

众包环境下的隐私保护研究

刘欢, 吴桂兴

(中国科学技术大学 软件学院 江苏 苏州 215123)

摘要: 自从众包的概念被提出以来,就日益受到学术界和工业界的广泛关注,而且随着移动互联网和智能设备的蓬勃发展,众包任务的执行也变得更加便利和高效。可以说,众包已经成为一个热门的研究方向。与此同时,人们也越来越重视个人的信息隐私,尤其是在执行众包任务时,人们不希望透露过多的个人信息,以免造成不必要的安全隐患。因此,需要在设计众包机制的时候,考虑到隐私保护,做到既能保证众包任务的顺利执行,又能保护参与者的隐私数据。为此,不仅介绍了众包的概念和 workflows,还详细介绍了差分隐私的数学定义以及实现机制。同时,也对差分隐私在不同众包场景下的应用做了分析,分别是:保护用户提供的数据、保护用户的位置信息以及保护用户的出价信息。通过总结现有的研究,最后对未来的研究方向进行展望。

关键词: 众包计算; 隐私保护; 差分隐私; 机制设计

中图分类号: TP391

文献标识码: A

文章编号: 1673-629X(2018)12-0111-05

doi: 10.3969/j.issn.1673-629X.2018.12.024

Research on Privacy Protection in Crowdsourcing

LIU Huan, WU Gui-xing

(School of Software Engineering, University of Science and Technology of China, Suzhou 215123, China)

Abstract: Since the concept of crowdsourcing has been raised, it has drawn increasing attention from academia and industry, and with the booming development of mobile internet and smart devices, the execution of crowdsourcing tasks has become more convenient and efficient. It can be said that crowdsourcing has become a very promising area of development. At the same time, people pay more and more attention to individual information privacy. Especially in the process of crowdsourcing, people do not want to disclose excessive personal information in order to avoid unnecessary security risks. Therefore, it is necessary to consider the privacy protection when designing the crowdsourcing mechanism so as to ensure the smooth execution of crowdsourcing tasks and protect the privacy data of participants. We not only introduce the concept and workflow of Crowdsourcing, but also introduce in detail the mathematical definition of differential privacy and implementation mechanisms. At the same time, we analyze the application of differential privacy in different crowdsourcing scenarios including protection of data provided by users, location information of users and bid information of users. By summarizing the existing research, we prospect the future research direction.

Key words: crowdsourcing calculation; privacy protection; differential privacy; mechanism design

0 引言

在经济全球化的发展进程中,企业将会面临越来越大的竞争压力。对于那些企业难以单独完成的任务,众包可以聚集海量的互联网用户,使它们通过合作或者独立的方式来完成,同时给予完成者一定的报酬奖励。然而在众包平台的运行中,平台、任务请求者和参与者之间常常需要进行信息交换,如果处理不当,可能会引起用户敏感信息的泄露,从而给用户的信息安

全造成意想不到的伤害。

为了解决这一问题,人们相继提出了多种隐私保护模型,如 k -anonymity^[1]、 (a, k) -anonymity^[2] 和 l -diversity^[3] 等。但是这些模型在面对一些新型的、有针对性的攻击手段时还需要对自身进行不断改进,才能继续提供安全保障。而且,这些模型也没有对其处理结果的隐私保护水平进行定量分析,即当模型的参数发生变化时,无从得知其隐私保护水平的变化情况。

收稿日期: 2018-01-12

修回日期: 2018-05-17

网络出版时间: 2018-09-21

基金项目: 江苏省科技项目-基础研究计划(BK20141209)

作者简介: 刘欢(1992-),男,硕士研究生,CCF会员(E200088349G),研究方向为移动众包和保护隐私的数据管理;吴桂兴,博士,高级工程师,研究方向为自然语言处理、多媒体信号处理与信息安全。

网络出版地址: <http://kns.cnki.net/kcms/detail/61.1450.TP.20180920.1535.006.html>

差分隐私(differential privacy, DP)的出现则完美解决了上述担忧。所以差分隐私一经推出便迅速得到了学术界的认可,并被广泛应用到了其他领域,产生了一系列新成果。

一方面,对众包的定义、工作流程以及差分隐私的数学定义和保护水平评价标准进行了介绍;另一方面,还对不同众包应用场景下,如何使用差分隐私保护参与者个人数据的安全,以及该技术路线如何实现等进行了阐述。接着在这两方面的基础上,对不同众包场景下差分隐私的应用情况进行了总结,寻找改进空间,并指出了未来的研究方向。

1 众包的基本定义和工作流程

1.1 基本定义

2006年6月,美国杂志《连线》的记者Jeff Howe提出了众包的概念^[4]。

定义1(众包):众包是一种公开面向互联网大众的分布式的问题解决机制,它通过整合计算机和互联网上未知的大众来完成计算机难以单独完成的任务^[5]。

1.2 工作流程

一般来说,在众包应用中通常存在两种角色:任务请求人(requester)和工人(worker)。

当任务请求人需要通过众包来完成自己的任务时,他可以执行以下操作:设计任务;将任务发布到众包平台上;接收/拒绝工人的答案;整合工人们提交的答案,获得最终的结果。而工人则需要执行以下步骤:选择自己感兴趣的任務;接受任务;执行任务;提交答案。根据每个步骤执行时间的先后顺序,可以将众包的工作流程大致分为三个阶段:任务准备、任务执行以及整合答案。

2 差分隐私

2.1 基本定义

2006年,为了解决统计数据库领域内的隐私泄露问题,Dwork提出了一种全新的隐私定义^[6]。

定义2(差分隐私):设有一随机算法 M , M 所有可能产生的结果组成的集合为 P_M , S_M 则是 P_M 的子集, D 和 D' 是两个邻近数据集。那么对于任意的 D 、 D' 和 S_M ,如果有 $\Pr[M(D) \in S_M] \leq \exp(\epsilon) \times \Pr[M(D') \in S_M]$,则称算法 M 提供 ϵ -差分隐私保护,其中参数 ϵ 被称为隐私保护预算。一般来说, ϵ 越小,意味着隐私保护水平越高。

在文献[7]中,通过对输出结果进行随机化处理,算法 M 对隐私提供了保护。与此同时,参数 ϵ 保证了当数据集中某一个记录发生变化时,算法输出相同结

果的概率差被控制在一个极小的范围内,而这个范围是可以被数据所有者接受的。

2.2 实现机制

差分隐私的实现机制主要有Laplace机制(Laplace mechanism)^[8]和指数机制(exponential mechanism)^[9]。其中,Laplace机制是对数值型结果进行保护,而指数机制则是对非数值型结果提供保护。

2.2.1 Laplace 机制

定义3(Laplace机制):对于给定的数据集 D ,假设有函数 $f: D \rightarrow R^d$,它的敏感度为 Δf ,则随机算法 $M(D) = f(D) + Y$ 提供 ϵ 水平的差分隐私保护。其中, $Y \sim \text{Lap}(\Delta f/\epsilon)$ 是随机噪声,并且服从尺度参数为 $\Delta f/\epsilon$ 的Laplace分布。其中 ϵ 越小,说明引入的噪声越大。

2.2.2 指数机制

在很多应用情况下,查询结果并不仅仅只是数值型的,也可能是实体对象型的(比如一种方案或选择),所以McSherry等提出了指数机制。

设Range是查询函数的输出域,对于任意 $r \in \text{Range}$,均为实体对象。在指数机制中,输出值 r 的可用性函数可表示为 $q(D, r) \rightarrow R$,可以用来对 r 的优劣程度进行评估。

定义4(指数机制):记数据集 D 为随机算法 M 的输入,实体对象 $r \in \text{Range}$ 为其输出, $q(D, r)$ 为可用性函数, Δq 作为函数 $q(D, r)$ 的敏感度。如果算法 M 以正比于 $\exp(\epsilon q(D, r)/(2\Delta q))$ 的概率从Range中选择并输出 r ,那么称算法 M 提供 ϵ 水平的差分隐私保护。

3 差分隐私在保护用户提供的数据中的应用

文献[10]提出了一个和以前的研究不同的模式,即假设数据收集者是不可信任的。这样一来,目标群体就可以完全掌控自己的数据隐私,只提交一个经过隐私处理的版本给数据收集者即可。同时,数据收集者也不再负有保护这些隐私数据的责任,消除了隐私泄露带来的风险。

3.1 基本设定

如图1所示,假设数据收集者最后想了解的信息是一个二进制随机变量 W 的状态,同时每个个体都对 W 有一个信号 S_i ,即隐私数据。由于知识水平的限制,每个个体的隐私数据 S_i 和 W 的真实数据之间相等的概率为 θ , θ 的范围一般为 $0.5 \leq \theta \leq 1$ 。为了保护隐私,每个个体只向数据收集者报告一个经过隐私处理的版本 X_i 。当使用差分隐私方法生成 S_i 时,产生的隐私损失是 ϵ ,个体的隐私损失成本是 ϵ 的函数,而 ϵ

水平隐私数据的价值则被记作 $V(\varepsilon)$ 。

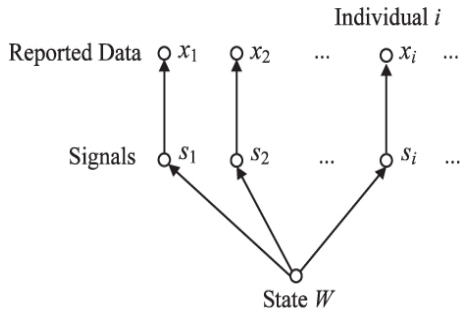


图 1 模型的信息结构

文献针对 $V(\varepsilon)$ 提出并证明了其下界和上界,还阐明了这些界限在报酬准确性问题里的应用,即数据收集者在保证 W 的准确性目标的情况下,如何最小化总报酬。

3.2 $V(\varepsilon)$ 的下界

首先,定义了 $V(\varepsilon)$ 的一个下界:

$$V_{LB}(\varepsilon) = g(\varepsilon) \frac{e^\varepsilon + 1}{e^\varepsilon} \left[\frac{\theta}{2\theta - 1} (e^\varepsilon + 1) - 1 \right] \quad (1)$$

其中, g 是个体隐私花费函数的一个变形。

经过证明,对任意 $\varepsilon > 0$,都有 ε 水平的隐私价值 $V(\varepsilon) \geq V_{LB}(\varepsilon)$ 。

3.3 $V(\varepsilon)$ 的上界

为了方便,定义了:

$$d = \frac{1}{2} \ln \frac{(e^\varepsilon + 1)^2}{4(\theta e^\varepsilon + 1 - \theta)((1 - \theta)e^\varepsilon + \theta)} \quad (2)$$

其中,对任意 $\varepsilon > 0$,都有 $d > 0$ 。

接着,给出了 $V(\varepsilon)$ 的上界: $V(\varepsilon) \leq V_{LB}(\varepsilon) + O(e^{-Nd})$,其中 $N \rightarrow \infty$ 。同样,文献给出了证明。

3.4 报酬-准确性问题

在得到上面关于隐私价值函数 $V(\varepsilon)$ 的上下界之后,还可以将其应用到报酬-准确性问题中去。

在给出准确性目标 τ 之后,可以用 $F(\tau)$ 来表示为了实现 τ 而要付出的最小总报酬。因为已经有了隐私价值函数 $V(\varepsilon)$ 的上下界,所以也能轻易得出 $F(\tau)$ 的上下界。

4 差分隐私在保护用户位置中的应用

在传统的移动群智应用中,有些任务是需要工人到达指定地点才能完成的,因此,为了实现最优的任务分配,比如使工人的移动距离最小化,平台往往需要事先知晓工人位置的准确位置。但是位置信息是工人个人隐私,暴露给平台会引起人们的安全担忧,同时,并不是每个工人最终都会被分配任务,这样他们的位置隐私也就白白泄露了。如何在任务分配的过程中既注意保护工人的位置隐私,又能使目标函数实现最优化就成为一个需要解决的难题。

文献[11]针对这个问题提出了一个解决框架,在此框架下,允许工人们将自己的位置信息模糊化之后再提供给众包平台。这样一来不仅满足了差分隐私,而且不用牵涉到任何第三方。接下来,将依次介绍该框架的工作流程、差分隐私保护、目标函数以及实验结果。

4.1 工作流程

该平台运行过程包含三个步骤:平台端生成地理位置模糊函数;用户端对位置信息进行模糊化;平台端分析模糊结果,并进行任务分配。

4.2 差分隐私保护

差分隐私是由 Andres 等介绍到位置保护中去的^[12]。其中,最为重要的就是模糊位置信息的概率函数 P ,它的基本思想是:对于模糊后的位置 l^* 将任意两个位置 l_1, l_2 映射到 l^* 的概率都是相同的。这样一来,即使攻击者知道了模糊函数 P ,并且观察到一个工人 u 处于位置 l^* ,也无法分辨出他的真实位置到底是 l_1 还是 l_2 ,从而实现了保护工人位置隐私的目标。

定义 5(差分隐私):假设目标区域包含一个位置集合 L ,那么当且仅当对 $\forall l_1, l_2, l^* \in L$,都有 $P(l^* | l_1) \leq e^{d(l_1, l_2)} P(l^* | l_2)$ 成立时,称概率函数 P 满足差分隐私保护。其中, $P(l^* | l)$ 是将 l 模糊到 l^* 的概率, $d(l_1, l_2)$ 是 l_1, l_2 之间的距离, ε 是隐私预算(ε 越小,表示隐私保护水平越高)。

4.3 目标函数

这部分将主要介绍的是位置模糊函数的求解和任务分配。

4.3.1 位置模糊函数

简单来说,在生成位置模糊函数 P 时,需要考虑的问题是:(1)使被选工人的移动距离最小化;(2) P 同时满足差分隐私。

在把位于 l_i 的一个任务分配给位于 l^* (模糊化后的位置)的工人之后,该工人预期需要移动的距离是:

$$d^*(l^* | l_i) = \frac{\sum_{l \in L} \pi(l) P(l^* | l) d(l, l_i)}{\sum_{l \in L} \pi(l) P(l^* | l)} \quad (3)$$

其中, π 是关注区域 L 中工人整体位置分布,满足 $\sum_{l \in L} \pi(l) = 1$ 。

假设 $x(l^* | l_i)$ 是把位于 l_i 的任务分配给位于 l^* 的工人的任务数量,那么就可以得出所有被选工人总的移动距离:

$$\sum_{l^* \in L} \sum_{l_i \in L} d^*(l^* | l_i) x(l^* | l_i) = \sum_{l^* \in L} \sum_{l_i \in L} \frac{\sum_{l \in L} \pi(l) P(l^* | l) d(l, l_i)}{\sum_{l \in L} \pi(l) P(l^* | l)} x(l^* | l_i) \quad (4)$$

4.3.2 任务分配

在工人们将自己模糊化后的位置信息上传到平台之后,接着就需要进行任务分配了,也就是求出 x 。

$$\begin{aligned} \min_x \quad & \sum_{l^* \in L} \sum_{l_i \in L} \frac{\sum_{l \in L} \pi(l) P(l^* | l) d(l, l_i)}{\sum_{l \in L} \pi(l) P(l^* | l)} x(l^*, l_i) \\ \text{s. t.} \quad & \sum_{l^* \in L} x(l^*, l_i) = N_i(l_i) \quad l_i \in L \\ & \sum_{l_i \in L} x(l^*, l_i) = N_c(l^*) \quad l^* \in L \\ & x(l^*, l_i) \in Z_{\geq 0} \quad l^*, l_i \in L \end{aligned} \quad (5)$$

其中, $N_c(l^*)$ 是处于位置 l^* 的实际工人数量。

接下来则可以使用 Benders 分解^[13]和遗传算法 (genetic algorithm, GA)^[14]对问题进行求解。

4.4 实验结果

文献[11]既在模拟数据集上,又在真实数据集 D4D 上进行了实验。从实验结果可以看出,在各种参数设定下,该文献算法的效果都要优于对比算法。

5 差分隐私在保护用户出价中的应用

文献[15]提出了一种差分隐私化的激励机制,以保护工人们的出价隐私。接下来,将详细介绍该机制的实现过程。

5.1 工作流程

(1) 众包平台首先向工人们公布任务集合 T ;

(2) 在拍卖阶段,工人需要向平台提交他的出价信息 $b_i = (\Gamma_i, \rho_i)$ 。其中, Γ_i 是工人感兴趣的任务集合, ρ_i 是他执行这些任务的出价;

(3) 根据工人们的出价信息,众包平台最终决定任务的分配,以及最终给工人们的报酬 p_i ;

(4) 当被分配到任务的工人执行完以后,需要将结果反馈给平台,平台对所有的结果进行处理整合,得到最终的任务结果,并付给工人报酬。

5.2 hSRC Auction

文献[15]中定义了一种 hSRC 拍卖机制。在该机制中,任意工人 w_i 都有 K_i 种可能的出价集合,记作

$T_i = \{\Gamma_{i,1}, \Gamma_{i,2}, \dots, \Gamma_{i,K_i}\}$, 而工人只对其中的一个真正感兴趣,记作 Γ_i^* , c_i^* 是工人为了完成这些任务所要付出的花费。

5.3 差分隐私保护

文献将提出的 hSRC 拍卖机制看作将出价信息 b 映射到报酬 P 的函数 $M(\cdot)$, 当且仅当 M 满足如下条件时,称 M 是符合 ε -差分隐私的:

$$\Pr[M(b) \in A] \leq \exp(\varepsilon) \Pr[M(b') \in A]$$

其中, A 是报酬集合; b 和 b' 是只相差一个出价的出价集合; ε 称为隐私预算,其值越小,表示隐私保护水平越高。

这样一来,当某个工人的出价集合 b 发生变化时,其最后对应的报酬集合并不会发生明显的变化,从而确保了某些好奇的工人无法从报酬集合上推断出该工人的出价信息,就实现了保护工人出价隐私的目标。

5.4 机制设计

在付给工人酬劳时,将使平台要支付的总酬劳最小化称为 TPM (total payment minimization) 问题。同时,文献证明了 TPM 是 NP 难问题。所以,无法在多项式时间内计算出使总酬劳最小的那个任务分配集。针对这种情况,文献设计了一种机制 DP-hSRC,能在多项式时间内得到一定比率的近似最优总酬劳,而且还考虑了保护工人出价隐私的目标。而且证明,该机制满足 ε -差分隐私保护、近似真实性和个体合理性等性质,时间复杂度被控制在了 $O(N^2 * K)$ 。

5.5 实验结果

5.5.1 对比算法

Optimal: 最优总报酬。

Baseline Auction: 实现方法略有不同的对比算法。

5.5.2 实验参数设定

实验设置了四种不同的参数设定,正如表 1 所示^[15]。

5.5.3 实验结果

最终的实验结果表明,提出的机制 DP-hSRC 在总酬劳上已经很接近最优算法,同时效果也比 Baseline Auction 好很多。

表 1 仿真参数设定

Setting	ε	c_{\min}	c_{\max}	Γ_i^*	$\theta_{i,j}$	δ_j	N	K
I	0.1	10	60	[10, 20]	[0.1, 0.9]	[0.1, 0.2]	[80, 140]	30
II	0.1	10	60	[10, 20]	[0.1, 0.9]	[0.1, 0.2]	120	[20, 50]
III	0.1	10	60	[50, 150]	[0.1, 0.9]	[0.1, 0.2]	[800, 1400]	200
IV	0.1	10	60	[50, 150]	[0.1, 0.9]	[0.1, 0.2]	1 000	[200, 500]

6 结束语

由于众包的工作流程牵涉到任务准备、任务执行和答案整合等一系列步骤,期间众包平台和参与者之间也需要进行多次的信息交换,所以有很多方面需要对参与者的隐私数据进行保护,也就是在应用场景方面还有很大的拓展空间。另一方面,在保护参与者隐私数据的同时,如何使系统的总体收益更加优化,也是一个可以开展的方向,比如在保护用户位置隐私的场景下,系统的总体目标函数还可以是传感数据质量最大化,也可以是总花费最小,也可以是工人的总体移动距离最小。

经过数年的发展,众包机制已经有了长足的进步,对参与者的隐私数据也有了相当的保护,但是仍然没有到达成熟的地步,所以未来还有大量的实际应用开发和优化工作有待完成。

参考文献:

- [1] SWEENEY L. K-anonymity: a model for protecting privacy [J]. International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems 2002, 10(5): 557-570.
- [2] WONG R C W, LI Jiuyong, FU A W C, et al. (α, k)-anonymity: an enhanced k -anonymity model for privacy preserving data publishing [C]//Proceedings of the 12th ACM SIGKDD international conference on knowledge discovery and data mining. Philadelphia, PA, USA: ACM 2006: 754-759.
- [3] MACHANAVAJJHALA A, GEHRKE J, KIFER D, et al. L-diversity: privacy beyond k -anonymity [C]//Proceedings of the 22nd international conference on data engineering. [s. l.]: IEEE 2006: 24.
- [4] HOWE J. The rise of crowdsourcing [J]. Wired Magazine, 2006, 14(6): 1-4.
- [5] 冯剑红, 李国良, 冯建华. 众包技术研究综述 [J]. 计算机学报 2015, 38(9): 1713-1726.
- [6] DWORK C. Differential privacy: a survey of results [C]//5th international conference on theory and applications of models of computation. Xi'an, China [s. n.]: 2008: 1-19.
- [7] 熊平, 朱天清, 王晓峰. 差分隐私保护及其应用 [J]. 计算机学报 2014, 37(1): 101-122.
- [8] DWORK C, MCSHERRY F, NISSIM K, et al. Calibrating noise to sensitivity in private data analysis [C]//Proceedings of the third conference on theory of cryptography. New York, NY: Springer 2006: 265-284.
- [9] MCSHERRY F, TALWAR K. Mechanism design via differential privacy [C]//48th annual IEEE symposium on foundations of computer science. Providence, RI, USA: IEEE 2007: 94-103.
- [10] WANG Weina, YING Lei, ZHANG Junshan. The value of privacy: strategic data subjects, incentive mechanisms and fundamental limits [J]. ACM SIGMETRICS Performance Evaluation Review 2016, 44(1): 249-260.
- [11] WANG Leye, YANG Dingqi, HAN Xiao, et al. Location privacy-preserving task allocation for mobile crowdsensing with differential geo-obfuscation [C]//Proceedings of the 26th international conference on world wide web. [s. l.]: International World Wide Web Conferences Steering Committee 2017: 627-636.
- [12] ANDRÉS M E, BORDENABE N E, CHATZIKOKOLAKIS K, et al. Geo-indistinguishability: differential privacy for location-based systems [C]//Proceedings of the 2013 ACM SIGSAC conference on computer & communications security. [s. l.]: ACM 2013: 901-914.
- [13] BENDERS J F. Partitioning procedures for solving mixed-variables programming problems [J]. Numerische Mathematik, 1962, 4(1): 238-252.
- [14] MITCHELL M. An introduction to genetic algorithms [M]. Cambridge: MIT Press 1998.
- [15] JIN Haiming, SU Lu, DING Bolin, et al. Enabling privacy-preserving incentives for mobile crowd sensing systems [C]//IEEE 36th international conference on distributed computing systems. Nara, Japan: IEEE 2016: 344-353.