

DU SUPPRESSIF AU PERTURBATIF, COMMENT PARAMÉTRER LES MÉTHODES DE BRUITAGE ? UNE PREMIÈRE DÉMARCHÉ UTILISANT LA MÉTHODE DES CLÉS ALÉATOIRES

Julien JAMME (*)

(*) Insee, Direction de la méthodologie et de la coordination statistique et internationale

julien.jamme@insee.fr

Mots-clés (6 maximum) : Confidentialité, Perturbation, Probabilités, Risque-Utilité.

Domaines : Confidentialité.

Résumé

À l'Insee, des méthodes suppressives (regroupement ou blanchiment de cases) sont majoritairement utilisées pour limiter les risques de divulgation d'informations confidentielles lors de la publication de tableaux statistiques. Aujourd'hui, ces méthodes atteignent certaines de leurs limites, en particulier lorsque les produits de diffusion deviennent complexes [1]. Certaines méthodes perturbatrices, en particulier la méthode des clés aléatoires, permettent de combler ces lacunes. Or, les règles du secret statistique sont particulièrement adhérentes aux méthodes suppressives: une case est supprimée si elle ne respecte pas le critère choisi. Comment, en utilisant des méthodes perturbatrices, s'assurer, par exemple, de respecter la règle de fréquence, celle qui fixe un nombre de contributeurs minimal à une case pour pouvoir la diffuser ? Le bruit injecté par ce type de méthode est en général piloté par un ensemble de paramètres. Comment assurer un choix optimal de ces paramètres, c'est-à-dire un choix qui offre les garanties nécessaires de protection tout en assurant un qualité suffisante de l'information transmise au public ?

Nous nous attacherons ici à étudier le cas de l'application de la méthode des clés aléatoires ([2], [3]) à des comptages issus de sources exhaustives. Cette méthode repose sur deux éléments fondamentaux: des clés individuelles aléatoires qui introduisent l'aléa dans le processus tout en assurant une certaine cohérence des requêtes entre elles et une matrice de probabilités de transition qui définit la façon dont les comptages sont perturbés. Pour définir ces distributions de probabilité, il est nécessaire de fixer - a minima - deux paramètres: la déviation maximale des comptages et la variance de la distribution de probabilité. Le choix de ces paramètres est réalisé pour assurer le meilleur compromis entre protection et utilité, les deux termes nécessitant chacun une métrique permettant de les objectiver.

Les probabilités de transition qui définissent le mécanisme de perturbation peuvent servir à estimer les deux plateaux de la balance. Elles fournissent directement une mesure d'utilité et, en les inversant avec la formule de Bayes, elles nous permettent d'estimer la capacité qu'aurait un attaquant d'inférer une valeur sensible à partir d'un comptage diffusé.

Notre démarche consiste ainsi à comparer l'effet sur le niveau de protection et sur la qualité de l'information transmise d'un ensemble de paramètres appliqués sur un ou plusieurs tableaux représentatifs de la source étudiée. Les résultats montrent en particulier la nécessité d'utiliser un

paramètre supplémentaire permettant d'interdire l'apparition de petits comptages et ainsi limiter la capacité d'inférence des attaquants. Le choix de ce seuil d'interdiction est d'autant plus important qu'il a un impact non négligeable sur le niveau de variance injecté dans les données et donc sur la qualité des informations diffusées auprès du public.

Abstract

At INSEE, suppressive methods (such as cell grouping or cell suppression) are predominantly used to limit the risks of disclosing confidential information when publishing statistical tables. Today, these methods are reaching some of their limitations, particularly when dissemination products become complex [1]. Some perturbative methods, particularly the cell key method, help address these shortcomings. However, the rules of statistical confidentiality are particularly tailored to suppressive methods: a cell is suppressed if it does not meet the chosen criterion. How, when using perturbative methods, can one ensure, for example, compliance with the frequency rule, which sets a minimum number of contributors to a cell for it to be disseminated? The noise injected by this type of method is generally controlled by a set of parameters. How can an optimal choice of these parameters be ensured, that is, a choice that provides the necessary protection guarantees while ensuring sufficient quality of the information transmitted to the public?

This study focuses on the application of the cell key method ([2], [3]) to counts derived from exhaustive sources. This method relies on two fundamental elements: individual random keys that introduce randomness into the process while ensuring a certain consistency across queries, and a transition probability matrix that defines how the counts are perturbed. To define these probability distributions, it is necessary to set, at least, two parameters: the maximum deviation of the counts and the variance of the probability distribution. The choice of these parameters is made to achieve the best trade-off between protection and utility, both of which require metrics to make them objective.

The transition probabilities that define the perturbation mechanism can be used to estimate the two sides of the balance. They directly provide a measure of utility and, by inverting them using Bayes' formula, allow estimation of an attacker's ability to infer a sensitive value from a disseminated count. Our approach thus consists of comparing the effects on the level of protection and on the quality of the transmitted information of a set of parameters applied to one or more representative tables from the studied source. The results particularly highlight the need to use an additional parameter to prohibit the appearance of small counts and thereby limit attackers' inference capabilities. The choice of this prohibition threshold is all the more important as it has a non-negligible impact on the level of variance injected into the data and thus on the quality of the information disseminated to the public.

Introduction

Les enjeux de protection de diffusions complexes

La protection des données tabulées contre les risques de divulgation est une tâche indispensable à assurer avant d'envisager toute diffusion. La complexité de cette tâche est très largement dépendante de la complexité de la diffusion elle-même, en particulier du nombre de tableaux à diffuser, des règles du secret statistique à respecter et de la pose du secret secondaire. En effet, cette étape de protection repose en général sur l'application de méthodes suppressives consistant à « blanchir » les cases ne respectant pas les règles du secret statistique (secret primaire) et celles qui permettraient de retrouver les premières (secret secondaire). Mais, il arrive que certaines diffusions rendent l'application des méthodes suppressives sinon inefficaces du moins inefficaces dans la protection des données.

En effet, lorsque les données sont ventilées sur de nombreux zonages géographiques non parfaitement emboîtés (par exemple une diffusion communale et une diffusion au carreau), on accroît les risques de divulgation par différenciation géographique [4]. S'il existe un outil efficace pour détecter ces risques entre deux zonages donnés avec le package R appelé `diffman` [5], le traitement de ces risques par les méthodes suppressives nécessite l'implémentation d'algorithme ad hoc et peut être pénible à gérer de manière efficiente.

Il arrive que des données issues d'une même source ne soient pas diffusées en même temps. On parlera de **diffusion synchrone** si le programme de diffusion est fixé à l'avance et parfaitement connu sans qu'aucun écart à ce programme ne soit possible. Au contraire, on parlera de **diffusion asynchrone** dès lors qu'une partie de la diffusion n'est pas connue ni fixée à l'avance. Les méthodes suppressives sont très bien adaptées pour gérer les premières. Il suffit pour cela de prendre en compte l'ensemble du programme de diffusion et de le traiter en une fois. En revanche, elles le sont moins pour les diffusions asynchrones. Dans ce cas, il faudrait pour publier un tableau au temps $t + 1$ prendre en compte les masques de secret posés sur l'ensemble des tableaux déjà publiés. Dans les faits, cette pratique est complexe à coordonner entre toutes les parties prenantes de la diffusion au sein d'un institut.

La pose du secret secondaire sur un ensemble de tableaux nécessite une analyse préalable des liens entre les tableaux. Or cette analyse, décrite dans [6], peut être complexe si les tableaux sont nombreux et font intervenir des variables hiérarchiques, et en particulier en présence de hiérarchies non-emboîtées [7]. Cette complexité est appréhendable par des experts rompus à l'exercice mais le niveau d'expertise requis pour traiter ce genre de demandes est plus difficile à exiger d'agents ne traitant qu'occasionnellement des problèmes de secret statistique.

Enfin, les outils classiques de pose du secret secondaire, en particulier τ -ARGUS [8], peuvent atteindre certaines de leurs limites dès lors qu'il faut traiter des tableaux de très grandes tailles, par exemple des tableaux à l'échelle communale. Basée sur des programmes d'optimisation sous contraintes, les calculs peuvent s'avérer très longs - jusqu'à plusieurs jours - voire ne pas converger dans un temps raisonnable - disons moins d'une semaine. L'utilisation d'une approche de pose du secret secondaire par réduction gaussienne [9] peut permettre parfois de réduire ces temps de calculs, mais ces traitements auront toujours une certaine lourdeur.

Les méthodes perturbatrices

Quand une demande cumule plusieurs de ces limites, il peut être tentant de se tourner vers des méthodes permettant de gérer sa protection de manière plus efficiente, si ce n'est plus efficace. Certaines méthodes perturbatrices relativement simples, telles que les méthodes basées sur des arrondis ou bien la méthode des clés aléatoires, consistent à brouter tous les comptages publiés. Ceci permet de réduire automatiquement les risques de divulgation par différenciation (en particulier géographique) sans détection et traitement *ad hoc* nécessaire. Leur relative simplicité rend ces méthodes faciles à implémenter et à appliquer, améliorant l'efficacité de la protection des données tabulées à l'échelle d'un institut lorsqu'une diffusion est complexe.

Opter pour de telles méthodes n'est pas sans conséquence sur les données : la perte d'additivité qu'elles peuvent engendrer dans les tableaux, la perception de la perturbation par le public et la manipulation des données sont à mettre en balance avec leurs avantages indéniables :

- Réduction du risque de divulgation par différenciation
- Implémentation peu coûteuse
- Centralisation de l'effort de protection
- Gestion des publications asynchrones

Ainsi, quand la production de données tabulées est industrielle, c'est-à-dire mêlant une diffusion de nombreux indicateurs sur de nombreux tableaux ayant des liens complexes, l'usage de méthodes perturbatrices comme les arrondis ou la méthode des clés aléatoires est à envisager.

Un exemple concret

La publication des données sur les Quartiers prioritaires de la politique de la ville (QPV) rassemble plusieurs des critères de complexité énoncés ci-dessus. Les quartiers sont des zones construites à partir de considération principalement statistiques et indépendamment des limites administratives, en particulier communales. Or, la diffusion des données sur les QPV s'accompagne toujours d'une diffusion concomitante sur des zonages administratifs tels que les communes, les EPCI, les départements et les régions, ainsi que sur un autre zonage statistique infracommunal: les IRIS. Autant d'occasions de risque de divulgation par différenciation. En outre, à l'occasion de la publication des données sur le nouveau zonage des QPV 2024, il a été souhaité, pour des raisons de continuité avec l'offre existante, de poursuivre la diffusion sur le précédent zonage des QPV 2015, augmentant encore les occasions de risque de divulgation en différenciant les deux zonages de QPV.

Les traitements originels de ce type de risque bien qu'assez lourds étaient relativement bien gérés par les équipes chargées de la protection des données, mais la diffusion concomitante des deux zonages a rendu cette tâche encore plus complexe et l'application de méthodes suppressives aurait diminué encore l'utilité des données diffusées en supprimant de nombreuses cases supplémentaires, d'autant que les cases concernées par les risques de différenciation ne concernent pas nécessairement que des petits comptages [1].

Ainsi, la gestion automatique des risques de divulgation par différenciation et la promesse d'une meilleure préservation de l'information a incité l'Insee, en accord avec les producteurs de données (CAF, CNAM, France Travail) et les principaux utilisateurs comme l'Agence nationale de la cohésion des territoires (ANCT), à se tourner vers la méthode des clés aléatoires.

Problématique

L'utilisation d'une méthode perturbatrice nécessite de relever plusieurs défis.

Le premier consiste à bien calibrer le bruit injecté afin d'obtenir le meilleur équilibre entre maîtrise des risques de divulgation et bonne conservation de l'utilité des données. Pour maîtriser cet équilibre, il faut pouvoir s'équiper de métriques de risque et d'utilité adaptées aux objectifs de protection et de diffusion.

Le second défi est intimement lié au premier: définir une métrique de risque de divulgation cohérente avec les règles du secret statistique en vigueur. Pour des tableaux de comptages, la principale règle du secret statistique consiste à définir comme sensible tout comptage inférieur à un certain seuil. L'application d'une méthode suppressive consiste à supprimer ces comptages et quelques autres en plus pour réduire les risques de réidentification des individus. Or, avec une méthode perturbatrice, aucun comptage n'est supprimé. Comment s'assurer alors qu'on a suffisamment le risque de réidentification, c'est-à-dire suffisamment réduit la capacité d'un attaquant à deviner quels comptages sont sensibles? C'est l'un des enjeux de cet article que de présenter la métrique sur laquelle repose l'arbitrage pour s'assurer d'une protection suffisante des données.

Le troisième défi concerne la perspective opposée: s'assurer que les données qui seront diffusées resteront suffisamment utiles pour les usagers. Si aucune information n'est plus supprimée, une injection de bruit dans les comptages mal maîtrisée peut rendre les données inutilisables. De plus, même maîtrisée, il est important de pouvoir conserver la confiance des usagers dans leur utilisation des données.

Plan

Les bonnes propriétés de la méthode des clés aléatoires ont conduit l'Insee à la privilégier à des méthodes d'arrondis par exemple. Les principes et le fonctionnement de cette méthode seront présentés dans un premier temps. Ensuite, les objectifs et étapes de la phase de calibration de la méthode seront explicités. En particulier, la métrique de risque de divulgation utilisée pour la calibration

de la méthode sera détaillée, ainsi que les différentes façons de mesurer la perte d'information engendrée par le mécanisme. Enfin, un exemple illustrera l'application de ces métriques pour calibrer la méthode dans un cas précis.

1 La méthode des clés aléatoires

La **méthode des clés aléatoires (cell key method, CKM)** est la technique de protection perturbative adoptée dans ce travail. Cette méthode a été initialement développée pour l'*Australian Bureau of Statistics* (ABS) en 2005 [2]. Cette implémentation originelle est décrite en détail dans [10]. Des développements ultérieurs fournis par des statisticiens du *Statistisches Bundesamt allemand* (*Destatis*) sont présentés dans [11] et [12]. Une présentation synthétique de la méthode est disponible dans une fiche méthodologique dédiée sur le site internet de l'Insee [3]. La CKM est l'une des principales méthodes choisies par les instituts nationaux de statistique européens pour protéger la diffusion des données du Censur européen de 2021 [13].

La méthode des clés aléatoires consiste à **dévier tous les comptages** d'un tableau tout en assurant, par l'utilisation de clés aléatoires individuelles, qu'**une même case est toujours perturbée de la même manière**. Chaque case est perturbée de manière indépendante, y compris les marges des tableaux. Ceci permet d'éviter les reconstructions des données à partir des liens naturels des cases dans un tableau, mais engendre nécessairement une perte d'additivité au sein des tableaux.

Les déviations

Considérons une base de données individuelles \mathcal{B} composée de N lignes (1 ligne = 1 individu), cette base ne contenant que des variables catégorielles. Pour une case \mathcal{C} issue du croisement des modalités de quelques unes des variables de \mathcal{B} , X dénotera la variable aléatoire correspondant au dénombrement original de \mathcal{C} et X' la variable aléatoire correspondant au dénombrement perturbé de \mathcal{C} . On notera Z la différence $X' - X$. Ainsi, X est le nombre d'individus contribuant à la case \mathcal{C} et $X' = X + Z$. La déviation Z appliquée doit répondre à plusieurs conditions:

- Elle est entière: $Z \in \mathbb{Z}$;
- Elle est bornée par un paramètre noté D : $|Z| \leq D$
- La déviation d'un comptage ne peut pas le rendre négatif: $X' \geq 0$.

De ces propriétés, on déduit que

$$\begin{cases} Z \in \llbracket -D; D \rrbracket \text{ si } X \geq D, \\ Z \in \llbracket 0; D \rrbracket \text{ sinon.} \end{cases}$$

Ce qui correspond aux propriétés suivantes pour le comptage après perturbation:

$$\begin{cases} X' \in \llbracket X - D; X + D \rrbracket \text{ si } X \geq D, \\ X' \in \llbracket 0; X + D \rrbracket \text{ sinon.} \end{cases}$$

La méthode se distingue des autres méthodes de perturbation par l'utilisation (i) d'un jeu de clés aléatoires et (ii) d'un mécanisme de perturbation piloté par une matrice de probabilités de transition.

Le rôle des clés aléatoires

La méthode tient son nom du fait qu'elle associe un nombre aléatoire à chaque comptage produit à partir de la source individuelle \mathcal{B} . Ce nombre aléatoire, appelé clé de la case ou *cell key*, est construit à partir de nombres aléatoires (ou clés individuelles) tirés, en début de processus, pour chaque individu de \mathcal{B} .

Notons RK une variable aléatoire distribuée selon une loi uniforme sur $[0; 1]$: $RK \sim \mathcal{U}([0; 1])$. Les clés individuelles $\{rk_1, \dots, rk_N\}$ sont des réalisations indépendantes de la variable aléatoire RK associées à chacun des individus de \mathcal{B} , comme dans le tableau 1. Ces clés individuelles, une fois posées, ne doivent plus être modifiées afin de conserver la cohérence globale des futurs résultats.

TAB. 1 – Ajout des clés individuelles dans la source de données

id	Commune de résidence	Âge	Clé (rk)
1	Amiens	25	0,9177275
2	Paris	20	0,8850062
3	Marseille	45	0,6266963
4	Amiens	45	0,1117820
5	Marseille	20	0,6496634
6	Marseille	20	0,2813433

Considérons la case \mathcal{C} d'un tableau agrégé construit à partir de \mathcal{B} . La clé aléatoire de la case est construite à partir des clés individuelles des individus qui y contribuent. Supposons, sans perte de généralité, que seuls les p premiers individus de \mathcal{B} contribuent à \mathcal{C} . Le comptage de la case vaut donc p et la clé associée, notons-la $ck(\mathcal{C})$, est définie comme la partie décimale de la somme des clés des p individus de \mathcal{C} : $ck(\mathcal{C}) = S_n - \lfloor S_n \rfloor$ où $S_n = \sum_{n=1}^p rk_n$. Cette définition permet en outre de s'assurer que la clé de chaque case est une réalisation d'une variable aléatoire uniforme sur $[0; 1]$ également.

Le tableau 2 fournit les comptages par commune à partir de la base individuelle présentée dans le tableau 1. En plus des comptages, le tableau fournit la composition individuelle des cases (colonne 2), la somme des clés individuelles, S_n (colonne 4) et la clé de la case $ck(\mathcal{C})$ (colonne 5).

TAB. 2 – Tableau agrégé par Commune de résidence et calcul de la clé associé à la case

Com.	ids	X	\sum clés	Clé de la case
Amiens	{1,4}	2	1,0295095	0,0295095
Marseille	{3,5,6}	3	1,5577030	0,5577030
Paris	{2}	1	0,8850062	0,8850062
Total	{1,2,3,4,5,6}	6	3,4722187	0,4722187

La perturbation qui sera appliqué aux comptages dépendra de leur clé associée, permettant ainsi d'assurer une cohérence globale et systématique des perturbations.

Le mécanisme de perturbation

L'association d'une perturbation à la clé d'un comptage dépend d'une matrice de probabilités de transition qui définit la distribution de probabilités du bruit injecté et d'une table de perturbation qui permet d'associer bruit et *cell key*.

Dans la suite, on notera i une valeur quelconque prise par X , z une valeur de Z et j une valeur prise par X' . On appellera probabilité de transition de i vers j la probabilité $p_{ij} = \mathbb{P}(X' = j | X = i) = \mathbb{P}(Z = j - i | X = i)$. Définir ces probabilités pour tous les cas possibles détermine entièrement la distribution de la variable aléatoire $X' | X = i$.

La même distribution est appliquée pour toutes les valeurs $i \geq D$, ce qui permet de limiter la taille de la matrice. Pour ces valeurs, le comptage original est dévié entre $\llbracket -D; D \rrbracket$. En revanche, une

distribution spécifique est nécessaire pour chacun des comptages $i \in]0; D[$, car le support de ces distributions ne sont pas identiques, la méthode des clés aléatoires - telle qu'elle est utilisée ici - n'autorisant pas la production de comptages négatifs. Pour $D = 2$, la matrice de transition M aura ainsi la forme suivante:

$$M = \begin{bmatrix} p_{00} & p_{01} & p_{02} & p_{03} & p_{04} \\ p_{10} & p_{11} & p_{12} & p_{13} & p_{14} \\ p_{20} & p_{21} & p_{22} & p_{23} & p_{24} \end{bmatrix}$$

La dernière ligne décrit la distribution de $X'|X = i$, pour $i \geq 2$. Sous plusieurs contraintes supplémentaires imposées (absence de biais, variance fixe, distribution symétrique quand c'est possible), un programme de maximisation de l'entropie permet de construire cette matrice de transition¹. Pour $D = 2$ et $V = 2$, on obtient la matrice suivante²:

$$M = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0.36649 & 0.36649 & 0.16757 & 0.09946 & 0 \\ 0.2 & 0.2 & 0.2 & 0.2 & 0.2 \end{bmatrix}$$

Cette matrice de transition est alors transformée en une table de perturbation permettant de lire la déviation à appliquer à un comptage selon la valeur prise par la clé. La table de perturbation associée à la matrice précédente est présentée dans le tableau 3. Les deux dernières colonnes correspondent aux intervalles cumulés des probabilités de transition. Ainsi, pour le comptage i d'une case C donnée, la procédure consiste à repérer l'intervalle $[p_{inf}; p_{sup}]$ dans laquelle se trouve la clé $ck(C)$. Ainsi, en prenant le comptage de la commune d'Amiens dans le tableau 2, qui vaut 2, et dont la clé associée vaut 0,0295095, la perturbation qui devra être appliquée sera de -2 (ligne surlignée en vert dans le tableau 3). Cette opération est effectuée pour tous les comptages construits à partir de la même source.

TAB. 3 – La table de perturbation pour $D = 2$ et $V = 2$

$X = i$	$X' = j$	p_{ij}	$Z = z$	p_{inf}	p_{sup}
0	0	1.0	0	0	1
1	0	0.36649	-1	0	0.36649
1	1	0.36649	0	0.36649	0.73297
1	2	0.16757	1	0.73297	0.90054
1	3	0.09946	2	0.90054	1
$i \geq 2$	0	0.2	-2	0	0.2
$i \geq 2$	1	0.2	-1	0.2	0.4
$i \geq 2$	2	0.2	0	0.4	0.6
$i \geq 2$	3	0.2	1	0.6	0.8
$i \geq 2$	4	0.2	2	0.8	1

Les caractéristiques de la méthode

Armé de ce mécanisme de perturbation, la méthode des clés aléatoires dispose donc des caractéristiques suivantes:

- Elle perturbe tous les comptages
 - Cela permet de gérer les risque de divulgation, en particulier par réidentification directe

1. C'est ce que fait le package R `ptable` [14]

2. On remarquera qu'ici les paramètres conduisent à produire une distribution uniforme pour $X'|X = 2$.

- et par différenciation géographique;
- En revanche, cela engendre une perte d'additivité dans les tableaux.
- La perturbation est sans biais
 - En moyenne, les déviations des comptages se compensent les unes des autres.
 - Cette propriété est probablement un désavantage sur les distributions appliquées aux petits comptages (voir par la suite).
- La perturbation d'une case est déterminée par les individus qui la composent:
 - Cela permet d'assurer, sans effort, une cohérence de toutes les diffusions ultérieures;
- Les zéros ne sont pas perturbés:
 - Aucun comptage non pertinent ne sera créé (par exemple des enfants à la retraite);
 - Pour les autres zéros, qu'on pourrait décrire comme conjoncturels, leur perturbation nécessiterait soit de permettre des comptages négatifs, soit d'accepter l'introduction de biais dans les données diffusées.

La méthode des clés aléatoires semble avoir beaucoup d'avantages dès lors qu'on cherche à produire des tableaux agrégés en quantité industrielle. Et si son implémentation ni sa prise en main par les producteurs ne présentent de difficultés importantes, reste à savoir comment il est possible de choisir un jeu de paramètres assurant un niveau de protection des données suffisant tout en préservant le mieux possible la qualité originale des agrégats.

2 Calibrer la méthode pour réaliser le meilleur arbitrage risque-utilité possible

Dans cette seconde partie, nous présentons les principales étapes et les ingrédients nécessaires pour mener cette phase de calibration de la méthode d'une manière la plus objective possible. La notion d'objectivité signifie ici de bien identifier les termes de notre arbitrage risque-utilité, les enjeux en termes de protection et de se doter de métriques qui puissent faciliter la prise de décision finale.

2.1 Les enjeux en termes de protection de la confidentialité

Pour des données de comptage, les risques de divulgation sont les suivants:

- le risque de réidentification: ce risque est géré principalement en considérant comme sensible tous les petits comptages, les règles du secret statistique fixant en général un seuil s , dit seuil de fréquence, en-dessous duquel une case ne peut être diffusée telle quelle.
- le risque de divulgation directe ou par inférence d'attributs sensibles. Ce risque nécessite de désigner les attributs sensibles diffusés et de repérer les situations de divulgation.
- le risque de divulgation par différenciation.

Notons que la perturbation systématique de tous les comptages réduit nécessairement le risque de divulgation par différenciation. La question sera alors de savoir si cette réduction du risque est suffisante ? Il sera tout à fait possible d'évaluer cette réduction.

Si des attributs sensibles sont présents dans les données agrégées diffusées, la perturbation systématique va là encore réduire le risque de divulgation automatique, et la question sera là aussi de savoir si cette réduction est suffisante ? Et là encore, la réduction du risque pourra être évaluée.

Dans cette partie, nous nous concentrerons uniquement sur la maîtrise du risque de réidentification, celui directement adossé à la règle de fréquence du secret statistique: tout comptage non nul $X < s$ est considéré comme sensible.

2.2 Les paramètres

Quelle décision a-t-on besoin de prendre ? Il s'agit de décider de la valeur prise par les paramètres de la méthode des clés aléatoires, avant tout la déviation maximale D qui pourra être appliqué à un comptage donné et la variance V du bruit injecté.

La forme de la distribution de probabilités du bruit appliqué aux grands comptages (derrière ligne de la matrice de transition) est dépendante de ces deux paramètres D et V .

Pour les petits comptages la distribution est sans biais mais n'est pas symétrique. Ceci conduit, en pratique, à peu perturber les petits comptages, en particulier les cas $X = 1$ et $X = 2$.

Par exemple, avec l'algorithme de construction des matrices de transition fourni par le package `R` `ptable` [14], pour $D = 10$ et $V = 5$, $\mathbb{P}(Z = 0|X = 1) = 0.39$ et $\mathbb{P}(|Z| \leq 1|X = 1) = 0.90$: dans 39% des cas, un unique ne sera pas perturbé et dans 90% des cas il sera perturbé uniquement de ± 1 au maximum. Cette métrique n'est pas en soi une métrique de risque: un comptage diffusé après perturbation valant $X' = 1$ peut provenir de comptages très différents ($X \in \{1, \dots, D + 1\}$, dans le cas standard). En revanche, ce manque de perturbation va probablement engendrer des risques de divulgation importants pour ces petits comptages.

Une première solution pourrait consister à autoriser l'apparition de comptages négatifs. Une deuxième solution consisterait à accepter un peu de biais dans certaines distributions de $X'|X = i$, en particulier pour les plus petits comptages. Une troisième solution consiste à se doter d'un troisième paramètre que nous noterons js ³, qui permet de définir un intervalle de comptages interdits dans la constitution des outputs. Si s désigne le seuil de la règle de fréquence, on choisira en général js dans l'intervalle $\llbracket 0; s - 1 \rrbracket$:

- Si $js = 0$, alors X' ne subit aucune restriction;
- Si $js > 0$, alors $X' = 0$ ou $X' > js$;
- $js = s - 1$ alors $X' = 0$ ou $X' \geq s$.

La dernière option est une façon d'imiter l'approche suppressive avec laquelle aucun output ne fait apparaître aucun comptage sensible dans les tableaux. Ici, il ne s'agit pas de supprimer les comptages mais de les dévier sur des comptages suffisamment grands pour réduire les risques de divulgation.

Le choix de js a des conséquences sur la forme de la distribution de probabilité, comme le montre la figure 1 pour la distribution de probabilité de la variable $Z|X = 1$. Elle transforme aussi la matrice de transition, celle-ci ayant un plus grand nombre de lignes pour une même valeur de D ($js + D + 2$ lignes exactement). Par exemple, si $js = 2$ et $D = 2$, la matrice n'aura plus 3 lignes comme dans le cas montré plus haut, mais 6 lignes. En effet, ce n'est que pour $i \geq js + D + 1$ que la distribution sera parfaitement symétrique. Dans la suite, nous appellerons « grands comptages » tous les comptages $i \geq js + D + 1$.

Muni de ces trois paramètres, la méthode des clés aléatoires doit permettre de limiter les risques de divulgation tout en conservant l'utilité des données diffusées. La phase de calibration est donc essentielle et doit pouvoir se reposer sur des critères objectifs que nous présentons dans la suite.

2.3 Mesurer le risque de divulgation

Les critères permettant d'établir une évaluation du risque adaptée aux tableaux de fréquences construits à partir d'un jeu de données exhaustif sont décrits de la manière suivante par [15]:

Une mesure du risque de divulgation pour un tableau de fréquences issu d'un recensement devrait satisfaire les propriétés suivantes : (a) les petites valeurs de case présentent un risque de divulgation plus élevé que les grandes valeurs ; (b) des fréquences

3. Par cohérence avec l'argument correspondant dans le package `ptable`.

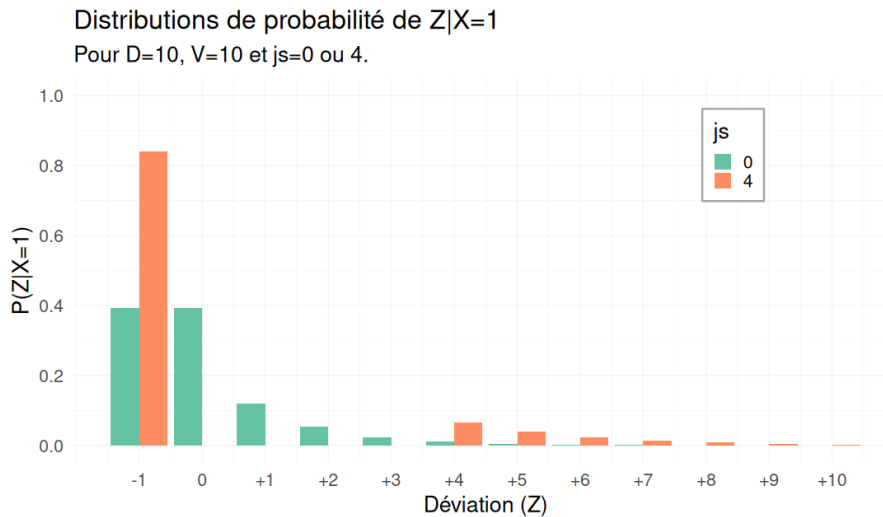


FIG. 1 – Impact de js sur les distributions de probabilité

uniformément réparties impliquent un faible risque de divulgation ; (c) plus le tableau de recensement contient de cases nulles, plus le risque de divulgation est élevé ; (d) la mesure de risque doit être bornée entre 0 et 1.

À la suite des travaux fondateurs de [16], notre mesure de risque repose sur un modèle d'inférence bayésien qui estime la probabilité qu'un attaquant détermine avec succès la véritable valeur d'une case protégée. En notant $(X = i)$ l'événement correspondant au fait que le comptage original soit égal à i , et $(X' = j)$ l'événement correspondant au comptage perturbé de la même case égal à j , cette mesure de risque s'appuie sur les probabilités a posteriori $q_{ij} = \mathbb{P}(X = i | X' = j)$, qui quantifient la confiance de l'intrus dans le fait que le comptage réel soit i étant donné le comptage perturbé observé j . Comme le mécanisme de bruit de la méthode des clés aléatoires repose sur des probabilités de transition $p_{ij} = \mathbb{P}(X' = j | X = i)$, les probabilités de transition inverses (q_{ij}) peuvent être dérivées à l'aide de la règle de Bayes :

$$q_{ij} = \frac{p_{ij} \times p_i}{q_j} \quad (1)$$

où $p_i = \mathbb{P}(X = i)$ et $q_j = \mathbb{P}(X' = j)$ correspondent respectivement aux distributions des comptages originaux et perturbés.

Affinements du scénario

Un aspect essentiel de l'évaluation du risque concerne le choix du domaine sur lequel le risque d'inférence est quantifié. Notre cadre se concentre principalement sur les *petits comptages sensibles*, en raison de leur fort potentiel de risque de divulgation. Deux stratégies complémentaires peuvent être adoptées :

- Conserver les probabilités *a posteriori* exactes q_{ij} , afin d'évaluer la capacité de l'intrus à inférer le comptage réel précis pour les cases initialement sensibles. Par exemple, se focaliser spécifiquement sur le risque d'inférer $X = 1$ pour tout comptage publié $X' = j$, noté q_{1j} .
- Considérer un risque d'inférence agrégé en regroupant les comptages en ensembles, comme mesurer la probabilité que le comptage réel appartienne à l'ensemble des comptages sensibles $I = \{1, \dots, s-1\}$, étant donné un ensemble de comptages perturbés J . Cette mesure

de risque est notée $q_{IJ} = \mathbb{P}(X \in I \mid X' \in J)$ ⁴.

Pour la suite, nous adoptons la seconde approche et fixons $J = \{1, \dots, s\}$ ⁵ afin d'évaluer la capacité de l'intrus à inférer qu'un comptage publié appartenant à J correspond à un comptage réel sensible.

Estimation des distributions a priori

Un autre aspect déterminant réside dans l'estimation des distributions des comptages originaux et perturbés, p_i et q_j . Comme les probabilités q_j peuvent être exprimées à partir des p_i et des p_{ij} , il est nécessaire d'estimer les p_i . Deux approches principales sont envisagées :

- Une estimation empirique fondée sur les fréquences observées dans les données, qui peut fournir une borne supérieure du risque en reflétant des distributions réalistes sur un ou plusieurs tableaux utilisés pour la calibration.
- L'hypothèse d'un a priori uniforme sur les comptages possibles, offrant une base de comparaison correspondant à l'estimation minimale du risque dans une situation d'incertitude maximale concernant les comptages.

Le choix de l'a priori influence fortement le risque évalué : les fréquences empiriques tendent à le maximiser, tandis que les a priori uniformes le minimisent. Considérer ces deux perspectives permet d'encadrer le risque et d'éclairer les décisions avec prudence et discernement.

2.4 Mesurer l'utilité

Il est important de se doter aussi de métriques d'utilité appropriées permettant de quantifier l'impact du mécanisme de perturbation sur les tableaux publiés. De nombreuses mesures d'utilité ont été proposées dans la littérature, en particulier pour les données de fréquences tabulaires. Parmi celles-ci, la *distance de Hellinger*, proposée par [15], s'impose comme un outil robuste pour évaluer les changements de distribution des cases induits par les méthodes de contrôle de la divulgation.

Mais on peut aussi choisir de profiter du fait que la méthode des clés aléatoires repose sur une matrice de transition décrivant les probabilités de perturbation de chaque case et qu'il est donc possible de déduire des mesures d'utilité ou de perte d'information directement à partir de ces distributions. Deux métriques illustratives sont proposées dans ce cadre :

- Une **mesure d'utilité** définie comme la probabilité que la déviation soit suffisamment contenue dans un intervalle choisi par l'utilisateur en calculant la proportion de comptages restant à une distance inférieure ou égale à d de leur valeur originale:

$$U_1 = \mathbb{P}(|Z| \leq d \mid X = i) = \sum_{j=i-d}^{i+d} p_{ij}$$

- Une **mesure de perte d'information** correspondant à la déviation absolue moyenne attendue produite par le mécanisme de perturbation:

$$U_2 = \mathbb{E}(|Z| \mid X = i) = \sum_{j=i-d}^{i+d} |j - i| p_{ij}$$

4. Cette métrique est implémentée dans un package R : <https://github.com/InseeFrLab/ckm>

5. La borne supérieure de l'ensemble est s et non $s - 1$, afin de permettre à la mesure d'être positive même lorsqu'aucun comptage publié n'est strictement « sensible », c'est-à-dire quand $j_s = s - 1$.

Ces mesures présentent deux avantages principaux. Premièrement, elles exploitent directement le mécanisme de perturbation sous-jacent via la matrice de transition, évitant ainsi de dépendre de tableaux perturbés simulés ou observés. Deuxièmement, ces métriques sont intuitivement interprétables et directement liées aux paramètres opérationnels de la méthode. Enfin, en se restreignant aux grands comptages ($i \geq js + D + 1$), la distribution de probabilités étant identique pour ceux-ci, les mesures présentées fournissent une information parfaitement aux données *a priori* les plus importantes du jeu de données diffusées.

2.5 Proposition d'une démarche pour la calibration des paramètres

Un « cadre d'évaluation du risque de divulgation et de l'utilité des données pour l'évaluation des méthodes de contrôle de la divulgation statistique » a déjà été proposé par [17]. Il s'agit d'un cadre très général, adapté à de nombreuses méthodes de protection des données. Dans le présent article, on se concentre spécifiquement sur la manière d'appliquer la méthode des clés aléatoires à partir des éléments présentés ci-dessus. Cette démarche se décompose en cinq étapes :

1. **Définition du domaine de recherche des paramètres:** Il s'agit d'établir des bornes raisonnables pour D , V et js à partir de considérations théoriques, d'études antérieures et d'une première exploration empirique. Par exemple, si le seuil de fréquence fixé par le secret statistique est $s = 5$, on s'orientera nécessairement vers $D > 5$. On pourra aussi prendre en considération des aspects de communication: il est peut-être plus facile de communiquer sur $D = 10$ que sur $D = 7$.
2. **Sélection de tableaux:** La calibration nécessite d'identifier un ou plusieurs tableaux, les plus détaillés parmi ceux qui sont envisagés à la diffusion. Dans le cas d'une diffusion synchrone, il peut s'agir de l'ensemble des tableaux diffusés.
3. **Évaluation du risque:** Estimer le risque de divulgation en évaluant la probabilité qu'un attaquant puisse inférer avec succès la véritable valeur des cases sensibles, à l'aide d'une approche bayésienne probabiliste.
4. **Évaluation de l'utilité:** Mesurer l'utilité des données à l'aide d'indicateurs globaux de perte d'information ou de validité statistique.
5. **Compromis risque–utilité:** Visualiser les résultats de risque et d'utilité sur un graphique croisant les deux dimensions, afin de faciliter l'identification des paramètres optimaux. Cet arbitrage nécessitera de fixer de manière empirique et consensuelle un seuil de tolérance au risque pour faciliter la prise de décision.

Ce cadre vise à transformer la sélection des paramètres, traditionnellement subjective, en un processus fondé sur des éléments objectifs, accompagné d'une documentation claire des arbitrages réalisés. Il est conçu pour faciliter la prise de décision des producteurs de données en la rendant aussi intuitive et transparente que possible. L'objectif d'une telle démarche est de permettre aux parties prenantes d'exprimer et de justifier clairement leurs décisions en s'appuyant sur des critères objectifs et compréhensibles, renforçant ainsi la responsabilisation des acteurs et la reproductibilité des choix. Chaque composante de la démarche doit reposer sur des procédures privilégiant la clarté, l'interprétabilité et la pertinence pratique. Ces principes garantissent que le processus de calibration soit non seulement rigoureux, mais aussi accessible, contribuant ainsi à une prise de décision éclairée.

3 Un exemple

Pour obtenir des résultats reproductibles et suffisamment réalistes⁶, nous utilisons des tableaux issus du recensement de la population publiés sur le site internet de l'Insee⁷ sans aucune restriction, en vertu d'un décret légal. Comme les données du recensement de la population sont en partie dérivées de réponses à des enquêtes, les tableaux de comptages ont été préalablement arrondis à des valeurs entières avant cette étude.

Le tableau est très ventilé: il contient plus de 5.3 millions de cases, dont 36.2% sont nulles et 37.7% ont une fréquence non nulle inférieure à 5. En supposant un seuil de confidentialité fictif de $s = 5$, plus d'un tiers des cases seraient sensibles, soit près de 60% de cases sensibles parmi les cases non nulles. Cette situation n'est pas nécessairement souhaitable, mais elle permet d'aborder un cas extrême de diffusion.

Les trois paramètres seront recherchés dans les domaines suivants:

- $D \in \{5; 10\}$ ⁸;
- $V \in \{2.5; 5; 10; 15\}$;
- $js \in \{0; 2; 4\}$.

Certains paramètres ne sont pas compatibles entre eux. Par exemple, il n'est pas possible de construire une matrice de transition avec $D = 10$, $js = 4$ et $V < 10$, avec l'ensemble des contraintes imposées sur cette matrice.

3.1 Le compromis risque–utilité

Comme mentionné en 2.3, la mesure du risque de divulgation est définie comme la capacité d'un intrus à inférer le caractère sensible d'un comptage original à partir d'un comptage perturbé. Plus précisément, suivant les notations précédentes, nous mesurons q_{IJ} , où $I = \{1, \dots, 4\}$ désigne l'ensemble des comptages sensibles originaux et $J = \{1, \dots, 5\}$ l'ensemble des comptages à partir desquels l'intrus formule son inférence.

La figure 2 illustre l'intervalle des mesures de risque de divulgation, entre un a priori uniforme (cercles) et un a priori empirique (carrés), pour chaque triplet de paramètres. On observe que la réduction du risque est plus efficace lorsque V ou js augmentent. D'un point de vue utilité, on peut s'appuyer dans un premier temps sur la figure 3 qui compare les distributions de probabilité du bruit injecté sur les grands comptages ($i \geq js + D + 1$) pour les différents jeux de paramètres testés. On peut observer, en particulier, qu'une variance de $V = 10$ pour $D = 5$ équivaut à appliquer une distribution uniforme et réduit nettement l'utilité. En outre, pour le cas $V = 5$ (deuxième colonne), l'utilité est mieux conservée avec une déviation maximale $D = 10$ qu'avec $D = 5$.

Le graphique risque–utilité (Figure 4) représente quant à lui, pour différentes combinaisons de paramètres (D, V, js) , le risque de divulgation en fonction de l'utilité des données, définie ici comme la probabilité qu'un comptage soit perturbé de moins de 3 unités par rapport à sa valeur originale. Le scénario idéal, c'est-à-dire celui qui réduirait à néant tout le risque de divulgation tout en conservant une utilité parfaite des données correspondrait au point de coordonnée $(1; 0)$ ⁹.

Comme attendu, les valeurs d'utilité plus élevées sont généralement associées à des niveaux de risque plus importants, tandis que les protections plus fortes (risque plus faible) entraînent une utilité réduite. Par exemple, les configurations avec $js = 0$ présentent l'utilité la plus élevée, atteignant

6. L'ensemble du matériel expérimental est disponible sur Github: https://github.com/julienjamme/ckm_risk_utility_experiments.

7. <https://www.insee.fr/fr/statistiques/8582668?sommaire=8582771>

8. Dans la pratique, il sera préférable de tester des valeurs nettement supérieures au seuil de confidentialité. Ce paramètre fixé trop bas peut avoir des effets négatifs sur l'utilité des données.

9. On fera attention au fait que l'échelle des ordonnées commence à 0,5 et non à 0.

jusqu'à 0.95, mais à des niveaux de risque d'environ 0.88. À l'inverse, augmenter j_s à 4 réduit nettement le risque (entre 0.60 et 0.7) – mais diminue également l'utilité (entre 0.72 et 0.76).

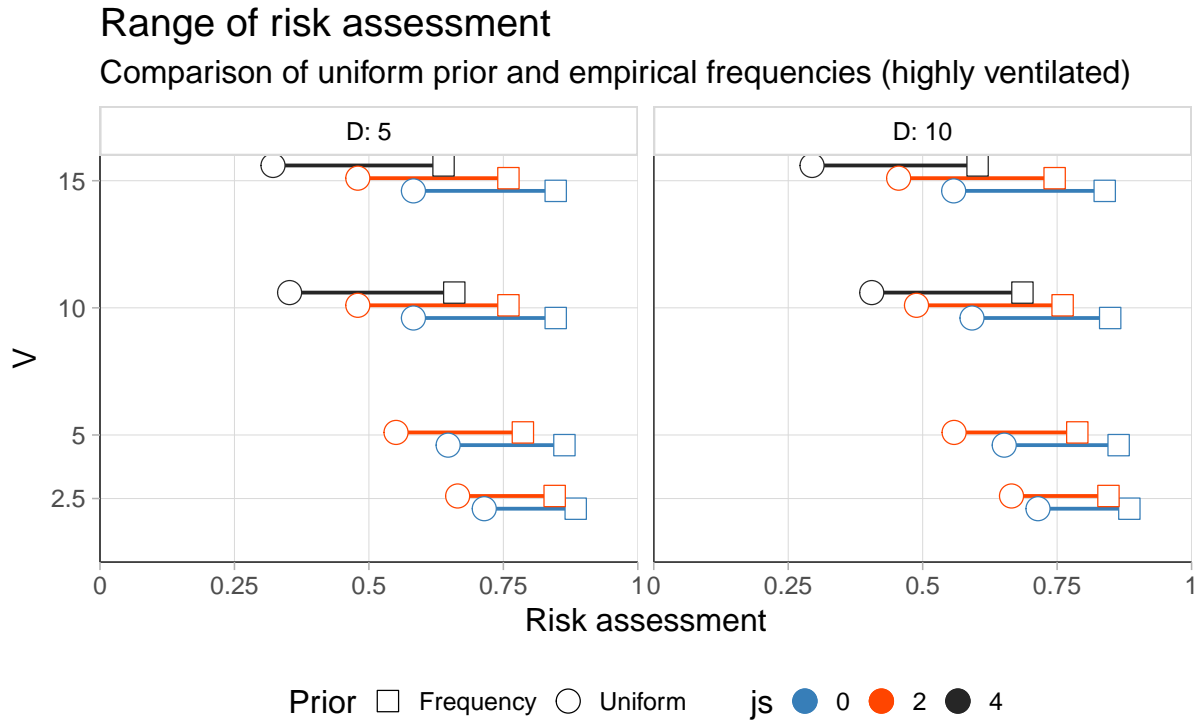


FIG. 2 – Intervalle des mesures de risque de divulgation selon D , V , j_s et le type de probabilités a priori (source: [18])

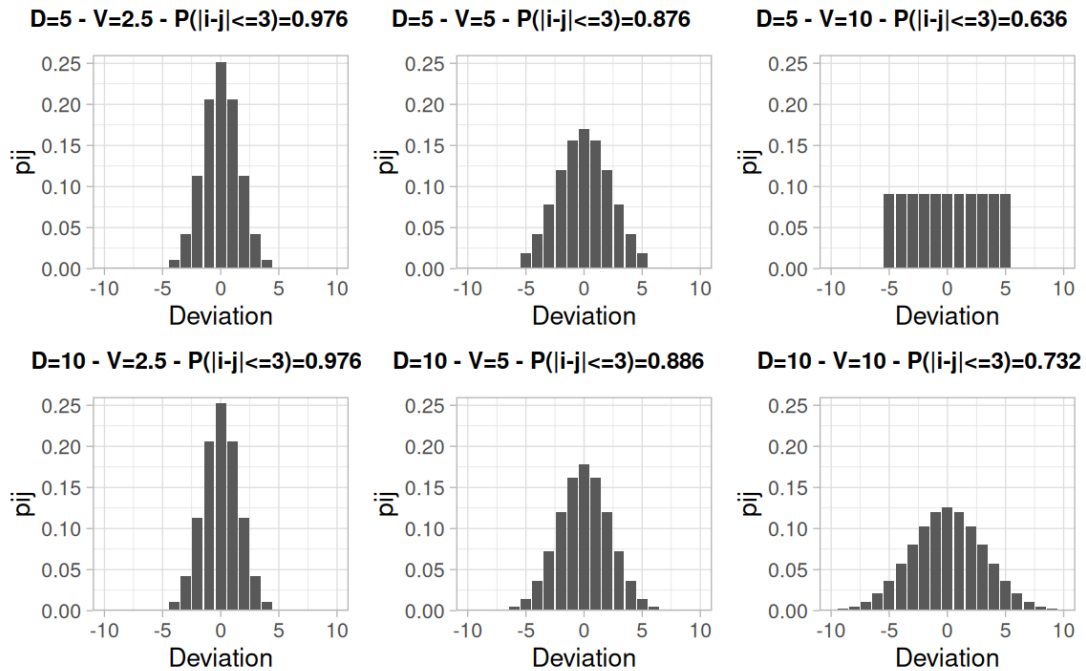


FIG. 3 – Distributions de probabilité pour différents paramètres D et V et leurs mesures d'utilité respectives (mesure U_1 , avec $d = 3$).

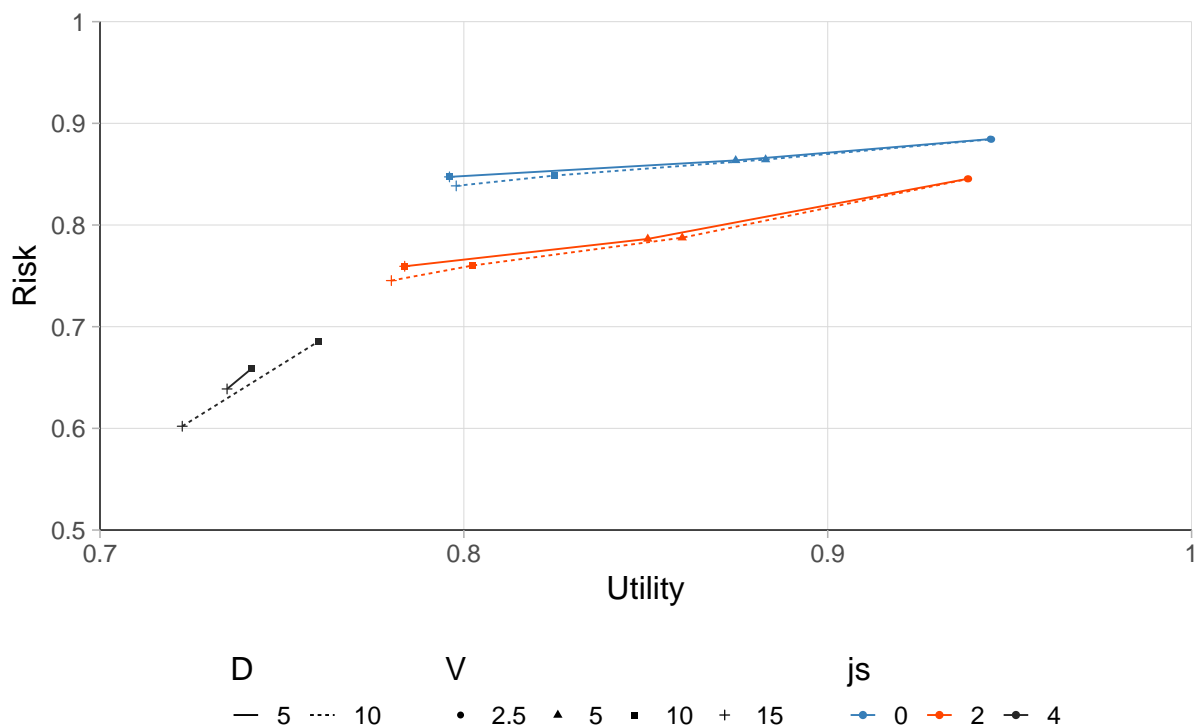


FIG. 4 – *Risque vs Utilité pour les données du recensement*

Pour des valeurs fixées de D et V , l'augmentation du seuil de petits comptages j_s abaisse systématiquement à la fois le risque et l'utilité. En observant les variations de V , l'augmentation de la variance déplace également les points vers un risque et une utilité plus faibles. En comparant les écarts maximums, les configurations avec $D = 10$ offrent généralement une utilité légèrement supérieure, à risque donné, comparativement à celles avec $D = 5$.

Pour procéder à l'arbitrage et prendre la décision des paramètres à retenir, l'approche la plus pragmatique consisterait à fixer un seuil maximal de risque de divulgation par inférence. Par exemple, fixer un seuil de 80% signifierait que nous acceptons que dans le pire des cas, l'attaquant se trompera de 20% dans son inférence pour déduire qu'un comptage perturbé donné est un comptage sensible. Définir ce seuil est une tâche consensuelle et on ne pourra pas ici dire quel seuil fixer. Quelques éléments pour éclairer ce choix :

- Pour les données entreprises, nous autorisons des inférences à 85%, puisque c'est le seuil de la règle de dominance utilisé dans toutes les productions statistiques sur les entreprises.
- Pour certaines données issues des sources fiscales, des seuils d'inférence à 80% sont utilisés (source Filosofi, par exemple).
- Il faut avoir en tête que la mesure de risque est par définition très conservatrice puisqu'en reposant les calculs sur les a priori empiriques (fréquences des comptages observées dans les vraies données), la qualité d'une inférence n'est pas connue précisément par un attaquant.
- Dans la majorité des cas, un attaquant se reposera sur un a priori non informatif (distribution uniforme des comptages) ce qui, comme on l'observe sur la figure 2, diminue beaucoup la qualité des inférences.

Supposons que l'institution fixe un risque de divulgation maximal acceptable à 0.8. Le graphique 4 suggère alors que la meilleure configuration est $(D = 10, V = 5, j_s = 2)$ (triangle orange sur la ligne en pointillés), atteignant un risque de 0.79 et une utilité de 0.86. Mais si l'institution souhaite publier des tableaux reproduisant les schémas de suppression des comptages sensibles, elle peut

décider de fixer $js = s - 1 = 4$. Dans ce cas, la meilleure configuration est ($D = 10, V = 10, js = 4$) (carré orange sur la ligne en pointillés), atteignant un risque de 0.69 et une utilité de 0.76.

Dans l'ensemble, le graphique Risque-vs-Utilité constitue un outil essentiel d'aide à la décision, permettant aux producteurs de données d'identifier les combinaisons de paramètres répondant à leur appétence au risque tout en maximisant l'utilité des données, dans une approche transparente et fondée sur des éléments probants. Le seul choix d'un seuil maximal du risque d'inférence rend alors le choix du paramétrage presque évident.

4 Conclusion

Les caractéristiques de la méthode des clés aléatoires rendent cette méthode particulièrement désirable dès lors qu'un producteur est confronté à la protection de données tabulées produites à une échelle industrielle, qui génèrent de nombreux risques de divulgation (directe ou par différenciation) et rendent les méthodes classiques coûteuses en perte d'information et dans leur implémentation. À des niveaux importants de complexité des diffusions, ces avantages peuvent largement compenser les inconvénients d'utiliser une méthode perturbatrice (additivité, perception du public, etc.).

Quand la décision est prise d'utiliser une telle méthode, le défi reste de la calibrer, c'est-à-dire de faire les bons choix de paramètres pour obtenir le meilleur compromis risque-utilité possible. La proposition faite ici est de se munir de mesures objectives du risque et de l'utilité basées sur les ingrédients de base de la méthode que sont les probabilités de transition. Doté de ces mesures, le producteur, accompagné d'un méthodologue, pourra baser son choix de paramètres à partir d'un graphique Risque-Utilité, choix qui pourra s'avérer rapidement évident si, par une démarche consensuelle, un seuil de risque maximal acceptable est défini préalablement.

Des travaux supplémentaires pourront être réalisés en particulier pour expérimenter une matrice de transition injectant légèrement du biais dans les perturbations des petits comptages, afin de voir si cela permet de mieux maîtriser la perte d'utilité sur les grands comptages par rapport à la solution actuelle d'utiliser un seuil d'interdiction de certains petits comptages.

Bibliographie

- [1] M. CHEVALIER, A. HACHID et J. JAMME, « Données sur les Quartiers prioritaires de la politique de la ville (QPV) : une nouvelle méthode pour protéger le secret statistique », jan. 2025. Blog de l'Insee - Section: Données et méthodes - <https://blog.insee.fr/qpv-nouvelle-methode-secret-statistique/>.
- [2] B. FRASER et J. WOOTON, « A Proposed Method for Confidentialising Tabular Output to Protect against Differencing », in *Monographs of Official Statistics: Work Session on Statistical Data Confidentiality*, p. 299–302, 2005.
- [3] J. JAMME, « La méthode des clés aléatoires (Cell Key Method) », fév. 2025.
- [4] M. MÖHLER, J. JAMME, E. DE JONGE, A. MŁODAK, J. GUSSENBAUER et P.-P. DE WOLF, *Guidelines for statistical disclosure control methods applied on geo-referenced data: 2025 edition*. Manuals and guidelines, LU: Publications Office, 2025.
- [5] V. COSTEMALLE, « Detecting geographical differencing problems in the context of spatial data dissemination », *Statistical Journal of the IAOS*, vol. 35, p. 559–568, déc. 2019.
- [6] C. BAUDRY, « Analyse automatique des métadonnées pour la protection des données tabulées », in *Journées de méthodologie statistique (JMS)*, (Paris), nov. 2025. Section: JMS 2025.
- [7] C. GUILLO, J. JAMME et N. RASTOUT, « Protection of linked tables with a suppressive approach. Method and Use Cases », in *NTTS 2023*, (Bruxelles), mars 2023.
- [8] P.-P. DE WOLF, A. HUNDEPOOL, S. GIESSING, J.-J. SALAZAR-GONZALEZ et J. CASTRO, « Tau-ARGUS », août 2024.

- [9] LANGSRUD, « Secondary Cell Suppression by Gaussian Elimination: An Algorithm Suitable for Handling Issues with Zeros and Singletons », in Privacy in Statistical Databases (J. DOMINGO-FERRER et M. ÖNEN, édés), vol. 14915, p. 87–101, Cham: Springer Nature Switzerland, 2024. Series Title: Lecture Notes in Computer Science.
- [10] G. THOMPSON, S. BROADFOOT et D. ELAZAR, « Methodology for the Automatic Confidentialisation of Statistical Outputs from Remote Servers at the Australian Bureau of Statistics », in Joint UNECE/Eurostat work session on statistical data confidentiality, (Ottawa), oct. 2013.
- [11] T. ENDERLE, S. GIESSING et R. TENT, « Designing Confidentiality on the Fly Methodology – Three Aspects », in Privacy in Statistical Databases (J. DOMINGO-FERRER et F. MONTES, édés), vol. 11126, p. 28–42, Cham: Springer International Publishing, 2018. Series Title: Lecture Notes in Computer Science.
- [12] S. GIESSING et R. TENT, « Concepts for generalising tools implementing the cell key method to the case of continuous variables », in UNECE - Expert Meeting on Statistical Data Confidentiality, (the Hague), oct. 2019.
- [13] E. SCHULTE NORDHOLT, R. TENT, S. GIESSING, M. GOLMAJER, F. BACH, M. de VRIES, R. van de LAAR, P.-P. de WOLF et N. KROL, Guidelines for statistical disclosure control methods for census and demographics data: 2024 edition. Manuals and guidelines, LU: Publications Office, 2024.
- [14] T. ENDERLE, « R package ptable : Generation of Perturbation Tables for the Cell-Key Method », mars 2023.
- [15] N. SHLOMO, L. ANTAL et M. ELLIOT, « Measuring Disclosure Risk and Data Utility for Flexible Table Generators », Journal of Official Statistics, vol. 31, p. 305–324, juin 2015.
- [16] T. ENDERLE, S. GIESSING et R. TENT, « Calculation of Risk Probabilities for the Cell Key Method », in Privacy in Statistical Databases (J. DOMINGO-FERRER et K. MURALIDHAR, édés), vol. 12276, p. 151–165, Cham: Springer International Publishing, 2020. Series Title: Lecture Notes in Computer Science.
- [17] N. SHLOMO et C. YOUNG, « Statistical Disclosure Control Methods Through a Risk-Utility Framework », in Privacy in Statistical Databases (D. HUTCHISON, T. KANADE, J. KITTLER, J. M. KLEINBERG, F. MATTERN, J. C. MITCHELL, M. NAOR, O. NIERSTRASZ, C. PANDU RANGAN, B. STEFFEN, M. SUDAN, D. TERZOPOULOS, D. TYGAR, M. Y. VARDI, G. WEIKUM, J. DOMINGO-FERRER et L. FRANCONI, édés), vol. 4302, p. 68–81, Berlin, Heidelberg: Springer Berlin Heidelberg, 2006. Series Title: Lecture Notes in Computer Science.
- [18] J. JAMME, « A Framework for Cell Key Method Parameters Calibration based on a Risk-Utility trade-off », in Expert Meeting on Statistical Data Confidentiality, (Barcelona), UNECE, oct. 2025.