

Safety Rules in Statistical Disclosure Control for Tabular Data

Giovanni M. Merola

Istituto Nazionale di Statistica, Dept. MPS

via C. Balbo 16, 00184, Roma, Italy

e-mail: gmmerola@istat.it

Abstract

We consider different safety rules currently used in Statistical Disclosure Control for tabular data. We generalize these rules and show how they can all be expressed as bounds on the relative error of estimation for partial sums. Thus we relate the Dominance rule to the p -rule and (p, q) -rule. Then we consider safety rules based on a different estimating procedure, showing that these are of the same kind of the previous ones. These different safety rules are then compared to each other and to the previous ones.

KEY WORDS: Disclosure risk; Dominance rule; p -rule; sensitivity measures.

Acknowledgements

We would like to thank Dr. Luisa Franconi and Mr. Giovanni Seri of ISTAT for some useful discussions. This work was partially supported by the European Union project IST-2000-25069 CASC on 'Computational Aspects of Statistical Confidentiality'. The views expressed are those of the author and do not necessarily reflect the policies of the Istituto Nazionale di Statistica.

1 Introduction

Statistical Disclosure Control (SDC) is used to protect confidential information that is to be released. It is mainly applied by official statistical institutes but also by other data providers, such as, for example, those dealing with clinical statistics. We consider SDC applied to tables with cells containing sums of nonnegative contributions, which are the values taken by a variable for each unit belonging to a cell. Such tables are the most common ones and also the easiest to violate. Cell totals are not considered confidential unless their knowledge allows the disclosure of confidential individual contributions. In most cases, contributions are not considered confidential unless they can be matched with the identity of the contributor. Often, large contributions can be matched more easily than others; therefore,

care must be taken with cells containing few dominating contributions. Generally, cells with few contributions require stricter protection, as also do cells containing enterprise data, for which the likelihood of a systematic attack is higher. Such considerations, however, are beyond the scope of this paper and we consider protecting any given latent information in a cell, regardless of its identifiability. We are concerned with the stage of SDC in which cells at risk are identified, hence not with the later protection stage.

Let C be a generic cell of a magnitude table with n nonnegative contributions z_i , for which the total

$$z = \sum_{i=1}^n z_i$$

is to be released. Without loss of generality, the contributions are indexed in nonincreasing order, so that $z_1 \geq z_2 \geq \dots \geq z_n \geq 0$. For generality we do not exclude the possibility that one or more, but not all, contributions could be zero. When not needed, we will drop the reference to the cell C , with the understanding that all quantities considered refer to the same cell. The sum of the $m \leq n$ largest contributions will be denoted as

$$t_m = \sum_{i=1}^m z_i$$

and the sum of the $(n - m)$ lowest contributions with $r_m = \sum_{i=m+1}^n z_i$, so that

$z = t_m + r_m = t_m + z_{m+1} + r_{m+1}$. It is always verified that

$$\frac{m}{n}z \leq t_m \leq z. \quad (1)$$

The risk of disclosure for a cell can be measured with respect to different criteria. Given a measure of risk, safety rules set a bound on the acceptable risk and a cell is considered at risk, or *sensitive*, if its contributions do not meet the safety conditions set by the rule. A safety rule almost always applied is the 3-rule, by which cells with less than three contributions are at risk. The rationale for this rule is that: if $n = 1$, releasing the total is equivalent to declaring the value to be protected and any intruder knowing n can estimate exactly z_1 ; if $n = 2$, anyone of the two respondents of that cell knowing n can estimate exactly the other contribution. Taking one of the respondents to be the intruder makes realistic the assumption of an intruder with knowledge of one contribution. Note that if n is not known the intruder cannot estimate exactly a contribution in either case. The 3-rule relies only on the number of respondents of a cell and it is often applied together with other rules based also on the values of the contributions, such as the widely accepted (m, k) -dominance rule (Cox 1981). By this rule, a cell is not sensitive if the sum of the first m contributions is less than $100k$ percent of the total. That is, if

$$\frac{t_m}{z} < k, \quad (2)$$

with k real in $(0, 1)$. The rationale for this rule, is the identifiability of large outliers, mentioned above. However, there does not seem to be agreement on which values of the parameters m and k should be chosen. In the current practice the choice is left to the protector's perception of concentration; common values are $m = 2, 3$ and $0.6 \leq k \leq 0.8$. de Wolf (2001) notes that a cell with minimal concentration, that is with n equal contributions, should not be considered at risk. In this case, letting $\bar{z} = z/n$ denote the mean contribution, $t_m = m\bar{z}$ and, therefore, the parameters should always be chosen such that $k \geq m/n$. This agrees with the lower bound in Equation (1). If the value of k had to be chosen with respect to some measures of concentration, it would vary with respect to the ratio m/n .

The 3-rule prevents the exact prediction of one contribution. In the p -rule (Cox 1981) a cell is safeguarded from predictive disclosure putting a lower bound on the absolute relative error of estimation of one contribution by one of the respondents of the cell. By this rule, the intruder knowing a contribution z_k would estimate a contribution z_i with $\hat{z}_i = z - z_k$. Denoting with $RE(z_i; z_k)$ the absolute relative error of estimation for such scenario, it can be shown that it is minimal when the second respondent estimates the largest contribution. Thus the p -rule is expressed as

$$RE(z_1; z_2) = \frac{|z - z_2 - z_1|}{z_1} < p, \quad (3)$$

where the threshold $p > 0$ is real. The p -rule was actually used for the protection of the 1992 US Economic Census data (Jewett 1993, cited in Willenborg and de Waal 2000). The *prior-posterior* rule (Cox 1981), or (p, q) -rule, widens the scenario of the p -rule assuming that the intruder estimates the remainder $(z - z_k - z_i)$ with $100q$ percent of error. Thus a contribution z_i is estimated by $\hat{z}_i = z - z_k - (1 + q)(z - z_k - z_i)$. Like for the p -rule, the lowest absolute relative error of estimation is obtained by the second contributor estimating z_1 . The prior-posterior rule is thus defined as

$$RE(z_1; z_2, q) = \frac{|z - z_2 - (1 + q)(z - z_2 - z_1) - z_1|}{z_1} < p, \quad (4)$$

where $p > |q|$, because it is required that the a-posteriori precision is not better than the a-priori precision.

In the next section we generalize and cast the currently used safety rules into the same predictive framework, establishing relations among them. In Section 3 we consider measures of risk for an intruder assuming uniform distribution of the quantities to be estimated, obtaining other safety rules. In Section 4 we compare the various safety rules and give some concluding remarks.

2 Dominance Rule and Sensitivity of Sums of Contributions

SDC for tabular data's operational tools are sensitivity measures. These are functions of the contributions of the cell that take positive value if the corresponding safety rule is not satisfied. Since they are only used as indicators of a safety rule being satisfied or not, proportional sensitivity measures are equivalent. A general class of such measures is that of *linear sensitivity measures* (LSM) (Cox 1981).

For our purposes a linear sensitivity measure is a function

$$S(C) = \sum_{i=1}^n \lambda_i z_i + d, \quad (5)$$

with λ_i , $i = 1, \dots, n$, and d real, for which a cell C is sensitive if $S(C) > 0$. Note that we can also write $S(C) = \sum_{i=1}^n [\lambda_i + d/(nz_i)] z_i$. For example, a LSM for the safety rule $n > m$, that generalizes the 3-rule stated above, can be built taking as weights $\lambda_i = -1/z_i$, $i = 1, \dots, n$; so that $S(C) = m + \sum_{i=1}^n -z_i/z_i = m - n$. A LSM is subadditive if $S(C_i \cup C_j) \leq S(C_i) + S(C_j)$, where C_i and C_j are two cells of the same table. Subadditivity ensures that the union of two or more nonsensitive cells is also nonsensitive, hence it is important because it ensures that the whole table is nonsensitive if all of its cells are nonsensitive. A safety rule with subadditive LSM is subadditive. It can be shown that a LSM with weights λ_i

is subadditive if and only if $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ (Cox 1981). Therefore, the LSM for the rule $n > m$ is subadditive. Equivalent forms of the LSM corresponding to the (m, k) -dominance rule in Equation (2) are

$$S = t_m - kz = (1 - k)t_m - kr_m \propto \frac{(1 - k)}{k}t_m - r_m. \quad (6)$$

This LSM, in the form (5), has weights proportional to $\lambda_j = (1 - k)$ for $j \leq m$ and to $\lambda_j = -k$ for $j > m$, hence it is subadditive.

In this section we consider predictive disclosure risk, which, like in the p -rule, is measured by the minimal achievable absolute relative error of estimation. Different rules are defined with respect to different predictive scenarios. Predictive scenarios differ for the intruder's prior knowledge, predictive aim, s , and estimate, \hat{s} . Safety rules bound the disclosure risk with a threshold $p > 0$ and will be denoted as

$$P_p(s; \dots) : \min_{\text{scenario}} RE(s; \dots) = \min_{\text{scenario}} \frac{|\hat{s} - s|}{s} > p,$$

where (\dots) indicate parameters of the predictive scenario. When not needed, we will drop the subscript for the lower bound, with the understanding that the quantities are defined for any $p > 0$. A rule P is implied by another rule P' if it is always satisfied when P' is satisfied. A sensitivity measure corresponding to the rule $P_p(s; \dots)$ will be denoted as $S_p(s; \dots)$. Henceforth, when not otherwise

specified, we will refer to sensitivity measures in the form: $S_p(s; \dots) = ps - |\hat{s} - s|$, as these forms referred to the same s enjoy the inverse ordering of the corresponding RE in different scenarios.

In this section we consider different predictive scenarios in which an intruder is interested in estimating a sum of $m < n$ contributions of a cell, denoted generically by s_m , and estimates this sum with its maximal possible value in the scenario; therefore, the corresponding rules will be generically referred to as M -rules. The prior knowledge of the intruder will be summarized by some parameters but the basic knowledge of the cell total, z , and that the number of contributions is larger than some integer m , is always assumed and will be omitted in the notation.

The simplest situation that we consider is that of an intruder with only basic knowledge who estimates s_m with $\hat{s}_m = z$. If the safety rule $P(s_m)$ is adopted the following proposition holds:

Proposition 2.1 *Safety rule $P(s_m) : \min_{\{s_m\}} RE(s_m) \geq p$, where s_m is the sum of any $m < n$ contributions of a cell and $\hat{s}_m = z$, has LSM $S(s_m) = pt_m - r_m$, which is subadditive and equivalent to the (m, k) -dominance in Equation (6) with $k = \frac{1}{1+p}$.*

Proof. The sums of m contributions are bounded such that $r_{n-m+1} \leq s_m \leq t_m$,

thus, by setting $\hat{s}_m = z$,

$$RE(s_m) = \frac{z - s_m}{s_m} = \frac{z}{s_m} - 1 \geq \frac{r_m}{t_m} = RE(t_m).$$

Hence, the sum of m contributions that gives lowest relative error of estimation in this scenario is the sum of the m largest ones, t_m . The LSM for this rule is obtained from the condition $RE(t_m) < p$ and it is

$$S(s_m) = pt_m - r_m, \tag{7}$$

which is equivalent to the LSM for the (m, k) -Dominance in (6) with $k = 1/(1 + p)$ and, thus, is subadditive.

q.e.d.

Since $P(s_m)$ is equivalent to the (m, k) -dominance we will simply call it Dominance. The choice of k in the (m, k) -dominance gives lower bound $p = (1 - k)/k$ in the Dominance. The values of p corresponding to different choices of k in the (m, k) -dominance rule, shown in Table 1, show that the values commonly chosen for k correspond to values of p between $\frac{1}{4}$ and $\frac{1}{3}$.

k	0.01	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
p	99	9	4	2.33	1.5	1	0.67	0.43	0.25	0.11	0

Table 1: Values of k in the (m, k) -dominance rule corresponding to different values of p in the Dominance rule.

The equivalence of the two rules is confirmed by another consideration. From Equation (1) follows that $RE(t_m) \leq (n - m)/m$. Hence, for this rule to make sense, p should always be chosen such that $p < (n - m)/m$. Substituting $k = 1/(1 + p)$ in the latter inequality we obtain the bound $n \geq \lfloor m/k \rfloor$, set for the Dominance. Elsewhere (e.g. Cox 1981 and Willenborg and de Waal 2000), the (m, k) -Dominance is interpreted in a predictive context in terms of the relative error of estimation for z_1 when a coalition of $m - 1$ contributors intrudes, however in a somewhat less convincing way.

Scenarios for stricter safety rules allow for the intruder to be one of the respondents. Consider generalizing the p -rule to estimating s_m with $\hat{s}_m = z - z_K$. Denoting with $\gamma_K = z_K/z$, the estimate can be written as $\hat{s}_m = (1 - \gamma_K)z$. Adopting the usual safety rule gives the following proposition.

Proposition 2.2 *Let the generalized p -rule be $P(s_m; z_K)$:*

$\min_{\{s_m, z_K\}} RE(s_m; z_K) \geq p$, where s_m is any sum of $m < n$ contributions for given m , z_K does not contribute to s_m and $\hat{s}_m = z - z_K$. This rule is equivalent to the Dominance rule (2) with bound $k = (1 - \gamma_K)/(1 + p)$; its SLM is subadditive and equal to $S(s_m; z_K) = pt_m - r_{m+1}$. The same rule with $m = K - 1$ would be obtained for given K and variable $m < K$.

Proof. Obviously, it holds that $s_m + z_\kappa \leq t_{m+1}$, $\forall \kappa$ and s_m , therefore

$$RE(s_m; z_\kappa) = \frac{z - z_\kappa - s_m}{s_m} \geq \frac{z - t_{m+1}}{t_{m+1} - z_\kappa} \geq \frac{z - z_{m+1} - t_m}{t_m} = RE(t_m; z_{m+1}).$$

$\min_{\{s_m, \kappa\}} RE(s_m; z_\kappa)$ is thus achieved for $s_m = t_m$ and $\kappa = m + 1$. Therefore, the safety rule is satisfied when $RE(t_m; z_{m+1}) > p$, that is when

$$\frac{t_m}{z} \leq (1 - \gamma_\kappa) \frac{1}{1 + p}. \quad (8)$$

It easily seen that the corresponding LSM is $S(s_m; z_\kappa) = pt_m - r_{m+1}$ which can be expressed in the form (5) with $d = 0$ and weights $\lambda_i = p$ for $i \leq m$, $\lambda_{m+1} = 0$ and $\lambda_i = -1$ for $i > m + 1$, from which follows the subadditivity. If we had, instead, fixed the rank κ of the known contribution and left undetermined $m < \kappa$, we would have obtained the same LSM, expressed as $pt_{\kappa-1} - r_\kappa$, by the same argument as above.

q.e.d.

The p -rule in (3) is a particular case of the above rule for $m = 1$.

The generalization of the (p, q) -rule to the estimation of a sum of m contributions is considered in the following proposition.

Proposition 2.3 *Let the generalized (p, q) -rule be $P(s_m; z_\kappa, q)$:*

$\min_{\{s_m, \kappa\}} RE(s_m; z_\kappa, q) \geq p$, where s_m is any sum of $m < n$ contributions for

given m , z_{κ} does not contribute to s_m and $\hat{s}_m = z - z_{\kappa} - (1 + q)(z - z_{\kappa} - s_m)$ with $0 < |q| < p$. The rule is equivalent to the Dominance rule (2) with bound $k = 1/[1 + (p/|q|)]$ and to the generalized p -rule with threshold $p' = p/|q|$; its LSM is subadditive and equal to $S_p(s_m; z_{\kappa}, q) = pt_m - |q|r_{m+1}$. The same rule is obtained for fixed κ and variable $m < \kappa$.

Proof. The absolute estimation error for this scenario is equal to $|\hat{s}_m - s_m| = |z - z_{\kappa} - (1 + q)(z - s_m - z_{\kappa}) - s_m| = |q|(z - s_m - z_{\kappa})$. From $s_m + z_{\kappa} \leq t_{m+1}$ follows that

$$RE(s_m; z_{\kappa}, q) = \frac{|q|(z - s_m - z_{\kappa})}{s_m} \geq \frac{|q|(z - t_{m+1})}{s_m} \geq \frac{|q|(z - t_{m+1})}{t_m} = RE(t_m; z_{m+1}, q).$$

Thus, also in this scenario, the minimal absolute relative error of estimation is achieved for $s_m = t_m$ and $\kappa = m + 1$, so that the safety rule is satisfied when $RE(t_m; z_{m+1}, q) > p$, that is when

$$\frac{t_m}{z} \leq (1 - \gamma_{\kappa}) \frac{1}{1 + p'}, \quad (9)$$

where $p' = p/|q|$. The equivalence with the generalized p -rule follows comparing the above safety condition with that in (8); the subadditivity of the rule follows from its equivalence with a subadditive rule. The LSM is easily obtained from (9)

and it is $S_p(s_m; z_\kappa, q) = pt_m - |q|r_{m+1}$. If we had, instead, fixed the rank κ of the known contribution and left undetermined $m < \kappa$, we would have obtained the same LSM, expressed as $pt_{\kappa-1} - |q|r_\kappa$, for the same argument as above.

q.e.d.

The scenario for the generalized (p, q) -rule may be restricted to considering only $q > 0$, because the LSM is symmetric about $q = 0$. When $q = 0$ the intruder knows exactly the remainder $(z - t_{m+1})$ and $RE(t_m; z_{m+1}, q = 0) = 0$. When $|q| = 1$, the rule reduces to the generalized p -rule; in this case the intruder estimates the remainder $z - s_m - z_\kappa$ with zero or with twice its value. Note that q can never be less than -1 because the intruder wouldn't estimate the remainder with a negative number. Also note that the scenario of the generalized (p, q) -rule is slightly different from the original *prior-posterior ambiguity rules* scenario (Cox 1981). In fact, we simply set q to be the relative error for the remainder, as opposed to being the relative error for each of its contributions. This scenario yields the same safety rule but, unlike the other, does not require the prior knowledge of n by the intruder.

The predictive scenarios we considered so far are nested into each other, for increasing prior knowledge is allowed to the intruder. The connections among the corresponding rules are stated in the following corollary.

Corollary 2.4 *The following properties for the safety rules in Propositions (2.1), (2.2) and (2.3) hold good:*

- a) $P_p(s_m)$ is implied by $P_p(s_m; z_k)$, which is implied by $P_p(s_m; z_k, q)$ if $|q| \leq 1$; each of these rules applied to s_m implies itself for any sum of less contributions, s_j , $j < m$, with same p ;
- b) $P_p(s_m)$ implies $P_p(s_j; z_k)$ and $P_p(s_j; z_k, q)$ if $|q| \geq t_{m-1}/t_m$, for all indices $j < m$;
- c) *The parameters of the predictive scenarios are such that: i) $P(s_m; z_k)$ is never satisfied if $p \geq (n - m - 1)/m$; ii) $P(s_m; z_k, q)$ is never satisfied if $p/|q| \geq (n - m - 1)/m$ or $2m > n - 1$ or $t_m > r_{m+1}$.*

Proof. Since $0 \leq \gamma_k < 1$, the safety bound for the generalized p -rule in (8) is always lower than that for the Dominance rule in(2) for equal p and m . Thus $P_p(s_m)$ is implied by $P_p(s_m; z_k)$. When $|q| < 1$ the value of p' in (9) is higher than p , hence the safety bound is lower than that in (8) for (8), which is thus implied. When $|q| > 1$ the converse is true. From this implication and the above follows that also the Dominance rule is implied by the generalized (p, q) -rule for same m and p and $|q| < 1$.

For the nonnegativity of the contributions it is always true that $t_m \geq t_{m-1}$ and $r_m \leq r_{m-1}$. Thus,

$$RE(t_m) = \frac{r_m}{t_m} \leq \frac{r_{m-1}}{t_{m-1}} = RE(t_{m-1})$$

for all $2 < m \leq n$, hence $P_p(s_j)$ is implied by $P_p(s_m)$ for all $j \leq (m - 1)$. The same argument applies to the generalized p -rule and (p, q) -rule, considering that $RE(t_m; z_{m+1}) = r_{(m+1)}/t_m$ and $RE(t_m; z_{m+1}, q) = |q|r_{(m+1)}/t_m$. Property (a) is thus proved.

Clearly, $RE(t_m) = r_m/t_m \leq r_m/t_{m-1} = RE(t_{m-1}; z_m)$ for all $2 < m \leq n$, thus $P(s_m)$ implies $P(s_j; z_k)$ for all j in $[1, (m - 1)]$. $RE(t_{m-1}; z_m, q) = |q|r_m/t_{m-1}$ is not lower than $RE(t_m)$ if and only if $|q| \geq t_{m-1}/t_m < 1$, in which case $P_p(s_{m-1}; z_m, q)$ is implied by $P_p(s_m)$. The implication for all sums of less contributions, s_j , $j < m$, follows from property (a), proving property (b).

Since the contributions are nonincreasing, it holds that $t_m \geq mz_{m+1}$ and $r_{m+1} \leq (n - m - 1)z_{m+1}$. Hence $S(s_m; z_k) \geq [pm - (n - m - 1)]z_{m+1}$, from which follows that $P(s_m; z_k)$ is never satisfied if $pm > n - m - 1$, which can be rearranged to give (i). Whence, since LSM $S_p(s_m; z_k, q)$ is equivalent to $S_{p/|q|}(s_m; z_k)$, it follows that $S_p(s_m; z_k, q)$ is never satisfied if $p/|q| > (n - m - 1)/m$. Since it must also be $|q| < p$, from inequality (1) follows that $S(t_m; z_{m+1}, q) > p(2m - n + 1)\bar{z}$. Therefore, $P(s_m; z_k, q)$ cannot be satisfied if $2m > (n - 1)$, which

proves the second of (ii). Also because $|q| \leq p$, $S(s_m; z_k, q) = pt_m - |q|r_{m+1} > 0$ if $t_m > r_{m+1}$ and the last of (ii) is thus proved.

q.e.d.

This proposition links the safety rules so far discussed. Since the generalized (p, q) -rule is obtained allowing more prior information to the intruder than in the generalized p -rule, it seems logical to always choose $|q| < 1$, so to obtain a tighter bound. The safety of all sums of less contributions, hence also of the individual values, gives motivation for the use of safety rules that consider the estimation of sums of contributions. Property b) gives a justification for the general preference for the Dominance rule over the other M -rules. Properties in c) show that the values of p in the M -rules cannot be chosen completely independently of m and n .

3 Safety Rules for Estimation Under Uniform Assumptions

In the M -rules any sum of m contributions is always estimated with its highest possible value. Hence, the minimal relative error is achieved for $s_m = t_m$. In this section we assume that the intruder is explicitly interested in estimating t_m

and that uses his prior knowledge to bound its value within the bounds $t_m^- \leq t_m^+$, such that $t_m^- \leq t_m \leq t_m^+$. Furthermore, we assume that the ignorance about the remaining quantities is modelled assuming that t_m is uniformly distributed within $[t_m^-, t_m^+]$ and that the estimate is obtained minimizing the mean squared error of estimation. In this way, repeating the procedure many times the intruder expects low squared error. The estimate \hat{t} that minimizes the mean squared error $\int_{t_m^-}^{t_m^+} (\hat{t}_m - t_m)^2 / (t_m^+ - t_m^-) dt_m$ is

$$\hat{t}_m = \frac{t_m^- + t_m^+}{2}, \quad (10)$$

for a well known property of the mean. We consider the usual safety rules for this estimating procedure, which will be generically referred to as U -rules and denoted as $P(t_m; U, \dots)$.

The safety conditions set by the U -rule for an intruder with only basic knowledge are stated in the following proposition.

Proposition 3.1 *Let an intruder with only basic knowledge estimate t_m as in (10).*

Then $\hat{t}_m = z/2$ and safety rule $P_p(t_m; U) : |\hat{t}_m - t_m| < pt_m$ is not satisfied when

$$\frac{1}{2(1+p)} \leq \frac{t_m}{z} \leq \frac{1}{2(1-p)}. \quad (11)$$

The rule cannot be satisfied if $t_m < r_m$ and $p > (n - 2m)/2m$ or if $t_m > r_m$ and $p > 1/2$.

Proof. The intruder knowing only z can only bound t_m such that $0 \leq t_m \leq z$.

Substituting these values in (10) gives

$$\hat{t}_m = \frac{z}{2}. \quad (12)$$

If $t_m < r_m$ then $\hat{t}_m = z/2 > t_m$. In this case safety rule $P_p(t_m; U)$ is satisfied if $z/2 - t_m > pt_m$, from which is easy to derive the safety condition

$$\frac{t_m}{z} < \frac{1}{2(1+p)}.$$

From the lower bound $t_m \geq m/n$ in Equation (1), follows that this condition cannot be satisfied if $p > (n - 2m)/2m$.

If $t_m > r_m$, which is always true if $2m > n$, then $\hat{t}_m = z/2 < t_m$. In this case, safety rule $P_p(t_m; U)$ is satisfied when $t_m - z/2 = (t_m - r_m)/2 > pt_m$, from which it is easy to derive the condition

$$\frac{t_m}{z} > \frac{1}{2(1-p)}.$$

From the upper bound $t_m \leq z$ in Equation (1), follows that this condition cannot be satisfied if $p > 1/2$. This means that if $t_m \geq r_m$ then $RE(t_m; U) < 1/2$. Applying both inequalities yields the required interval.

q.e.d.

The following proposition sets the safety conditions against an intruder knowing the number of contributions in the cell.

Proposition 3.2 *Let an intruder knowing the number of contributions, n , estimate t_m by (10). Then the estimate is $\hat{t}_m = (z + m\bar{z})/2$ and safety rule $P_p(t_m; U, n)$ is not satisfied when*

$$\frac{n+m}{2n(1+p)} \leq \frac{t_m}{z} \leq \frac{n+m}{2n(1-p)}. \quad (13)$$

The rule can neither be satisfied if $(n-m)t_m < (n+m)r_m$ and $p > (n-m)/2m$ nor if $(n-m)t_m > (n+m)r_m$ and $p > (1-m/n)/2$.

Proof. From the knowledge of the cell mean \bar{z} , applying inequality (1), t_m can be bounded as $m\bar{z} \leq t_m \leq z$. Substituting these values in (10) gives

$$\hat{t}_m = \frac{z}{2} + \frac{m\bar{z}}{2} = \left(\frac{n+m}{2n} \right) z. \quad (14)$$

If $(n+m)z > 2nt_m$ then $\hat{t}_m > t_m$ and safety rule $RE(t_m; U, n) > p$ is satisfied if $(n+m)z > 2n(1+p)t_m$, that is when

$$\frac{t_m}{z} < \frac{m+n}{2n(1+p)}.$$

From the lower bound $t_m/z \geq m/n$ in Equation (1) follows that this condition cannot be satisfied if $p > (n-m)/2m$.

If $(n+m)z \leq 2nt_m$ then $\hat{t}_m < t_m$ and the safety rule is satisfied if $(n+m)z > 2n(1-p)t_m$, that is when

$$\frac{t_m}{z} > \frac{m+n}{2n(1-p)}.$$

From the upper bound $t_m \leq z$ in Equation (1) follows that this condition cannot be satisfied if $p > (1 - m/n)/2$, hence never if $p > 1/2$. Applying both inequalities yields the required interval.

q.e.d.

The scenario in which an intruder knows the contribution z_{m+1} is considered in the next proposition.

Proposition 3.3 *Let an intruder knowing $z_k = z_{m+1}$ estimate t_m by (10). Then the estimate is $\hat{t}_m = (z + (m - 1)z_k)/2$ and safety rule $P_p(t_m; U, z_{m+1})$ is not satisfied if*

$$\frac{1 + (m - 1)\gamma_k}{2(1 + p)} \leq \frac{t_m}{z} \leq \frac{1 + (m - 1)\gamma_k}{2(1 - p)}, \quad (15)$$

where $\gamma_k = z_k/z$. The safety rule cannot be satisfied if $(t_m - r_m)/2 \leq (m - 1)z_k$ and $p > (n - 2m + n(m - 1)\gamma_k)/2m$ or if $(t_m - r_m)/2 \leq (m - 1)z_k$ and $p > 1/2 - [\gamma_k(m - 1)/2]$.

Proof. From the knowledge of z_k , t_m can be bounded as $mz_k \leq t_m \leq z - z_k$, because $z_i \geq z_{m+1}$ for $i < m + 1$. Substituting these values in (10) gives

$$\hat{t}_m = \frac{z}{2} + \frac{(m - 1)}{2} z_k = \left(\frac{1 + (m - 1)\gamma_k}{2} \right) z \quad (16)$$

If $t_m/z \leq [1 - (m - 1)\gamma_k]/2$ then $\hat{t}_m > t_m$. In this case, safety rule $P(t_m; U, z_{m+1}) >$

p is satisfied when

$$\frac{t_m}{z} < \frac{1 + (m-1)\gamma_{\mathbb{K}}}{2(1+p)}.$$

From the lower bound $t_m \geq m/n$ in Equation (1) follows that this condition cannot be satisfied if $p > [n - 2m + n(m-1)\gamma_{\mathbb{K}}]/2m$.

When $\hat{t}_m < t_m$ the safety rule is satisfied if

$$\frac{t_m}{z} > \frac{1 + (m-1)\gamma_{\mathbb{K}}}{2(1-p)}.$$

From the upper bound $t_m \leq z$ in Equation (1) follows that this condition cannot be satisfied if $p > 1/2 - [(m-1)\gamma_{\mathbb{K}}/2]$, hence never if $p > 1/2$. Applying both inequalities yields the required interval.

q.e.d.

The union of the previous two scenarios, that is an intruder knowing both n and z_{m+1} , is considered in next proposition.

Proposition 3.4 *Let an intruder knowing $z_{\mathbb{K}} = z_{m+1}$ and n estimate t_m by (10).*

Then if $z_{\mathbb{K}} > \bar{z}$ the estimate \hat{t}_m and the safety conditions are the same as in Proposition (3.3). If $z_{\mathbb{K}} \leq \bar{z}$ then the estimate is $\hat{t}_m = z - (n - m + 1)z_{\mathbb{K}}/2$ and safety rule $P_p(t_m; U, \mathbb{K}, n)$ is not satisfied if

$$\frac{2 - (n - m + 1)\gamma_{\mathbb{K}}}{2(1+p)} \leq \frac{t_m}{z} \leq \frac{2 - (n - m + 1)\gamma_{\mathbb{K}}}{2(1-p)}. \quad (17)$$

The safety rule cannot be satisfied when $2r_m > (n - m + 1)\gamma_K$ and $p > [2(n - m) - \gamma_K(n - m + 1)]/2m$ or when $2r_m < (n - m + 1)\gamma_K$ and $p > \gamma_K(n - m + 1)/2n$ or $m > (n - 1)/2$.

Proof. If $z_K > \bar{z}$ then t_m can be bounded as $mz_K \leq t_m \leq z - z_K$ and the safety conditions must be the same as for $P_p(t_m; U, K)$. In this case the knowledge of n is redundant. If $z_K \leq \bar{z}$ then t_m can be bounded as $z - (n - m)z_K \leq t_m \leq z - z_K$, where the lower bound is derived from $t_m \geq z - z_K - \max\{r_K\}$. In this case both pieces of information contribute to the lower bound.

When $z_K \leq \bar{z}$ the estimate of t_m is

$$\hat{t}_m = z - \frac{(n - m + 1)}{2} z_K = \left[\frac{2 - (n - m + 1)\gamma_K}{2} \right] z. \quad (18)$$

If $2r_m > (n - m + 1)\gamma_K$ then $\hat{t}_m > t_m$ and safety rule $P_p(t_m; U, K, n)$ is satisfied if

$$\frac{t_m}{z} < \frac{2 - (n - m + 1)\gamma_K}{2(1 + p)}.$$

From the lower bound $t_m \geq m/n$ in Equation (1) follows that this condition cannot be satisfied if $p > [2(n - m) - n(n - m + 1)\gamma_K]/2m$. If $2r_m > (n - m + 1)\gamma_K$ then $\hat{t}_m < t_m$ and safety rule $RE(t_m; U, n) > p$ is satisfied if

$$\frac{t_m}{z} > \frac{2 - (n - m + 1)\gamma_K}{2(1 - p)}.$$

From the upper bound $t_m \leq z$ in Equation (1) follows that this condition cannot be satisfied if $p > \gamma_k(n - m + 1)/2$, hence never if $p > 1/2$ or if $m > (n - 1)/2$. Applying both inequalities yields the required interval.

q.e.d.

The above propositions show that for the U -rules cells are at risk if the ratio t_m/z is inside an interval. As it can be easily seen, this implies that the rules are not subadditive.

We propose to adopt the U -rules restricted to being satisfied only if t_m/z is lower than the lower bound. In so doing we obtain subadditive rules of the Dominance kind. Observing that the lower bound for acceptably large t_m/z ratios increases rapidly with p and m to the maximal value 1, such restriction will often be implied by the conditions themselves. Furthermore, this restriction agrees with the idea of protecting large contributions because they are more easily identifiable. Thus, for the restricted U -rules, that will be denoted as $P^R(t_m; U, \dots)$, cells are safe only if t_m is overestimated and the absolute relative error of estimation is less than the threshold p .

The safety conditions satisfying the M -rules and the restricted U -rules are shown in Table 2, together with the restrictions on the values of the parameters.

rule	bound	restrictions
$P(s_m)$	$\frac{t_m}{z} < \frac{1}{(1+p)}$	$p < \frac{(n-m-1)}{m}$
$P(s_m; z_{\kappa})$	$\frac{t_m}{z} \leq \frac{(1-\gamma_{\kappa})}{1+p}$	$p < \frac{(n-m-1)}{m}$
$P(s_m; z_{\kappa}, q)$	$\frac{t_m}{z} \leq \frac{(1-\gamma_{\kappa})}{1+\frac{p}{ q }}$	$\frac{p}{ q } < \frac{(n-m-1)}{m}, 2m > n - 1$
$PR(t_m; U)$	$\frac{t_m}{z} < \frac{1}{2(1+p)}$	$p < \frac{(n-2m)}{2m}$
$PR(t_m; U, n)$	$\frac{t_m}{z} < \frac{n+m}{2n(1+p)}$	$p < \frac{(n-m)}{2n}$
$PR(t_m; U, \kappa)$	$\frac{t_m}{z} < \frac{1+(m-1)\gamma_{\kappa}}{2(1+p)}$	$p < \frac{(n-2m)+(m-1)\gamma_{\kappa}}{2m}$
$PR(t_m; U, \kappa, n)^*$	$\frac{t_m}{z} < \frac{2-(n-m+1)\gamma_{\kappa}}{2(1+p)}$	$p < \frac{2(n-m)-n(n-m+1)\gamma_{\kappa}}{2m}$

Table 2: Comparison of different safety rules, safety bounds and restrictions on the parameters. The asterisk indicates that the bound applies if $\gamma_{\kappa} < 1/n$, otherwise the bound above applies.

4 CONCLUSIONS

The Dominance and the generalized p -rule assume that the intruder always estimates s_m with its maximum possible value, setting the remainder equal to zero. In passing, note that this is equivalent to setting the absolute relative error of estimation of the remainder constantly equal to one. The generalized prior-posterior rule allows for a positive estimate of the remainder and the value of q can be used to lower the bound on t_m/z set by the other rules. However, the difficulty of

quantifying this prior knowledge for a generic intruder and the restrictions on the parameters given in Corollary (2.4), may restrict considerably the application of the generalized prior-posterior rule. Furthermore, the restrictions indicate that in the M -rules the value of p and q cannot be chosen completely arbitrarily but are bounded by the values of m and n . The restricted U -rules give upper bounds for the ratio t_m/z that depend on the prior information the intruder is allowed to have. For an intruder with only basic knowledge the bound on t_m/z set by $P(t_m; U)$ rule will be half of that for the corresponding M -rule. The bounds set by the restricted U -rules will not always be lower than those set by the corresponding M -rules, a comparison is given in Table 3. The safety conditions set by the U -rules are not always stricter the more prior knowledge is allowed to the intruder. In particular, it turns out that in many cases the lowest bound is that of the $P^R(t_m)$ rule. This fact is explained by the estimation criterium adopted by the rules. It can be observed that the more prior knowledge is allowed the wider are rejection intervals for the complete U -rules.

Finally, we would like to stress the fact that we have only considered rules protecting from the disclosure of sums of observations, without considering the further use that the intruder could make of these estimates, such as estimating the single contributions. In fact, safety conditions could be obtained assuming one or

$>$	\cdot	z_K	$U,$	U, n	U, z_K
z_K	N				
U	N	$\gamma_K > \frac{1}{2}$			
U, n	N	$\gamma_K > \frac{(n-m)}{n}$	A		
U, z_K	$\gamma_K > \frac{1}{(m-1)}$	$\gamma_K > \frac{1}{(m+1)}$	A	$\gamma_K > \frac{m}{n(m-1)}$	
U, n, z_K^*	N	N	N	$\gamma_K < \frac{(n-m)}{n(n-m+1)}$	A

Table 3: Conditions for which the M -rules and restricted U -rules with parameters as in the row labels yield higher upper bounds for the ratio t_m/z than the rules with parameters as in the column labels. N stands for *never*, A for *always* and the asterisk denotes that $\gamma_K < 1/n$.

another estimating procedure for the single contributions. It would be difficult to decide upon which estimating procedure to assume and, in many cases, one would obtain very strict safety conditions. The U -rules, applied with an appropriate value of p , should provide a reasonable protection against disclosure risk.

References

- [1] Cox, L. H. (1981), "Linear Sensitivity Measures in Statistical Disclosure Control," *Journal of Statistical Planning and Inference*, 5, 153-164.

- [2] de Wolf, P. P. (Sept. 2001), *Notes for the TES course on SDC* held at CBS, Voorburg (NL).
- [3] Jewett, R. (1993), Disclosure analysis for the 1992 Economic Census. *Unpublished manuscript*. Economic Programming Division, US Bureau of the Census, Washington, DC. [142, 183]
- [4] Willenborg, L. and de Waal, T. (2000) *Elements of Statistical Disclosure Control*, New York: Springer-Verlag.