

**UNITED NATIONS STATISTICAL COMMISSION  
and  
ECONOMIC COMMISSION FOR EUROPE**

**CONFERENCE OF EUROPEAN STATISTICIANS  
METHODOLOGICAL MATERIAL**

**EVALUATING EFFICIENCY OF STATISTICAL  
DATA EDITING:**

**GENERAL FRAMEWORK**



**UNITED NATIONS  
Geneva, 2000**



---

**CONTENTS**

PREFACE.....iv

1 WHY IS DATA EDITING A FOCAL POINT .....1

2 STATISTICAL QUALITY IN A MARKET PERSPECTIVE .....2

3 STATISTICAL EDITING .....4

4 MEASURING STATISTICAL QUALITY AND EDITING PERFORMANCE .....5

    4.1 *Measuring quality* .....5

    4.2 *Measuring process and cost data* .....6

5 ANALYSIS .....8

6 NEEDS FOR FURTHER RESEARCH .....11

7 REFERENCES .....13

## **PREFACE**

The methodological material, "Evaluating Efficiency of Statistical Data Editing: General Framework", was prepared based on the request of countries participating in the activities on statistical data editing organised by the UN/ECE Statistical Division within the framework of the programme of work of the Conference of European Statisticians.

The document was reviewed at the Work Session on Statistical Data Editing in June 1999. National Statistical Offices of the UN/ECE member countries and the Food and Agriculture Organisation (FAO) participated in this meeting. The material reflects the outcome of the discussion on the document.

At its 1999 plenary session, the Conference of European Statisticians agreed to reproduce this document and to distribute it to the interested statistical offices as a methodological material.

The document was prepared by Professor Svein Nordbotten.

## 1. Why is data editing a focal point

Data editing is a step in the preparation of statistics, the goal of which is to improve the quality of the statistical information. International research indicates that in a typical statistical survey, the editing may consume up to 40% of all costs. The following questions have been raised:

- Is the use of these resources spent on editing justified?
- Can more effective editing strategies be applied? or
- Can the quality perhaps be improved by allocating some of the editing resources to other statistical production processes to prevent errors [Granquist 1996 and 1997]?

Large statistical organisations and national statistical offices regard their activities as processes producing many statistical *products* in parallel. Each production can be considered as a thread through a sequence of special processes. The overall task for a statistical organisation is to specify, tune and run each thread of processes to deliver a product with as high a quality as possible taking the available resources into account.

We assume that quality can be conceived as a measure of how well the statistical producer succeeds in serving his users. The success will depend on the market demand for statistical products and how the producer allocates his resources to the production of each product and to each process in the production. The better knowledge the statistical producer can acquire about the market for statistical products and the production processes, the better his chances will be for a successful and efficient statistical production. Editing has a particular role in statistical production because its only aim is to improve the quality of the statistical products.

The purpose of this paper is to present a general framework for evaluating the efficiency of statistical data editing in improving the quality of statistical products. The paper includes discussion of:

- the market for statistical products,
- the statistical quality in a market perspective,
- how the quality depends on editing process variables,
- how to measure quality and editing process performance data,
- model tools to support the design of editing processes.

Further work is needed, and the presentation is concluded by suggestions for some important tasks for future research.

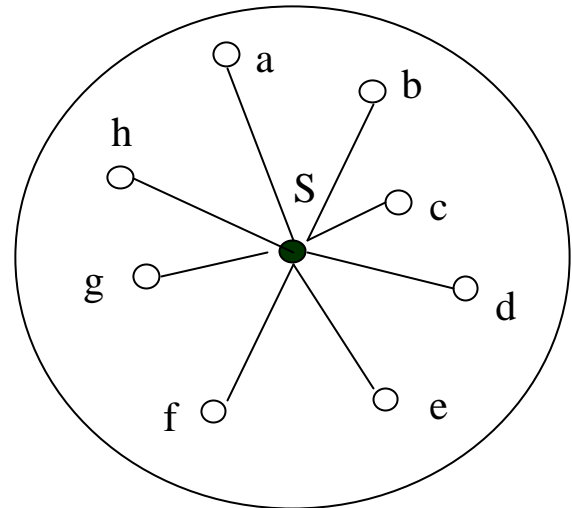
## 2. Statistical quality in a market perspective

The different needs for statistical information are likely to be so numerous that it would be prohibitive for a statistical producer to serve them all. Different users will therefore frequently have to use the same statistical product as a substitute for their varying user needs.

Consider a conceptual definition  $S$  of a statistical product as the centre of a circle symbolising all conceptual definitions for which the definition of  $S$  can be considered to be a feasible substitute. *Figure 1* illustrates a situation in which the applications  $a-h$  have different conceptual needs symbolised by a different spatial location in a circle. As long as the conceptual distances from the centre are within an acceptable length represented by the radius of the circle, the users can be served by the statistical concept  $S$ . For example, users needing a certain population estimate for different points of time, may all use statistics from a census as long as the census was taken within an acceptable time distance.

As a simplification, we ignore the multiplicity of user conceptual needs for statistical products and assume that all needs in the circle can be served by the statistical concept symbolised in the figure by the black circle in the centre. When measured by a perfect process, the statistical concept will be referred to as the *target* product and the attribute value of the product will be referred to as the *target size*<sup>1</sup>. Needs outside the circle cannot be served satisfactorily by this product.

The quality related to a statistical product, is determined by a number of factors including product *relevance* (correspondence between the concept measured and the concept required by the



**Figure 1: Individual target needs served by one statistical target**

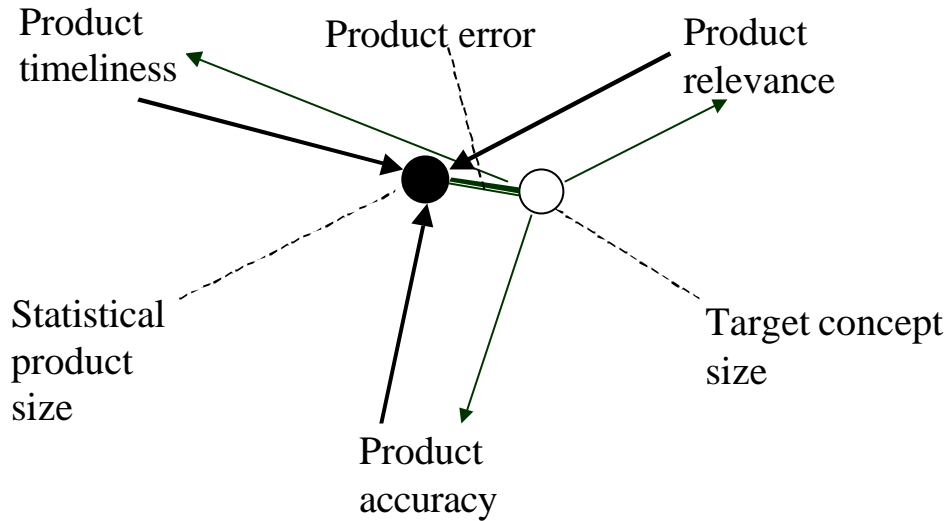
application), *timeliness* (the period between the time of the observations and the time to which the application refers), and *accuracy* (the deviation between the target size determined by a perfect process and the product size determined by the imperfect process) [Depoutot 1998]. Wider quality concepts, as used for example by Statistics Canada, include also accessibility, interpretability and coherence [Statistics Canada 1998].

Figure 2 symbolises by arrows how the 3 factors may pull the statistical product size (the black circle) away from target size (the white circle). The deviation between the actual product size and the ideal target size, is an inverse indicator of quality and frequently referred to as the error of the statistical product.

To justify the preparation of statistics, the users must benefit from the products. We can imagine a market place in which the statistical producers and users trade. We assume that any statistical product has a certain economic value for each user determined by the product quality.

The market value for a statistical product can be described by a sum of all user values. This sum may be considered as a function of the product quality. The cost of production can also be conceived as a function of the product quality.

<sup>1</sup> We use product *size* as a general term for the measurement to avoid confusion with the utility value of the product for a user. A measured population total, an average income, a percentage, etc. are examples of different product sizes.



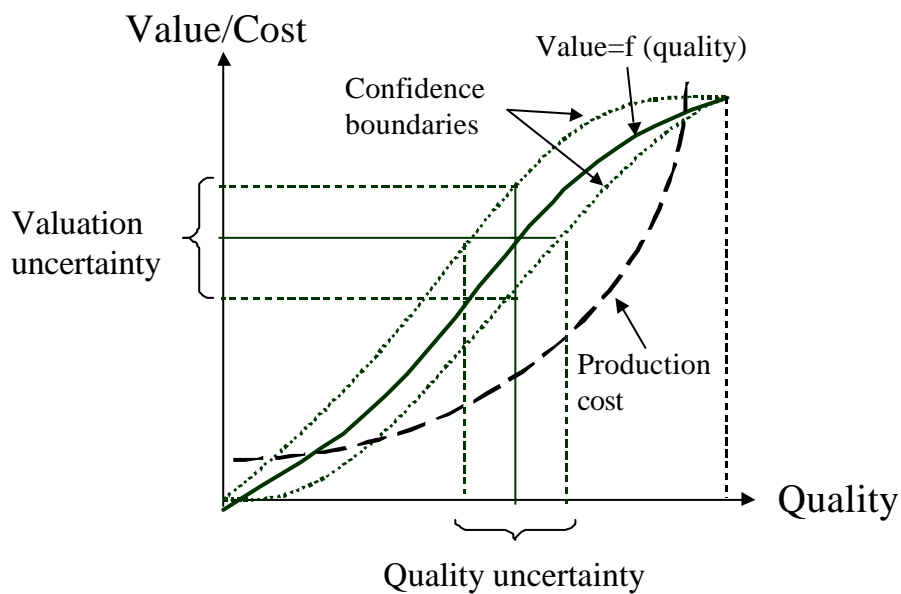
**Figure 2: Factors affecting statistical quality**

The market value for a statistical product can be described by a sum of all user values. This sum may be considered as a function of the product quality. The cost of production can also be conceived as a function of the product quality.

Figure 3 presents a simple graphical model of such a market. According to elementary theory of production, the statistical producer should aim at a quality level, which justifies the costs, i.e. at a quality level for which the product value curve is above the cost curve. The market would

theoretically be in optimal economic balance when the marginal value and cost are equal.

The *users* want data about quality to decide if the supplied statistics are suitable for their needs, while the *producers* need data on quality to analyse alternative production strategies and to allocate resources for improving overall production performance. However, quality can never be precise. One obvious reason is that the precise quality of a statistical product presumes knowledge of the target size, and then there would be no need



**Figure 3: A statistical market mechanism**

for measuring the fact. Another reason is, as mentioned above, that the desired target concept may vary among the users. While a quality statement expresses uncertainty about a statistical product, uncertainty will also be a part of the quality measurement itself. This is illustrated by the stippled curves in *Figure 3* indicating a confidence interval for the value-quality curve.

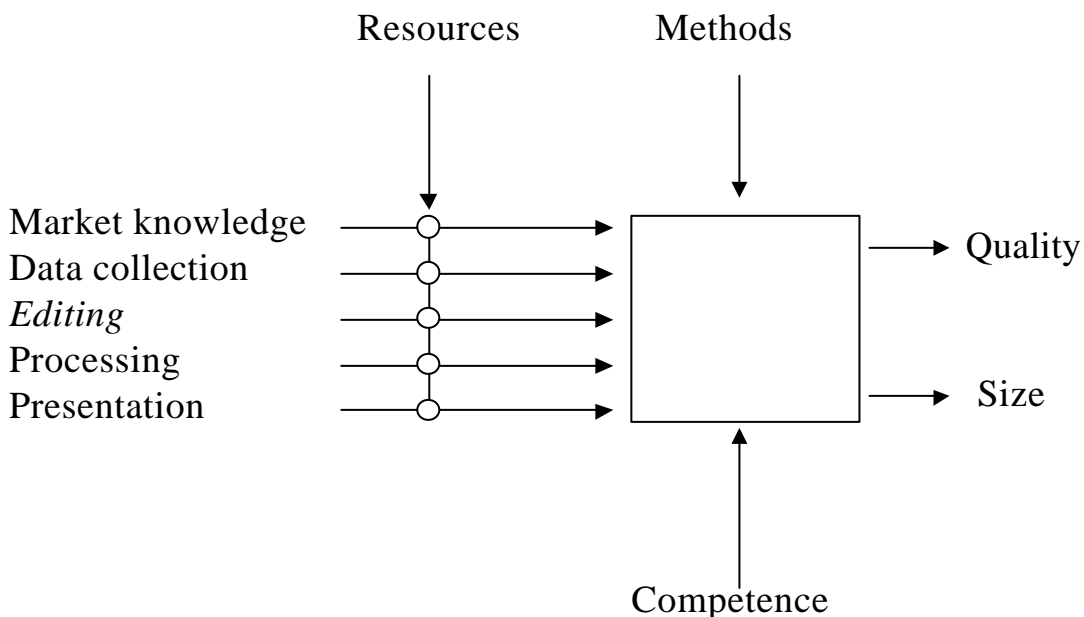
### 3. Statistical editing

The preparation of statistics can be presented as a number of threads, each representing a separate product passing through a sequence of processes. In this paper, we shall limit the discussion to a single product and thread even though we are aware of the interdependence among the threads competing for available resources. *Figure 4* illustrates how each statistical product can be considered as being described by two main variables, the *size* representing the product demanded by the users and the *quality* expressing how reliable the product is. Both variables will depend on how the different processes are designed and how the resources including the available professional competence, are allocated to each process.

Each of the processes can be described by several variables. We would like to identify how the quality of a product is determined by these variables. If we succeed, we shall be able to discuss if the resources are efficiently allocated and if the editing process for the product considered should have 10%, 20% or 40% of the resources allocated to the product considered.

The target size can only be measured correctly by means of a statistical process carried out according to an ideal procedure without any resource restrictions. But because an ideal production usually will be prohibitive for the producer, errors appear in connection with the observation and processing of data for the individual units. In designing a statistical product, resources should first be allocated to statistical processes, which can effectively prevent errors to be generated. The aim of the editing process is to catch individual errors, which are too expensive to prevent efficiently by other processes.

While the result of the editing process is described by data on *quality*, the execution of the process is the source for *performance* data. Description of an editing process requires four types of data:



*Figure 4: Allocations to statistical processes*



- data about the editing *architecture* describing how the process is set up using different methods. For example, this may be a control algorithm for detecting errors, an algorithm for imputing one category of errors and instructions for manual actions of another category of rejected data. The architecture data inform us how the designer wanted the process to be constructed and are obtained during the design of the process.
- data about the *implementation* of the editing process with numerical characteristics, such as specific bounds for edit ratios, imputation functions, etc. These data describe how the process was implemented with all its detailed specifications.
- *performance* data, which document the operational characteristics of the process applied on a specific set of data. They include data on the quality of the editing process.
- *cost* data on which kind of resources were used and how they were spent on different activities.

The first two types of data are obtained during the preparation of the process while the last two types are collected during the execution of the process.

Description of the editing process by means of these four types of data will, in addition to being useful information for evaluating and trimming the process, also provide indications about alternative ways to improve the quality of statistical products, and about opportunities for studying the relationship between the editing process and the statistical product quality.

## 4. Measuring statistical quality and editing performance

In the two previous sections, the statistical product quality and the editing process were discussed from a rather abstract perspective. To be useful, these theoretical notions must be replaced by operational variables, which can be measured and processed. In this section we associate the abstract variables from the previous sections with operational variables which can be observed.

### 4.1 Measuring quality

Quality cannot usually be observed by an exact measurement, but can, subject to a specified risk, be *predicted* as an upper bound for the product error, i.e. for the deviation of the product size from the target size.

Consider the expression:

$$\Pr (|Y'-Y|>D)=1-p \quad 4.1$$

which implies that the probability or risk is  $(1-p)$  that the product size  $Y'$  deviates from its target size  $Y$  by more than an amount  $D$ . We shall denote  $D$  as a *quality predictor* even though it decreases by increasing quality and in fact is an error indicator. Because  $D$  is unknown, we shall substitute it with the prediction  $D'$  [Nordbotten 1998]. It can be demonstrated that  $D' = \mathbf{a}(p)*\text{var } Y'$  where the value of  $\mathbf{a}$  is determined by the probability distribution of  $Y'$  and the assumed value of  $p$ . Assuming that  $Y'$  has a normal distribution,  $\mathbf{a}$  is easily available in statistical tables. The variance of the product size  $Y'$  can be derived from a small sample as described below.

To compute a prediction  $D'$ , we need a small sample of individual records with edited as well as raw data. If the raw records for these units can be re-edited in as ideal a manner as possible to obtain a third set of records containing individual target data, we can compute  $Y$  as well as  $D'$  for different

confidence levels  $p$ . It can be shown that as expected, a smaller risk  $(1-p)$  is related to a larger  $D'$  for the same product and sample.

Because  $D'$  is itself subject to errors, the prediction may or may not provide satisfactory credibility. It is therefore important to test the prediction empirically. In experiments with individual data for which both edited and target data versions exist, we can perform statistical tests comparing predicted quality  $D'$  and actual quality  $D$  of the edited data [Nordbotten 1998 and Weir 1997].

		Actual		Total
		$D \leq 5$	$D > 5$	
Predicted deviation $D'$	$D' \leq 5$	363	67	430
	$D' > 5$	51	23	74
Total		414	90	504

**Table 1: Testing 504 product estimates requiring  $|Y'-Y| \leq 5$  assuming  $p=0.75$ .**

Table 1 illustrates how 504 products or estimates were classified in an experiment to evaluate accuracy predictions [Nordbotten 1999]. The figures, which are based on 1990 Norwegian Population Census data, refer to imputed population totals compared with the corresponding target totals. Only estimates with a deviation from the target with 5 or less people were assumed acceptable. The quality prediction algorithm classified 430 product estimates (first row sum) as satisfactory while 414 (first column sum) were within the pre-set requirement. 51 estimates were predicted as outside the boundary while they in fact were acceptable, a typical Type 1 classification error. On the other hand, 67 values were predicted acceptable while their deviations were greater than 5, misclassifications of Type 2.

With a normal distribution and a  $p=0.75$ , we should expect that 25 percent of the values (i.e. 126 product estimates) would be subjected to a Type 1 misclassification. As the table shows, the

number of Type 1 errors (51) is well within the expected limit. The explanation of this unexpected result is that the distribution of  $D'$  does not approximate closely the normal distribution.

Manzari and Della Rocca distinguish *between output oriented approaches and input oriented approaches* to evaluation of editing and imputation procedures [Manzari and Della Rocca 1999]. In an output oriented approach they focus on the effect of the editing on resulting products, while in an input oriented approach they concentrate on the effect of the editing on the individual data items. Because they evaluate editing processes by means of data with synthetic errors introduced, they are able to follow an input oriented approach. The quality indicator  $D'$  presented in this section is a typical example of an output oriented approach. In the next section, we will also discuss an input oriented approach.

**4.2 Measuring process and cost data**

In section 3, we stressed the need for identifying the variables of the editing process, which determined the quality of a statistical product. Two logical steps constitute the editing process:

- *classification* of an observation as acceptable or suspicious, and
- *correction* of components believed to be wrong.

Before the advent of automation in statistical production, subject matter experts carried out editing, frequently with few formal editing guidelines. Later, computers were introduced and provided new possibilities for more efficient editing, but required also a formalisation of the process [Nordbotten 1963]. Editing principles were developed and implemented in a number of tools for practical application. Today, a wide spectrum of different methods and tools exists. An editing architecture adjusted to a particular survey can be designed by a combination of available tools [UN/ECE 1997].

While the quality evaluation focused on the *final effects* of the editing process on the statistical products, the objective of the process evaluation is to describe what is happening with data *during* the editing process [Engström 1996 and 1997]. But because of the close relationship between the performance of the process and the quality of the results, properties of the editing process can also be useful quality indicators.

The measurement of the quality effects of editing is based on comparisons between edited data and target data. The process measurement on the other hand, is based on comparison between raw (unedited) and edited data. Process data are generated during the process itself and can therefore be frequently used for continuous monitoring of the process. Continuous monitoring of the process permits changes during the editing process execution based on the operational variables observed.

Some typical variables, which can be recorded during the process, are shown in *List 1*. These basic variables give us important facts about the editing process. They are descriptive facts, and can be useful if we learn how to combine and interpret them correctly. Since we have no theoretical system guiding us with respect to selecting which variables to observe, the approach of this section is explorative.

If the number of observations rejected as suspicious

N:	Total number of observations
$N_C$ :	Number of observations rejected as suspicious
$N_I$ :	Number of imputed observations
X:	Raw value sum for all observations
$X_C$ :	Raw value sum for rejected observations
$Y_I$ :	Imputed value sum of rejected observations
Y:	Edited value sum of all observations
$K_C$ :	Cost of editing controls
$K_I$ :	Cost of imputations

**List 1: Typical operational and cost**

in a periodic survey increased from one period to another, it can for example be interpreted as an indication that the raw data have decreasing quality. On the other hand, it can also be regarded as indication of increased quality of the final results, because more units are rejected for careful inspection. A correct conclusion may require that several of the variables be studied simultaneously. As a first step toward a better understanding of the editing process, the basic variables can be combined in different ways. *List 2* gives examples of a few composite variables frequently being used for monitoring and evaluating the editing process.

Frequencies:

$$F_C = N_C/N \quad (\text{Reject frequency})$$

$$F_I = N_I/N \quad (\text{Impute frequency})$$

Ratios:

$$R_C = X_C/X \quad (\text{Reject ratio})$$

$$R_I = Y_I/X \quad (\text{Impute ratio})$$

Per unit values:

$$\underline{K}_C = K_C/N \quad (\text{Cost per rejected unit})$$

$$\underline{K}_I = K_I/N \quad (\text{Cost per imputed})$$

**List 2: Some typical operational and cost ratios**

The *reject frequency*,  $F_C$ , indicates the relative extent of the control work performed. This variable gives a measure of the workload a certain control method implies, and is used to tune the control criteria according to available resources. In an experimental design stage, the reject frequency is used to compare and choose between alternative methods.

The imputation effects on the rejected set of  $N_C$  observations are the second group of variables. The *impute frequency*,  $F_I$ , indicates the relative number of observations which have their values changed during the process.  $F_I$  should obviously not be larger than  $F_C$ . If the difference  $F_C - F_I$  is significant, it may be an indication that the rejection criteria are too narrow, or perhaps that more resources should be allocated to make the

inspection and imputation of rejected observations more effective.

The *rejected value ratio*,  $R_C$ , measures the impact of the rejected values relative to the raw value sum for all observations. A small rejected value ratio may indicate that the suspicious values are an insignificant part of the total of values. If combined with a high  $F_C$ , a review of the process may conclude that the resources spent on inspection of rejected values cannot be justified and are in fact better used for some other process.  $R_C$  may show that even though the  $F_C$  is large, the  $R_C$  may be small which may be another indication that the current editing procedure is not well balanced.

The *impute ratio*,  $R_I$ , indicates the overall effect of the editing and imputation on the raw observations. If  $R_I$  is small, we may suspect that resources may be wasted on editing.

*Costs per rejected unit*,  $K_C$ , and *cost per imputed unit*,  $K_I$ , add up to the total editing cost per unit. The costs per product (item) have to be computed based on a cost distribution scheme since only totals will be available from the accounting system.

The process data are computed from both raw and edited micro data. The importance of preserving also the original raw data has now become obvious and it should become usual practice that the files of raw and edited micro data are carefully stored.

As already pointed out, we have yet no theory for the operational aspects of the editing process. The process variables computed are often used independently of each other. The editing process can easily be evaluated differently depending on which variables are used. The purpose of the next section is to investigate how the process can be described by a set of interrelated variables which may provide further knowledge about the nature of the editing process and a basis for improved future designs.

## 5. Analysis

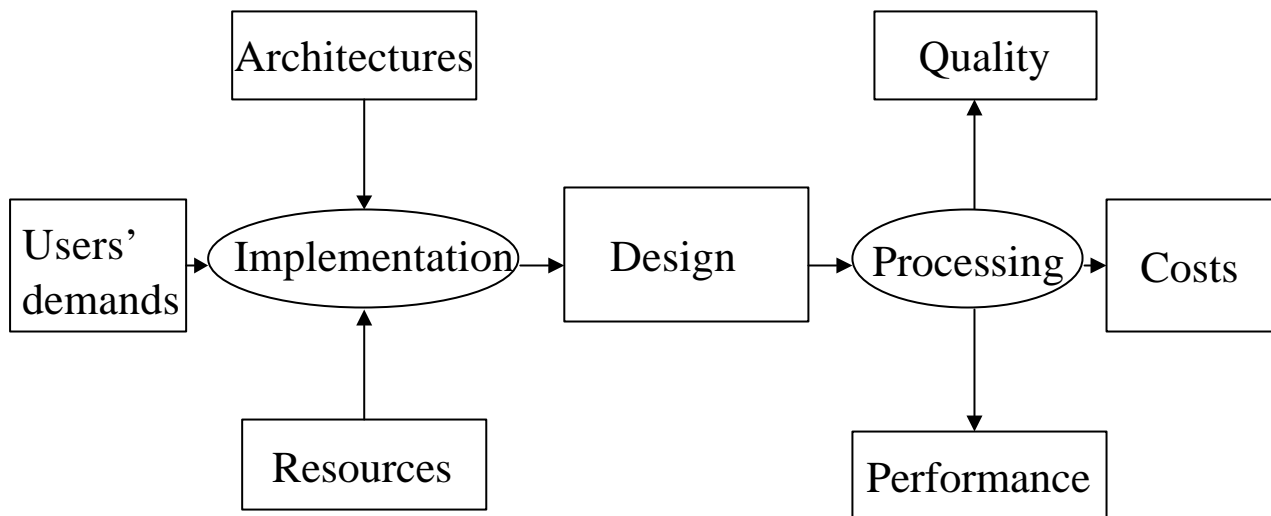
*Metadata* of the type outlined in section 4 offer opportunities for systematic exploration and evaluation of relationships among the statistical product quality and the editing process variables considered. The research objective is to develop a model of the editing process, which describes the *causal* relationships among the variables discussed and can serve as a tool for designing efficient editing processes.

The set of editing architectures, the users' demands and the available resources including mental and stored knowledge and experience, are the environmental conditions within which the implementation for a specific statistical product can be selected. The selected implementation is assumed to determine the operational performance and finally the quality and cost of the editing associated with the product.

*Figure 5* outlines the general structure of a model in mind. On the left are the 3 classes of input variables available for the implementation design. Each class may comprise several variables, which in turn can take a set of values representing different alternatives. There may be, for example, variables identifying alternative control and correction methods, variables representing different edit criteria and imputation parameters, etc. For these variables, values must be selected, designed or estimated. The selection of an implementation design is usually done based on knowledge, which may be mental or represented in a metadata system maintained systematically. When executed, the implementation design is assumed to determine the product quality, cost and performance levels.

The causal relations among the different sets of variables are symbolised by arrows in the figure. Using the notation already introduced, we can write down the model in symbolic form by:

$$I=f(A, D, R), \quad 5.1$$



*Figure 5: Causal model*

where  $A$ ,  $D$  and  $R$  are variables representing architectures, users' demands and levels of available resources, respectively, from the available sets  $\mathbf{A}$ ,  $\mathbf{D}$ ,  $\mathbf{R}$ . The design variables are represented by the implementation variable  $I$  from the set  $\mathbf{I}$ . The mapping  $f$  represents the mappings from the elements in  $\mathbf{A}$ ,  $\mathbf{D}$  and  $\mathbf{R}$  to the elements in  $\mathbf{I}$  and corresponds to the implementation design activity. The selected implementation design  $I$  determines the editing process represented by the mappings  $q$ ,  $p$  and  $k$ , the quality  $Q$ , the performance variables  $P$  and the costs variables  $K$  from the available sets  $\mathbf{Q}$ ,  $\mathbf{P}$  and  $\mathbf{K}$ :

$$Q = q(I), \quad 5.2$$

$$P = p(I) \quad 5.3$$

and

$$K = k(I). \quad 5.4$$

The expressions 5.2 and 5.3 illustrate that the performance variables may be considered as indicators of quality. If the relations 5.2 and 5.3

exist, combination of the two relations may give a new relation:

$$Q = q'(P) \quad 5.5$$

where the new mapping  $q'$  symbolising the mapping from  $P$  to  $Q$ . This expression indicates that  $P$  can be used as an indicator of  $Q$  if the  $q'$  can be determined.

When the quality  $Q$  and costs  $K$  both are determined,  $Q$  needs to be compared with  $K$ . A model corresponding to the market *Figure 3* is needed, i.e. an equation reflecting the relationship between the quality  $Q$  and the market value  $V$ :

$$V = v(Q). \quad 5.6$$

The market value  $V$  can be compared and evaluated with the associated costs  $K$ . Alternative designs, i.e. different implementations, could also be evaluated, compared and ranked.

Exploring these relations empirically will be an important challenge and long-term objective for research in editing of statistical data. It will require

the collection of data from several surveys as well as observations from the statistical market.

The aim stated above was the development of a model to support the producer in finding answers to the questions about which is the ‘best’ architecture and design of an editing process for a given market situation and available architectures and resources. To create a tool for improving the editing strategy, we ‘turn around’ the causal model discussed in the last paragraphs to a decision support model. This transformed model is outlined in Figure 6. The variables on the left side are available architectures and resource alternatives, while at the upper right side we have required quality. On the lower right side, the output of the model are design specifications and cost estimates.

When a statistical market demands a statistical product, the decision support model should assist the statistical producer to investigate:

- if a feasible architecture exists given the repository of editing methods/techniques and the financial and human resources he commands;
- which editing process design can be implemented within the input constraints;
- what the cost will be of the designed process.

There may be several editing designs, which satisfy the input conditions. Based on the discussion in section 2. and the causal model, we search the implementation design  $I$  that gives the highest non-negative solution to the expression:

$$H=V(q(I))-K(I) \tag{5.7}$$

In the long run, research must be extended to study also the impact of other statistical processes, e.g. data acquisition, estimation and presentation of statistical products, on statistical quality and how other statistical processes interact with editing [Jong 1996, Nordbotten 1957]. Only such research may give the necessary tools for tuning resource allocations across all processes in order to obtain the best quality statistics in a general sense.

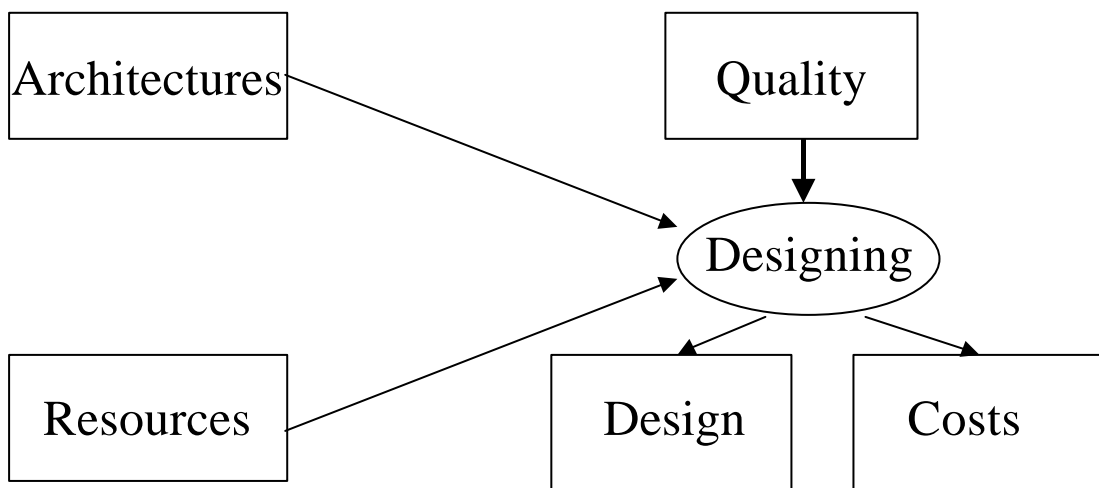


Figure 6: Strategy model

## 6. Needs for further research

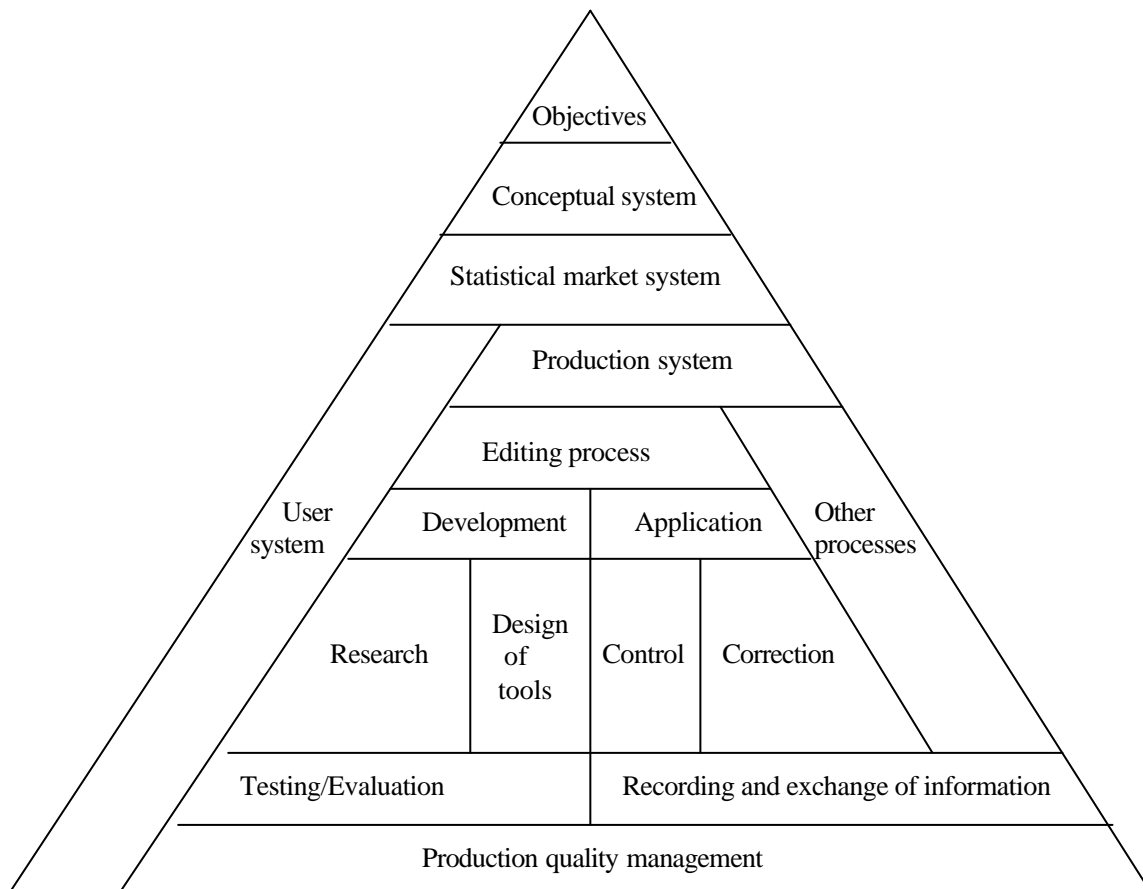
There exists no overall theory from which the producers of official statistics can obtain support for their design and management decisions. So far, producers have relied on theories from several disciplines for separate statistical processes. For some processes and tasks, well-developed theory is available. The theory of sample survey design and estimation is an example of a theory which is an important foundation for design decisions. It illustrates how errors due to random sampling can be taken into account in designing effective samples and evaluating the quality in results due to the sample design. In the last couple of decades, the theory of databases is another example of a theory which has become an important basis for the producers.

Control and correction of non-random errors, on

the other hand, have a less strong theoretical basis. Up to now, a large set of editing methods has been developed [Winkler 1999]. Little progress has been seen so far, however, in integrating the different theories into a general theory of statistical quality.

We can characterise much of the research and methodological development in official statistics as fragmented and explorative associated with different processes. One explanation may be that the purpose of producing statistical information has not been clearly stated and widely understood and that different types of specialists have been involved in different statistical production processes.

There is therefore a need for a general theory of statistical systems and their environments permitting systematic, empirical research and co-operation among the different groups of specialists. *Figure 7* illustrates how the research in editing can be envisaged as a part in a wider scheme for research



*Figure 7: A framework for research in statistical editing*

in statistical production.

This paper focuses on the editing process, but also takes into account other statistical production processes and their environment as outlined in section 2. It emphasises the interactions between the editing process, the other statistical processes and the statistical production environment.

The proposed framework does not intend to be the missing theory for official production of statistics. Its purpose is limited to proposing some relevant research activities connected to editing and aimed at improving the statistical product qualities, and suggesting some research topics and an infrastructure to work within.

Some tasks for future research in statistical data editing may be:

- *Development of a conceptual framework for description of editing processes.*

A common language for communication is needed. The UN/ECE Draft Glossary of Terms Used in Data Editing has recently been updated [Winkler 1999a]. There are terms still missing and the glossary should be continuously updated with terms used, for example, in papers contributed in the UN/ECE Work Sessions on Statistical Data Editing.

- *Collecting empirical data sets suitable for experimentation.*

From statements given at the Work Session in Rome, available data sets suitable for testing new editing methods and architectures are missing and would be appreciated by many researchers. Such sets should be stored in a standard form in a repository and made accessible to researchers working with statistical editing method development and evaluation. Both raw and edited microdata should be stored. When existing, a sample of re-edited ('target') data will be very useful for quality evaluations. Mainly because of confidentiality rules, it is very difficult today to obtain access to data sets used by colleagues in their research. If real data sets cannot be made available, an

alternative is data sets with synthetic data as discussed by Manzari and Della Rocca [Manzari and Della Rocca 1999].

- *Comparison and evaluation of relative merits of available editing tools.*

Useful information can be made available by systematic comparison of the functionality of editing methods based on their description [Poirier 1999]. However, the essential condition for comparison and evaluation of editing architectures is access to empirical microdata. So far, few examples of data sets exchanged and used for comparison of methods have been reported [Kovar and Winkler 1996].

- *Research on causal model description of the editing process.*

Research on causal models will require detailed data from the editing process of the type pointed out above. Data from simulations, can in many situations be a substitute for real data. In some situations, synthetic data can even be superior for studying in detail how different editing methods handle special error types. How to construct useful generators for synthetic data and errors, is therefore also a relevant research task in connection with evaluation of editing methods.

- *Exchange of information on research*

UN/ECE made an important contribution to the exchange of information on editing research work by compiling the statistical editing bibliography [UN/ECE 1996]. In a field like statistical data editing it is important that this bibliography be kept up-to-date. Internet can be exploited more effectively for dissemination of research ideas, experience, references, general problems with answers, etc. Internet can also be used as a highway for researchers to data sets released for comparative research, to stored editing methods, programs and systems made available by authors and developers who wish to share their products with colleagues for comments and applications.



## 7. References

- Depoutot, R. (1998): QUALITY OF INTERNATIONAL STATISTICS: COMPARABILITY AND COHERENCE. Conference on Methodological Issues in Official Statistics. Stockholm.
- Engström, P. (1996): MONITORING THE EDITING PROCESS. Working Paper No. 9, UN/ECE Work Session on Statistical Data Editing. Voorburg.
- Engström, P. (1997): A SMALL STUDY ON USING EDITING PROCESS DATA FOR EVALUATION OF THE EUROPEAN STRUCTURE OF EARNINGS SURVEY. Working Paper No. 19, UN/ECE Work Session on Statistical Data Editing, Prague.
- Engström, P. and Granquist, L. (1999): IMPROVING QUALITY BY MODERN EDITING. Working Paper No. 23, UN/ECE Work Session on Statistical Data Editing, Rome.
- Granquist, L. (1996): THE NEW VIEW ON EDITING. UN/ECE Work Session on Statistical Data Editing, Voorburg. Also published in the International Statistical Review, Vol. 65, No. 3, pp.381-387.
- Granquist, L (1997): AN OVERVIEW OF METHODS OF EVALUATING DATA EDITING PROCEDURES. Statistical Data Editing. Methods and Techniques, Vol. 2. Statistical Standards and Studies No 48. UN/ECE. pp. 112 122.
- Jong, W.A.M. de (1996): DESIGNING A COMPLETE EDIT STRATEGY - COMBINING TECHNIQUES. Working Paper No. 29, UN/ECE Work Session on Statistical Data Editing, Voorburg.
- Kovar, J. and Winkler, E.W. (1996): EDITING ECONOMIC DATA. Working Paper No. 12, UN/ECE Work Session on Statistical Data Editing. Voorburg.
- Manzari, A. and Della Rocca, G. (1999): A GENERALIZED SYSTEM BASED ON SIMULATION APPROACH TO TEST THE QUALITY OF EDITING AND IMPUTATION PROCEDURES. Working Paper No. 13, UN/ECE Work Session on Statistical Data Editing, Rome.
- Nordbotten, S. (1957): ON ERRORS AND OPTIMAL ALLOCATION IN A CENSUS. Skandinavisk Aktuarietidskrift. pp. 1-10.
- Nordbotten, S. (1963): AUTOMATIC EDITING OF INDIVIDUAL STATISTICAL OBSERVATIONS. Statistical Standards and Studies. Handbook No. 2. United Nations, N.Y.
- Nordbotten, S. (1995): EDITING STATISTICAL RECORDS BY NEURAL NETWORKS. Journal of Official Statistics, Vol.11, No. 4, pp. 391-411.
- Nordbotten, S. (1998): ESTIMATING POPULATION PROPORTIONS FROM IMPUTED DATA. Computational Statistics & Data Analysis, Vol. 27, pp. 291-309.
- Nordbotten, S. (1999): SMALL AREA STATISTICS FROM SURVEY AND IMPUTED DATA. To be published.
- Poirier, C. (1999): A FUNCTIONAL EVALUATION OF EDIT AND IMPUTATION TOOLS. Working Paper No. 12, UN/ECE Work Session on Statistical Data Editing. Rome.
- Statistics Canada (1998): QUALITY GUIDELINES. Third Edition, Ottawa.
- UN/ECE (1996): BIBLIOGRAPHY ON STATISTICAL DATA EDITING. Working Paper No. 5, UN/ECE Work Session on Statistical Data Editing. Voorburg.

UN/ECE (1997): STATISTICAL DATA EDITING METHODS AND TECHNIQUES. Vol. 2. Statistical Standards and Studies No. 48. UN/ECE, N.Y. and Geneva.

Weir, P. (1997): DATA EDITING AND PERFORMANCE MEASURES. Working Paper No. 38, UN/ECE Work Session on Statistical Data Editing, Prague.

Winkler, W.E. (1999a): DRAFT GLOSSARY OF TERMS USED IN DATA EDITING. Working Paper No. 2, UN/ECE Work Session on Statistical Data Editing. Rome.

Winkler, W.E. (1999b): STATE OF STATISTICAL DATA EDITING AND CURRENT RESEARCH PROBLEMS. Working Paper No. 29, UN/ECE Work Session on Statistical Data Editing, Rome.