# ON ESTIMATION AND OTHER PROBLEMS OF STATISTICAL INFERENCE IN THE MICRO SIMULATION APPROACH

Anders Klevmarken

The micro simulation approach to economic analysis is still in the beginning of its development. Although "numbers" are involved in the simulations much work is largely of a theoretical character one step away from empirical applications. This is so partly because of data shortage but also because there is a need to use the simulation approach to learn about the properties of ones theoretical constructs. The ultimate goal must, however, be to make an inference to the economy, whether on a macro or a micro level. To do this adequate micro data are needed as well as a basis for the inference.

The general principles of statistical inference apply to the micro simulation approach as well as to other research in econometrics. As a matter of fact, it is hard to find any useful alternative. This does not exclude, however, that there are methodological problems which are more or less specific to this approach. In the following I will first give a few comments on the analysis of micro data in general and then turn to some problems more specific to the micro simulation approach.

## Analysis of micro data, some common problems

Micro data, and in particular longitudinal micro data, certainly offer new possibilities to obtain a better understanding of micro and macro behaviour, but nothing is for free. The use of micro data makes it necessary to solve problems we tend to neglect at the macro level.

1.   There is usually a large individual
     variability in micro data which show up
     in low $R^2$:s. To explain this variability
     we will probably have to use models which
     involve more parameters than is typically
     the case at an aggregate level. For
     instance, an analysis of household
     consumption would not only involve
     household income and lagged consumption
     but also measures of household charac-
     teristics.

2.   Partly because of the large range of
     variability micro relations are frequently
     non-linear which makes the statistical
     inference difficult.

3.   Measurement errors become relatively
     important. Sometimes we will work with
     proxy or indicator variables which
     "suggest" models with latent structures,
     (c.f. Aigner & Goldberger (1977), Wold
     (1973, 1974, 1975)).

4.   There are selectivity problems in micro
     data which may be difficult to handle.
     In panel data in particular self-
     selectivity may demand a separate
     treatment. One promising approach is
     to incorporate the selection mechanism
     into the basic model and estimate both
     at the same time, (c.f. Heckman (1976),
     Maddala (1977)).

5.   Although micro data are expected to be
     a rich source of information there will
     most certainly remain unmeasurable
     individual characteristics. In panel
     data these have sometimes been taken
     care of by a variance-components approach.

6.   The relationships between cross-section,
     cohort and time  series data deserve
     more attention. We do not only need to
     know how macro activities influence micro
     units and how micro units should be
     aggregated to macro. Because the in-
     creasing demand for personal integrity
     will limit our possibilities to obtain
     micro data, and in particular panel
     data, we will often also have to in-
     vestigate if cross-sectional data could
     be used for an inference about longi-
     tudinal behaviour.

We already have statistical methods which can
be used to treat some of these problems, but the
new emphasis on micro data will have to "generate"
new methods. To indicate the nature of these
methods I would like to give a few <u>key words</u>:

a) Although macro theory usually has a micro
theoretical foundation it is not always good
enough for empirical studies of micro be-
haviour. Our methods will thus have to be
<u>exploratory</u>.

b) Because the sample size will be rela-
tively large it is possible to emphasise
<u>consistency</u> rather than efficiency. In
traditional macro econometrics consistency
is a completely uninteresting property
because of the short time-series usually
available. Frequently, however, we only
know the asymptotic properties of our
estimators. For this reason I agree with
those who claim that one should not give
much credence to confidence intervals
computed in macro econometric models. On
the other hand, from this does not
follow that statistical inference is
useless.

c) One should also emphasize <u>robustness</u> of
methods. There is usually a conflict
between our desire to have robust and
efficient methods. With large samples of
micro data, however, we will not have
to be overly concerned about the loss in
efficiency.

d) In traditional econometrics we concentrate
on mean relationships, while with micro
data the <u>distributional aspects</u> will be
more emphasized. For this purpose we
will probably have to develop better
statistical methods than those available
now.

e) There will be a need for methods which
require neither linearity nor assumptions
of particular non-linear forms, but
rather admit <u>data to determine the
functional form</u> of the relationships
estimated.

## Problems in the micro simulation approach

Next I would like to comment on a few problems
which are more specific to the micro simulation
method. The <u>size</u> of the models contributes to
many of the practical difficulties. It is

important to know the properties of an estimated
model and the predictions produced by this model.
It has been suggested that these properties
could be explored by tracing out "reaction
surfaces" by alternative assumptions about model
structure and parameter values (sensitivity
analysis). This is a good idea for small or medium
sized models or for exploring particular features but
cannot be used to evaluate a large micro simulation
model. The sources of uncertainty in the pre-
dictions are the same as in most other econometric
predictions. There will be genuine residual
variation as well as measurement errors. Par-
ameters will be unknown but estimated. Exogenous
variables are not known but predicted. There
will be specification errors, etc. The multiple
of these errors cannot be explored in "reaction
surfaces" because it would be unmanageable to
analyse the large amount of computer printout
required. With these large models it is not
feasible to simulate all possible implications
of a model and discover unrealistic features.
Also, such an approach would not give the
probability of the occurance of a simulated
event. For these reasons it is very important
that each detail (assumption) in the model be
tested by statistical methods. It is also
important to test the model carefully to balance
what I would like to call the "size law", namely
that the vested interest in our own model is
proportional to its size.

Large size models also make simulations expensive.
Methods have to be found which quickly trace out
the distributions for strategic variables.
Although the simulation methods will depend on
the model structure there are general, efficient
Monte Carlo methods and there are also powerful
computer languages for simulations like for
instance SIMULA. Experts on numerical methods
and computer simulations could undoubtedly
contribute to a more efficient use of the
computer.

Another major problem in micro simulation studies
is the lack of data. A typical feature of some
micro analytic studies is that the objective
function which is maximized (or minimized) to
obtain estimates of the micro parameters is
formulated in macro variables because micro data
are not available. For instance, with respect to
the micro parameters one might attempt to minimize
some quadratic function of the residuals between
observed and predicted GNP, consumption expendi-
tures, investment expenditures, rate of un-
employment, rate of increase in consumer prices
etc. This procedure might easily lead into

identification problems. To illustrate by a
simple example, if we only know the sum of two
variables each of which are linearly related to
two other variables, it is not possible to
identify the two intercepts. In a more complex
model it might be difficult to see if the model
is identified or not. If not, the search for a
maximum (minimum) may go on for ever. Even if the
model is formally identified there may be cases
analogous to multicollinearity in ordinary linear
models, i.e. the surface of the objective function
in the neighbourhood of the extremum is flat.
It might then be possible to change some parameter
values with but a very small change in the value
of the objective function.

Gunnar Eliasson in his paper "How does inflation
affect growth - Experiments on the Swedish Model"[1]
presented a slightly different data problem. He
wanted to investigate if the "over shooting"
response of his model to an external shock is a
realistic feature. The problem is that so far we
have not observed such an "over shooting" in the
economy which makes it difficult to put this
property of the model to a direct test.

First we would like to know if this particular
property is the result of the general model
structure or the particular parameter estimates
obtained. Suppose we can write the model

$$M1: F(y, \theta) = 0; \quad \theta \varepsilon S$$

where y is a vector of variables and $\theta$ a vector
of unknown parameters which belong to the set S.
These relations define our maintained hypothesis.
If F has the over shooting property for every $\theta$
in S no sample would be able to reject this
property, i.e. no test is possible. In this case
there is no support for the property and one
would like to consider a more general model
which would include M1.

Even if there are $\theta$:s in S which do not imply
"over shooting" one might think of cases when
this property is "almost" untestable. Suppose
our data are generated by another (stochastic)
model M2 which does not have the "over shooting"
property and that the distribution of y is such
that we with a probability close to 1 will
obtain estimates of $\theta$ in M1 which give over
shooting, then the probability to reject this

---

[1] See pp. 277 ff in this conference volume.

property will be close to 0. To obtain some protection against this possibility one would like to investigate if theoretically plausible models different from Ml with about the same fit would also give the over shooting property. If they do, some support for overshooting is obtained.

In general I can see no other way to solve the testing problem than to test each part of the model against micro data by statistical methods. If micro data are unavailable we will most certainly encounter difficulties in discriminating between model structures. Suppose our data are generated by Ml but there are many parameter vectors $\hat{\theta}$ which give almost the same fit to the observed (macro) data and some give "over shooting" while others do not. This result neither give support to the over shooting property, nor rejects it. Equivalently, if one estimate $\hat{\theta}$ implies overshooting but it is possible to find another $\hat{\theta}$ which gives almost the same fit but no overshooting, then there is no support.

Eliasson discovered the over shooting property of his model by deterministic simulation. But assigning the value zero to the random errors does not always give unbiased predictions, c.f. the case of log-normally distributed errors. Depending on the structure of the model it might also generate random shocks which would counteract the over shooting. If the random errors implicit in the behavioural relations are taken into account by stochastic simulations one might thus obtain different results vis à vis over shooting.

Finally I would like to comment on what is called "the dynamic approach" to estimation. Let us take the following simple example:

$$y_t = \alpha + \beta y_{t-1} + \varepsilon_t \; ; \; \varepsilon_t \text{ is } NID(0, \sigma^2_\varepsilon)$$

Minimization of

$$\sum_1^T (y_t - \hat{y}|y_{t-1})^2 = \sum_1^T (y_t - \alpha - \beta y_{t-1})^2 \; ,$$

gives the Ordinary Least Squares estimates which are maximum likelihood estimates and they are consistent, asymptotically unbiased and asymptotically efficient. In the dynamic approach the following residual sum of squares is minimized

$$\sum_{1}^{T} (y_t - \hat{y}_t|\hat{y}_{t-1})^2 = \sum_{t=1}^{T} (y_t - \alpha \sum_{i=1}^{t} \beta^{i-2} - \beta^{t-1}y_1)^2 \, ,$$

where $y_1$ is the first y-observation. It remains
to be shown that the estimates obtained have any
desirable properties.

If the OLS estimates are used for "dynamic
predictions", i.e. only the first y-observation
is used to start the forecasting, and if all $\varepsilon_t$
are set equal to zero, one would probably obtain
a sequence of y-predictions which deviates from
the observed series in a seemingly non-random way.
Is this result an indication of a bad model? Not
necessarily! In a mean-square sense the prediction
was the best possible given that we only knew
the first y-value. The random number generator
which we call the economy will generate a y-series
with all $\varepsilon$ set equal to zero only with a probability
close to zero. The probability that our random
number generator would be able to generate the
same series of $\varepsilon$ values as generated by the
economy is also almost zero. To simulate only
one future y-path thus is almost useless. What
is of interest is to simulate the whole distri-
bution of y-paths. Our interest must then be
concentrated on building models which yield
distributions with small variances.

268

References

Aigner, D J & Goldberger, A S (ed) Latent
        variables in Socio-Economic Models,
        North-Holland Publ. Co. Amsterdam 1977

Heckman, J D, The common structure of statistical
        models of truncation, sample selection
        and limited dependent variables and a
        simple estimator for such models,
        Annals of Economic and Social Measurement
        5(4), 475-492, 1976

Maddala, G S, Self-Selectivity problems in
        Econometric Models in Applications of
        Statistics (ed) by P R Krishniah,
        North-Holland Publishing Co 1977

Wold, H, Nonlinear Interative Partial Least
        Squares (NIPALS) Modelling: Some current
        developments in Multivariate Analysis III
        ed by P R Krishnaiah, Academic Press,
        New York 1973

"        Casual flows with latent variables.
        Partings of the ways in the light of
        NIPALS modelling, European Economic
        Review, No 5, 67-86 (1974)

"        Path models with latent variables:
        The NIPALS approach, Quantitative
        Sociology, Academic Press, 1975