

---

# Concepts for the Evaluation of Anonymized Data

Josep Domingo-Ferrer, Josep M. Mateo and Àngel Torres \*)

**Abstract:** We present in this paper some criteria for empirical comparison of SDC methods for continuous microdata. Based on re-identification experiments, we try to optimize the tradeoff between information loss and disclosure risk. SDC methods compared include additive noise, distortion by probability distribution, microaggregation, resampling, rank swapping and a novel approach based on lossy compression. Generic information loss measures (not targeted to specific data uses) are defined, and two approaches to empirical re-identification are used: Euclidean record linkage and probabilistic record linkage. Some weighting schemes to aggregate information loss and disclosure risk measures are discussed and empirical results are given for one of them.

**Keywords:** Statistical disclosure control, Continuous microdata, Record linkage, Re-identification experiments, Information loss measures.

## 1. Introduction

This paper describes an empirical approach to evaluating the protection of statistical microdata. This work was started in the context of the U.S. Census Bureau contract OTTILIE-R (Optimizing the Tradeoff between Information Loss and disclosure risk for continuous microdata) and continued during the European project CASC (Computational Aspects of Statistical Confidentiality). The idea is to define information loss and disclosure risk measures and then construct a score which combines both types of measures. Specifically, the following steps have been taken in this work:

- *Literature analysis.* Literature on SDC for microdata has been analyzed to identify those methods which are relevant for protecting continuous data. In addition, SDC of continuous microdata based on lossy compression has been introduced.
- *Test data.* Test data have been obtained from publicly available microdata files.
- *Disclosure risk assessment.* Two record linkage algorithms have been used to establish the disclosure risk associated to a particular SDC method. In addition, an interval disclosure measure has been defined.
- *Metrics definition.* Information loss actually depends on the data uses to be supported by masked data. Since data uses did not fall within the scope of OTTILIE-R, we have defined a battery of

---

\* ) Universitat Rovira i Virgili, Dept. of Computer Engineering and Mathematics  
E-43006 Tarragona, Catalonia, Spain, E-mail jdomingo@etse.urv.es

---

generic, robust information loss metrics which try to capture structural differences between the original and masked data files.

- *Empirical work*. Experiments carried out are directed to obtaining t-uples of the form (*method, parms, risk, loss*), where *parms* are the input parameters to *method*, *risk* is the percent of re-identified records in the test data set and *loss* is the information loss. The obtained t-uples can be aggregated to rank methods; depending on the aggregation used, different method rankings are conceivable.

Section 2 reviews relevant SDC methods for the protection of continuous microdata. Section 3 lists information loss measures which have been taken into account in experimentation. Section 4 describes record linkage approaches to assessing disclosure risk. Section 5 reports on actual comparison results. Section 6 discusses alternative score constructions for aggregating information loss and disclosure risk. Section 7 is a conclusion.

## 2. Relevant SDC methods for continuous microdata

Sampling methods consist of publishing a sample of records from the original microdata set instead of publishing the whole original microdata set. Sampling methods are suitable for categorical microdata, but their adequacy for continuous microdata is less clear in a general disclosure scenario. The reason is that such methods leave a continuous variable  $V$  unperturbed for all individuals in the sample. Thus, if variable  $V$  is present in an external administrative public file, unique matches with  $V'$  (the version of variable  $V$  restricted to the published sample) are likely, since it is unlikely that a continuous variable (even one truncated due to digital representation) takes exactly the same value for more than one individual. Thus, we will concentrate in what follows on perturbative methods, which have the additional advantage of allowing the entire microdata set to be released.

Perturbative methods distort the microdata set before publication. Perturbative methods considered in our work are a subset of those making sense for continuous microdata:

- *Additive noise* (*Noise $p$*  for short). Gaussian noise is added to the original data to get the masked data (Kim 1986). If the standard deviation of the original variable is  $s$ , noise is generated using a  $N(0, ps)$ . Values of  $p$  considered in the experiments below are 0.01, 0.02, 0.04, 0.06, 0.08 up to 0.2 with 0.02 increments.
- *Data distortion by probability distribution* (*Distr* for short, (Liew/Choi/Liew 1985)). For each variable in the original variable, the best fitted distribution is found; then the fitted distribution is used to generate the masked data set. There are no parameters.
- *Resampling*. Originally proposed for protecting tabular data (Domingo-Ferrer/Mateo-Sanz 1999; Heer 1993), resampling can also be used for microdata. Let  $V$  be an original variable in a dataset with  $n$  records. Take with replacement  $t$  independent samples  $X_1, \dots, X_t$  of size  $n$  of the values of  $V$ .

---

Independently rank each sample (using the same ranking criterion for all samples). Finally, for  $j=1$  to  $n$ , compute the  $j$ -th value  $v'_j$  of the masked variable  $V'$  as the average of the  $j$ -th ranked values in  $X_1, \dots, X_t$ . Resampling has been tested for  $t=1$  (Resamp1) and  $t=3$  (Resamp3).

- *Microaggregation*. Records are clustered into small aggregates or groups of size at least  $k$  (Defays/Nanopoulos 1993; Domingo-Ferrer/Mateo-Sanz 2002). Rather than publishing a variable for a given individual, the average of the values of the variable over the group to which the individual belongs is published. Variants of microaggregation considered include: individual ranking (MicIRk); microaggregation on projected data using z-scores projection (MicZk) and principal components projection (MicPCPk); microaggregation on unprojected multivariate data considering two variables at a time (Mic2mulk), three variables at a time (Mic3mulk), four variables at a time (Mic4mulk) or all variables at a time (Micmulk). Values of  $k$  between 3 and 10 have been considered.
- *Lossy compression (JPEGq)*. This method is new and proposed by these authors for continuous data. The idea is to regard a numerical microdata file as an image (with rows being records and columns being variables). Lossy compression, and more specifically the JPEG algorithm (Joint Photographic Experts Group, Standard IS 10918-1 (ITU-T T.81), <http://www.jpeg.org>), is then used on the image, and the compressed image is interpreted as a masked microdata file. Depending on the lossy compression algorithm used, appropriate mappings between variable ranges and color scales will be needed. The JPEG quality  $q$  has been taken as a parameter with values from 5% up to 100% with 5% increments.
- *Rank swapping (Rankp)*. Although originally described only for ordinal variables, this method can be used for any numerical variable (Moor 1996). First values of variable  $V_i$  are ranked in ascending order; then each ranked value of  $V_i$  is swapped with another ranked value *randomly* chosen within a restricted range (e.g. the rank of two swapped values cannot differ by more than  $p\%$  of the total number of records). The following values of  $p$  have been considered in experimentation: 1, 2, 3, 4, 5, 6, 7 and 10.

### 3. Information loss measures

To evaluate the information loss caused by an SDC method on a continuous microdata set, we want to assess how different the masked data set is from the original data set. We will say there is little information loss if the structure of the masked data set is very similar to the structure of the original data set. In fact, the motivation for preserving the structure of the data set is to ensure that the masked data set will be analytically valid and interesting. We can actually try several complementary ways to assess the preservation of the structure of the original data set:

1. Compare the data in the original and the masked data sets. The more similar the SDC method to the identity function, the less impact (but the higher the disclosure risk!).

---

2. Compare some statistics computed on the original and the masked data sets.

Let  $X$  and  $X'$  be the original and the masked data set. Let  $V$  and  $V'$  be the covariance matrices of  $X$  and  $X'$ , respectively; similarly, let  $R$  and  $R'$  be the correlation matrices. Table 1 summarizes the measures proposed. In this table,  $p$  is the number of variables,  $n$  the number of records, and components of matrices are represented by the corresponding lowercase letters (e.g.  $x_{ij}$  is a component of matrix  $X$ ). Regarding  $X - X'$  measures, it also makes sense to compute those on the averages of variables rather than on all data (see the  $\bar{X} - \bar{X}'$  row in Table 1). Similarly, for  $V - V'$  measures, it is also sensible to compare only the variances of the variables, *i.e.* to compare the diagonals of the covariance matrices rather than the whole matrices (see the  $S - S'$  row in Table 1).

**Table 1. Information loss measures**

	Mean square error	Mean abs. Error	Mean variation
X-X'	$\frac{\sum_{j=1}^p \sum_{i=1}^n (x_{ij} - x'_{ij})^2}{np}$	$\frac{\sum_{j=1}^p \sum_{i=1}^n  x_{ij} - x'_{ij} }{np}$	$\frac{\sum_{j=1}^p \sum_{i=1}^n \frac{ x_{ij} - x'_{ij} }{ x_{ij} }}{np}$
$\bar{X} - \bar{X}'$	$\frac{\sum_{j=1}^p (\bar{x}_j - \bar{x}'_j)^2}{p}$	$\frac{\sum_{j=1}^p \bar{c}_j  \bar{x}_j - \bar{x}'_j ^2}{p \bar{c}_j}$	$\frac{\sum_{j=1}^p \frac{ \bar{x}_j - \bar{x}'_j }{ \bar{x}_j }}{p}$
V-V'	$\frac{\sum_{j=1}^p \sum_{1 \leq i \leq j} (v_{ij} - v'_{ij})^2}{\frac{p(p+1)}{2}}$	$\frac{\sum_{j=1}^p \sum_{1 \leq i \leq j}  v_{ij} - v'_{ij} }{\frac{p(p+1)}{2}}$	$\frac{\sum_{j=1}^p \sum_{1 \leq i \leq j} \frac{ v_{ij} - v'_{ij} }{ v_{ij} }}{\frac{p(p+1)}{2}}$
S-S'	$\frac{\sum_{j=1}^p (v_{jj} - v'_{jj})^2}{p}$	$\frac{\sum_{j=1}^p  v_{jj} - v'_{jj} }{p}$	$\frac{\sum_{j=1}^p \frac{ v_{jj} - v'_{jj} }{ v_{jj} }}{p}$
R-R'	$\frac{\sum_{j=1}^p \sum_{1 \leq i < j} (r_{ij} - r'_{ij})^2}{\frac{p(p-1)}{2}}$	$\frac{\sum_{j=1}^p \sum_{1 \leq i < j}  r_{ij} - r'_{ij} }{\frac{p(p-1)}{2}}$	$\frac{\sum_{j=1}^p \sum_{1 \leq i < j} \frac{ r_{ij} - r'_{ij} }{ r_{ij} }}{\frac{p(p-1)}{2}}$

#### 4. Disclosure risk measures

The assessment of the quality of an SDC method cannot be limited to information loss; disclosure risk is another magnitude that should be measured. The method that optimizes the tradeoff between both magnitudes subject to some user requirements turns out to be the best option.

Literature on disclosure risk is basically related to sampling methods, in which a sample of the original data set is published. Disclosure risk here is measured as the probability that a sample unique is a population unique (Skinner/Marsh/Openshaw/Wymer 1994). If the size of the sample is similar to the size of the whole population, such a probability can be dangerously high; in that case, an intruder who locates a unique value in the released sample could be almost sure that there is a single individual in the population with that value. This could lead to identification of that individual.

---

The uniqueness property as stated above is no longer relevant for perturbative methods, since in this case the whole microdata set is published, but with some distortion. There is not much literature on disclosure risk that can be used for a broad class of perturbative methods; disclosure risk measures tend to be method-specific (measures described in (Adam/Wortmann 1989) are still up-to-date). Empirical methods, like record linkage techniques, provide a more unified approach to disclosure risk assessment for perturbative methods. We briefly describe below two approaches to record linkage and one measure of interval disclosure.

#### **4.1 Distance-based record linkage**

This approach to record linkage is described in (Pagliuca/Seri 1998) for the specific case of microaggregation masking and using the Euclidean distance. However, it can be generalized for any perturbative method provided that a distance between the original and the masked value can be defined. As in any record linkage context, it is assumed that an intruder has an external data set containing as key variables some of the same variables present in the released masked data set. The intruder is assumed to try to link the masked data set with the external data set.

Linkage then proceeds by computing the distances between records in the original and the masked data sets. The distances used are standardized to avoid scaling problems. For each record in the masked data set, the distance to every record in the original data set is computed. Then the "nearest" and "second nearest" records in the original data set are considered. A record in the masked data set is labelled as "linked" when the nearest record in the original data set has the same record number is the corresponding original record). A record in the masked data set is labelled as "linked to 2nd nearest" when the second nearest record in the original data set has the same record number. In all other cases, a record in the masked data set is labelled as "not linked". The percent of "linked" and "linked to 2nd nearest" is a measure of disclosure risk.

#### **4.2 Probabilistic record linkage**

In (Jaro 1989), a probabilistic record linkage method was described and illustrated on the 1985 Census of Tampa, Florida. The matching algorithm uses the linear sum assignment model to "pair" records in the two files to be matched (the original file and the masked file in our case). The percent of correctly paired records is a measure of disclosure risk.

Although less simple than the Euclidean method described in the previous section, this approach is attractive because it only requires the user to provide two probabilities as input: one is an upper bound of the probability of a false match, and the other an upper bound of the probability of false non-match. The Euclidean method above requires rescaling variables as well as an assumption on the weight of

---

variables when computing a distance: for instance, in the proposal of (Pagliuca/Seri 1998), all variables have the same weight.

The U.S. Census Bureau implementation of probabilistic record linkage provided by W. Winkler (U. S. Bureau of Census 2000; Winkler 1998) has been used (with some additions) in the experimentation.

### 4.3 Interval disclosure

For a record in the masked data set, take a rank interval centered on the values of that record as follows: each variable is independently ranked and a rank interval is defined around the value the variable takes on each record; the ranks of values within the interval for a variable around record  $r$  should differ less than  $p\%$  of the total number of records and the rank in the center of the interval should correspond to the value of the variable in record  $r$ . Then the measure is the proportion of original values which fall into the interval centered around their corresponding masked value. A 100% proportion means that an intruder is completely sure that the original value lies in the interval around the masked value (interval disclosure). Values of  $p$  ranging between 1% and 10% have been considered for experimentation.

## 5. Comparison results

A microdata set was constructed using the Data Extraction System (DES) of the U.S. Census Bureau (<http://www.census.gov/DES>). 13 continuous variables were chosen and 1080 records were selected so that there were not many repeated values for any of the variables (in principle, one would not expect repeated values for a continuous variable, but there were repetitions in the data set). Table 2 contains a ranking of methods described in Section 2 (the parameter values described in that section were tried for each method). The Information Loss column (IL) is computed by averaging the mean variations of  $X-X'$ ,  $\bar{X}-\bar{X}'$ ,  $V-V'$ ,  $S-S'$  and the mean absolute error of  $R-R'$ ; the resulting average has been multiplied by 100. The Distance Linkage Disclosure risk column (DLD) contains the average percent of linked records using distance-based record linkage; the average is computed over the number of key variables that the intruder is assumed to know (we have considered knowledge of 1 up to 7 variables). Similarly, the Probabilistic Linkage Disclosure risk column (PLD) is the average percent of correctly paired records using probabilistic linkage. The Interval Disclosure (ID) column contains the average percent of original values falling in the intervals around their corresponding masked values (averages have been computed over all parameter values, i.e. 1% to 10% with 1% increments). Finally, the column Score has been used to rank Table 2 and has been computed as

$$\text{Score} = 0.5 \cdot \text{IL} + 0.125 \cdot \text{DLD} + 0.125 \cdot \text{PLD} + 0.25 \cdot \text{ID}.$$

---

The rationale of the above weighting is to give equal weight to information loss (0.5) and to disclosure risk. The 0.5 weight of disclosure risk is equally divided among ID (0.25) and record linkage. The 0.25 weight of record linkage is equally divided among both approaches to record linkage. The correlation between DLD and PLD is actually 0.962, so both approaches are very similar. The (IL,DLD), (IL,PLD) and (IL,ID) correlations are -0.605, -0.551 and -0.807; thus, the lower the information loss, the higher the disclosure risk, as one would expect. The IL Rank, DLD Rank, PLD Rank and ID Rank columns contain the ranking of each method with respect to IL, DLD, PLD and ID; the lower the rank, the better a method performs (i.e. lower information loss and disclosure risk).



**Table 2. Comparison results**

Method	IL	DLD	PLD	ID	Score	IL Rank	DLD Rank	PLD Rank	ID Rank
Rank15	<b>19.01</b>	1.19	<b>0.15</b>	35.05	<b>18.44</b>	53	<b>6</b>	7	<b>21</b>
Rank19	<b>22.95</b>	0.93	<b>0.08</b>	28.04	<b>18.61</b>	59	<b>2</b>	2	<b>2</b>
Rank16	<b>20.91</b>	1.39	<b>0.11</b>	32.18	<b>18.69</b>	56	<b>8</b>	5	<b>16</b>
Rank13	<b>16.77</b>	2.17	<b>0.12</b>	40.35	<b>18.76</b>	48	<b>12</b>	6	<b>28</b>
Rank14	<b>19.72</b>	1.92	<b>0.07</b>	37.00	<b>19.36</b>	55	<b>10</b>	1	<b>25</b>
Rank11	<b>14.32</b>	2.43	<b>0.25</b>	47.81	<b>19.45</b>	44	<b>13</b>	14	<b>39</b>
Rank12	<b>16.37</b>	2.50	<b>0.25</b>	43.73	<b>19.46</b>	47	<b>14</b>	11	<b>35</b>
iRank20	<b>25.81</b>	0.69	<b>0.09</b>	26.83	<b>19.71</b>	64	<b>1</b>	3	<b>1</b>
Rank18	<b>25.74</b>	0.95	<b>0.09</b>	29.25	<b>20.31</b>	63	<b>4</b>	4	<b>6</b>
Rank10	<b>13.37</b>	3.90	<b>0.38</b>	53.17	<b>20.51</b>	41	<b>24</b>	17	<b>45</b>
Rank17	<b>25.12</b>	1.52	<b>0.20</b>	30.95	<b>20.51</b>	61	<b>9</b>	9	<b>10</b>
Rank09	<b>11.66</b>	5.01	<b>0.52</b>	57.58	<b>20.91</b>	38	<b>37</b>	29	<b>49</b>
Rank08	<b>11.60</b>	6.07	<b>0.85</b>	63.37	<b>22.51</b>	37	<b>39</b>	39	<b>56</b>
Rank07	<b>9.25</b>	7.51	<b>1.08</b>	68.71	<b>22.87</b>	30	<b>41</b>	43	<b>63</b>
Rank06	<b>7.87</b>	9.02	<b>2.79</b>	73.80	<b>23.86</b>	26	<b>43</b>	56	<b>71</b>
Mic3mul07	<b>11.06</b>	19.34	<b>4.70</b>	72.34	<b>26.62</b>	36	<b>68</b>	65	<b>69</b>
Rank05	<b>6.78</b>	16.80	<b>13.60</b>	78.89	<b>26.91</b>	22	<b>58</b>	70	<b>77</b>
Mic3mul09	<b>13.46</b>	19.22	<b>3.44</b>	69.91	<b>27.04</b>	42	<b>67</b>	60	<b>65</b>
Mic3mul10	<b>14.84</b>	17.99	<b>3.44</b>	68.61	<b>27.25</b>	46	<b>64</b>	59	<b>62</b>
Mic4mul04	<b>12.14</b>	19.76	<b>6.67</b>	71.85	<b>27.33</b>	39	<b>69</b>	68	<b>68</b>
Mic4mul05	<b>14.50</b>	17.43	<b>5.45</b>	69.09	<b>27.39</b>	45	<b>61</b>	66	<b>64</b>
Mic3mul08	<b>13.51</b>	20.81	<b>4.15</b>	70.68	<b>27.54</b>	43	<b>71</b>	63	<b>66</b>
Mic4mul08	<b>18.89</b>	17.78	<b>3.35</b>	62.84	<b>27.80</b>	52	<b>62</b>	58	<b>55</b>
Mic3mul06	<b>10.24</b>	20.41	<b>13.90</b>	74.00	<b>27.91</b>	33	<b>70</b>	71	<b>72</b>
Mic4mul07	<b>19.36</b>	17.10	<b>2.08</b>	64.41	<b>28.18</b>	54	<b>60</b>	53	<b>58</b>
Mic4mul06	<b>17.91</b>	17.82	<b>3.98</b>	66.41	<b>28.28</b>	50	<b>63</b>	62	<b>60</b>
Mic4mul09	<b>21.35</b>	15.93	<b>2.00</b>	61.66	<b>28.33</b>	58	<b>57</b>	52	<b>54</b>
Mic4mul10	<b>22.98</b>	16.85	<b>2.37</b>	60.56	<b>29.03</b>	60	<b>59</b>	55	<b>51</b>
Mic3mul05	<b>9.73</b>	23.78	<b>18.29</b>	76.59	<b>29.27</b>	31	<b>76</b>	73	<b>74</b>
Mic3mul04	<b>7.45</b>	23.49	<b>22.75</b>	79.14	<b>29.29</b>	24	<b>75</b>	75	<b>79</b>
Mic4mul03	<b>10.69</b>	22.88	<b>16.69</b>	76.89	<b>29.51</b>	35	<b>74</b>	72	<b>75</b>
Rank04	<b>5.90</b>	22.77	<b>22.78</b>	84.12	<b>29.67</b>	20	<b>73</b>	76	<b>86</b>
Micmul03	<b>27.67</b>	14.26	<b>1.88</b>	57.23	<b>30.16</b>	65	<b>54</b>	50	<b>47</b>
Micmul04	<b>31.74</b>	13.72	<b>1.38</b>	52.44	<b>30.86</b>	67	<b>53</b>	48	<b>44</b>
Mic3mul03	<b>6.29</b>	29.70	<b>29.06</b>	82.95	<b>31.23</b>	21	<b>79</b>	80	<b>85</b>
Micmul05	<b>35.12</b>	11.73	<b>1.14</b>	48.43	<b>31.27</b>	70	<b>46</b>	44	<b>41</b>
Micmul07	<b>37.68</b>	13.20	<b>1.20</b>	43.46	<b>31.50</b>	72	<b>52</b>	45	<b>34</b>
Micmul06	<b>38.77</b>	13.00	<b>1.22</b>	45.76	<b>32.60</b>	73	<b>50</b>	46	<b>37</b>
Micmul08	<b>41.53</b>	13.12	<b>0.99</b>	42.66	<b>33.19</b>	75	<b>51</b>	42	<b>32</b>
Rank03	<b>5.07</b>	31.73	<b>36.92</b>	89.53	<b>33.50</b>	18	<b>80</b>	83	<b>93</b>
Mic2mul10	<b>10.68</b>	49.38	<b>27.29</b>	77.43	<b>34.28</b>	34	<b>86</b>	78	<b>76</b>
Micmul10	<b>44.69</b>	14.66	<b>0.50</b>	40.41	<b>34.34</b>	76	<b>55</b>	27	<b>29</b>
Noise0.16	<b>32.56</b>	15.65	<b>4.66</b>	64.39	<b>34.91</b>	68	<b>56</b>	64	<b>57</b>
Micmul09	<b>45.98</b>	12.82	<b>0.85</b>	40.99	<b>34.95</b>	79	<b>49</b>	40	<b>30</b>

Table 2. Comparison results

Method	IL	DLD	PLD	ID	Score	IL Rank	DLD Rank	PLD Rank	ID Rank
Mic2mul09	<b>9.93</b>	51.03	<b>33.04</b>	78.94	<b>35.21</b>	32	<b>87</b>	81	<b>78</b>

Mic2mul08	<b>8.55</b>	54.31	<b>33.70</b>	79.77	<b>35.22</b>	27	<b>88</b>	82	<b>80</b>
Mic2mul07	<b>7.53</b>	54.72	<b>37.41</b>	81.40	<b>35.63</b>	25	<b>89</b>	84	<b>83</b>
Noise0.12	<b>25.24</b>	22.21	<b>22.39</b>	71.58	<b>36.09</b>	62	<b>72</b>	74	<b>67</b>
Noise0.1	<b>21.14</b>	27.70	<b>29.03</b>	75.20	<b>36.46</b>	57	<b>78</b>	79	<b>73</b>
Mic2mul06	<b>7.03</b>	56.38	<b>42.00</b>	82.89	<b>36.54</b>	23	<b>90</b>	86	<b>84</b>
JPEG080	<b>33.97</b>	19.13	<b>6.93</b>	66.35	<b>36.83</b>	69	<b>65</b>	69	<b>59</b>
Noise0.14	<b>35.13</b>	19.21	<b>6.24</b>	67.62	<b>37.65</b>	71	<b>66</b>	67	<b>61</b>
Noise0.18	<b>41.12</b>	11.96	<b>3.52</b>	60.95	<b>37.73</b>	74	<b>47</b>	61	<b>52</b>
Noise0.08	<b>17.43</b>	36.06	<b>39.76</b>	79.84	<b>38.15</b>	49	<b>82</b>	85	<b>81</b>
Rank02	<b>2.90</b>	47.26	<b>57.47</b>	94.56	<b>38.18</b>	11	<b>85</b>	90	<b>96</b>
JPEG070	<b>44.92</b>	9.66	<b>2.34</b>	57.28	<b>38.28</b>	77	<b>44</b>	54	<b>48</b>
Noise0.2	<b>45.97</b>	10.01	<b>0.97</b>	57.63	<b>38.77</b>	78	<b>45</b>	41	<b>50</b>
Mic2mul05	<b>5.88</b>	58.97	<b>56.84</b>	85.40	<b>38.77</b>	19	<b>92</b>	89	<b>88</b>
JPEG085	<b>29.47</b>	23.85	<b>24.48</b>	72.80	<b>38.98</b>	66	<b>77</b>	77	<b>70</b>
Mic2mul04	<b>4.90</b>	61.53	<b>60.69</b>	87.26	<b>39.54</b>	17	<b>94</b>	91	<b>89</b>
JPEG090	<b>18.17</b>	35.37	<b>46.98</b>	80.87	<b>39.60</b>	51	<b>81</b>	87	<b>82</b>
Noise0.06	<b>13.03</b>	45.54	<b>56.22</b>	84.16	<b>40.28</b>	40	<b>84</b>	88	<b>87</b>
Mic2mul03	<b>3.28</b>	66.97	<b>64.79</b>	90.51	<b>40.74</b>	15	<b>95</b>	92	<b>94</b>
Noise0.04	<b>8.93</b>	58.51	<b>65.28</b>	88.95	<b>42.18</b>	28	<b>91</b>	94	<b>90</b>
JPEG075	<b>50.45</b>	12.67	<b>2.90</b>	61.27	<b>42.49</b>	80	<b>48</b>	57	<b>53</b>
JPEG095	<b>9.06</b>	60.11	<b>66.56</b>	89.23	<b>42.67</b>	29	<b>93</b>	96	<b>92</b>
Resamp3	<b>3.15</b>	67.90	<b>67.63</b>	96.81	<b>42.72</b>	14	<b>96</b>	97	<b>97</b>
Rank01	<b>2.34</b>	69.19	<b>66.35</b>	99.54	<b>43.00</b>	9	<b>97</b>	95	<b>106</b>
JPEG065	<b>57.77</b>	7.02	<b>1.90</b>	53.87	<b>43.47</b>	81	<b>40</b>	51	<b>46</b>
Noise0.02	<b>4.24</b>	77.34	<b>71.32</b>	94.42	<b>44.31</b>	16	<b>99</b>	98	<b>95</b>
Resamp1	<b>3.11</b>	75.42	<b>71.85</b>	98.36	<b>44.56</b>	13	<b>98</b>	99	<b>99</b>
MicPCP03	<b>69.62</b>	3.16	<b>0.77</b>	38.41	<b>44.90</b>	84	<b>17</b>	38	<b>26</b>
JPEG055	<b>63.70</b>	5.57	<b>1.26</b>	49.70	<b>45.13</b>	83	<b>38</b>	47	<b>42</b>
Noise0.01	<b>2.57</b>	85.19	<b>74.13</b>	97.03	<b>45.46</b>	10	<b>100</b>	103	<b>98</b>
JPEG100	<b>3.06</b>	87.14	<b>73.03</b>	99.14	<b>46.34</b>	12	<b>101</b>	100	<b>101</b>
MicIR10	<b>1.19</b>	97.37	<b>74.07</b>	99.12	<b>46.81</b>	8	<b>102</b>	102	<b>100</b>
MicIR08	<b>1.03</b>	97.84	<b>74.07</b>	99.29	<b>46.83</b>	6	<b>108</b>	101	<b>103</b>
MicIR09	<b>1.14</b>	97.96	<b>74.40</b>	99.24	<b>46.93</b>	7	<b>109</b>	104	<b>102</b>
MicIR06	<b>0.87</b>	97.66	<b>75.28</b>	99.51	<b>46.93</b>	5	<b>106</b>	105	<b>105</b>
MicIR05	<b>0.69</b>	97.58	<b>75.99</b>	99.58	<b>46.94</b>	3	<b>104</b>	106	<b>107</b>
MicIR03	<b>0.45</b>	97.39	<b>78.96</b>	99.79	<b>47.22</b>	1	<b>103</b>	107	<b>109</b>
MicIR04	<b>0.64</b>	97.63	<b>79.78</b>	99.67	<b>47.41</b>	2	<b>105</b>	108	<b>108</b>
MicIR07	<b>0.81</b>	97.79	<b>88.06</b>	99.42	<b>48.49</b>	4	<b>107</b>	109	<b>104</b>
MicPCP04	<b>78.84</b>	3.43	<b>0.62</b>	36.00	<b>48.92</b>	87	<b>19</b>	32	<b>23</b>
JPEG050	<b>73.20</b>	4.26	<b>0.67</b>	47.96	<b>49.21</b>	86	<b>31</b>	36	<b>40</b>
JPEG060	<b>71.24</b>	7.66	<b>1.52</b>	51.71	<b>49.69</b>	85	<b>42</b>	49	<b>43</b>
MicPCP05	<b>82.55</b>	3.94	<b>0.69</b>	34.10	<b>50.38</b>	88	<b>25</b>	37	<b>20</b>
MicPCP07	<b>89.28</b>	4.02	<b>0.62</b>	32.56	<b>53.36</b>	91	<b>27</b>	33	<b>17</b>

Table 2. Comparison results

Method	IL	DLD	PLD	ID	Score	IL Rank	DLD Rank	PLD Rank	ID Rank
MicPCP09	<b>90.78</b>	4.54	<b>0.25</b>	31.40	<b>53.84</b>	94	<b>34</b>	12	<b>13</b>
MicPCP06	<b>90.26</b>	3.37	<b>0.50</b>	33.42	<b>53.97</b>	93	<b>18</b>	26	<b>19</b>
MicZ03	<b>90.25</b>	3.16	<b>0.61</b>	35.71	<b>54.52</b>	92	<b>16</b>	31	<b>22</b>
JPEG035	<b>88.80</b>	3.65	<b>0.44</b>	43.20	<b>55.71</b>	90	<b>20</b>	23	<b>33</b>

JPEG045	<b>87.55</b>	4.15	<b>0.67</b>	46.78	<b>56.07</b>	89	<b>30</b>	35	<b>38</b>
MicZ04	<b>94.94</b>	3.70	<b>0.53</b>	33.04	<b>56.26</b>	96	<b>21</b>	30	<b>18</b>
MicPCP08	<b>96.93</b>	3.97	<b>0.34</b>	32.04	<b>57.02</b>	97	<b>26</b>	16	<b>14</b>
MicPCP10	<b>97.82</b>	4.13	<b>0.46</b>	31.19	<b>57.28</b>	98	<b>29</b>	24	<b>11</b>
JPEG040	<b>90.99</b>	3.72	<b>0.66</b>	44.98	<b>57.29</b>	95	<b>22</b>	34	<b>36</b>
MicZ07	<b>102.87</b>	4.27	<b>0.38</b>	30.53	<b>59.65</b>	99	<b>32</b>	20	<b>9</b>
MicZ06	<b>103.92</b>	3.88	<b>0.41</b>	30.43	<b>60.10</b>	100	<b>23</b>	21	<b>8</b>
MicZ05	<b>104.06</b>	4.03	<b>0.42</b>	31.30	<b>60.41</b>	101	<b>28</b>	22	<b>12</b>
MicZ08	<b>107.92</b>	4.55	<b>0.52</b>	29.60	<b>61.99</b>	102	<b>35</b>	28	<b>7</b>
MicZ10	<b>109.79</b>	4.83	<b>0.38</b>	28.20	<b>62.59</b>	103	<b>36</b>	18	<b>3</b>
MicZ09	<b>110.91</b>	4.35	<b>0.38</b>	28.36	<b>63.14</b>	105	<b>33</b>	19	<b>4</b>
Distr	<b>58.62</b>	43.05	<b>64.88</b>	88.98	<b>65.04</b>	82	<b>83</b>	93	<b>91</b>
JPEG030	<b>110.48</b>	3.02	<b>0.48</b>	41.79	<b>66.12</b>	104	<b>15</b>	25	<b>31</b>
JPEG025	<b>155.15</b>	2.13	<b>0.25</b>	38.76	<b>87.56</b>	106	<b>11</b>	13	<b>27</b>
JPEG020	<b>164.91</b>	1.36	<b>0.29</b>	36.11	<b>91.69</b>	107	<b>7</b>	15	<b>24</b>
JPEG015	<b>202.66</b>	1.10	<b>0.15</b>	32.06	<b>109.50</b>	108	<b>5</b>	8	<b>15</b>
JPEG010	<b>269.38</b>	0.93	<b>0.22</b>	28.44	<b>141.94</b>	109	<b>3</b>	10	<b>5</b>

## 6. Alternative score constructions

In (Torres 2003), an alternative score construction is proposed which can be used instead of the one proposed in the previous section. Like in the previous score, information loss and disclosure risk are both assigned 0.5 weight. The differences are in the way of information loss and disclosure risk are computed.

IL is computed as the 100 times the average of IL1, IL2 and IL3, where IL1 is the mean variation of  $X-X'$ ; IL2 is the average of the mean variation of  $\bar{X}-\bar{X}'$  and the mean variation of  $S-S'$ ; and IL3 is the average of the mean variation of  $V-V'$  and the mean absolute error of  $R-R'$ . Thus, IL1 represents the discrepancy between raw original and raw protected data, IL2 represents the discrepancy between univariate statistics for original and protected data, while IL3 represents the discrepancy between bivariate statistics for original and protected data. The philosophy here is to give equal weight to those three discrepancies.

Disclosure risk assessment does not use probabilistic record linkage (only distance-based record linkage DLD is used), which allows larger experiments to be carried out in less time. Two kinds of interval disclosure are distinguished: ID1, which is based on ranks and exactly corresponds to ID as defined in Sections 4.3 and used in Section 5, and ID2, which is based on standard deviations. ID2 is defined as the average percent of original values falling in the intervals around their corresponding masked values, where intervals are centered on the actual masked values (not on their rank) and interval widths are  $p\%$  of the standard deviation of the variable.

Thus, this score can be formalized as:

$$\text{Score2} = 0.5 \cdot ((\text{IL1} + \text{IL2} + \text{IL3}) / 3) + 0.25 \cdot \text{DLD} + 0.125 \cdot \text{ID1} + 0.125 \cdot \text{ID2}$$

---

Results obtained with this new score are similar to those reported in Section 5 and Table 2 regarding the relative ranking of families of methods. Rank swapping scores best, closely followed by microaggregation. The differences show up in the best scoring parameterizations: while in Table 2, Rank15 was best, the best method with Score2 is Rank07. Also, with Score2 microaggregation Micmul16 and Micmul10 appear as the fourth and eighth best scoring methods, respectively.

## 7. Conclusions

There is a rich array of methods for microdata disclosure limitation. A set of proposals for continuous microdata have been identified and described in this paper. Measures for assessing information loss have also been described. Experimental results presented in Table 2 are self-explanatory. One thing that stands out is that rankswapping with parameter around 10% is a very good option; next follows multivariate microaggregation taking groups of three or four variables at a time; for microaggregation, the group size has no significant effect. Data distortion by probability distribution turns out to perform very poorly. For most methods, performance depends on parameter choice, even if some methods are more parameter-dependent than other.

Changing the way the score is constructed can change the experimental results to a limited extent: the relative ranking between families of methods (rank swapping, microaggregation, etc.) is likely to stay the same as long as the score strikes a balance between information loss and disclosure risk.

## Acknowledgments

This work was partly funded by the U.S. Bureau of the Census under contract OTTILIE-R (ref. no OBLIG-2000-29158-0-0) and by the European Commission under project CASC (ref. no. IST-2000-25069). Thanks go to Francesc Seb  for his help in automating the probabilistic record linkage software and running the experiments. The contribution of Sarah Giessing as Ko-Referent is also gratefully acknowledged.

## References

- Adam, N. R., Wortmann, J. C., (1989), Security-control methods for statistical databases: a comparative study, *ACM Computing Surveys*, vol. 21(4):515-556.
- Defays, D., Nanopoulos, P., (1993), Panels of enterprises and confidentiality: the small aggregates method, in *Proc. of 92 Symposium on Design and Analysis of Longitudinal Surveys*, Ottawa: Statistics Canada, 195-204.
- Domingo-Ferrer, J., Mateo-Sanz, J. M., (1999), On resampling for statistical confidentiality in contingency tables, *Computers & Mathematics with Applications*, vol 38: 13-32.

- 
- Domingo-Ferrer, J., Mateo-Sanz, J.M., (2002), Practical Data-Oriented Microaggregation for Statistical Disclosure Control, IEEE Transactions on Knowledge and Data Engineering, vol 14(1): 189-201.
- Heer, G. R., (1993), A bootstrap procedure to preserve statistical confidentiality in contingency tables, in Proceedings of the International Seminar on Statistical Confidentiality (ed. D. Lievesley), Luxemburg: Office for Official Publications of the European Communities, 261-271.
- Jaro, M. A., (1989), Advances in record-linkage methodology as applied to matching the 1985 Census of Tampa, Florida, Journal of the American Statistical Association, vol. 84:414-420.
- Joint Photographic Experts Group, Standard IS 10918-1 (ITU-T T.81) <http://www.jpeg.org>.
- Kim, J. J., (1986), A method for limiting disclosure in microdata based on random noise and transformation, in Proc. of the ASA Sect. on Survey Res. Meth., pp. 303-308.
- Liew, C. K., Choi, U. J., Liew, C. J., (1985), A data distortion by probability distribution, ACM Transactions on Database Systems, vol. 10: 395-411.
- Moore, R., (1996), Controlled data swapping techniques for masking public use microdata sets, U. S. Bureau of the Census (unpublished manuscript).
- Pagliuca, D., Seri, G., (1998), Some Results of Individual Ranking Method on the System of Enterprise Accounts Annual Survey, Esprit SDC Project, Deliverable MI-3/D2.
- Skinner, C., Marsh, C., Openshaw, S., Wymer, C., (1994), Disclosure Control for Census Microdata, Journal of Official Statistics, vol. 10:31-51.
- Torres, A. (2003), Contributions to Microaggregation for Statistical Data Protection, Ph.D. Dissertation (in Catalan), Polytechnical University of Catalonia, Barcelona, 2003. Advisors: J. Domingo-Ferrer and J. M. Mateo-Sanz.
- U. S. Bureau of the Census, (2000), Record Linkage Software: User Documentation. Available from U. S. Bureau of the Census.
- Winkler, W., (1998), Re-identification methods for evaluating the confidentiality of analytically valid microdata, in Statistical Data Protection, Luxembourg: Office for Official Publications of the European Communities, 1999. Journal version in Research in Official Statistics, vol. 1(2): 50-69, 1998.