

---

Amos Golan — Moshe Kim\*

Research Department

11.9.1992

Maximum Entropy and the Existence  
of Consistent Aggregates:  
An Application to U.S. Banking

\* Department of Economics, University of Haifa, Haifa, 31905, ISRAEL.  
The paper was finalized while Professor Kim visited the Bank of Finland  
Research Department in August–September 1992.

ISBN 951-686-339-6  
ISSN 0785-3572

Suomen Pankin monistuskeskus  
Helsinki 1992

## Abstract

A non-parametric framework for testing the existence of aggregates is developed. The framework is based on the notion of the maximum entropy formalism which is applied to model the size distribution of firms. The theory, once presented, is applied to the US banking industry in order to test for the existence of a consistent output-aggregate. The existence of an output-aggregate in banking can facilitate much of the empirical research in this industry and shed some light on the industry's long-run market structure.



# Contents

	Page
Abstract	3
1 Introduction	7
2 Methodology	8
2.1 Background	8
2.2 The Multi-Variable Size Distribution Model	8
2.3 The Aggregate Test	12
3 The Banking Industry – Empirical Analysis	13
3.1 The Data	13
3.2 The Non-Aggregated Analysis	14
3.3 The Output Aggregate Analysis	16
4 Conclusions	20
References	21



# 1 Introduction

Determination of the existence and validity of economic aggregates is a first step in estimating aggregated data. Theoretical and empirical methods have been developed to test for the existence and validity of these economic aggregates. Most studies concentrate on investigating the separability of the aggregate from the other economic variables of interest. Then, after specifying some functional form, test the parametric restrictions implied by various aggregation theorems. This line of research was initiated by Berndt and Christiansen (1974) and was followed by Denny and Fuss (1977). Later research has applied these concepts to the examination of the existence of aggregates in trucking, Wang-Chang and Friedlaender (1985), and in the banking industry, Kim (1986).

However, as has been noted recently (Aizcorbe, 1990), there may exist several problematic issues with the above line of research. First, a functional form that is consistent with the aggregate of interest must be specified. The second problem that arises involves the locality of tests in the (most) flexible functional forms, such as the widely used translog. The third and main pitfall in the above studies, however, is the fact that once a parametric approach is used, a functional form has to be identified *a-priori*. This, in turn, poses the problem of whether the applied tests for aggregation are indeed aggregation tests or whether they are tests of the functional form.

The objective of this paper is to suggest a non-parametric model capable of testing for the existence of aggregates. This model is based on the notion of the maximum entropy formalism (e.g., Jeffreys, 1939, Jaynes, 1979, Zellner, 1990). This theory was applied by Golan (1989, 1991b) to model the size distribution of firms and also shows the unique relationship between the size distribution and the general properties of the production function. Specifically, the size distribution of firms reveals the general properties of the industry's production technology, such as, returns to scale and curvature properties. The theory, once presented, will be applied to the US banking industry in order to test for the existence of a consistent output-aggregate in this industry. The existence of an output-aggregate in banking can facilitate much of the empirical research in this industry and shed some light on the industry's long-run market structure.

The framework developed here is based on the theory of information developed by Shannon (1948) together with the maximum entropy (ME) formalism. This framework is used to describe the long-term market structure of the banking industry based on the currently available information as represented by the data. The framework presented here is a generalization of a large class of the size distribution of firms theories that analyzes the market structure of an industry (e.g., Ijiri and Simon, 1977, Lucas, 1978). This generalization is developed in Golan (1989, 1991b). The present work applies the above theory to investigate the existence and validity of consistent output-aggregates in the U.S. banking industry. The theory is constructed in two steps. First, based on a definition of bank size depending on more than just one variable, and using non-aggregated data, the industry's size distribution is calculated for each year. Next, using the same data set and using output-aggregates the industry's size distribution is calculated again, for each year. If both cases (with and without

output-aggregates) yield the same size distribution, output-aggregates can be used.

Section 2 introduces the theoretical model and the methodology. Section 3 discusses the data used, the empirical results, and the nonparametric tests for the validity of output-aggregates. Finally, Section 4 summarizes the research.

## 2 Methodology

### 2.1 Background

The analytical technique used in this paper is known as the maximum entropy (ME) formalism (Jeffreys (1939), Jaynes (1957, 1979), Montroll (1981) and Reiss et. al. (1986)). The ME formalism is related to Shannon's information theory in the sense that all possible economic states (i.e., banks with certain characteristics such as given level of inputs and outputs, or labor, computers, demand deposits and time deposits) are equally probable with the exception of those excluded by the constraints (the data). These constraints represent the only information available, they are the only prior knowledge one has about the banking industry. Given the constraints (prior information), all possibilities are equally probable. Like Bayes rule, the ME approach is 100 % efficient (Jaynes (1979), Zellner (1986)). A detailed description of the ME in general and its application to chemical/physical systems is given in Jaynes (1957, 1979) and Levine and Tribus (1979). Application of the ME to traffic systems is given in Reiss et. al. (1986), and its economic implications and applications are given in Zellner (1990) and Golan (1991b). Since the first step toward analyzing (non-parametrically) the existence and validity of economic aggregates, say output-aggregates, requires estimating the multi-variable size distribution of the industry investigated, this size distribution theory is summarized here briefly. This is done next.

### 2.2 The Multi-Variable Size Distribution Model

Consider an industry with  $N$  banks. Each bank is represented by a vector of outputs,  $Y$ , and a vector of inputs (or resources),  $I$ . It is assumed that over a given period of time resources are limited, i.e., that  $I$  is fixed (but large). It then follows that the average output per bank, over the same period of time, is fixed. These assumptions place constraints on the distribution of outputs and inputs and represent the researcher's only available information about this industry. In other words, one determines the banking structure at some period  $t$  based only on the inputs used and outputs produced by all the banks, in that time period.

Since data are discrete it is possible to view the industry in discrete terms in the sense that outputs and inputs are expressed as multiples of a definite quantity, and both quantities change in arbitrarily small jumps, say one dollar

(or a unit of an input).<sup>1</sup> For simplicity, the model is defined in terms of one output and one input. The extension to the multi-input multi-output case is immediate and is given in Golan (1991b) and is summarized at the end of this section. The output of bank  $i$ , say demand deposit (DD), is defined as

$$y_{n_i} = (n_i)x \quad (1)$$

where  $(n_i)$  is an integer and  $x$  is the smallest amount of production capable of existing independently. In the current work this amount is one dollar, referred to as essential economic unit (EU). Similarly, the input of each bank is

$$r_{k_i} = (k_i)f \quad (2)$$

where  $(k_i)$  is an integer and  $f$  is the EU of inputs, say a unit of labor (L). For simplicity of notations (1) and (2) will be written as  $y_n = nx$  and  $r_k = kf$ . The upper and lower bounds for  $n$  are defined as  $n(+)=\text{Max}\{n_i\}$  and  $n(-)=\text{Min}\{n_i\}$ . The upper and lower bounds for  $k$  are similarly defined.

Define  $q_{nk}$  as the number of banks having an output DD in the range  $nx$  to  $(n+1)x$  and input  $L$  in the range  $kf$  to  $(k+1)f$  where  $n$  and  $k$  are not independent. There is a relationship between production level and the firm's level of resources. This fact is expressed through a function  $\eta_k(n)$  which measures the number of banks having demand deposits in the range  $nx$  to  $(n+1)x$ , when the labor is  $kf$ . Call this function the *production index function*. It is emphasized here that this  $\eta_k(n)$  function is not known a priori and hence cannot be specified before analyzing the data. The production index function represents the relative weight of each technology used in that banking industry. In view of the above definitions,  $q_{nk}$  is the joint distribution function of output/outputs and inputs while  $\eta_k(n)$  is the conditional distribution of output ( $n$ ) given the input ( $k$ ).

It is now possible to express the constraints imposed on the model by the assumptions of constant  $N$ , constant  $I$ , and constant average production as follows:

$$\sum_n \sum_k q_{nk} \cdot (kf) = I \quad (3a)$$

$$\sum_n \sum_k q_{nk} \cdot (nx) = Ny \quad (3b)$$

where  $y$  is constant average production. The sums in (3a–3b) *include* the implicit dependence (technology) of  $n$  upon  $k$ . This dependence of  $n$  upon  $k$  is expressed through the production index function of the economy being modeled. Constraint (3a) is trivial and represents the limited resources (inputs)

---

<sup>1</sup> Most of the relationships in this model are insensitive to the size of the EUs as long as the size is fixed. Thus, it is possible to retain the discrete description of the economy without having to be too precise about the magnitude of  $x$  and  $f$ . This is considered in Golan (1991a).

in that economy. It then follows that constraint (3b) represents the notion of conservation of total output.

To estimate the joint distribution of the economy ( $q_{nk}$ ) one may use the ME formalism. (See for example Reiss et. al., 1986, Jaynes, 1979, and Zellner, 1990 for a thorough discussion of the ME approach.) Using the ME formalism to find  $q_{nk}^{\wedge}$ , the estimated  $q_{nk}$ , one maximizes

$$\Delta = \binom{N}{q_{n_1}, \dots, q_{n_k}} = \frac{N!}{\prod_{n,k} q_{nk}!}$$

subject to (3a-3b) where  $\Delta$  is the multinomial coefficient and represents the number of ways of partitioning the  $N$  distinct banks into  $(n(+)+1-n(-)) \cdot (k(+)+1-k(-))$  subsets (types) with  $q_{nk}$  banks in each subset. The maximization result is

$$q_{nk}^{\wedge} = \frac{N e^{-\alpha k f} e^{-\beta n x}}{\sum_n \sum_k e^{-\alpha k f} e^{-\beta n x}} \quad (4)$$

where  $\alpha$  and  $\beta$  are the Lagrange multipliers which will be determined by substituting (4) into (3a) and (3b). Equation (4) gives the distribution of the different states (sequences) of this banking industry where the sum goes only over those integers that are allowed by the constraints. It is the long run *steady state multi-variable size distribution* of banks in the industry. Using (4) it is possible to develop the behavior and dynamics of the banking industry.

Next, defining  $q_{nk}^{\wedge}$  in probability terms yields

$$P_{nk} \equiv \frac{q_{nk}^{\wedge}}{N} = \frac{e^{-\alpha k f} e^{-\beta n x}}{\Omega} \quad (5)$$

Equation (5) is the probability of outcome  $q_{nk}^{\wedge}$ . (For simplicity of notation  $q_{nk}^{\wedge}$  will be written as  $q_{nk}$  in the rest of this paper.)

Before discussing the model further, a generalization of the above framework to a multi-input multi-output system is in place. Following the previous development of a single-input single-output model, the multivariable case is presented below. This generalized model is constructed in terms of an input vector  $X$ , and an output vector  $Y$ , where one can define the input vector generally enough to include assets, endowments, and so on. All data (constraints) consist of specifying mean values of certain functions  $\{g_1(y), g_2(y), \dots, g_{m_1}(y), f_1(x), f_2(x), f_{m_2}(x)\}$ :

$$\sum_{i=1}^n p_i f_1(x_i) = F_1 \quad 1 \leq i \leq m_1$$

$$\sum_{i=1}^n p_i g_j(y_i) = G_j \quad 1 \leq j \leq m_2$$
(6)

where  $\{F_1\}$  and  $\{G_j\}$  are numbers given in the statement of the problem and  $i$  is an index representing the number of sizes (groups) of an industry. This set of constraints represents the data (e.g., eqs. (3a-3b)) by  $m_1$  input constraints and  $m_2$  output constraints. Let  $m = m_1 + m_2$ , then, if  $m < n$ , the entropy maximization method is a standard variational problem solvable by using the Lagrange multipliers technique. It has the formal solution:

$$p_i = \frac{1}{\Omega(\alpha_1, \dots, \alpha_{m_1}, \beta_1, \dots, \beta_{m_2})} \exp[-\alpha_1 f_1(x_i) - \dots - \alpha_{m_1} f_{m_1}(x_i) - \beta_1 g_1(y_i) - \dots - \beta_{m_2} g_{m_2}(y_i)]$$
(6a)

where

$$\Omega(\alpha_1, \dots, \alpha_{m_1}, \beta_1, \dots, \beta_{m_2}) = \sum_{i=1}^n \exp[-\alpha_1 f_1(x_i) - \dots - \alpha_{m_1} f_{m_1}(x_i) - \beta_1 g_1(y_i) - \dots - \beta_{m_2} g_{m_2}(y_i)].$$
(6b)

Alternately, in terms of the discrete (EU) version,

$$\Omega = \sum_{n_1} \dots \sum_{n_{m_2}} \sum_{k_1} \dots \sum_{k_{m_1}} e^{-\alpha_1 k_1 f_1} \dots e^{-\alpha_{m_1} k_{m_1} f_{m_1}} e^{-\beta_1 n_1 g_1} \dots e^{-\beta_{m_2} n_{m_2} g_{m_2}}$$
(6c)

where  $\{\alpha_i\}$  and  $\{\beta_j\}$  are the Lagrange multipliers, which are chosen as to satisfy constraints (6); and  $\{f_1\}$  ( $l = 1, \dots, m_1$ ) and  $\{X_j\}$  ( $j = 1, \dots, m_2$ ) are the EU's associated with each input or output. Similarly  $\{k_1\}$  and  $\{n_j\}$  are the

associated integers. This solution holds if the data can be represented by a set of simultaneous equations for  $m$  unknowns given by:

$$F_1 = \frac{\partial}{\partial \alpha_1} \log \Omega \quad 1 \leq i \leq m_1$$
(7a)

and

$$G_j = \frac{\partial}{\partial \beta_j} \log \Omega \quad 1 \leq j \leq m_2 \quad (m_1 + m_2 = m).$$
(7b)

The value of the entropy-maximum is then a function only of the given data:

$$S(F_1, \dots, F_{m_1}, G_1, \dots, G_{m_2}) = \log \Omega + \sum_i \sum_j \alpha_i F_i \beta_j G_j. \quad (7c)$$

If this function is known, the explicit solution of (7a-7b) is  $\alpha_1 = \frac{\partial S}{\partial F_1}$   $1 \leq i \leq m_1$  and  $\beta_j = \frac{\partial S}{\partial G_j}$   $1 \leq j \leq m_2$ . Given this distribution, the best prediction one can make (i.e., minimizing the expected square of the errors) of any quantity, say  $w(x)$ , is

$$w(x) = \sum_{i=1}^n p_i w(x_i). \quad (8)$$

This completes the generalization of the above simple model to a multi-input multi-output theory where the constraints are the data (such as, in the present application, labor, computers, number of offices, demand deposit, time deposit, mortgage loans, other loans, ATM machines, etc.) one wishes to analyze. It is due to this framework that one does not have to specify a-priori the technological relations (i.e., production function) of the analyzed industry. Furthermore, as can be seen from the above equations, there is no need to specify the causality relations specifically. That is, unlike the common econometric estimation where one has to specify both the dependent and the independent variables, the proposed model avoids this problem by calculating the multi-variable size distribution based on data, and only then the technological parameters are inferred.

The solution to the ME formalism, eq. (4), yields the values of the parameters  $\alpha$  and  $\beta$  (or similarly,  $p$  and  $w$ )<sup>2</sup> together with the estimated probabilities of each multi-size group, i.e., the relative weight of each group in the total population. These parameters and probabilities, together with the production index function, characterize the banking industry. This estimation procedure, requiring a solution of non-linear equations, is done with a Fortran program based on an algorithm developed by Agmon, Alhassid and Levine (1979a, 1979b) and Alhassid et. al. (1978). A test for the existence of aggregates is discussed next.

## 2.3 The Aggregate Test

Having estimated the multi-variable size distribution of the firms one can finally investigate the question of whether or not it is accurate to specify output-aggregates when investigating the banking industry. It is emphasized that the aggregate test developed here can be performed on variables other than outputs. In fact, this technique enables one to test any economic aggregate of interest. This is due to the fact that in essence, with this procedure, one

---

<sup>2</sup>  $p$  and  $W$  are defined as  $W = 1/\beta$  and  $p = \alpha/\beta$ . See Golan (1989) for further details regarding  $p$  and  $W$ .

investigates the value of information. The question under investigation is whether two sets of information (constraints) are different (in terms of their effect on the size distribution) from a set of information consisting of an aggregation of the first two sets. The rest of this paper analyzes output-aggregates.

The investigation of output-aggregates is done by comparing the non-aggregated multi-variable size (ME) distribution of banks with the output-aggregated multi-variable size (ME) distribution. If both multi-variable size distributions are similar one can use output-aggregates. For example, consider using the above approach to estimate the size distribution of banks in the U.S. based on two inputs: labor and computers and two outputs: DD and TD (time deposits). Then, the same analysis is repeated for the aggregated case which is based on the same two inputs and an aggregated measure of DD and TD, say DD+TD where both are measured in dollar terms. If both (non-aggregated and aggregated) estimations yield the same size (ME) distribution, an output-aggregate can be used when analyzing this particular banking industry. This procedure holds for any number of inputs and outputs and for any type of output-aggregates. This concludes the theoretical background describing the application of the ME formalism to test the validity of using economic aggregates.

### 3 The Banking Industry Empirical Analysis

#### 3.1 The Data

The data used include information on 175 U.S. banks during the years 1979 to 1986 generated from the Federal Reserve Functional Cost Analysis (FCA) which has been widely used in empirical research in banking. The main variables in the data set are demand deposit (DD), time deposit (TD), mortgage loans (ML), installment loans (IL), agricultural and construction loans, labor (L), computers (C), number of offices, number of ATM machines, and capital (book value of buildings and equipments). Most outputs are measured in dollar terms where the number of offices, number of ATMs, labor and computers are measured in physical quantities. The analysis is presented in the following way. First, the (non-aggregated) long-run multi-variable size (ME) distribution of the banking industry is calculated for each year. This analysis is based on 4 variables (two inputs and two outputs): L, C, DD and ML. The reason for working with only these four variables is twofold. First, it was found that, in most cases, these variables are sufficient to correctly represent the banking industry,<sup>3</sup> and second, in order to show the strength of the theory it is sufficient to present a four variables' analysis. In the next step of the analysis, the long-run multi-variable size (ME) distribution will be calculated for

---

<sup>3</sup> In other words, adding another variable (constraint) does not change significantly the multi-variable size distribution of the industry, and hence, a sufficient representation of the industry does not require adding this additional variable into the information set analyzed.

different output-aggregates and compared to the first stage (the nonaggregated) of the analysis.

### 3.2 The Non-Aggregated Analysis

Based on the ME formalism described earlier, the size distribution of the banking industry (as given by the available data) is calculated where size is defined over DD, ML, L and C simultaneously. The resulting size distribution, for each year, is presented in Fig. 1. The horizontal axis represents the group (from the smaller group size to the largest where in this analysis the industry is divided arbitrarily into 15 groups), and the vertical axis is the relative weight (probability) of each group in the whole industry. Table 1 shows the weight of each group for the eight analyzed years.

Figure 1. **The size (ME) distribution of the 8 years for the U.S. banking industry (based on Table 1).**

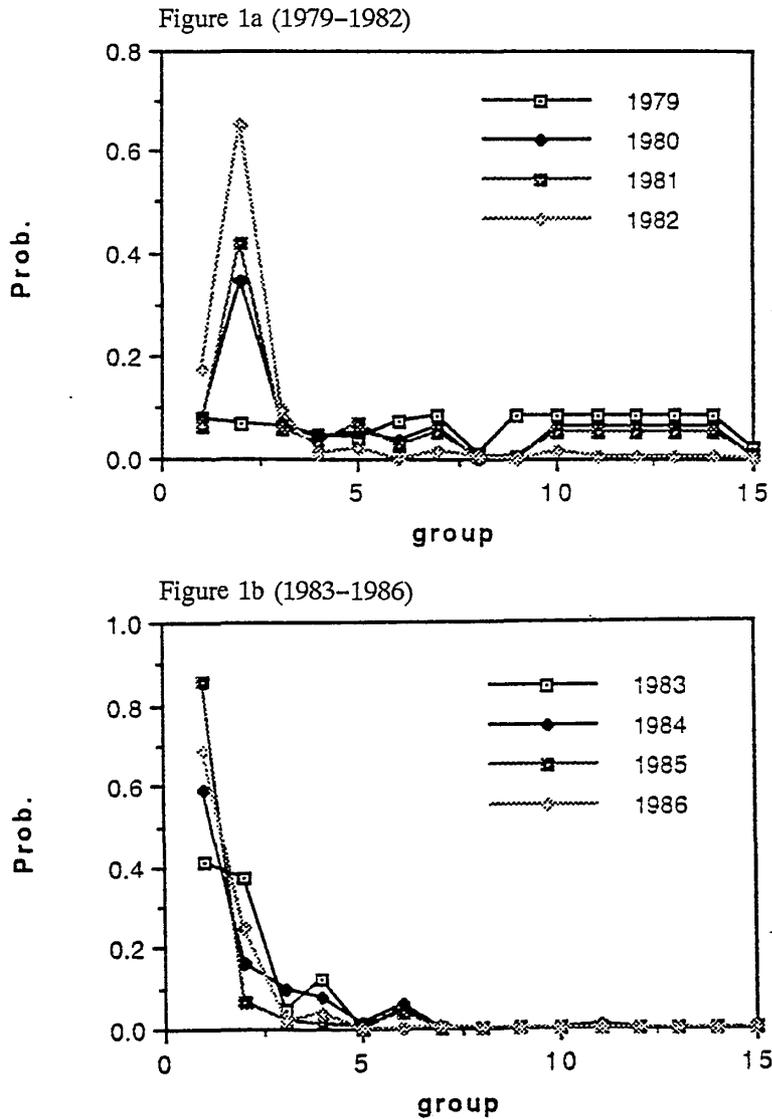


Table 1.

**Size (ME) distributions of the U.S. banking industry for the years 1979–1986. Each column represents one year in increasing order of the multi-variable size.**

1979	1980	1981	1982	1983	1984	1985	1986
7.9732e-2	7.3969e-2	6.3904e-2	1.70351e-1	4.14157e-1	5.85986e-1	8.55282e-1	6.83878e-1
7.0511e-2	3.46388e-1	4.16542e-1	6.55066e-1	3.70038e-1	1.60225e-1	6.3316e-2	2.50909e-1
6.2831e-2	7.0053e-2	5.6281e-2	9.2697e-2	4.5888e-2	9.5432e-2	1.6745e-2	2.1563e-2
4.8646e-2	4.1056e-2	3.2766e-2	1.2792e-2	1.22058e-1	7.5584e-2	1.0362e-2	3.7227e-2
4.0455e-2	5.2192e-2	6.8073e-2	1.8773e-2	3.08e-4	1.5982e-2	6.626e-3	0e+0
7.5104e-2	3.536e-2	2.8537e-2	1.1e-5	4.6886e-2	5.6691e-2	3.7569e-2	0e+0
8.4345e-2	6.1956e-2	5.305e-2	1.4914e-2	1e-6	1.93e-4	1.667e-3	5.539e-3
1.1101e-2	1.709e-3	2.772e-3	4.406e-3	0e+0	2.3e-4	1.174e-3	1.26e-4
8.4345e-2	3.03e-4	2.832e-3	1e-5	1.33e-4	2.51e-4	1.174e-3	1.26e-4
8.4345e-2	6.1956e-2	5.305e-2	1.3358e-2	1.33e-4	2.51e-4	1.174e-3	1.26e-4
8.4345e-2	6.1956e-2	5.305e-2	4.406e-3	0e+0	8.423e-3	1.174e-3	1.26e-4
8.4345e-2	6.1956e-2	5.305e-2	4.406e-3	1.33e-4	2.51e-4	1.174e-3	1.26e-4
8.4345e-2	6.1956e-2	5.305e-2	4.406e-3	1.33e-4	2.51e-4	1.174e-3	1.26e-4
8.4345e-2	6.1956e-2	5.305e-2	4.406e-3	1.33e-4	2.51e-4	1.174e-3	1.26e-4
2.1203e-2	7.237e-3	9.995e-3	0e+0	0e+0	0e+0	2.13e-4	0e+0

Examination of Fig. 1 together with Table 1 gives a notion of the amount of information one receives just by using the ME formalism to analyze the size distribution of firms. Figure 1 and Table 1 show that the banking industry went through a continuing change from 1979 to 1986. In 1979 the industry was almost evenly distributed with about 4–8.5 % in each group except for groups number 8 (relatively big) and 15 (the largest banks). A structural change that effected the small banks (groups 1 and 2) occurred in 1980. The smallest banks remained as in 1979, but the next group increased its weight in the industry to almost 40 % on account of all the other larger groups (6 to 15). In other words, the weight of the smaller banks in the industry increased; the banking industry became more competitive. This result is supported by the diversity-concentration (i.e., the entropy) measure discussed in Golan (1989, 1991b). During 1981 the same trend continued. The second group increased further to about 42 % on account of the bigger banks (groups 10 to 15). This trend continued further in 1982. As can be seen easily in Fig. 1 and Table 1, during this period, group number two increased to almost 66 % together with an increase of group number 1 (the smallest group) to about 17 % of the market. This increase came on account of groups 4 to 15 (the middle to large banks). Note also that the largest banks "disappeared" from the industry. If one considers the first 3 groups (the smallest banks) together it is easy to verify that in 1982 about 90 % of the market consisted of small banks. In essence, one sees that the industry went through a complete change and became much more competitive (from almost a uniform distribution in 1979 to almost a fully competitive market in 1982). The next period of years is characterized by moving into a different market structure. Specifically, in 1983 the smallest banks increased their weight to about 40 % while the second group decreased its weight to about 37 %. Group number 3 almost "disappeared" and group number 4 increased its weight to more than 12 %. This change continued in the next year, 1984, where this time the smallest group captured about 59 % (a 20 % increase) of the market on account of a similar decrease of the second

group. This same trend continued in 1985, where in that period the smallest group increased its weight to almost 86 % of the market while the second group reduced its weight to less than 10 %. Finally, even though the same structure remained during 1986, a different trend is starting. The relative weight of the smallest group size reduces (after a 3 year upward trend) to about 69 % and the second group increases its weight back to about 25 %. In general, the structure of the middle and large size banks did not change significantly during this four year period. To conclude, it is clear that the banking industry went through a major change during the analyzed period; it went from a relatively concentrated (non-competitive) market to a significantly more competitive structure. This analysis is consistent with the observed banking industry during this period.

For purpose of comparison, additional ME estimates were constructed. First, and as was discussed earlier, an analysis similar to the above was done for the case where size is defined over DD, TD, ML, IL, L, C and capital. The results, however, are consistent with the above analysis, and therefore are not presented here (i.e., the additional information does not change significantly the size, ME, distribution). Second, the same analysis was carried for the case where the market is divided into 6 groups only. Again, the results are consistent with the above results but, are less refined. That is, the multi-size measure is defined over a larger range, making the size (ME) distribution less accurate compared with the 15 group case. Cases of more than 15 groups, however, are not analyzed due to lack of data, i.e., the number of banks in each multi-size group is too small. Third, a similar analysis can be performed to investigate the distribution within each group. That is, refining the results even further. Due to insufficient number of banks (except for the first two, or three groups) this analysis is not done here.

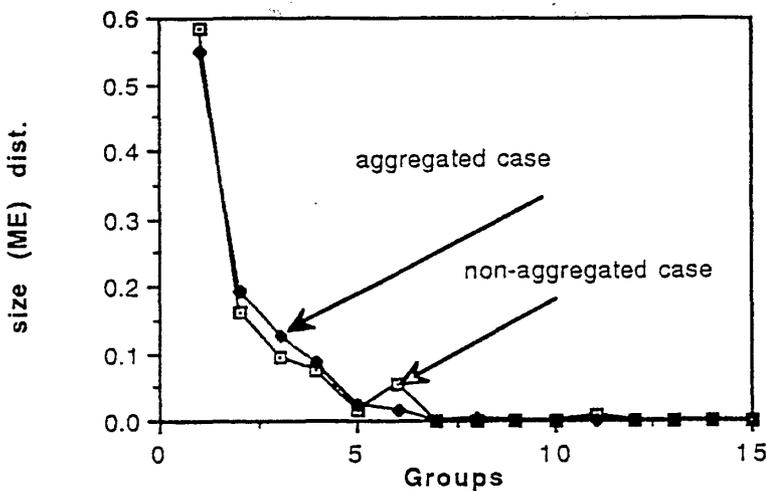
### 3.3 The Output Aggregate Analysis

To test the validity of output-aggregates the ME analysis is performed on the same data as before (i.e., same inputs, same number of groups, same years, and same outputs) with only one difference; the outputs are now aggregated. Different types of aggregations are tested (e.g., all the outputs are aggregated into one measure, the outputs are aggregated according to subgroups, say TD+DD and ML+IL, etc.) with the main conclusion that, in a large number of cases, one can use outputaggregates in the banking industry. These results are now discussed.

It was shown before that in most cases the main variables necessary to analyze the size distribution of firms, and hence its technological characteristics, are the two inputs, labor (L) and computers (C), and the two outputs, demand deposits (DD) and mortgage loans (ML), referred to as the "Non-Aggregated Base Case". Therefore, the output-aggregate needed to be tested here is DD+ML, call it "Case A". In terms of the ME formalism this aggregation implies one less constraint on the system since the two output "constraints" reduce to one. The comparison of the two cases for 1984 is given in Fig. 2 where the horizontal axis represents the groups (from small to large) and the vertical axis represents the relative weight of each group in the total industry, i.e., the multi-variable size (ME) distribution of banks.

Figure 2.

**The Base Case (non-aggregated size distribution) compared with the aggregated Case A for 1984.**



Even though the two cases do not coincide perfectly one can see that the two distributions are qualitatively similar except for a 5 % change in the relative weight of group number 6, the middle size banks. More precisely, using the Kolmogorov-Smirnov two-sample test,<sup>4</sup> the two distributions are proved to be equal at a significance level of  $\alpha = .05$ . Using the MannWhitney U test<sup>5</sup> the two distributions are proved to be equal at a significance level of  $\alpha = .002$  and are not equal for a significance level of  $\alpha = .02$ . Table 2 presents these results together with some other aggregates. Investigating the results (Table 2) carefully reveals that an analysis based on output-aggregates yields (in most cases) a "smoother" multi-variable size (ME) distribution. Furthermore, as is shown in Golan (1988, 1989) the technology inferred in both cases is the same. Since in this paper we are interested solely in the aggregation question, the exact production parameters are not calculated here. Using the same inputs L

<sup>4</sup> The Komogorov-Smirnov two-sample two-tailed test is as follows. Let  $S_{n_1}(X)$  be the observed cumulative step function of one of the samples, say the Base Case. Let  $S_{n_2}(X)$  be the observed cumulative step function of the other sample, say Case A. The test is, then, focuses on

$$D = \text{maximum} |S_{n_1}(X) - S_{n_2}(X)|.$$

The results of this test are given in Tables 2, 3, and 4.

<sup>5</sup> The Mann-Whitney U test is as follows. Let  $n_1$  be the number of cases in one distribution, and  $n_2$  in the second. To apply the U test one has to combine the observations, or the weights in the present paper, and rank these weights in an order of increasing size. Next, one focuses on one of the two distributions, say Case A with  $n_1$  different size groups. The value of U is given by the number of times that a weight in the distribution with  $n_2$  cases precedes a weight in the distribution with  $n_1$  cases in the ranking. The results of this test, for different significant levels, are given in Tables 2, 3, and 4. Note that other tests to compare the two distributions (aggregated and non-aggregated) are possible. For example the Chisquared or the minimum cross entropy estimator. However, the Komogorov-Smirnov and U tests are chosen due to their complete independence from the ME formalism.

and C, other output-aggregates (and a mix of output-aggregates) are analyzed. The first is DD+TD+ML+IL, referred to as "Case B". The second is DD+TD, ML+IL (referred to as "Case C"), and the third is DD+TD, ML+IL+Other Loans (OL), referred to as "Case D". In most of these cases, the multi-variable size (ME) distribution generated by the output-aggregate, or from the mix of outputaggregates, proved to coincide with the multi-variable size (ME) distribution generated by the non-aggregated analysis: output-aggregates proved to be a valid measure (in most cases) for the analyzed data. All of the above experiments, together with (non-parametric) significance tests, are summarized in Table 2.

Table 2. **Comparison of the Base Case (non-aggregated) to other outputaggregate cases for 1984.**

The last 3 rows show the significant test results of comparing each case with the base case, where K-S stands for the Kolmogorov-Smirnov test, and U stands for the Mann-Whitney U test discussed in footnotes 4 and 5. Note the different  $\alpha$  values for the different cases.

Non-Agg.	Case A	Case B	Case C	Case D
5.85986e-1	5.50859e-1	5.72554e-1	5.71526e-1	5.42955e-1
1.60225e-1	1.9123e-1	1.81423e-1	1.85888e-1	1.72565e-1
9.5432e-2	1.24643e-1	8.7004e-2	8.2783e-2	1.17389e-1
7.5584e-2	8.652e-2	1.10523e-1	1.02249e-1	1.17167e-1
1.5982e-2	2.1755e-2	1.5198e-2	1.2089e-2	1.83e-2
5.6691e-2	1.5663e-2	1.9354e-2	2.507e-2	1.9441e-2
1.93e-4	3.8e-5	2.31e-4	1.65e-4	1.22e-4
2.3e-4	4.603e-3	6.112e-3	6.881e-3	1.155e-3
2.51e-4	6.41e-4	6.89e-4	6.75e-4	4.65e-4
2.51e-4	6.41e-4	6.89e-4	6.75e-4	4.65e-4
8.423e-3	1.484e-3	4.154e-3	9.976e-3	6e-6
2.51e-4	6.41e-4	6.89e-4	6.75e-4	4.65e-4
2.51e-4	6.41e-4	6.89e-4	6.75e-4	4.65e-4
2.51e-4	6.41e-4	6.89e-4	6.75e-4	4.65e-4
0e+0	0e+0	0e+0	0e+0	8.577e-3
	K-S ( $\alpha=.5$ ):accept			
	U ( $\alpha=.02$ ):reject	U ( $\alpha=.1$ ):reject	U ( $\alpha=.1$ ):reject	U ( $\alpha=.1$ ):reject
	U ( $\alpha=.002$ ):accept	U ( $\alpha=.05$ ):accept	U ( $\alpha=.05$ ):accept	U ( $\alpha=.05$ ):accept

Table 3 shows a similar analysis for the year 1985 where the base case (L, C, DD, ML) is compared with three types of output-aggregates. Except for group number 1 (smallest banks) where the base case has about 12 % more banks than all the other cases, the multi-variable size (ME) distributions are almost similar. According to the Kolmogorov-Smirnov test all the distributions are similar to the Base Case at a significance level of  $\alpha = .05$ . That is, output-aggregates can be used correctly. According to the Mann-Whitney U test, the distributions are all alike at  $\alpha = .002$ , and are different than the base Case at  $\alpha = .02$ . A similar analysis is done for 1986 (Table 4) where in this case output-aggregates are proved to perform similar to the non-aggregated case for the Kolmogorov-Smirnov test, and are significantly different than the Base Case for the Mann-Whitney U test.

Table 3. Comparison of the Base Case (non-aggregated) to other outputaggregate cases for 1985.

Note the different  $\alpha$  values for the different cases.

Non-Agg.	Case A	Case C	Case D
8.55282e-1	7.41298e-1	7.37411e-1	7.4313e-1
6.3316e-2	7.9126e-2	6.7541e-2	9.1418e-2
1.6745e-2	3.1986e-2	7.7683e-2	4.4162e-2
1.0362e-2	4.0058e-2	3.3315e-2	4.3766e-2
6.626e-3	6.911e-3	1.7805e-2	4.502e-3
3.7569e-2	2.8685e-2	2.2085e-2	1.6367e-2
1.667e-3	2.3771e-2	5.602e-3	8.47e-4
1.174e-3	6.859e-3	5.409e-3	6.667e-3
1.174e-3	6.859e-3	5.409e-3	6.667e-3
1.174e-3	6.859e-3	5.409e-3	6.667e-3
1.174e-3	6.859e-3	5.409e-3	6.667e-3
1.174e-3	6.859e-3	5.409e-3	6.667e-3
1.174e-3	6.859e-3	5.409e-3	6.667e-3
1.174e-3	6.859e-3	5.409e-3	6.667e-3
2.13e-4	1.51e-4	6.94e-4	9.138e-3
	K-S ( $\alpha=.5$ ):accept	K-S ( $\alpha=.5$ ):accept	K-S ( $\alpha=.5$ ):accept
	U ( $\alpha=.02$ ):reject	U ( $\alpha=.02$ ):reject	U ( $\alpha=.02$ ):reject
	U ( $\alpha=.002$ ):accept	U ( $\alpha=.002$ ):accept	U ( $\alpha=.002$ ):accept

Table 4. Comparison of the Base Case (non-aggregated) to other outputaggregate cases for 1986.

Note the different  $\alpha$  values for the different cases.

Non-Agg.	Case A	Case B	Case C
6.83878e-1	6.66302e-1	6.746e-1	6.74835e-1
2.50909e-1	2.38836e-1	2.61546e-1	2.61166e-1
2.1563e-2	6.5982e-2	2.8098e-2	2.8387e-2
3.7227e-2	2.0129e-2	7.323e-3	7.534e-3
0e+0	0e+0	1.3408e-2	1.3294e-2
0e+0	0e+0	0e+0	0e+0
5.539e-3	0e+0	0e+0	0e+0
1.26e-4	1.25e-3	2.146e-3	2.112e-3
1.26e-4	1.25e-3	2.146e-3	2.112e-3
1.26e-4	1.25e-3	2.146e-3	2.112e-3
1.26e-4	1.25e-3	2.146e-3	2.112e-3
1.26e-4	1.25e-3	2.146e-3	2.112e-3
1.26e-4	1.25e-3	2.146e-3	2.112e-3
1.26e-4	1.25e-3	2.146e-3	2.112e-3
0e+0	0e+0	0e+0	0e+0
	K-S ( $\alpha=.5$ ):accept	K-S ( $\alpha=.5$ ):accept	K-S ( $\alpha=.5$ ):accept
	U ( $\alpha=.02$ ):reject	U ( $\alpha=.02$ ):reject	U ( $\alpha=.02$ ):reject
	U ( $\alpha=.002$ ):reject	U ( $\alpha=.002$ ):reject	U ( $\alpha=.002$ ):reject

## 4 Conclusions

A non-parametric test for the validity of economic aggregates has been introduced. This test is performed by applying the ME formalism, on non-aggregated and aggregated data, to derive the multi-variable size (ME) distribution of an industry. Comparison of the multi-variable size (ME) distribution of the aggregated case with the non-aggregated case, gives the desired validity test. If both distributions are similar, aggregates can be used in the estimation procedure. This theory is general enough to investigate the validity of all economic aggregates and was applied to test the validity of output-aggregates in the U.S. banking industry. It is a nonparametric test since no functional form is given a-priori for the multiproduct production function. The ME formalism requires no prior production relation assumptions; it only requires the (ME) non-parametric analysis of the data. Moreover, since the technology can be inferred from the size distribution directly, and is not assumed a-priori, in those cases where the use of output-aggregates proves to be correct, one can use the more common econometric procedures to estimate the production technology using an output-aggregate instead of a cumbersome (and sometimes intractable) multi-output production/cost functions estimation.

## References

- Agmon, N., Alhassid, Y. and Levine, R.D. (1979a) An algorithm for finding the distribution of maximal entropy, *J. of Computational Physics*, 30, 250–259.
- Agmon, N., Alhassid, Y. and Levine, R.D. (1979b) An algorithm for determining the Lagrange parameters in the maximal entropy formalism, in: R.D. Levine and M. Tribus Eds. *The Maximum Entropy Formalism* (MIT Press, Cambridge, MA).
- Aizcorbe, A.M. (1990) Testing the Validity of Aggregates, *Journal of Business & Economic Statistics*, vol. 8, no. 4, Oct., 373–383.
- Alhassid, Y., Agmon, N. and Levine, R.D. (1978) An upper bound for the entropy and its applications to the maximal entropy problem, *Chemical Physics Letters*, 53, 22–26.
- Berndt, E.R. and Christiansen, L.R. (1974) Testing for the Existence of a consistent Aggregate Index of Labor Inputs, *American Economic Review*, vol. 64, no. 3, (June) 391–404.
- Denny, M. and Fuss, M. (1977) The Use of Approximation Analysis to Test for Separability and the Existence of Consistent Aggregates, *American Economic Review*, vol. 67, 404–418.
- Fienberg, S.E. and Zellner, A. (1975) *Studies in Bayesian Econometrics and Statistics* (North-Holland).
- Golan, A. (1988) *A Discrete Stochastic Model of Economic Production and a Model of Fluctuations in Production – Theory and Empirical Evidence*, Ph.D. Thesis, Univ. of California, Berkeley, April.
- Golan, A. (1989) *A Discrete Stochastic Theory of Size Distribution of Firms*, manuscript, Department of Economics, Univ. of Haifa, Haifa, Israel.
- Golan, A. (1991a) On discrete continuous choice of economic modeling or quantum economic chaos, *Mathematical Social Sciences*, (April).
- Golan, A. (1991b) *A Multi-Variable Stochastic Theory of Size Distribution of Firms With Empirical Evidence*, manuscript, Department of Economics, Univ. of Haifa, Haifa, Israel.
- Jaynes, E.T. (1957) Information theory and statistical mechanics, *Physical Review*, 106, 620–630.
- Jaynes, E.T. (1979) Where do we stand on maximum entropy? in: R.D. Levine and M. Tribus Eds., *The Maximum Entropy Formalism* (MIT Press, Cambridge, MA).
- Jeffreys, H. (1967) *Theory of Probability* (3rd ed., Oxford Univ. Press; first edition 1939).
- Ijiri, Y. and Simon, H.A. (1977) *Skew Distributions and the Sizes of Business Firms* (North Holland).
- Kim, M. (1986) Banking Technology and the Existence of a Consistent-Output Aggregate, *Journal of Monetary Economics*, 18 (Sept.) 181–195.
- Levine, R.D. and Tribus, M. Eds. (1979) *The Maximum Entropy Formalism* (MIT Press, Cambridge, MA).
- Lucas Jr. R.E. (1978) On the Size Distribution of Business Firms, *Bell Journal of Economics*, 508–523.

- Montroll, E.W. (1981) On the entropy function in sociotechnical systems, *Proc. Natl. Acad. Sci., USA* 78, 7839–7843.
- Reiss, H., Hammerich Dell A. and Montroll, E.W. (1986) Thermodynamic treatment of nonphysical systems: Formalism and example, *J. Statistical Physics*, 42, 647–687.
- Shannon, C.E. (1949) *Bell Systems Tech. J.*, 27, 379, 623 (1948). Reprinted in: C.E. Shannon and W. Weaver, *The Mathematical Theory of Communication* (Univ. of Illinois Press, Urbana, 1949).
- Wang-Chang, S.J. and Friedlaender, A.F. (1985) Output Aggregation, Network Effects, and the Measurement of Trucking Technology, *Review of Economics and Statistics*, 67, 267–276.
- Zellner, A. (1986) Optimal information-processing and Bayes' theorem, Working Paper, Graduate School of Business, Univ. of Chicago, March.
- Zellner, A. (1990) Bayesian Methods and Entropy in Economics and Econometrics, Working Paper, Graduate School of Business, Chicago.
- Zellner, A. and Highfield, R.A. (1987) Calculation of maximum entropy distributions and approximation of marginal posterior distributions, Working Paper, Graduate School of Business, Univ. of Chicago.

**BANK OF FINLAND DISCUSSION PAPERS**

ISSN 0785-3572

- 1/92 Jaakko Autio **Valuuttakurssit Suomessa 1864–1991**. Katsaus ja tilastosarjat. (Exchange Rates in Finland 1864–1991. Survey and Statistical Series). 1992. 36 + 245 p. ISBN 951-686-309-4. (TU)
- 2/92 Juha Tarkka – Johnny Åkerholm **Fiscal Federalism and European Monetary Integration**. 1992. 29 p. ISBN 951-686-310-8. (KP)
- 3/92 Päivikki Lehto-Sinisalo **The History of Exchange Control in Finland**. 1992. 92 p. ISBN 951-686-311-6. (TO)
- 4/92 Erkki Koskela – Matti Virén **Inflation, Capital Markets and Household Saving in Nordic Countries**. 1992. 21 p. ISBN 951-686-312-4. (TU)
- 5/92 Arto Kovanen **International Capital Flows in Finland 1975–1990: The Role of Expectations**. 1992. 18 p. ISBN 951-686-313-2. (KT)
- 6/92 Ilmo Pyyhtiä **Investointien kohdentuminen Suomessa** (Allocation of investments in Finland). 1992. 54 p. ISBN 951-686-314-0. (KT)
- 7/92 Margus Hanson **Eesti Pank ja Viron rahajärjestelmä 1920- ja 1930-luvulla** (Eesti Pank and the Monetary System of Estonia in the 1920s and 1930s). 1992. 58 p. ISBN 951-686-315-9. (TU)
- 8/92 Markku Malkamäki **Estimating Conditional Betas and the Price of Risk for a Thin Stock Market**. 1992. 36 p. ISBN 951-686-317-5. (TU)
- 9/92 Markku Malkamäki **Conditional Betas and the Price of Risk in a Thin Asset Market: A Sensitivity Analysis**. 1992. 39 p. ISBN 951-686-318-3. (TU)
- 10/92 Christian Starck **Keskuspankkien riippumattomuus – kansainvälinen vertailu** (The Independence of Central Banks – an International Comparison). 1992. 43 p. ISBN 951-686-319-1. (KP)
- 11/92 Juha Tarkka **Tax on Interest and the Pricing of Personal Demand Deposits**. 1992. 21 p. ISBN 951-686-320-5. (TU)
- 12/92 Terhi Kivilahti – Jyri Svanborg – Merja Tekoniemi **Itä-Euroopan maiden valuuttojen vaihdettavuudesta** (Currency convertibility in Eastern Europe). 45 p. ISBN 951-686-322-1. (IT)
- 13/92 Heikki Koskenkylä **Norjan pankkikriisi ja vertailua Suomen pankkeihin** (The Bank Crisis in Norway and a Comparison with Finnish Banks). 1992. 37 p. ISBN 951-686-323-X. (TU)
- 14/92 Tom Kokkola **An International Bibliography of Payment Instruments and Systems Literature for the Years 1985–1991**. 1992. 94 p. ISBN 951-686-325-6. (RM)
- 15/92 Paavo Peisa – Kari Takala **Selviämmekö lamasta ehjin nahoin? Bruttokansantuotteen rakennemallien estimointituloksia** (Are We Going to Survive of the Recession Unmarred? Estimation Results of Structural Models of GDP). 1992. 32 p. ISBN 951-686-326-4. (KT)

- 16/92 Markku Malkamäki **Cointegration and Causality of Stock Markets in Two Small Open Economies and Their Major Trading Partner Nations** 1992. 37 p. ISBN 951-686-328-0. (TU)
- 17/92 T.R.G. Bingham **The Ecu and Reserve Management**. 1992. 22 p. ISBN 951-686-329-9. (KP)
- 18/92 Kari Takala **Työttömyyden ennustaminen lyhyellä aikavälillä (Forecasting Unemployment in the Short Term)**. 1992. 40 p. ISBN 95-686-330-2. (KT)
- 19/92 Anne Turkkila **Suomalaisten pankkien kansainvälistyminen (Internationalization of Finnish Banks)**. 1992. 47 p. ISBN 951-686-331-0. (TU)
- 20/92 Bent Christiansen – Már Gudmundsson – Olli-Pekka Lehmuusaari – Christina Lindenius – Sigurd Simonsen – Christian Starck – Johnny Åkerholm **De nordiska centralbankerna i det framtida Europa (The Nordic Central Banks in Future Europe)**. 1992. 18 p. ISBN 951-686-332-9. (KP)
- 21/92 Paavo Peisa – Heikki Solttila **Pankkikriisi yritysaineiston valossa (The Banking Crisis on the Basis of Panel Data at Firm Level)**. 1992. 18 p. ISBN 951-686-334-5. (RM)
- 22/92 Timo Tyrväinen **Wage Setting, Taxes and Demand for Labour: Multivariate Analysis of the Cointegrating Relations**. 1992. 37 p. ISBN 951-686-335-3. (TU)
- 23/92 Timo Tyrväinen **Tax Incidence in Union Models**. 1992. 37 p. ISBN 951-686-336-1. (TU)
- 24/92 Bent Christiansen – Már Gudmundsson – Olli-Pekka Lehmuusaari – Christina Lindenius – Sigurd Simonsen – Christian Starck – Johnny Åkerholm **Pohjoismaiden keskuspankit tulevassa Euroopassa (The Nordic Central Banks in Future Europe)**. 1992. 24 p. ISBN 951-686-337-X. (KP)
- 25/92 Christian Starck – Matti Virén **Bankruptcies and Aggregate Economic Fluctuations**. 1992. 20 p. ISBN 951-686-338-8. (TU)
- 26/92 Amos Golan – Moshe Kim **Maximum Entropy and the Existence of Consistent Aggregates: An Application to U.S. Banking**. 1992. 22 p. ISBN 951-686-339-6. (TU)