

计算 $\mathcal{R}_i = \mathcal{R}_i - \bigcup_{j=0}^{i-1} \mathcal{R}_j$; //去掉 \mathcal{R}_i 中与前面集合重复点集;
 }
 //设 $\mathcal{R}_n = \{[-a_n, -\Delta f, -a_n], [a_n, a_n + \Delta f]\}$
 步骤 4: 计算 $sum1 = \sum_{i=0}^n \exp(-i\varepsilon)\mu(\mathcal{R}_i)$ 及 $sum2 = 2\Delta f^* \exp(-\varepsilon(n+1))/(1-\exp(-\varepsilon))$;
 步骤 5: 计算 $\alpha = sum1 + sum2$;
 步骤 6: 设置 $p_i(r) = \exp(-i\varepsilon)/\alpha$, 其中, $r \in \mathcal{R}_i, i \in \mathbb{N}$;
 步骤 7: Return $\{(p_i(r), \mathcal{R}_i): i \in \mathbb{N}\}$.

算法 1 需要下面的集合序列收敛性质:

定义 2. 设 $\mathcal{R} = \mathbb{R}$. 对数据集 $x \in \mathcal{D}$, 称集合序列 $\{\mathcal{R}_i: i \in \mathbb{N}\}$ 在第 $n \in \mathbb{N}$ 步收敛, 如果存在 $a_n \in \mathcal{R}$, 有:

$$\mathcal{R}_n = \pm(a_n + [0, \Delta f]) \text{ 且 } \mathcal{R}_{n+1} = \pm(a_n + [\Delta f, 2\Delta f]).$$

我们的实验设计如下: 首先, 计算算法 1 生成的密度函数的数学期望值 $mean = \sum_{i=0}^{\infty} p_i \times \int_{r \in \mathcal{R}_i} r \mu(dr)$; 然后, 与 Laplace 机制^[27]的期望值 $\Delta f/\varepsilon$ 及 Staircase 机制^[18]的期望值 $\Delta f \times \exp(\varepsilon/2)/(\exp(\varepsilon)-1)$ 进行比较, 越小的期望值代表噪声复杂度越低. 由于发现 Staircase 机制的期望值比 Laplace 机制的期望值小, 因此, 本实验只比较 $mean$ 与 Staircase 机制期望值 $\Delta f \times \exp(\varepsilon/2)/(\exp(\varepsilon)-1)$ 的大小, 即 $rate = mean \times (\exp(\varepsilon)-1)/(\Delta f \times \exp(\varepsilon/2))$. $rate$ 值小于 1, 代表算法 1 对应机制优于 Staircase 机制; 反之, 则差于 Staircase 机制.

我们构造两组实验: 第 1 组使用相同的敏感度及不同的函数邻居集测度, 第 2 组使用不同的敏感度及相同的函数邻居集测度. 第 1 组实验有 4 个不同的查询函数 f_1, f_2, f_3, f_4 , 其邻居集分别为 $\mathcal{V}_1 = [0, 1] \cup [1000, 1001], \mathcal{V}_2 = [0, 100] \cup [1000, 1001], \mathcal{V}_3 = [0, 500] \cup [1000, 1001], \mathcal{V}_4 = [0, 1001]$. 4 个不同集合可以理解为 4 个不同单位的工资分布情况: \mathcal{V}_1 代表了工资两级分化极其严重的单位; \mathcal{V}_4 代表了该单位的工资虽然最高工资和最低工资差距很大, 但是高中低档工资都可以出现; 其他两个集合代表了折中的情形. 实验结果如图 1 所示.

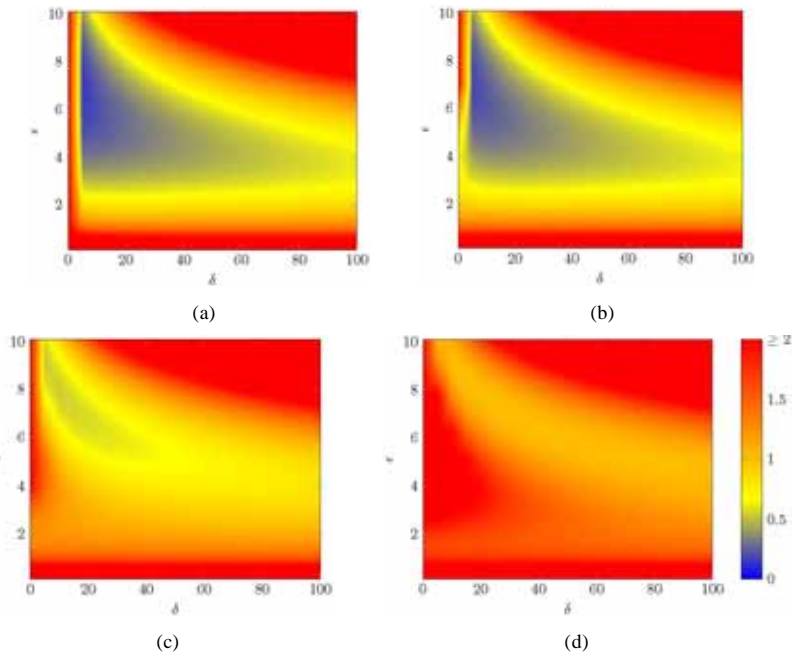


Fig.1
图 1

图 1 的 4 个子图分别代表了 $\mathcal{V}_1, \mathcal{V}_2, \mathcal{V}_3, \mathcal{V}_4$ 的实现效果图.其中,横坐标表示 δ 的取值,纵坐标表示 ϵ 的取值,坐标点 (δ, ϵ) 对应的值表示相应的 *rate* 的值.从图 1 可以发现:随着邻居集测度 $\mu(\mathcal{V}_i)$ 的增大, *rate* 的值相应增大,从而说明算法 1 中机制的优势是随 $\mu(\mathcal{V}_i)$ 的增大而递减的.

第 2 组实验有 3 个查询函数,其邻居集分别为 $\mathcal{V}_5=[0,1] \cup [100,101], \mathcal{V}_6=[0,1] \cup [1000,1001], \mathcal{V}_7=[0,1] \cup [2000,2001]$.显然,3 个函数的全局(及局部)敏感度分别为 101,1001,2001,但有相同的邻居集测度 $\mu(\mathcal{V}_i)=2$.其实验结果如图 2 所示.

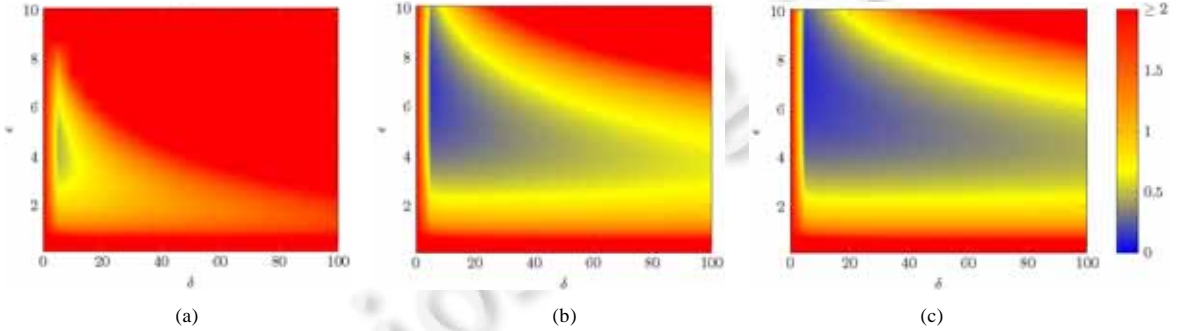


Fig.2
图 2

图 2 的 3 个子图分别代表了 $\mathcal{V}_5, \mathcal{V}_6, \mathcal{V}_7$ 的实现效果图.与图 1 一样,每个子图中坐标点 (δ, ϵ) 对应的值表示相应的 *rate* 的值.从图 2 可以发现:当敏感度增加的时候,算法 1 中机制的噪声复杂度要比 Staircase 机制的复杂度相对越来越小.也就是说:在函数邻居集的测度不变化的条件下,随着敏感度的增加,算法 1 中机制的噪声复杂度要比 Staircase 等敏感度方法的噪声复杂度越来越小.

综合图 1、图 2 的结果,我们可以得到如下的结论:若查询函数的邻居集 \mathcal{V} 的敏感度 Δf 与其测度 $\mu(\mathcal{V})$ 的比率 $\Delta f/\mu(\mathcal{V})$ 越来越大时,敏感度方法将添加越来越大的噪声.而本文的方法会极大地降低这个比率带来的噪声复杂度升高的问题,这也是本文的方法优于敏感度方法的最显著特征.

Laplace 机制、Staircase 机制及算法 1 中机制的密度函数如图 3 所示.图 3(c)的密度函数与前两个密度函数的显著区别是:密度函数并不是从中心向两边递减的,而是总体递减,但是局部有增的密度函数.这种密度函数更适合于 $\Delta f/\mu(\mathcal{V})$ 很大的情形,因为其波浪状形态可以不受全局敏感度的限制而更适合于邻居集的(不同查询函数的)多变特性.

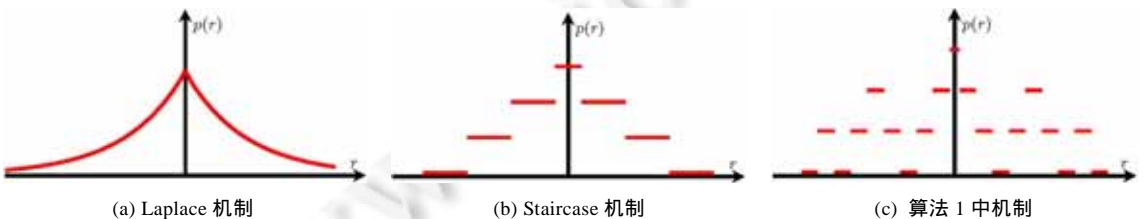


Fig.3 The density functions of Laplace mechanism, Staircase mechanism and the mechanism in algorithm 1
图 3 Laplace 机制、Staircase 机制、算法 1 中机制的密度函数

算法 2 是算法 1 生成的密度函数的随机变量生成算法.

算法 2. 生成随机变量.

输入:敏感度 Δf 、迭代次数 n 、 $\{\mathcal{R}_i; i \in \{0, 1, \dots, n\}\}$ 、 a_n 、 $\alpha = \text{sum}1 + \text{sum}2$;

//设 $\mathcal{R}_n = \{-a_n, -\Delta f, -a_n\}, [a_n, a_n + \Delta f]$, 即,算法从第 n 步后收敛;

输出:随机变量 X .

步骤 1:以 $\exp(-i\varepsilon)\mu(\mathcal{R}_i)/\alpha$ 的概率输出 $i, i \in \{0, 1, \dots, n\}$; 以概率 $\text{sum}2/\alpha$ 输出 $n+1$;

步骤 2:If 步骤 1 输出值 $i \in \{0, 1, \dots, n\}$ then

在 \mathcal{R}_i 中均匀抽样出一个数 x ;

Return x ;

Else

从参数为 $1-\exp(-\varepsilon)$ 的几何分布中抽样出一个整数 j ;

从 $[-a_n-(j+1)\Delta f, -a_n-j\Delta f] \cup [a_n+j\Delta f, a_n+(j+1)\Delta f]$ 中均匀随机抽样出一个数 y ;

Return y ;

6.1 时间复杂度

算法 1 和算法 2 的时间复杂度由集合序列 $\{\mathcal{R}_i; i \in \mathbb{N}\}$ 的收敛性决定. 本节的例子中, 该序列都在 2 200 步内达到了收敛. 但是对一般线性查询函数, 我们还没能有方法证明其一定在有限步内收敛. 不过, 我们相信这个结论是对的.

若 $\{\mathcal{R}_i; i \in \mathbb{N}\}$ 在第 n 步收敛, 且 $\max_{i \in \{0, 1, \dots, n\}} N_i = N$, 其中, N_i 是 \mathcal{R}_i 中区间的个数, 则算法 1 的时间复杂度为 $O(n^2 \times N^2)$, 算法 2 的时间复杂度为 $O(n)$.

7 结论

传统的观点认为, 查询函数的(全局或局部)敏感度是查询函数噪声复杂度的(最主要)标志. 本文的结论发现: 敏感度其实只是查询函数的邻居集 \mathcal{V}^x 的一个极值特征, 还有很多敏感度无法刻画的特征(如该集中点的分布). 本文的方法可根据 \mathcal{V}^x 中点的分布情形, 构造适合于该分布的差分隐私机制, 相应的密度函数类似于图 3(c) 所示. 与敏感度方法(如图 3(a)、图 3(b)所示)不同, 本文机制的密度函数不是(向两边)全局递减的, 而是以局部有起伏的方式(向两边)递减的, 这种递减方式更具灵活性, 更适用于 \mathcal{V}^x 中点的不规则分布情形.

第 6 节的实际例子也说明, 本文的方法一般要比 Laplace 机制及 Staircase 机制更精确. 但是本文的方法具有很高的时间复杂度, 不适合于直接使用在大数据集上. 如何通过本文的方法和结论构造低时间复杂度的非敏感度方法, 将作为未来的一项工作.

另外, 本文的很多分析方法和结论(如第 3 节、第 4 节中的内容都没有设定函数类型)也可对非线性查询问题^[15, 16, 34-36]进行分析, 因为非线性查询问题具有更加复杂的邻居集. 至于其分析有效程度如何, 还需要进一步探索. 非线性查询问题将具有更加复杂的最优机制, 不同数据集将对对应形状各异的密度函数, 对非线性函数的非敏感度机制研究将作为未来的另一项工作.

致谢 本文作者感谢匿名审稿专家对本文初稿提出的宝贵建议和意见, 这些建议和意见对本文的完整性及易读性有很大的帮助. 同时, 这些建议和意见促使我们发现了初稿中的一个证明错误.

References:

- [1] Ganta SR, Kasiviswanathan SP, Smith A. Composition attacks and auxiliary information in data privacy. In: Proc. of the 14th ACM SIGKDD Int'l Conf. on Knowledge Discovery and Data Mining. 2008. 265-273. [doi: 10.1145/1401890.1401926]
- [2] Narayanan A, Shmatikov V. Robust de-anonymization of large sparse datasets. In: Proc. of the 2008 IEEE Symp. on Security and Privacy (S&P 2008). 2008. 111-125. [doi: 10.1109/SP.2008.33]
- [3] Aggarwal CC, Yu PS. Privacy-Preserving data mining—Models and algorithms. In: Proc. of the Advances in Database Systems, Vol.34. Springer-Verlag, 2008. [doi: 10.1007/978-0-387-70992-5]
- [4] Dwork C. A firm foundation for private data analysis. Communications of the ACM, 2011, 54(1):86-95. [doi: 10.1145/1866739.1866758]

- [5] Dwork C, Roth A. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 2014,9(3-4):211–407. [doi: 10.1561/04000000042]
- [6] Jain P, Thakurta A. Differentially private learning with kernels. In: *Proc. of the 30th Int'l Conf. on Machine Learning*. 2013. 118–126.
- [7] Zhang J, Zhang ZJ, Xiao XK, Yang Y, Winslett M. Functional mechanism: Regression analysis under differential privacy. *Proc. of the VLDB Endowment*, 2012,5(11):1364–1375. [doi: 10.14778/2350229.2350253]
- [8] Zhang J, Cormode G, Procopiuc CM, Srivastava D, Xiao XK. Privbayses: Private data release via bayesian networks. In: *Proc. of the 2014 ACM SIGMOD Int'l Conf. on Management of Data*. 2014. 1423–1434. [doi: 10.1145/2588555.2588573]
- [9] Li NH, Qardaji W, Su D. Provably private data anonymization: Or, k -anonymity meets differential privacy. *CERIAS Tech Report*, 2010.
- [10] Soria-Comas J, Domingo-Ferrer J, Sánchez D, Martínez S. Enhancing data utility in differential privacy via microaggregation-based k -anonymity. *The Int'l Journal on Very Large Data Bases*, 2014,23(5):771–794. [doi: 10.1007/s00778-014-0351-4]
- [11] Kasiviswanathan SP, Lee HK, Nissim K, Raskhodnikova S, Smith AD. What can we learn privately? In: *Proc. of the IEEE 49th Annual IEEE Symp. on Foundations of Computer Science*. 2008,40(3):793–826. [doi: 10.1109/FOCS.2008.27]
- [12] Wasserman L, Zhou S. A statistical framework for differential privacy. *Journal of the American Statistical Association*, 2010, 105(489):375–389. [doi: 10.1198/jasa.2009.tm08651]
- [13] Chaudhuri K, Hsu D. Convergence rates for differentially private statistical estimation. In: *Proc. of the Int'l Conf. on Machine Learning*. 2012. 1327–1334.
- [14] Kellaris G, Papadopoulos S, Xiao XK, Papadias D. Differentially private event sequences over infinite streams. *Proc. of the VLDB Endowment*, 2014,7(12):1155–1166. [doi: 10.14778/2732977.2732989]
- [15] Karwa V, Raskhodnikova S, Smith AD, Yaroslavtsev G. Private analysis of graph structure. *ACM Trans. on Database Systems*, 2014,39(3):22:1–22:33. [doi: 10.1145/2611523]
- [16] Chaudhuri D, Sarwate AD, Sinha K. A near-optimal algorithm for differentially-private principal components. *Journal of Machine Learning Research*, 2013,14(1):2905–2943.
- [17] Nikolov A, Talwar D, Zhang L. The geometry of differential privacy: The sparse and approximate cases. In: *Proc. of the Annual ACM Symp. on Theory of Computing*. 2013. 351–360. [doi: 10.1145/2488608.2488652]
- [18] Geng Q, Viswanath P. The optimal noise-adding mechanism in differential privacy. *IEEE Trans. on Information Theory*, 2016, 62(2):925–951. [doi: 10.1109/TIT.2015.2504967]
- [19] Gupte M, Sundararajan M. Universally optimal privacy mechanisms for minimax agents. In: *Proc. of the 29th ACM SIGMOD-SIGACT-SIGART Symp. on Principles of Database Systems*. 2010. 135–146. [doi: 10.1145/1807085.1807105]
- [20] Li C, Miklau G. Optimal error of query sets under the differentially-private matrix mechanism. In: *Proc. of the 16th Int'l Conf. on Database Theory*. 2013. 272–283. [doi: 10.1145/2448496.2448529]
- [21] Yaroslavtsev G, Cormode G, Procopiuc CM, Srivastava D. Accurate and efficient private release of datacubes and contingency tables. In: *Proc. of the IEEE 29th Int'l Conf. on Data Engineering (ICDE)*. 2013. 745–756. [doi: 10.1109/ICDE.2013.6544871]
- [22] Wang ZT, Fan K, Zhang JQ, Wang LW. Efficient algorithm for privately releasing smooth queries. In: *Proc. of the Advances in Neural Information Processing Systems*. 2013. 782–790.
- [23] Geng Q, Viswanath P. Optimal noise adding mechanisms for approximate differential privacy. *IEEE Trans. on Information Theory*, 2016,62(2):952–969. [doi: 10.1109/TIT.2015.2504972]
- [24] Brenner H, Nissim K. Impossibility of differentially private universally optimal mechanisms. In: *Proc. of the 51st IEEE Annual Symp. on Foundations of Computer Science*. 2010. 71–80[doi: 10.1109/FOCS.2010.13]
- [25] Zhang J, Cormode G, Procopiuc CM, Srivastava D, Xiao XK. Private release of graph statistics using ladder functions. In: *Proc. of the 2015 ACM SIGMOD Int'l Conf. on Management of Data*. 2015. 731–745. [doi: 10.1145/2723372.2737785]
- [26] Raskhodnikova S, Smith AD. Smooth sensitivity and sampling in private data analysis. In: *Proc. of the 39th Annual ACM Symp. on Theory of Computing*. 2007. 75–84. [doi: 10.1145/1250790.1250803]
- [27] Dwork C. Differential privacy. In: *Proc. of the Int'l Colloquium on Automata, Languages, and Programming*, 2006. 1–12. [doi: 10.1007/978-3-540-79228-4_1]

- [28] Hardt M, Talwar K. On the geometry of differential privacy. In: Proc. of the 42nd ACM Symp. on Theory of Computing. 2010. 705–714. [doi: 10.1145/1806689.1806786]
- [29] Ghosh A, Roughgarden T, Sundararajan M. Universally utility-maximizing privacy mechanisms. In: Proc. of the 41st ACM Symp. on Theory of Computing. 2009. 351–360. [doi: 10.1145/1536414.1536464]
- [30] Gupta A, Roth A, Ullman J. Iterative constructions and private data release. In: Proc. of the 9th Int'l Conf. on Theory of Cryptography. 2012. 339–356. [doi: 10.1007/978-3-642-28914-9_19]
- [31] Roth A, Roughgarden T. Interactive privacy via the median mechanism. In: Proc. of the 42nd ACM Symp. on Theory of Computing. 2010. 765–774. [doi: 10.1145/1806689.1806794]
- [32] Hardt M, Ligett K, McSherry F. A simple and practical algorithm for differentially private data release. In: Proc. of the Advances in Neural Information Processing Systems. 2010. 2339–2347.
- [33] Collette Y, Siarry P. Multiobjective Optimization: Principles and Case Studies. Springer Science & Business Media, 2013. 15–43.
- [34] Lu GQ, Zhang XJ, Ding LP, Li YF, Liao X. Frequent sequential pattern mining under differential privacy. Journal of Computer Research and Development, 2015,52(12):2789–2801 (in Chinese with English abstract) [doi: 10.7544/issn1000-1239.2015.20140516]
- [35] Ouyang J, Yin J, Liu SP. Differential privacy publishing strategy for distributed transaction data. Ruan Jian Xue Bao/Journal of Software, 2015,26(6):1457–1472 (in Chinese with English abstract). <http://www.jos.org.cn/1000-9825/4526.htm> [doi: 10.13328/j.cnki.jos.004576]
- [36] Chen SX. New methods for linear queries in providing differential privacy protection [Ph.D. Thesis]. Shanghai: Fudan University, 2012 (in Chinese with English abstract).

附中文参考文献:

- [34] 卢国庆,张啸剑,丁丽萍,李彦峰,廖鑫.差分隐私下的一种频繁序列模式挖掘方法.计算机研究与发展,2015,52(12):2789–2801. [doi: 10.7544/issn1000-1239.2015.20140516]
- [35] 欧阳佳,印鉴,刘少鹏.一种分布式事务数据的差分隐私发布策略.软件学报,2015,26(6):1457–1472. <http://www.jos.org.cn/1000-9825/4576.htm> [doi: 10.13328/j.cnki.jos.004576]
- [36] 陈世熹.提供差分隐私保护的线性查询新方法[博士学位论文].上海:复旦大学,2012.



武跟强(1980 -),男,甘肃天水人,博士生, CCF 学生会会员,主要研究领域为安全多方计算与隐私保护,差分隐私理论.



夏娴瑶(1986 -),女,博士生,CCF 学生会会员,主要研究领域为隐私保护.



贺也平(1962 -),男,博士,研究员,博士生导师,主要研究领域为系统安全,隐私保护.