

Quadratic interior-point methods in statistical disclosure control

Jordi Castro*

Department of Statistics and Operations Research, Universitat Politècnica de Catalunya, Pau Gargallo 5, 08028 Barcelona, Spain (e-mail: jordi.castro@upc.edu)

Abstract. The safe dissemination of statistical tabular data is one of the main concerns of National Statistical Institutes (NSIs). Although each cell of the tables is made up of the aggregated information of several individuals, the statistical confidentiality can be violated. NSIs must guarantee that no individual information can be derived from the released tables. One widely used type of methods to reduce the disclosure risk is based on the perturbation of the cell values. We consider a new controlled perturbation method which, given a set of tables to be protected, finds the closest safe ones – thus reducing the information loss while preserving confidentiality. This approach means solving a quadratic optimization problem with a much larger number of variables than constraints. Real instances can provide problems with millions of variables. We show that interior-point methods are an effective choice for that model, and, also, that specialized algorithms which exploit the problem structure can be faster than state-of-the art general solvers. Computational results are presented for instances of up to 1000000 variables.

Keywords: Interior-point methods, Quadratic Programming, Large-scale programming, Statistical confidentiality, Controlled perturbation methods

AMS subject classification 90C06, 90C20, 90C51, 90C90

1 Introduction

In the current Information Society, National Statistical Institutes (NSIs) play a fundamental role, routinely releasing large volumes of data for their further exploitation. The released data are usually classified as aggregated or disaggregated.

* Partially supported by the EU IST-2000-25069 CASC project and by the Spanish MCyT project TIC2003-00997.

		z_1	z_2	
⋮
51–55	...	38000 €	40000 €	...
56–60	...	39000 €	42000 €	...
⋮

a

		z_1	z_2	
⋮
51–55	...	20	1 or 2	...
56–60	...	30	35	...
⋮

b

Fig. 1a,b. Example of disclosure in tabular data. **a** Average salary per age and ZIP code. **b** Number of individuals per age and ZIP code

Disaggregated data (a.k.a. microdata or microfiles) correspond to files of records, each record providing the values for a set of variables of an individual. Aggregated data (a.k.a. tabular data) are obtained from microdata crossing two or more variables, which results in sets of tables with a likely large number of cells. In this paper we focus on tabular data protection (see, e.g., [22, 25] for a comprehensive introduction to this field). Although each cell shows aggregated data for several individuals, there is a risk of disclosing individual data. This is clearly shown in the example of Fig. 1. Table (a) in that Figure gives the average salary for age interval and ZIP code, while table (b) shows the number of individuals for the same variables. If there was only one individual in ZIP code z_2 and age interval 51–55, then any external attacker would know the salary of this single person is 40000 €. For two individuals, any of them could deduce the salary of the other, becoming an internal attacker. Usually, cells showing information about few individuals are considered sensitive, although other rules can be used in practice. See, for instance, [9] for a general description and [11, 24] for a recent discussion about sensitivity rules.

In the example of Fig. 1 we should protect table (a), a single two-dimensional table. This can be considered the simplest case. However, in practice we must deal with more complex situations. The full range can be classified as:

- Single two-dimensional, three-dimensional, and, in general, multidimensional tables. Those tables can be individually protected.
- Hierarchical tables, i.e., sets of tables with variables that have a hierarchical relation (e.g., ZIP code and city). In that case, the total or marginal cells for some table can be the internal cells for the others. They have to be protected together, to avoid the disclosure of sensitive data.
- Linked tables. It is a generalization of the previous situation, where several tables are made from the same microdata, thus sharing information or cells, either hierarchical or not. Again, they have to be protected together.

Eventually, we could consider the whole set of linked tables that can be produced from some microfiles (e.g., a population census). Clearly, the number of cells involved in that case could be of several millions. All the above situations can both refer to frequency tables (i.e., cell values are integer and are usually associated to

the number of individuals in that cell) or magnitude tables (i.e., cell values are real, and, for instance, can show the mean for some other variable of all the individuals in that cell).

The methods for protection of tabular data can be classified as perturbative (they change the cell values) or nonperturbative (no change is performed). The most widely used nonperturbative method is *cell suppression*, where some *secondary* cells are removed to avoid the disclosure of some sensitive *primary* cells (which are removed as well). That results in a difficult combinatorial optimization problem, which finds the pattern of secondary suppressions that makes the table safe with a minimum number of cells or information loss. Some heuristics [3, 8, 18] and exact methods [12, 13] have been suggested for the cell suppression problem.

Among the perturbative methods, one of the techniques that received more attention due to its simplicity was *rounding*. This method rounds cell values to a multiple of a fixed integer rounding base. *Controlled rounding* is a variant where the additivity of the table is preserved (i.e., rounded marginal values are the sum of the corresponding slice of rounded cells) [2]. However, preserving (integer) additivity is not easy in a multidimensional table or a set of linked tables. Moreover, in practice it can be necessary to maintain the original marginal values, instead of rounding them.

To avoid the above shortcomings of rounding, we suggest a new perturbation method that finds the minimum- L_2 -distance (or closest) tables to those to be protected, preserving marginal values, as well as any set of additional linear constraints. (If Δ is the vector of deviations between the cells of the original and perturbed tables, we define the L_2 distance as $\|\Delta\|_2 = \sqrt{\sum_i \Delta_i^2}$.) Finding the minimum- L_2 -distance tables means we try to minimize the information loss when delivering the perturbed table. In this work we focus on tables of magnitudes (i.e., cell values $\in \mathbb{R}$). For tables of frequencies (integer cell values) this procedure could still be applied followed by some heuristic post-processing. The main drawback of this approach is the solution of a very large quadratic optimization problem. We will show that interior-point methods can solve this problem very efficiently. Recently, and independently of this work, a similar idea was suggested in [10], using the L_1 distance for the perturbed table. (Again, if Δ is the vector of deviations between the cells of the original and perturbed tables, the L_1 distance is $\|\Delta\|_1 = \sum_i |\Delta_i|$.) In practice, however, the optimization problem obtained with the L_2 distance can be solved more efficiently, as shown in [6].

The structure of the document is as follows. In Sect. 2 we introduce the *minimum- L_2 -distance* perturbation method. In Sect. 3 we show some computational experience using a general state-of-the-art interior-point solver. From those results we conclude that it is worth using a specialized interior-point algorithm that exploits the problem structure. This is the subject of Sect. 4. Finally in Sect. 5 we present some computational results in the solution of three-dimensional tables of up to 1000000 cells using a specialized interior-point algorithm.

2 The *minimum- L_2 -distance* perturbation method

Any problem instance, either with one table or a number of (linked or hierarchical) tables, can be represented by the following elements:

- A set of cells $a_i, i = 1, \dots, n$, that satisfy some linear relations $Aa = b$ (a being the vector of a_i 's). For instance, for a two-dimensional table of $r + 1$ rows and $c + 1$ columns (last row and column correspond to marginal values) we have

$$\sum_{i=1}^r a_{ij} = a_{(r+1)j} \quad j = 1 \dots c \quad (1)$$

$$\sum_{j=1}^c a_{ij} = a_{i(c+1)} \quad i = 1 \dots r. \quad (2)$$

For a three-dimensional table with $l + 1$ levels (levels correspond to third dimension, last level is marginal), the relations are

$$\sum_{i=1}^r a_{ijk} = a_{(r+1)jk} \quad j = 1 \dots c, \quad k = 1 \dots l \quad (3)$$

$$\sum_{j=1}^c a_{ijk} = a_{i(c+1)k} \quad i = 1 \dots r, \quad k = 1 \dots l \quad (4)$$

$$\sum_{k=1}^l a_{ijk} = a_{ij(l+1)} \quad i = 1 \dots r, \quad j = 1 \dots c. \quad (5)$$

In practice most tables have positive cell values, and bounds $a \geq 0$ have to be considered.

- A lower and upper bound for each cell $i = 1, \dots, n$, respectively \underline{a}_i and \bar{a}_i , which are considered to be known by any attacker. If no previous knowledge is assumed for cell i , we would simply set $\underline{a}_i = 0$ ($\underline{a}_i = -\infty$ if bounds $a \geq 0$ were not assumed) and $\bar{a}_i = +\infty$.
- A set $\mathcal{P} = \{i_1, i_2, \dots, i_p\}$ of indices of confidential cells.
- A lower and upper protection level for each confidential cell $i \in \mathcal{P}$, respectively lpl_i and upl_i , such that the released value should be greater or equal than $a_i + upl_i$ or less or equal than $a_i - lpl_i$. Modelling these “or” constraints would need an extra binary variable for each confidential cell, resulting in a large combinatorial optimization problem which would constrain the effectiveness of the approach to small and medium sized problems. To avoid it, we will assume the user (e.g., the NSI) fixes in advance the sense of the protection for each confidential cell. In practice tabular data protection is the last stage of the “data cycle”, and, in an attempt to meet publication deadlines, NSIs require fast solutions to large and complex tables. That justifies the above simplifying

assumption. If the particular choice of protection senses results in an infeasible problem, we can solve an alternative one considering marginal cells as nonfixed, but with a large penalization for possible perturbations.

The *minimum- L_2 -distance* method finds (using the L_2 distance) the closest set of *perturbed* values x_i to a_i , $i = 1, \dots, n$, such that the tables relations, lower and upper cell bounds, and sensitive cells protection levels are satisfied. The optimization problem can be written as

$$\begin{aligned} \min_x \quad & \|x - a\|_2^2 \\ \text{subject to} \quad & Ax = b \\ & \underline{a}_i \leq x_i \leq \bar{a}_i \quad i = 1 \dots n \\ & x_i \leq a_i - lpl_i \text{ or } x_i \geq a_i + upl_i \quad i \in \mathcal{P}, \end{aligned} \tag{6}$$

x being the vector of the cell values x_i , $i = 1 \dots n$ of the perturbed table.

Problem (6) can be applied to any kind of table, since it does not constrain the structure of the cell relations $Ax = b$. Any other set of linear conditions can also be added to (6). For instance, we could impose that, in the perturbed table, values x_i related to some non-confidential cells must be close enough to the original values a_i , e.g., $(1 - \alpha)a_i \leq x_i \leq (1 + \beta)a_i$ for some small α and β . For cells corresponding to national or regional totals, or for cells with a zero value, $\alpha = \beta = 0$ can be a good choice (i.e., we do not perturb the original cell value). This is the usual practice in those situations. We may also want to affect the distance by any positive semidefinite diagonal metric matrix $W = \text{diag}(w_1, \dots, w_n)$. For instance, we can set $w_i = 1$ or $w_i = 1/a_i$ to deal with, respectively, absolute or relative perturbations in the objective function. In the computational results of Sects. 3 and 5 we used $w_i = 1$. The more general model can be written as:

$$\min_x (x - a)^T W(x - a) \tag{7}$$

$$\text{subject to } Ax = b \tag{8}$$

$$c \leq Tx \leq d \tag{9}$$

$$l \leq x \leq u, \tag{10}$$

(9) being any set of additional linear equality or inequality constraints defined by T , c and d , if necessary; and l and u in (10) the final lower and upper bounds of the perturbed cell values.

3 Computational experience using a general interior-point solver

(7–10) is a large (possibly with millions of variables), separable convex quadratic problem. It is known that the computational cost of quadratic separable problems is the same of linear ones, if solved through an interior-point algorithm [27]. Therefore, in principle, that seems to be the best choice.

Table 1. Dimensions of some small instances and results with four solvers

r	c	l	$ \mathcal{P} $	n	m	Cplex 8.0			Minos
						Barrier	Dual	Primal	
25	25	–	10	625	50	0.01	0.02	0.08	1.63
50	50	–	20	2500	100	0.03	0.09	1.02	11.14
100	100	–	20	10000	200	0.21	0.66	22.50	>700*
10	10	10	20	1000	300	0.05	0.19	0.66	3.15
15	15	15	20	3375	675	0.29	2.92	16.3	164.71
25	25	25	20	15625	1875	4	160	868	>4600*

* Execution was stopped.

To confirm the good behaviour of interior-point methods in this problem, we performed some preliminary limited computational experience with some small two-dimensional and three-dimensional tables. These instances, and those used below in this Section and in Sect. 5, were obtained with two different generators that have been used in the literature. The first generator follows the description of [18]. Cell values are randomly obtained from an integer uniform distribution $[1 \dots 1000]$ with probability 0.8 and are 0 with probability 0.2. The second one is similar to the first generator of [12]. Cell values are randomly obtained from integer uniform distributions $[1 \dots 4]$ for confidential cells and $\{0\} \cup [5 \dots 500]$ for the remaining entries. Cells to be protected are randomly chosen from the internal (i.e., nonmarginal) cells in both generators. We extended the original generators for three-dimensional tables using the same distributions. They can be obtained from the author on request. Note that the performance of the method does not depend on any particular distribution. Therefore, for real tables, the method is expected to behave as for random ones.

We generated three small two-dimensional and three-dimensional tables, whose sizes are shown in Table 1. Columns r , c , l and $|\mathcal{P}|$ give, respectively, the number of rows, columns, levels and number of sensitive cells for each instance. The first three rows correspond to two-dimensional tables, and thus column l is empty. Columns n and m show, respectively, the number of variables (number of cells) and constraints of the resulting quadratic optimization problem. Note that, for two-dimensional problems and from (1,2), $n = rc$ and $m = r + c$, while for three-dimensional instances and from (3–5), $n = rcl$ and $m = (r + c)l + rc$. As for the rest of instances of the paper, we considered that the lower and upper bounds known for each cell are $0.9a_i$ and $1.1a_i$, respectively, and the value of the sensitive cells was set to $x_i = 1.1a_i$ in the perturbed table. We solved each problem with four solvers: the interior-point barrier algorithm of Cplex 8.0 [17], the dual and primal simplex algorithms for quadratic problems of Cplex 8.0 (see, e.g. [26] for a description of this simplex variant), and Minos 5.5 [20, 21] – one of the most efficient reduced-gradient type solvers. The CPU time in seconds required by each solver is shown, respectively, in columns “Barrier”, “Dual”, “Primal” and “Minos”. The executions were carried on a Compaq Evo N610c notebook, with a Pentium Mobile 4 at 1.8

Table 2. Dimensions of some large instances and results with the barrier Cplex 8.0 solver

r	c	l	$ \mathcal{P} $	n	m	CPU
1000	500	–	10000	500000	1500	47.1
1000	750	–	10000	750000	1750	72.9
1000	1000	–	10000	1000000	2000	136.0
100	100	25	10000	250000	15000	198.5
100	100	50	10000	500000	20000	896.7
100	100	100	10000	1000000	30000	Not enough memory

GHz and 512 MB of RAM. The problem (7–10) was implemented using the AMPL modelling language [14]. Looking at Table 1, it is clear that the interior-point solver is the best option, mainly when the size of the instances increases.

Unfortunately, real problems are much larger than those used for Table 1. And for large instances even general interior-point solvers can be computationally expensive. This is clearly shown in Table 2, which reports the CPU time in seconds (column “CPU”) required by the interior-point barrier algorithm of Cplex 8.0 for three two-dimensional and three-dimensional large instances. The meaning of the other columns is the same as in Table 1. The largest three-dimensional problem, involving one million of cells, could not be solved with the available memory. The two-dimensional problem with the same number of variables did not have such limitation. In fact, two-dimensional problems could be solved more efficiently. This is mainly due to the lesser number of constraints they involve. It can be concluded that general interior-point solvers are too expensive, both in memory and execution time, when applied to large instances. Next section shows that using a specialized interior-point solver which exploits the problem structure is instrumental.

4 Exploiting the problem structure

4.1 Two-dimensional tables

The linear relations (1,2) of a $(r + 1) \times (c + 1)$ table can be modeled as the network of Fig. 2. Arcs are associated to cells and nodes to equations. Injections correspond to marginal row $r + 1$ and marginal column $c + 1$. Thus (7,8,10) – assuming no extra constraints (9) are considered – is a large convex separable quadratic minimum-cost network flows problem. Some effective specialized interior-point methods have been devised for linear network problems. They solve the normal equations at each iteration of the interior-point method (discussed below) by a preconditioned conjugate gradient. The most effective preconditioner is the “maximum spanning tree preconditioner” [23], and its variants [15, 19]. As far as we know, it has not been applied to quadratic network flows problems. In that case, the preconditioner would probably not preserve its good properties, mainly when we are close to the optimal solution: in a quadratic problem the optimizer does no longer have to be located in a vertex, and thus the maximum spanning tree does not correspond to any

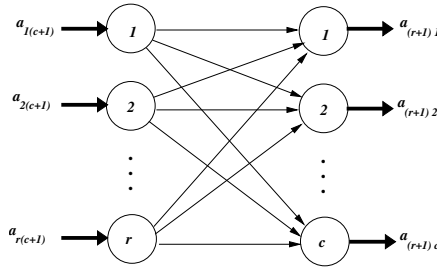


Fig. 2. Network representation of a $(r + 1) \times (c + 1)$ table

optimal basis. However, in the first and intermediate iterations of the interior-point method, it could be more efficient than using a Cholesky factorization, as a general solver does. Therefore, in principle, there is room for improving the execution times shown in Table 2. This is still part of the further work to be done.

4.2 Three-dimensional tables

As we observed in Table 2, three-dimensional tables are computationally more challenging than two-dimensional ones with a similar number of cells. As shown in [5], the relations (3–5) of a three-dimensional table are equivalent to those of a multicommodity network flows problem with equality mutual capacity constraints (see, e.g., Chapter 17 of [1] for an introduction to multicommodity problems). Indeed, the structure of the constraints (8) defined by (3–5) is

$$\begin{bmatrix} N & & & \\ & N & & \\ & & \ddots & \\ & & & N \\ I & I & \dots & I \end{bmatrix} \begin{bmatrix} x^1 \\ x^2 \\ \vdots \\ x^l \end{bmatrix} = \begin{bmatrix} b^1 \\ b^2 \\ \vdots \\ b^l \\ a^{l+1} \end{bmatrix}, \quad (11)$$

$N \in \mathbb{R}^{(r+c) \times (rc)}$ being the network linear relations of the two-dimensional table associated to each level (depicted in Fig. 2), $x^k \in \mathbb{R}^{rc}$, $k = 1 \dots l$, the cells (flows) of level k , $b^k \in \mathbb{R}^{r+c}$, $k = 1 \dots l$, the row and column marginals (injections) of level k , $a^{l+1} \in \mathbb{R}^{rc}$ the level marginal values (mutual arc capacities), and $I \in \mathbb{R}^{(rc) \times (rc)}$ the identity matrix. Constraints involving the network matrix N correspond to (3–4), while (5) are the linking constraints. If marginal values (i.e., the right-hand-side term of (11)) are considered also variables (for instance, with a large penalization for deviations from the original values) the structure of (8) would be an extension

of the standard multicommodity network flows problem:

$$\begin{bmatrix} N & & & & & & & & -I \\ & \ddots & & & & & & & \\ & & N & & & & & & -I \\ I & \dots & I & & & & & & -I \end{bmatrix} \begin{bmatrix} x^1 \\ \vdots \\ x^l \\ b^1 \\ \vdots \\ b^l \\ x^{l+1} \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}. \tag{12}$$

The optimization problem (7,8,10) to be solved for a three-dimensional table is thus a convex separable quadratic multicommodity network flows problem, which can be written in standard form as

$$\begin{aligned} \min_x & \sum_{k=1}^l \left((c^k)^T x^k + (x^k)^T Q^k x^k \right) \\ \text{subject to} & Ax = b \\ & 0 \leq x^k \leq u^k, k = 1 \dots l, \end{aligned} \tag{13}$$

where $Ax = b$ is either (11) or (12) (in the latter case, marginal cells are also variables, and then the two l 's in (13) should be replaced by $l + 1$, and the variables b^k of (12) should be considered included in the associated x^k of (13)). c^k and Q^k are respectively the vector of linear costs and the diagonal matrix of quadratic costs for each x^k .

We extended the specialized interior-point method of [4] (initially developed for linear multicommodity problems) for the solution of (13). As far as we know, this is the only specialized method for general quadratic multicommodity flows, unlike the linear case, where there are several available algorithms (see, e.g, Chapter 17 of [1]). In fact, the method developed in [4] is not restricted to multicommodity problems. It can also be used for a wide range of block diagonal structured problems. In particular, it solves the extended formulation (12). Next, we justify the above assertions. Most details about the original specialized interior-point algorithm for linear multicommodity problems are omitted; they can be found in [4].

The most expensive computation of a primal-dual interior-point method is the solution of the normal equations

$$(A\Theta A^T)\Delta y = \bar{b} \tag{14}$$

at each iteration. A is the constraints matrix of (13), Δy is the direction for the dual variables, and Θ is a positive definite matrix, which can be partitioned in l (or $l + 1$ if formulation (12) is considered) blocks, one for each commodity. The expression of each block is

$$\Theta^k = ((\Theta_{lin}^k)^{-1} + Q^k)^{-1}, \tag{15}$$

Θ_{lin}^k being the positive definite diagonal matrix of the linear problem, and Q^k the matrix of quadratic costs of commodity k in (13). In our context, Q^k is diagonal; therefore Θ^k is also diagonal and easily computable.

Exploiting the structures of A in (11) and Θ , we obtain

$$A\Theta A^T = \begin{bmatrix} B & C \\ C^T & D \end{bmatrix} = \left[\begin{array}{ccc|c} N\Theta^1 N^T & \dots & 0 & N\Theta^1 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & N\Theta^l N^T & N\Theta^l \\ \hline \Theta^1 N^T & \dots & \Theta^l N^T & \sum_{k=1}^l \Theta^k \end{array} \right]. \quad (16)$$

If, instead, the extended formulation (12) is considered we have

$$A\Theta A^T = \begin{bmatrix} B & C \\ C^T & D \end{bmatrix} = \left[\begin{array}{ccc|c} N\Theta_I^1 N^T + \Theta_M^1 & \dots & 0 & N\Theta_I^1 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & N\Theta_I^l N^T + \Theta_M^l & N\Theta_I^l \\ \hline \Theta_I^1 N^T & \dots & \Theta_I^l N^T & \sum_{k=1}^{l+1} \Theta_I^k \end{array} \right], \quad (17)$$

where the Θ_I^k and Θ_M^k matrices correspond, respectively, to the internal and marginal cells of the table.

Either using (16) or (17), and appropriately partitioning Δy and \bar{b} , we can write (14) as

$$\begin{bmatrix} B & C \\ C^T & D \end{bmatrix} \begin{bmatrix} \Delta y_1 \\ \Delta y_2 \end{bmatrix} = \begin{bmatrix} \bar{b}_1 \\ \bar{b}_2 \end{bmatrix}. \quad (18)$$

By block multiplication, we can reduce (18) to

$$(D - C^T B^{-1} C) \Delta y_2 = (\bar{b}_2 - C^T B^{-1} \bar{b}_1) \quad (19)$$

$$B \Delta y_1 = (\bar{b}_1 - C \Delta y_2). \quad (20)$$

Following [4], (20) is solved by performing a Cholesky factorization of each diagonal block of B , while the system with matrix

$$H = D - C^T B^{-1} C, \quad (21)$$

the Schur complement of (18), is solved by a preconditioned conjugate gradient method. Note that H is symmetric and positive definite, since both B and $A\Theta A^T$ are symmetric and positive definite. A good preconditioner is instrumental for the performance of the method. The preconditioner of [4] was developed for linear multicommodity problems. However, it can be applied to any problem where the following result holds:

Proposition 1 *Let $H = D - C^T B^{-1} C$, the Schur complement of (18), be a positive definite matrix. If D and $D + C^T B^{-1} C$ are respectively nonsingular and positive definite then the inverse of H can be computed as*

$$H^{-1} = \left(\sum_{i=0}^{\infty} (D^{-1} (C^T B^{-1} C))^i \right) D^{-1}. \tag{22}$$

Proof. See pp. 860–861 of [4].

Both (16) and (17) satisfy the premises of Proposition 1. Their D^{-1} matrices can be easily computed, since

$$D = D_{lin} + \sum_{k=1}^l (Q^k)^{-1},$$

D_{lin} being the diagonal positive definite matrix for the linear problem without the quadratic term. The preconditioner for those problems, as for the linear ones, is thus obtained by truncating the infinite power series (22) at some term h . In practice $h = 0$ or $h = 1$ are good choices. Note that for $h = 0$ the preconditioner is equal to $\hat{H}^{-1} = D^{-1}$, which, since Q^k are diagonal, is also diagonal. This is instrumental in the overall performance of the algorithm. All the computational results of this work were obtained with $h = 0$. The quality of the preconditioner depends of the spectral radius of $D^{-1} (C^T B^{-1} C)$, and, in practice, it was observed to work better for quadratic than for linear problems [7]. A thorough study of the behaviour of the spectral radius for quadratic problems is part of the further work to be done.

5 Computational experience using a specialized interior-point algorithm

We implemented the specialized interior-point method described in the last Section in a code named QIPM. It is a quadratic extension of the package IPM, originally developed in [4] for linear multicommodity problems. IPM can be found in <http://www-eio.upc.es/~jcastro>. QIPM can be obtained from the author on request. We compared the performance of QIPM against Cplex 8.0 using a set of 162 three-dimensional instances, obtained with the two generators described in Sect. 3. The 81 instances for each generator were produced considering all the combinations for $r, c, l \in \{25, 50, 100\}$ and $|\mathcal{P}| \in \{1000, 5000, 10000\}$ (r, c, l and $|\mathcal{P}|$ with the same meaning as in Tables 1 and 2). The lower and upper bounds, and the values of the sensitive cells were computed as in Sect. 3.

Table 3 shows the results obtained with Cplex 8.0 and QIPM for some of the largest instances. Column “g” gives the generator used. Columns r, c, l and $|\mathcal{P}|$ have the same meaning as in previous tables. Columns K_n and K_m show, respectively, the number of thousands of variables and constraints of the resulting

Table 3. Dimensions and results with Cplex 8.0 and QIPM for some of the largest instances

g	r	c	l	\mathcal{P}	Kn	Km	Cplex 8.0		QIPM		
							it.	CPU	it.	$\overline{\text{CG}}$	CPU
1	100	50	100	1000	500	20	8	893	9	1.6	6
2	100	50	100	1000	500	20	7	924	7	3.0	6
1	100	50	100	5000	500	20	8	1284	9	1.6	7
2	100	50	100	5000	500	20	7	909	7	2.6	5
1	100	50	100	10000	500	20	8	885	9	1.6	6
2	100	50	100	10000	500	20	7	913	7	2.6	5
1	100	100	25	1000	250	15	8	185	9	2.7	4
2	100	100	25	1000	250	15	8	206	7	3.9	3
1	100	100	25	5000	250	15	9	205	10	2.4	4
2	100	100	25	5000	250	15	8	205	7	4.1	4
1	100	100	25	10000	250	15	8	179	11	2.4	5
2	100	100	25	10000	250	15	8	199	7	4.1	3
1	100	100	50	1000	500	20	8	875	9	1.8	8
2	100	100	50	1000	500	20	7	899	7	3.4	7
1	100	100	50	5000	500	20	8	792	9	1.8	8
2	100	100	50	5000	500	20	7	910	7	3.0	7
1	100	100	50	10000	500	20	8	866	9	1.9	8
2	100	100	50	10000	500	20	7	897	7	3.7	7
1	100	100	100	1000	1000	30		*	8	1.6	14
2	100	100	100	1000	1000	30		*	7	3.4	13
1	100	100	100	5000	1000	30		*	9	1.6	16
2	100	100	100	5000	1000	30		*	7	3.0	13
1	100	100	100	10000	1000	30		*	9	1.6	16
2	100	100	100	10000	1000	30		*	7	2.6	13

* Not enough memory

quadratic optimization problem. Columns “it.” and “CPU” show the number of interior-point iterations and the execution times (in seconds) for both codes. Column “ $\overline{\text{CG}}$ ” shows the average number of conjugate gradient iterations per interior-point iteration performed by QIPM. The execution environment was described in Sect. 3. The largest instances could not be solved with Cplex 8.0 due to a lack of memory. Clearly, the specialized interior-point method is about two orders of magnitude faster than the general solver, although the number of iterations is similar. Also, the execution times of QIPM smoothly increase with the size of the problem, and almost proportionally to the number of variables. It is worth noting that QIPM uses standard Cholesky factorization routines, whereas Cplex 8.0 includes a highly tuned and optimized factorization code. Then, in principle, there is still room for improving the QIPM performance.

The results obtained for all the 162 instances are summarized in Figs. 3–5. The holes observed in those figures correspond to infeasible problems due to the tight bounds considered. Figure 3 shows the ratio between the CPU times of Cplex 8.0 and QIPM with respect to the number of variables of the problem. As shown by the Figure, the ratios increase with the problem size, which makes QIPM a very efficient tool for large instances. For instances with more than 250000 variables,

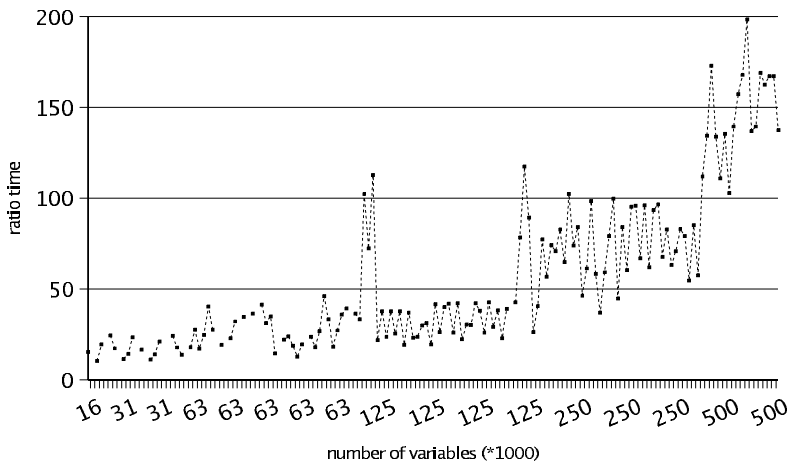


Fig. 3. Ratio of the Cplex 8.0 and QIPM CPU times

QIPM was at least 50 times faster than Cplex 8.0. For the largest instances it was about 175 times faster.

Figure 4 shows the objective function relative error – in absolute value – of the QIPM solution, $\left| (f_{\text{Cplex}}^* - f_{\text{QIPM}}^*) / (1 + f_{\text{Cplex}}^*) \right|$, assuming Cplex 8.0 provides the exact one. The default optimality tolerance used in [4] for linear problems was 10^{-6} . Although for quadratic problems this tolerance could likely be decreased, since the preconditioner works better than in the linear case, we preserved that default value. This explains why most of the relative errors are around 10^{-6} in Fig. 4. Although there is no a clear tendency, the largest relative errors were obtained in some of the largest instances.

Finally, Fig. 5 shows the difference between the number of interior-point iterations of QIPM and Cplex 8.0, with respect to the problem size. All the differences are in the range $-1, \dots, 3$. QIPM, at most, saved one iteration, while Cplex 8.0 saved two in several instances. The values are quite uniformly distributed and independent of the number of variables of the instances.

6 Conclusions

The results obtained with the minimum- L_2 -distance method show that, first, it can be a promising tool for the protection of statistical tabular data; and second, that specialized interior-point methods can solve large instances in few seconds, more efficiently than general solvers. However, this work can be improved in several ways. First, the minimum- L_2 -distance method can be adjusted to fit the real needs of NSIs, which would likely mean the inclusion of additional constraints or terms in the objective function. Second, we have to consider the general situation, that

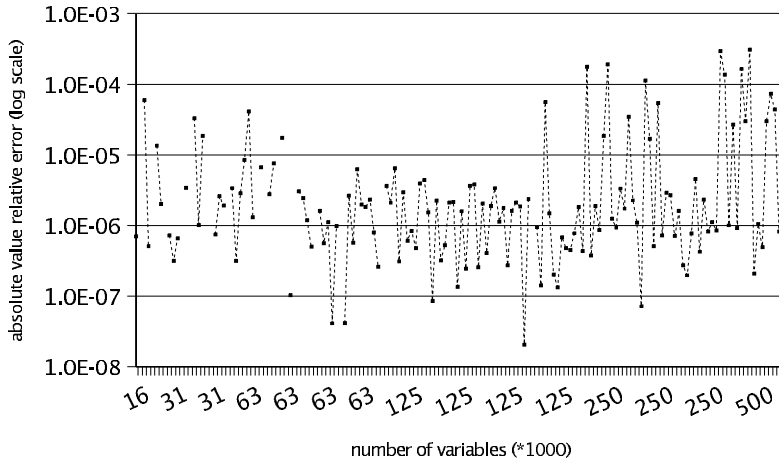


Fig. 4. Objective function relative error of the QIPM solution, in absolute value

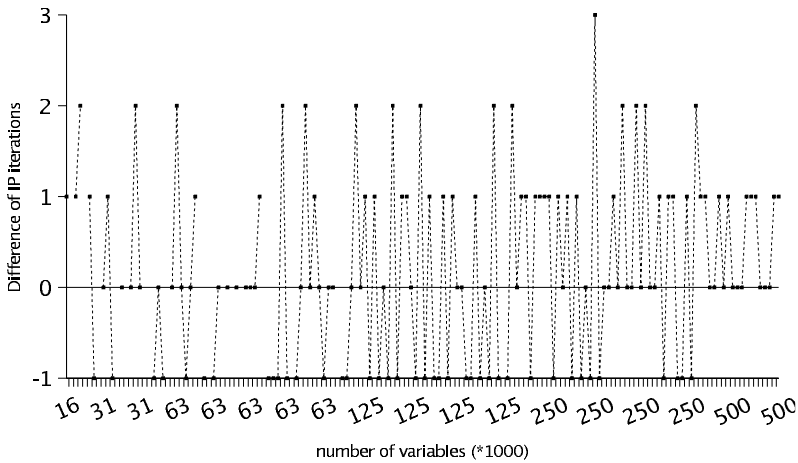


Fig. 5. Difference between the number of interior-point iterations performed by QIPM and Cplex 8.0

involves multidimensional, hierarchical and linked tables. That means the development of specialized interior-point solvers for the resulting structured problems [16]. Extending the method for frequency tables, through some type of heuristic post-process, is another of the future tasks to be done.

References

- [1] Ahuja, R.K, Magnanti, T.L, Orlin, J.B. (1993) *Network Flows*. Prentice Hall, Upper Saddle River
- [2] Bacharach, M. (1966) Matrix rounding problems. *Management Science* **9**: 732–742
- [3] Carvalho, F.D., Dellaert, N.P., Osório, M.D. (1994) Statistical disclosure in two-dimensional tables: general tables. *J. of the American Statistical Association* **89**: 1547–1557

- [4] Castro, J. (2000) A specialized interior-point algorithm for multicommodity network flows. *SIAM J. on Optimization* **10**(3): 852–877
- [5] Castro, J. (2002) Network flows heuristics for complementary cell suppression: an empirical evaluation and extensions. In: Domingo-Ferrer (ed.) *Inference Control in Statistical Databases*. Lecture Notes in Computer Science **2316**, pp. 59–73
- [6] Castro, J. (in press) Minimum-distance controlled perturbation methods for tabular data protection. European Journal of Operational Research. Available from <http://www-eio.upc.es/~jcastro> as research report DR2003-14, Dept. of Statistics and Operations Research, Universitat Politècnica de Catalunya
- [7] Castro, J. (2003) Solving quadratic multicommodity problems through an interior-point algorithm. In: Sachs, E.W., Tichatschke, R. (eds.) *System Modelling and Optimization XX*, pp. 199–212. Kluwer, Boston
- [8] Cox, L.H. (1995) Network models for complementary cell suppression. *J. of the American Statistical Association* **90**: 1453–1462
- [9] Cox, L.H. (2002) Disclosure risk for tabular economic data. In: Doyle, P., Theeuwes, J., Lane, J., Zayatz, L. (eds.) *Confidentiality, Disclosure and Data Access: Theory and Practical Applications for Statistical Agencies*, pp. 167–183. North Holland, Amsterdam
- [10] Dandekar, R.A., Cox, L.H. (2002) Synthetic tabular data: an alternative to complementary cell suppression. Manuscript, Energy Information Administration, U.S. Dept. of Energy. Available from the first author on request (Ramesh.Dandekar@eia.doe.gov)
- [11] Domingo-Ferrer, J., Torra, V. (2002) A critique of the sensitivity rules usually employed for statistical table protection. *Int. J. of Uncertainty Fuzziness and Knowledge-based Systems* **10**(5): 545–556
- [12] Fischetti, M., Salazar, J.J. (1999) Models and algorithms for the 2-dimensional cell suppression problem in statistical disclosure control. *Mathematical Programming* **84**: 283–312
- [13] Fischetti, M., Salazar, J.J. (2000) Models and algorithms for optimizing cell suppression in tabular data with linear constraints. *J. of the American Statistical Association* **95**: 916–928
- [14] Fourer, R., Gay, D.M., Kernighan, B.W. (1993) *AMPL: A Modeling Language for Mathematical Programming*. Boyd & Fraser, Danvers
- [15] Frangioni, A., Gentile, C. (2001) New preconditioners for KKT systems of network flow problems. *SIAM J. on Optimization* **14**(3): 894–913
- [16] Gondzio, J., Sarkissian, R. (2003) Parallel interior point solver for structured linear programs. *Mathematical Programming* **96**: 561–584
- [17] ILOG CPLEX (2002) *ILOG CPLEX 8.0 Reference Manual Library*. ILOG, Gentilly
- [18] Kelly, J.P., Golden, B.L., Assad, A.A. (1992) Cell Suppression: disclosure protection for sensitive tabular data. *Networks* **22**: 28–55
- [19] Mehrotra, S., Wang, J. (1995) Conjugate gradient based implementation of interior point methods for network flow problems. In: Adams, L., Nazareth, J. (eds.) *AMS Summer Conference Proceedings*, pp. 124–142. SIAM, Philadelphia
- [20] Murtagh, B.A., Saunders, M.A. (1978) Large-scale linearly constrained optimization. *Mathematical Programming* **14**: 41–72
- [21] Murtagh, B.A., Saunders, M.A. (1983) MINOS 5.0. User’s guide, Dept. of Operations Research, Stanford University
- [22] Oganian, A. (2002) *Security and Information Loss in Statistical Database Protection*. PhD thesis, Dept. of Applied Mathematics 4, Universitat Politècnica de Catalunya
- [23] Resende, M.G.C., Veiga, G. (1993) An implementation of the dual affine scaling algorithm for minimum-cost flow on bipartite uncapacitated networks. *SIAM J. on Optimization* **3**: 516–537
- [24] Robertson, D.A., Ethier, R. (2002) Cell suppression: experience and theory. In: Domingo-Ferrer, J. (ed.) *Inference Control in Statistical Databases*. Lecture Notes in Computer Science **2316**, pp. 8–20
- [25] Willenborg, L., de Waal, T. (2000) *Elements of Statistical Disclosure Control*. Lecture Notes in Statistics, vol. 155 Springer, New York
- [26] Wolfe, P. (1959) The simplex method for quadratic programming. *Econometrica* **27**: 382–398
- [27] Wright, S.J. (1997) *Primal-Dual Interior-Point Methods*. SIAM, Philadelphia