

$$L^* = \sum_i \sum_j (-T_{ij}' + T_{ij} \ln T_{ij}' - \ln T_{ij}!) \quad (3.35)$$

Since  $T_{ij}$  is given,  $\ln T_{ij}!$  can be ignored in the maximisation and  $L^*$  will be a maximum when

$$Z = \sum_i \sum_j (T_{ij} \ln T_{ij}' - T_{ij}') \quad (3.36)$$

is at a maximum. Hence, the parameter estimates associated with  $T_{ij}'$  that maximise  $Z$  are required. These are the estimates of the parameters in equation (3.2) that maximise the expression for  $Z$  in equation (3.36). Using calculus, these estimates are obtained when

$$\partial Z / \partial \xi = \sum_i \sum_j T_{ij}' \ln x_{ij} - \sum_i \sum_j T_{ij} \ln x_{ij} = 0 \quad (3.37)$$

where  $x_{ij}$  is an independent variable in equation (3.2) and  $\xi$  is the parameter associated with that variable. For example, if  $k$  in equation (3.2) is defined as  $e^{\mu}$ , the constraint equation for  $\mu$  is

$$\partial Z / \partial \mu = \sum_i \sum_j T_{ij}' - \sum_i \sum_j T_{ij} = 0 \quad (3.38)$$

since  $e$ , the variable associated with  $\mu$  is a constant. Equation (3.38) indicates that the sum of the predicted interactions ( $\sum_i \sum_j T_{ij}'$ ) will be equal to the sum of the observed interactions ( $\sum_i \sum_j T_{ij}$ ). This means that a total flow constraint is met automatically in the maximum likelihood calibration of the unconstrained gravity model.

In a similar manner, the maximum-likelihood equation for  $\beta$ , the distance-decay parameter, will be (from 3.37):

$$\sum_i \sum_j T_{ij}' \ln d_{ij} - \sum_i \sum_j T_{ij} \ln d_{ij} = 0 \quad (3.39)$$

which can be interpreted as a cost constraint.

### 3.4.2 The Production-Constrained (Attraction-Constrained) Model

Interactions can be assumed to have a multinomial distribution (Batty and Mackie, 1972), in which case the log-likelihood of observing a set of flows is

$$L^* = \sum_i \sum_j T_{ij} \ln p_{ij} \quad (3.40)$$

where  $p_{ij}$  represents the predicted probability of moving between  $i$  and  $j$  and is defined as

$$p_{ij} = T_{ij}' / \sum_j T_{ij}' \quad (3.41)$$

Alternatively,  $p_{ij}$  can be defined as the product of two other probabilities:

$$p_{ij} = p_{ji} \cdot p_i \quad (3.42)$$

where  $p_{ji}$  is the conditional probability of interacting with  $j$  given one originates at  $i$ , and  $p_i$  is the probability of an interaction originating in  $i$ . In a production-constrained model,

$$\sum_j p_{ji} = 1 \quad \text{for all } i \quad (3.43)$$

and, clearly,

$$\sum_i p_i = 1 \quad (3.44)$$

The probability  $p_{ji}$  is given by a production-constrained model,

$$p_{ji} = A_i \exp[\sum_h \alpha_h f_h(x_{ijh})] \quad (3.45)$$

so that the objective is to determine the estimates of the  $\alpha_h$ s which maximise equation (3.40) subject to the constraints in equations (3.43) and (3.44). It is relatively straightforward to demonstrate that the solution to this is a series of equations of the form:

$$\sum_i \sum_j T_{ij}' f_h(x_{ijh}) = \sum_i \sum_j T_{ij} f_h(x_{ijh}) \quad \text{for all } h. \quad (3.46)$$

In each case the MLE of  $\alpha_h$  is therefore obtained when the constraint in equation (3.46) is met. For example, in the case of a power distance function, the estimate of  $\beta$  is obtained when

$$\sum_i \sum_j T_{ij}' \ln d_{ij} = \sum_i \sum_j T_{ij} \ln d_{ij} \quad (3.47)$$

which is the same constraint as in the maximum-likelihood calibration of the unconstrained model. The only difference in the two models is that the latter has only a single constant to be estimated whereas the production-constrained model has a separate constant,  $A_i$ , to be estimated for each origin. The estimate of  $A_i$  is obtained when

$$\sum_j T_{ij}' = \sum_j T_{ij} \quad \text{for all } i. \quad (3.48)$$

### 3.4.3 The Doubly Constrained Model

The maximum likelihood estimator of the cost parameter in the doubly constrained model is obtained through the same form of constraint equation as given in (3.46) above. The only difference in the calibration of this model from that of the production-constrained model being that an extra set of parameters, the  $B_j$ s, is estimated from the destination constraint set

$$\sum_i T_{ij}' = \sum_i T_{ij} \quad \text{for all } j. \quad (3.49)$$

### 3.5 An Algorithm for Maximum-Likelihood Calibration

As a demonstration of the general maximum-likelihood calibration of an interaction model, consider the estimation of  $\beta$ , the distance(cost)-decay parameter, in a doubly constrained model with a power function of distance. The constraint equation is that given in (3.47) and the general procedure for obtaining an estimate of  $\beta$  using this equation is outlined in Figure 3.1. Starting values for  $\beta'$  (usually 1.0 in a negative power function) and the set of  $B_i$  values (usually 1.0 for each  $B_i$ ) are chosen. The initial values of  $A_i$  are then calculated and are input into the updated calculation of the  $B_i$ s. This iterative calculation of the balancing factors continues until all the values are stable under further iteration. The stable  $A_i$  and  $B_i$  values are then input into the model and a set of predicted interactions is obtained (with  $\beta = 1.0$ ). If the constraint equation is met at this stage, the value of  $\beta'$  is retained; if it is not, then  $\beta'$  is changed and the whole cycle repeated until the constraint equation for  $\beta$  is met.

Variations in the above procedure only occur in the way in which  $\beta'$  is changed and algorithms of varying degrees of speed and complexity exist for this purpose. The simplest and slowest, for example, is a straightforward iteration where  $\beta'$  is changed by a set amount on each iteration. Other, faster procedures include first order iteration, Newton-Raphson techniques and the Secant method on which further discussion can be found in Batty and Mackie (1972), Batty (1976a) and Cheney and Kincaid (1980). Here, we will briefly outline Newton's method because it is probably the most widely applied procedure for solving a set of nonlinear equations and because it forms the basis of the SIMODEL algorithm which is discussed below.

#### 3.5.1 Newton's Method for Solving One Nonlinear Equation

For any given model, the maximum-likelihood constraint equation for each parameter in that model can be represented in a general form by,

$$f(\beta) = 0 \quad (3.50)$$

where  $\beta$  represents the parameter to be estimated in the model. The calibration of the model takes place by finding the value of  $\beta$  generating  $f(\beta) = 0$  and this is denoted by  $\beta'$  in Figure 3.2. The specific form of  $f(\beta)$  is the nonlinear equation

$$\sum_i \sum_j T_{ij}' f(x_{ij}) - \sum_i \sum_j T_{ij} f(x_{ij}) = 0 \quad (3.51)$$

where  $T_{ij}'$  is a function of  $\beta$ . The function is differentiable over all  $\beta$  and so the graph of  $f(\beta)$  against  $\beta$  has a definite slope at each point and hence has a unique tangent at that point. Consider the tangent to the curve at the point  $\beta_0$ ,  $f(\beta_0)$ , denoted by the line  $P\beta_1$  in Figure 3.2. From simple trigonometry, the line  $P\beta_1$  at  $\beta_0$  is defined in terms of the linear function,

$$L(\beta) = f'(\beta_0)(\beta - \beta_0) + f(\beta_0) \quad (3.52)$$

Figure 3.1 Maximum Likelihood Calibration Procedure

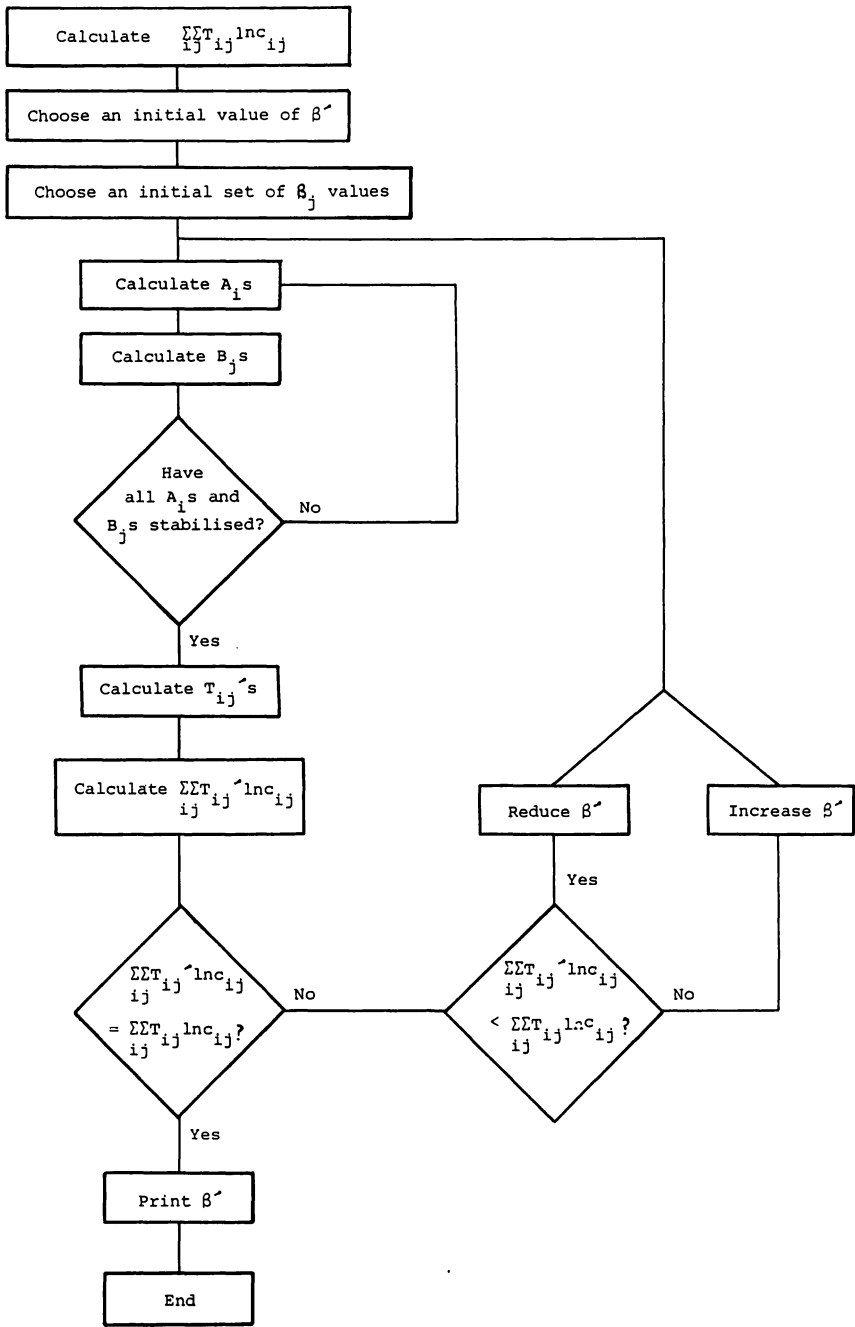
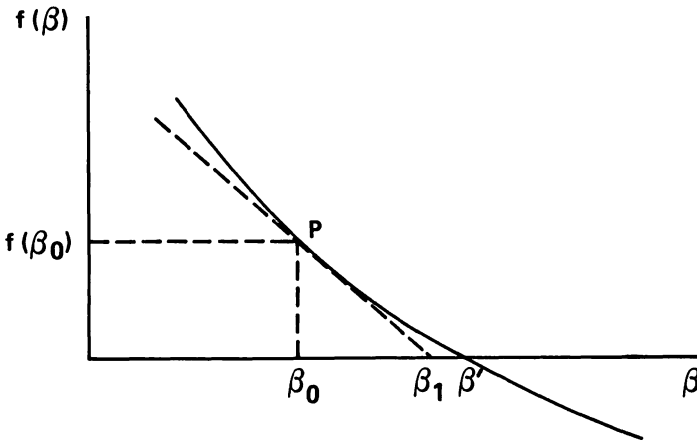


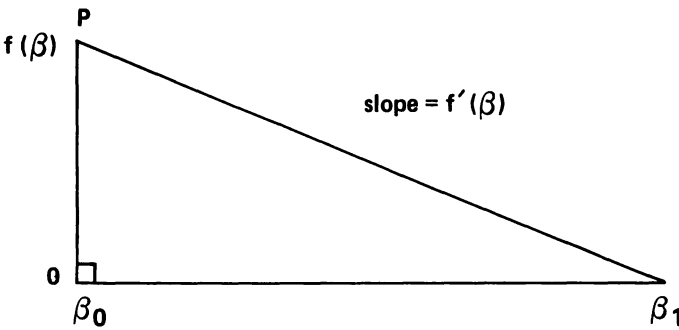
Figure 3.2 Graph of  $f(\beta)$  against  $\beta$



where  $f'(\beta_0)$  is the slope of the curve at  $\beta_0$  (and so is the slope of the tangent to this point) and  $f(\beta_0)$  is the value of  $f(\beta)$  at  $\beta_0$ . Thus, as  $\beta$  increases beyond  $\beta_0$ , the value of  $L(\beta)$  decreases by  $f'(\beta)(\beta - \beta_0)$ , as shown in Figure 3.3. The zero of  $L(\beta)$  is an approximation to the zero of  $f(\beta)$  and from equation (3.52) this occurs when

$$\beta_1 = \beta_0 - [f(\beta_0)/f'(\beta_0)] . \quad (3.53)$$

Figure 3.3 Determining  $\beta_1$



Thus, starting with a point  $\beta_0$  (an initial guess), a more accurate estimate,  $\beta_1$ , is obtained from the formula in equation (3.53). This can be used to generate a more accurate estimate,  $\beta_2$ , by

$$\beta_2 = \beta_1 - [f(\beta_1)/f'(\beta_1)] \quad (3.54)$$

and so on until  $\beta$  converges. In general, the formula that is iterated is,

$$\beta_{n+1} = \beta_n - [f(\beta_n)/f'(\beta_n)] \quad \checkmark \quad (3.55)$$

The function  $f'(\beta)$  is simply the derivative of the nonlinear function with respect to the parameter being estimated.

### 3.5.2 Newton's Method for Solving a Set of Nonlinear Equations

The following discussion is based largely on that of Johnston (1982) in which an extension of the previous one-parameter estimation procedure is made so that it can be used to estimate any number of parameters.

Single parameter estimation deals with a one-dimensional curve as in Figure 3.2. For  $n$  parameters it is necessary to deal with an  $n$ -dimensional hypersurface. The derivative  $f'(\beta)$  at  $\beta = \beta_k$  is a number  $f'(\beta_k)$  that satisfies the condition

$$[f(\beta_k) - f(\beta_p)]/(\beta_k - \beta_p) - f'(\beta_k) \rightarrow 0 \text{ as } \beta_p \rightarrow \beta_k \quad (3.56)$$

which is equivalent to

$$f(\beta_k) - f(\beta_p) - f'(\beta_k) \cdot (\beta_k - \beta_p) \rightarrow 0 \text{ as } \beta_p \rightarrow \beta_k. \quad (3.57)$$

The analogue of  $f'(\beta_k)$  in  $n$  dimensions is an  $n \times n$  matrix,  $J(\beta_k)$ , that satisfies the condition

$$f(\beta_k) - f(\beta_p) - J(\beta_k) \cdot (\beta_k - \beta_p) \rightarrow 0 \text{ as } \beta_p \rightarrow \beta_k \quad (3.58)$$

where  $J(\beta_k)$  is a Jacobian matrix whose entries are the partial derivatives,

$$J_{ab} = \partial f_a(\beta)/\partial \beta_b \quad (3.59)$$

where the indices  $a$  and  $b$  refer to parameters to be estimated. Then, given a set of initial estimates of  $\beta$ , a more accurate set of estimates,  $\beta_k$ , is given by

$$\beta_{k+1} = \beta_k - [J(\beta_k)]^{-1} \cdot f(\beta_k). \quad (3.60)$$

Computationally, the above procedure involves the following steps;

- (i) Compute  $f(\beta_k)$  and  $J(\beta_k)$
- (ii) Invert  $J(\beta_k)$
- (iii) Compute  $\beta_k - [J(\beta_k)]^{-1} \cdot f(\beta_k)$
- (iv) Using the value in (iii) as  $k$  return to (i)

### 3.5.3 Standard Errors of Maximum Likelihood Parameter Estimates

Consistent estimators of the asymptotic variances of the parameter estimates by maximum likelihood are given by the diagonals of the matrix  $-H^{-1}$  where  $H$  represents an  $n \times n$  matrix of second derivatives;  $n$  being the number of parameters in the model. The typical element of the matrix  $H$ ,  $H_{ab}$ , is defined as

$$H_{ab} = \partial^2 L^*(\beta')/(\partial \beta_a \partial \beta_b) \quad (3.61)$$

where  $\beta'$  represents the vector of parameters estimated in the model and  $L^*(\beta')$  represents the logarithm of the likelihood function that is maximised in order to obtain the values in

$\beta'$ . That is,

$$L'(\beta') = \sum_i \sum_j T_{ij} \ln p_{ij}(\beta') \quad (3.62)$$

where  $T_{ij}$  represents the observed number of interactions between  $i$  and  $j$  and  $p_{ij}(\beta')$  represents the model prediction of the probability of interaction between  $i$  and  $j$ . In the case of interaction models with only one parameter estimated by maximum likelihood, the computation of the variance of the estimate simplifies to:

$$\sigma^2(\beta') = - \partial^2 L'(\beta') / \partial \beta'^2 \quad (3.63)$$

More details on the variance of maximum likelihood parameters are given by Mood and Graybill (1963) and by Cox (1970). Applications to spatial interaction modelling are provided by Giles and Hampton (1981) and computational notes are provided in Williams and Fotheringham (1984).

Because the MLEs are asymptotically normally distributed, their variances can be used to examine the significance of individual parameters in the usual manner. With maximum likelihood estimation, however, an alternative significance testing procedure exists through the construction of log-likelihoods. The statistic,  $\tau$ , defined as

$$\tau = 2.T.[L'(\beta') - L'(\beta'; \beta_h = 0)] \quad (3.64)$$

where  $T$  is the total volume of interaction is asymptotically chi-square distributed with 1 degree of freedom. The notation  $L'(\beta'; \beta_h = 0)$  represents the log-likelihood when the parameter whose significance is being examined is set to zero. A value of  $\tau$  significantly different from zero allows one to reject the hypothesis that  $\beta_h = 0$ . Further details on the use of relative likelihood statistics in significance testing can be found in Hathaway (1975), Horowitz (1981) and Stopher and Meyburg (1979).

### 3.6 An Empirical Comparison of OLS and ML Parameter Estimates

Tobler (1983) reports data on 1970-1980 interregional migration between the nine major census divisions of the United States. These data, along with relevant distances and populations, are reported in Tables 3.1 and 3.2. Each of the four interaction models defined in Chapter 2 is calibrated by both OLS and MLE using these data. In all cases the intrazonal flows are ignored since the emphasis of the model calibration is on explaining interregional migration. Consequently, the doubly constrained model is calibrated by regression using equation (3.25) rather than (3.24). In this example, population is used as the sole measure of destination attractiveness and the effect of spatial separation is represented as a negative power function of distance. Thus, a maximum of three parameters can be estimated: an origin propulsiveness parameter,  $\mu$ ; a destination attractiveness parameter,  $\alpha$ ; and a distance-decay parameter,  $\beta$ . The estimates for each model are reported in Table 3.3.

It is clear that the OLS and MLE parameter estimates are not identical for any of the four models. This is due to different criteria involved in the OLS and MLE procedures. In OLS the criterion is to minimise the sum of squared differences between the predicted and actual independent variables; in MLE the criterion is to meet a non-linear constraint. The parameter estimates are similar, however, because when the ML constraints are met, it is likely that the predicted and observed interactions will be similar and that the sum of squared differences between them will be near the minimum.

There is a noticeable trend in the estimated distance-decay parameters as constraints are added to the interaction model. This can be explained in terms of the ability of each model



to replicate the set of migration flows. This ability is measured in Table 3.3 by the value of the goodness-of-fit statistic, SRMSE, defined in the subsequent section. The statistic has a value of zero when the observed flows are replicated perfectly and increasing values are indicative of increasingly poor model accuracy. Consequently, from Table 3.3 it is clear that adding constraints to the interaction model formulation increases accuracy, as would be expected.

Table 3.1 Interregional Migration between Census Regions\*

	1	2	3	4	5	6	7	8	9
1	0	180048	79223	26887	198144	17995	35563	30528	110792
2	283049	0	300345	67280	718673	55094	93434	87987	268458
3	87267	237229	0	281791	551483	230788	178517	172711	394481
4	29877	60681	286580	0	143860	49892	185618	181868	274629
5	130830	382565	346407	92308	0	252189	192223	89389	279739
6	21434	53772	287340	49828	316650	0	141679	27409	87938
7	30287	64645	161645	144980	199466	121366	0	134229	289880
8	21450	43749	97808	113683	89806	25574	158006	0	437255
9	72114	133122	229764	165405	266305	66324	252039	342948	0

\* The census regions are: 1-New England, 2-Mid Atlantic, 3-East North-Central, 4-West North-Central, 5-South Atlantic, 6-East South-Central, 7-West South-Central, 8-Mountain, and 9-Pacific.

Table 3.2 Populations and Distances in Miles between Region Centroids

Region	Population	1	2	3	4	5	6	7	8	9
1	11848000	-								
2	37056000	219	-							
3	40266000	1009	831	-						
4	16327000	1514	1336	505	-					
5	29920000	974	755	1019	1370	-				
6	13096000	1268	1049	662	888	482	-			
7	19025000	1795	1576	933	654	1144	662	-		
8	8289000	2420	2242	1451	946	2278	1795	1278	-	
9	25476000	3174	2996	2205	1700	2862	2380	1779	754	-

The variation in the estimated distance-decay parameter between the models is thus probably related to the variation in the accuracy of the models. Because the doubly constrained model is most accurate it could be assumed that the estimate derived from this model is the most accurate representation of the true relationship between interaction and distance. The difference between the distance-decay parameter estimates of the two singly-constrained models can also be explained in terms of the accuracy of the two models. The attraction-constrained model is more accurate in replicating migration flows than the production-constrained model and hence the parameter estimate of the former is more similar to that of the doubly constrained model. A possible reason for the superior performance of the attraction-constrained model is that the population size of a region is probably a more accurate measure of propulsiveness than of attractiveness. Many other variables such as climate, economic factors and quality-of-life determine the attractiveness of a region for migration. Such variables do not seem to be as important in determining the number of people leaving particular regions.

Table 3.3 OLS and ML Parameter Estimates for Four Interaction Models

Model	OLS Results				MLE Results			
	$\mu$	$\alpha$	$\beta$	SRMSE	$\mu$	$\alpha$	$\beta$	SRMSE
Unconstrained	.828	.742	.452	.596	.692	.635	.367	.583
Production-constrained	*	.639	.572	.563	*	.658	.494	.560
Attraction-constrained	.703	*	.709	.343	.737	*	.718	.342
Doubly-constrained	*	*	.994	.245	*	*	.905	.234

Through the use of equation (3.15), it is possible to employ the data presented in Tables 3.1 and 3.2 to derive estimates of the relative attractiveness of each of the nine regions for migrants in the United States. These results are reported in Table 3.4 where each value is an estimate of attractiveness relative to the East North-Central region which was set to 1.0 (see discussion above). Clearly, by far the most attractive regions are the Pacific region and the South Atlantic region while the least attractive regions are the North-East and the West North-Central.

Table 3.4 Estimates of Relative Attractiveness

Number	Region	Estimated Relative Attractiveness
1	North-East	0.39
2	Mid-Atlantic	0.82
3	East North-Central	1.00
4	West North-Central	0.60
5	South Atlantic	2.25
6	East South-Central	0.42
7	West South-Central	1.03
8	Mountain	0.98
9	Pacific	3.90

Other empirical comparisons of OLS and MLE parameter estimates in interaction modelling can be found in Stetzer (1976), van Est and van Setten (1979), Openshaw (1979) and Fotheringham and Williams (1983).

### 3.7 Goodness-of-Fit Statistics and Spatial Interaction Models

An important component of spatial interaction modelling is the assessment of a model's ability to replicate a known set of flows. Accurate replication supports the theoretical propositions on which the model is based: that is, it supports one particular model form over others. It also lends confidence in the accuracy of parameter estimates and in the ability of a model to predict flows in systems other than that in which the model was calibrated. Many goodness-of-fit statistics have been employed in spatial interaction modelling to assess a model's ability to replicate a data set and reviews are presented by Knudsen and Fotheringham (1986) and by Fotheringham and Knudsen (1987). All such statistics involve a quantitative description of some aspect of the difference between  $T'$ , the matrix of predicted flows, and  $T$ , the matrix of observed flows, but there has been little consistency in the use of a particular statistic(s) which hinders comparison of model performance across studies (see Hathaway, 1975; Thomas, 1977; Southworth, 1983; and Fotheringham and Williams, 1983). It is also unfortunate that the use of different goodness-of-fit statistics can lead to different conclusions being reached regarding model performance. For example, Willmott (1984) employs three goodness-of-fit measures to evaluate model performance (although not in a spatial interaction context), and each statistic indicates a different model as the most accurate.

In the light of the research on goodness-of-fit by Knudsen and Fotheringham (1986) and Fotheringham and Knudsen (1987), a reasonable strategy to employ in evaluating spatial interaction models would appear to be to employ a combination of two of the following three statistics:  $R^2$ , Information Gain and SRMSE, the Standardised Root Mean Square Error. The statistic  $R^2$  and its properties are well-known, and its significance can be

examined through a t-test. Information Gain is slightly less well-known and is calculated as:

$$I = \sum_i \sum_j T_{ij} \ln(T_{ij}/T_{ij}') \quad (3.65)$$

It has a lower value of zero corresponding to a perfect set of predictions and upper limit of infinity. Its significance can be found through its relationship to the minimum discrimination information (MDI) statistic,

$$MDI = 2TI \quad (3.66)$$

which is asymptotically chi-square distributed (Bishop *et al.*, 1975; Phillips, 1981) with  $mn-k$  degrees of freedom, where  $mn$  represents the number of origin-destination pairs and  $k$  represents the number of parameters in the model. SRMSE is calculated as

$$SRMSE = (1/T)[\sum_i \sum_j (T_{ij} - T_{ij}')^2/n] \quad (3.67)$$

It has a lower limit of zero, indicating a completely accurate set of predictions, and an upper limit that, although variable and dependent on the distribution of the observed flows, is usually 1.0. Values of SRMSE greater than 1.0 only occur when the mean error in predicting a set of flows is greater than the mean flow. Unfortunately, SRMSE does not have a known theoretical distribution so that significance testing can only be carried out through elaborate experimental procedures (see Fotheringham and Knudsen, 1987). However, its value has a strongly linear relationship with our general concept of error which makes it useful as a comparative measure of model performance.

On theoretical grounds  $R^2$  is more applicable to assessing the goodness-of-fit of linear models and so can be employed to examine the performance of models made linear by the transformations described above in Section 3.2. Information Gain is most applicable to models calibrated by maximum-likelihood and SRMSE can be applied to any models calibrated by any means. Thus, a useful combination of goodness-of-fit statistics to report for linear interaction models is  $R^2$  and SRMSE; while a useful combination for models calibrated by MLE is Information Gain and SRMSE.

Examples of the use of the above statistics in spatial interaction modelling include Fotheringham's (1983a) use of  $R^2$  to assess differences in linear interaction model specification; Lewis's (1975) use of  $R^2$  to assess differences in model performance between gravity and Heckscher-Ohlin models of interregional trade; the use of SRMSE by Pitfield (1978) to discriminate between a linear programming model and a doubly-constrained interaction model of freight movements in Britain; and the use of Information Gain by Thomas (1977) in a study of journeys-to-work on Merseyside. Fotheringham and Williams (1983) report all three of the above statistics to compare the calibration results of a production-constrained model by OLS and MLE in four different data sets. Reassuringly, the three statistics had the same order of relative magnitude in all four data sets.

Many goodness-of-fit statistics, other than those reported above, can be employed in spatial interaction modelling and the reader is referred to Fotheringham and Knudsen (1987) for a survey of such statistics. However, in the interests of comparability and consistency, we recommend that a combination of the above three statistics always be reported.

### 3.8 The Use of SIMODEL to Calibrate Spatial Interaction Models

While the calibration of interaction models by least squares regression can easily be undertaken with one of a multitude of computer software packages, maximum-likelihood calibration packages are less frequently encountered. As a consequence, an overview is now

presented of a package, SIMODEL, that has been written for the exclusive purpose of calibrating spatial interaction models by MLE.

SIMODEL (Williams and Fotheringham, 1984) is a FORTRAN program which has been extensively tested in classroom and research situations and has proven to be a robust and efficient algorithm. To give an idea of the execution time of SIMODEL, the calibration of a two-parameter model with a  $25 \times 25$  interaction matrix takes 1.76 seconds of CPU time on a CDC Cyber 170/855. SIMODEL is now implemented at research sites throughout the world.

**Table 3.5 Options Available for Calibration in SIMODEL**

OPTION	MODEL DESCRIPTION	FORMAT	NO. OF PARAMETERS
1	Production-constrained	$T_{ij} = A_i O_i D_j f(d_{ij})$	1
2	Production-constrained	$T_{ij} = A_i O_i W_j f(d_{ij})$	2
3	Production-constrained	$T_{ij} = A_i O_i D_j f(d_{ij})c_j^\delta$	2
4	Production-constrained	$T_{ij} = A_i O_i W_j f(d_{ij})c_j^\delta$	3
5	Destination-constrained	$T_{ij} = B_j O_i D_j f(d_{ij})$	1
6	Destination-constrained	$T_{ij} = B_j v_i^\mu D_j f(d_{ij})$	2
7	Poisson unconstrained	$T_{ij} = v_i^\mu w_j^\alpha f(d_{ij})$	3
8	Poisson unconstrained	$T_{ij} = v_i^\mu w_j^\alpha f(d_{ij})c_j^\delta$	4
9	Doubly constrained	$T_{ij} = A_i B_j O_i D_j f(d_{ij})$	1

**Notes on Table 3.5**

In the above notation,  $T_{ij}$  represents the interaction between origin  $i$  and destination  $j$ ;  $O_i$  represents the known total outflow from  $i$ ;  $D_j$  represents the known total inflow into  $j$ ;  $v_i$  represents origin propulsiveness;  $w_j$  represents destination attractiveness;  $d_{ij}$  represents the separation between  $i$  and  $j$ ; and  $c_j$  represents the competition destination  $j$  faces from all other destinations and is usually measured by the accessibility of  $j$  to all other destinations (see Chapter 4 for more details).  $A_i$  and  $B_j$  are balancing factors. The parameters  $\mu$ ,  $\alpha$ , and  $\delta$  are estimated in the model calibration. A distance-decay parameter,  $\beta$ , is also calibrated in the separation function,  $f(d_{ij})$ .

## 62 Chapter 3

The data input to SIMODEL has a very simple structure consisting of one logical record describing the model option chosen and the structure of the data to be read followed by the data themselves. The logical record consists of the following eight elements:

1. Model option to be calibrated (see Table 3.5)
2. Distance/Cost function to be used (power or exponential)
3. Number of origin-destination pairs for which data are available
4. Number of origins
5. Number of destinations
6. Minimum distance declaration. Only interactions over distances greater than this value are used in the model calibration
7. Maximum iteration limit. A value of 100 is suggested
8. Probability data declaration. A value of one is entered if the interaction data are in the form of probabilities; 0 otherwise.

After reading the above record, SIMODEL then reads interaction and distance data which are entered on one logical record per origin-destination pair. Following this, data on origins and destinations are read in as blocks. As an example, the following SIMODEL program would calibrate a Poisson unconstrained model with a power distance function for a  $2 \times 2$  interaction model using all flows over a distance greater than zero units. The interaction and distance data are read in first, followed by an origin propulsiveness measure and two destination variables. Four parameters would then be estimated in the model calibration.

```

6POWER      4      2      2  0.0  100      0
1999.0      23.0      1      1
423.0       124.1      1      2
627.0       99.2       2      1
3241.0      11.2       2      2
110682.0
342891.0
221264.0
684682.0
834.89
253.71

```

In terms of output, SIMODEL produces a series of calibration results plus a range of diagnostic information. These are demonstrated in Table 3.6 for a production-constrained model calibrated with data on airline passenger flows from Atlanta, Georgia to 25 other cities in the United States. The model has two parameters to be estimated: a destination attractiveness parameter and a distance-decay parameter. These are reported with confirmation that their respective constraint equations are met. This is followed by some general information on the interaction matrix and then a series of goodness-of-fit measures. The standard errors of the parameter estimates are then reported plus a series of log-likelihood values that can also be used to assess the significance of each parameter estimate. Following that, the balancing factors (only one in this case) are printed plus the observed and predicted flows and absolute and percentage errors. Further details on SIMODEL output can be found in Williams and Fotheringham, 1984.

**Table 3.6 An Example of SIMODEL Output for a Production-Constrained Interaction Model**

SIMODEL VERSION 1.0 UNDER NOS 2.2 - RELEASED JULY 1984

PRODUCTION CONSTRAINED GRAVITY MODEL RESULTS WITH OPTION 2  
FOR 1 ORIGIN(S) AND 25 DESTINATION(S) CONSIDERING  
ONLY INTERACTIONS AT DISTANCES OR COSTS GREATER THAN 160.0

THE OBSERVED MEAN TRIP BENEFIT = 14.9831  
THE PREDICTED MEAN TRIP BENEFIT = 14.9831  
MAXIMUM LIKELIHOOD MASS PARAMETER = .7818262

THE OBSERVED MEAN TRIP LENGTH = 6.5003  
THE PREDICTED MEAN TRIP LENGTH = 6.5003  
MAXIMUM LIKELIHOOD DISTANCE PARAMETER = -.7365098

AFTER 10 ITERATIONS OF THE CALIBRATION ROUTINE  
WITH A POWER DISTANCE OR COST FUNCTION

THE NUMBER OF ORIGIN-DESTINATION PAIRS CONSIDERED = 24  
THE TOTAL INTERACTIONS OBSERVED = 242873.0  
THE TOTAL INTERACTIONS PREDICTED = 242873.0  
THE ASYMMETRY INDEX FOR THIS INTERACTION DATA = 0

REGRESSING THE OBSERVED INTERACTIONS ON THE PREDICTED  
INTERACTIONS YIELDS AN R SQUARED VALUE OF .605

TIJ (OBS) =  $-1235.6 + 1.122 \text{ TIJ (PRED)}$   
(.195)

T STATISTIC FOR REGRESSION (SIG. DIFF. PARA. FROM 1) = .6262

PERCENTAGE DEVIATION OF OBSERVED INTS. FROM THE MEAN ( 10119.7) = 57.965

PERCENTAGE DEVIATION OF PREDICTED INTS. FROM THE OBSERVED INTS. = 42.189

PERCENTAGE REDUCTION IN DEVIATION = 27.216

AYENI S INFORMATION STATISTIC (PSI) = .617062E-01

MINIMUM DISCRIMINANT INFORMATION STATISTIC = .299735E+05

THE STANDARDIZED ROOT MEAN SQUARE ERROR STAT.= .5787

## 64 Chapter 3

```
THE MAXIMUM ENTROPY FOR      24 CASES      = 3.1781
THE ENTROPY OF THE PREDICTED INTERACTIONS = 3.0063
THE ENTROPY OF THE OBSERVED INTERACTIONS = 2.8736

MAXIMUM ENTROPY - ENTROPY OF PREDICTED INTS. = .1718

ENTROPY OF PRED. INTS.- ENTROPY OF OBS. INTS.= .1327

ENTROPY RATIO STATISTIC = .5641

VARIANCE OF THE ENTROPY OF PREDICTED INTERACTIONS = .165272E-05
VARIANCE OF THE ENTROPY OF OBSERVED INTERACTIONS = .255475E-05
T STATISTIC FOR THE ABSOLUTE ENTROPY DIFFERENCE = 64.6963

THE INFORMATION GAIN STATISTIC = .13270573

AVERAGE DISTANCE TRAVELED IN SYSTEM = 736.5283
AVERAGE ORIGIN-DESTIN. SEPARATION = 851.2500

STANDARD ERROR OF MLE DISTANCE PARAMETER = .527456E-02
STANDARD ERROR OF MLE MASS PARAMETER = .276766E-02

THE LOG-LIKELIHOOD VALUE OF THE FITTED MODEL WITH ALL
PARAMETERS = -.300629E+01

THE LOG-LIKELIHOOD VALUE OF THE FITTED MODEL WITHOUT
THE PARAMETER .7818262 = -.317735E+01

THE RELATIVE LIKELIHOOD (LAMBDA) STATISTIC FOR THE
PARAMETER .7818262 = .830900E+05

THE LOG-LIKELIHOOD VALUE OF THE FITTED MODEL WITHOUT
THE PARAMETER -.7365098 = -.305742E+01

THE RELATIVE LIKELIHOOD (LAMBDA) STATISTIC FOR THE
PARAMETER -.7365098 = .248337E+05

THE LOG-LIKELIHOOD VALUE OF THE FITTED MODEL WITHOUT
ALL THE MODEL PARAMETERS = -.317805E+01

THE RHO-SQUARED STATISTIC FOR THE MODEL = .540466E-01

THE ADJUSTED RHO-SQUARED STATISTIC = -.319492E-01

THE LOG-LIKELIHOOD VALUE OF THE MEAN MODEL= -.317805E+01
```



I A(I) BALANCING FACTOR      O(I) OBSERVED OUTFLOW

1      .485498E-04                      242873.0

J      OBSERVED INFLOWS              PREDICTED INFLOWS

1	0	0
2	6469.0	9476.4
3	7629.0	9635.4
4	20036.0	23858.4
5	4690.0	9312.6
6	6194.0	9662.2
7	11688.0	9078.8
8	2243.0	3675.5
9	8857.0	16637.3
10	7248.0	8052.7
11	3559.0	5729.4
12	9221.0	10110.4
13	10099.0	6983.1
14	22866.0	6305.0
15	3388.0	6515.2
16	9986.0	6949.6
17	46618.0	30924.7
18	11639.0	16384.9
19	1380.0	2480.7
20	5261.0	11373.4
21	5985.0	12127.7
22	6731.0	4621.1
23	2704.0	2652.7
24	12250.0	7364.3
25	16132.0	12961.5

ORIGIN DESTINATION OBSERVED FLOW PREDICTED FLOW ABSOLUTE ERROR PERCENT ERROR

1	2	6469.0	9476.4	3007.4	46.49
1	3	7629.0	9635.4	2006.4	26.30
1	4	20036.0	23858.4	3822.4	19.08
1	5	4690.0	9312.6	4622.6	98.56
1	6	6194.0	9662.2	3468.2	55.99
1	7	11688.0	9078.8	-2609.2	-22.32
1	8	2243.0	3675.5	1432.5	63.87
1	9	8857.0	16637.3	7780.3	87.84
1	10	7248.0	8052.7	804.7	11.10
1	11	3559.0	5729.4	2170.4	60.98
1	12	9221.0	10110.4	889.4	9.65
1	13	10099.0	6983.1	-3115.9	-30.85
1	14	22866.0	6305.0	-16561.0	-72.43
1	15	3388.0	6515.2	3127.2	92.30
1	16	9986.0	6949.6	-3036.4	-30.41
1	17	46618.0	30924.7	-15693.3	-33.66
1	18	11639.0	16384.9	4745.9	40.78
1	19	1380.0	2480.7	1100.7	79.76
1	20	5261.0	11373.4	6112.4	116.18
1	21	5985.0	12127.7	6142.7	102.64
1	22	6731.0	4621.1	-2109.9	-31.35
1	23	2704.0	2652.7	-51.3	-1.90
1	24	12250.0	7364.3	-4885.7	-39.88
1	25	16132.0	12961.5	-3170.5	-19.65