**UNITED NATIONS STATISTICAL COMMISSION and
ECONOMIC COMMISSION FOR EUROPE
CONFERENCE OF EUROPEAN STATISTICIANS**

**EUROPEAN COMMISSION
STATISTICAL OFFICE OF THE
EUROPEAN COMMUNITIES (EUROSTAT)**

**Joint UNECE/Eurostat work session on statistical data confidentiality**
(Geneva, Switzerland, 9-11 November 2005)

Topic (ii): Disclosure risk, information loss and usability of data

# A COMBINED METHODOLOGY FOR ASSESSING IDENTITY AND
# VALUE DISCLOSURE RISK FOR NUMERICAL MICRODATA

**Supporting Paper**

Submitted by the University of Kentucky and Oklahoma State University, United States of America[1]

---

[1] Prepared by Krish Muralidhar and Rathindra Sarathy.

# A Combined Methodology for Assessing Identity and Value Disclosure Risk for Numerical Microdata

Krish Muralidhar
University Of Kentucky
Lexington KY 40506 USA
krishm@uky.edu

Rathindra Sarathy
Oklahoma State University
Stillwater OK 74073 USA
sarathy@okstate.edu

## Abstract

Numerical microdata are often masked prior to release to prevent disclosure of confidential information.  One key aspect of assessing the effectiveness of masking techniques is their ability to prevent *identity* and *value* disclosure.  Identity disclosure refers to the ability of an intruder to match a particular released record as belonging to an individual.  Value disclosure refers to the ability of an intruder to predict the true value of confidential variable(s) using the released microdata.  In this study, we establish a *theoretical basis* to assess whether a masking technique is capable of minimizing (both types of) disclosure risk.  For masking techniques that do not provide minimum disclosure risk, we provide a *common basis* for assessing (both types of) disclosure risk.

## Introduction

Government agencies and other organizations often mask numerical microdata prior to releasing or sharing it with other entities.  The *primary* purpose for such masking is to prevent the disclosure of sensitive or confidential information contained in the data.  Hence, it is critical that we have the ability measure the actual or the potential risk of disclosure resulting from a particular masking technique.

Prior to actually describing this study, we would first like to define disclosure risk since these definitions play an important role in the evaluation of disclosure risk.  Dalenius (1977) provides a general description of disclosure risk as having occurred if an intruder is able to determine the value of a microdata point more accurately with the release of information (than without that information).  Similar generic definitions of disclosure risk have been provided by Duncan and Lambert (1986).  Our interest in disclosure risk is more specific than the generic (but useful) definitions provided above.  We are interested in the specific ability to evaluate a set of masking techniques and disclosure risk resulting from these techniques.  Hence, it is necessary to define disclosure risk in more concrete terms.

In practice, there are two types of disclosure, namely, identity disclosure and value disclosure.  Identity disclosure refers to the case where, using the released data, an

intruder is able to identify a particular released record as belonging to a particular individual. Clearly, this type of disclosure is relevant in situations where the identity of an individual is in itself considered sensitive. Value disclosure occurs if an intruder is able to estimate the value of a confidential variable for a particular record. Whether identity disclosure or value disclosure (or both) are important depends on the particular context. In some situations, being able to identify an individual as belonging to a particular record alone could constitute disclosure (such as when the released data consists of a set of individual with a disease). In others, that an individual belongs to the released data set alone does not constitute disclosure. In these cases, disclosure occurs only when an individual is able to estimate the value of a confidential variable. This situation occurs in the case of organizational databases where that an individual is the employee of the organization does not in itself constitute disclosure. However, if an intruder is able to estimate the value of a confidential variable for this particular individual, then such estimation constitutes disclosure. It is also easy to see that in some situations, it may be necessary for an intruder to first identify the record as belonging to an individual and then estimate the value of a confidential variable in order for disclosure to occur.

In practice, disclosure could also be deterministic (exact) or probabilistic (partial). In exact disclosure, an intruder is able to either identify a particular record as belonging to an individual with certainty or is able to compute the exact value of a confidential variable. As the name implies, in probabilistic disclosure, the intruder is not certain that a identity or value has been disclosed. In terms of identity disclosure, the intruder is able to identify that a record belongs to a particular individual with a high probability and/or estimate the true value of a confidential variable with a greater degree of accuracy. It is very clear that any masking technique that results in deterministic disclosure is unlikely to be used in practice. Hence, in the remainder of the paper, we will use the term "disclosure" to represent "probabilistic" or "partial value" disclosure.

The objective of this paper is to develop a methodology for assessing disclosure risk from two perspectives. First, consistent with the definitions of Dalenius (1977) and Duncan and Lambert (1986), we establish a *theoretical basis* to assess whether a masking technique is capable of minimizing (both types of) disclosure risk. For masking techniques that do not provide minimum disclosure risk, we provide a *common basis* for assessing (both types of) disclosure risk.