

WSS President's Invited Seminar
May 1st, 2019

**Statistical Disclosure Limitation and
Differential Privacy**

Natalie Shlomo
Professor of Social Statistics



Topics

- Overview of types of disclosure risk in traditional forms of statistical data
- Common statistical disclosure limitation methods
- Disclosure risk-data utility paradigm
- Inferential disclosure and differential privacy
- New dissemination strategies:
 - Online flexible table builder
 - Microdata and Synthetic Data
 - Remote access and remote analysis
- Discussion

Traditional Statistical Outputs

- **Survey Microdata**
 - Social surveys (census/register and business survey microdata generally not released)
 - Available from data archives for registered users
- **Tabular Data**

Frequency Tables

Census/register
(whole population) counts

Weighted sample counts

Magnitude Tables

Business Statistics,
eg., total turnover

Types of Disclosure Risks

Identity Disclosure

Identification is widely referred to in confidentiality pledges and code of practice

Individual Attribute Disclosure

Confidential information about a data subject is revealed and can be attributed to the subject (Identity disclosure a necessary pre- condition)

Group Attribute Disclosure

Confidential information is learnt about a group and may cause harm

Common SDL Methods

Social Survey Microdata

Identity Disclosure (assume no response knowledge)-
rare categories of identifying variables (population unique)

Attribute disclosure - individual(s) identified and survey target variables learnt, eg. health, income

Recoding/grouping identifying variables, eg. k-anonymity

Suppressing variables such as high level geographies

Sub-sampling, eg. census samples

Top-coding sensitive variables

Recoding / Microaggregation, eg. l-diversity

Common SDL Methods

Frequency Tables (whole population counts)

Identity Disclosure –small cells

Table design, eg. spanning variables and grouped categories

Minimum population thresholds

Attribute disclosure - zeros in row/column and one populated cell

Pre-tabular and/or post-tabular perturbation to introduce ambiguity in zero cells

Nested tables to avoid disclosure by differencing

Common SDL Methods

Magnitude Tables (Business statistics)

Assumptions:

- Intruders are competitors in the cell and can form coalitions
- Businesses in a cell are known
- The ranking of the businesses with respect to their size is known

Attribute disclosure - What can a competitor learn with sufficient precision

Table design

Minimum population thresholds

Cell suppression: primary and secondary

(mathematical programming and optimization)

Disclosure Risk and Data Utility

Disclosure risk

Frequency tables:

Whole population counts and disclosure risk is visible: small cells, placement of zero cells

Let $F = \{F_1, F_2, \dots, F_K\}$

$$H\left(\frac{F}{N}\right) = -\sum_k \frac{F_k}{N} \log\left(\frac{F_k}{N}\right) \text{ and } 1 - \left[\frac{H\left(\frac{F}{N}\right)}{\log K}\right]$$

Microdata:

Set of cross-classified quasi-identifiers defined by $k=1, \dots, K$

$$\sum_k I(f_k = 1, F_k = 1)$$

where

f_k sample count

F_k population count

Probabilistic modelling for estimation: Poisson-log linear modelling

Magnitude tables (Business statistics):

Let $T_k = \sum_{i \in k} x_i$ in cell k

(n, p) Dominance Rule classifies cell as disclosive if

$$x_{(1)} + \dots + x_{(n)} \geq (p/100) \times T_k$$

Disclosure Risk and Data Utility

Utility

- Impact on variance
- Impact on bias

Distortions to distributions:
distance metrics, eg.
Hellinger's Distance*,
variation in propensity scores

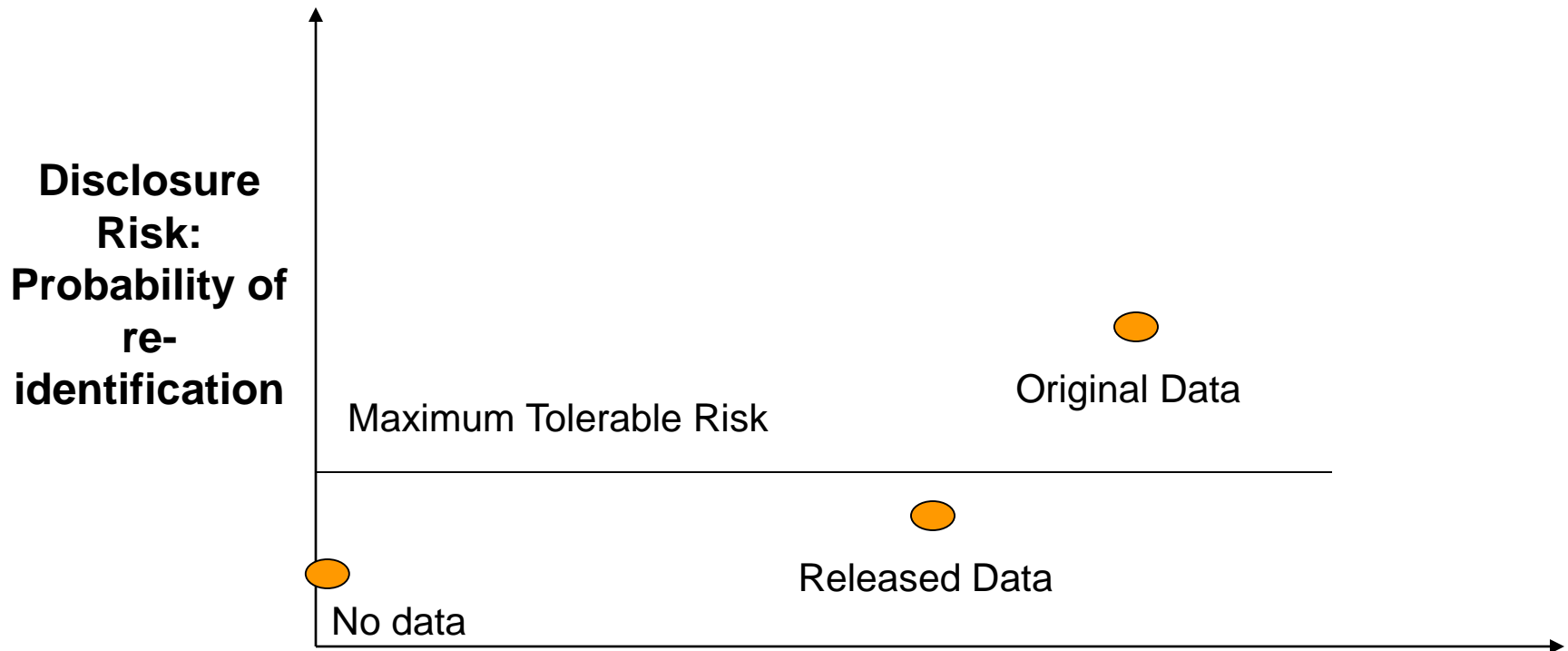
Changes in inference:
confidence interval overlap,
change in χ^2 or R^2

Changes in associations:
change in correlations and
rankings, Cramer's V

$$* HD(F, F') = \sqrt{\frac{1}{2} \sum_{k=1}^K (\sqrt{F_k} - \sqrt{F'_k})^2}$$

Disclosure Risk and Data Utility

R-U Confidentiality Map (Duncan, et.al. 2001)



Data Utility: Quantitative measure on the statistical quality

Inferential Disclosure

Confidential information may be revealed exactly or to a close approximation with high confidence from statistical properties of released and combined data

Examples:

Survey microdata – a good prediction model with very high R^2

Census tables – disclosure by differencing and linking tables

This type of disclosure has largely been ignored and dealt with through strict control of data that is released

- Microdata deposited in archives for registered users
- Strict control of tabular data, eg. review boards for special request tabulations

Where do we go from here?

- Traditional forms of statistical data and their confidentiality protection rely heavily on assumptions that may no longer be relevant

Digitalization of all aspects of our society leading to new and linked data sources offering opportunities for research and evidence-based policies



With detailed personal information easily accessible from the internet, traditional SDL may no longer be sufficient and agencies relying more on restricting and licensing data

- Growing demand for more open and accessible data via web-based applications
- Need for more rigorous data protection mechanisms with stricter privacy guarantees
- Collaborations with computer scientists through scientific programs

Differential Privacy

- Computer Science **differential privacy** (Dwork and Roth 2014): the intruder has knowledge of entire database except for one target unit (“worst case” scenario)

Definition: Mechanism M satisfies (ϵ, δ) -differential privacy if for all neighbouring databases D, D' differing by one individual, all possible queries q and $S \subseteq \text{Range}(M)$ all possible outputs:

$$P(M(q(D)) \in S) \leq e^\epsilon P(M(q(D')) \in S) + \delta$$

and the probability is taken over the randomness of the mechanism

If $\delta = 0$ then we have ϵ -differential privacy

Example of Differential Privacy Mechanism

Laplace Mechanism

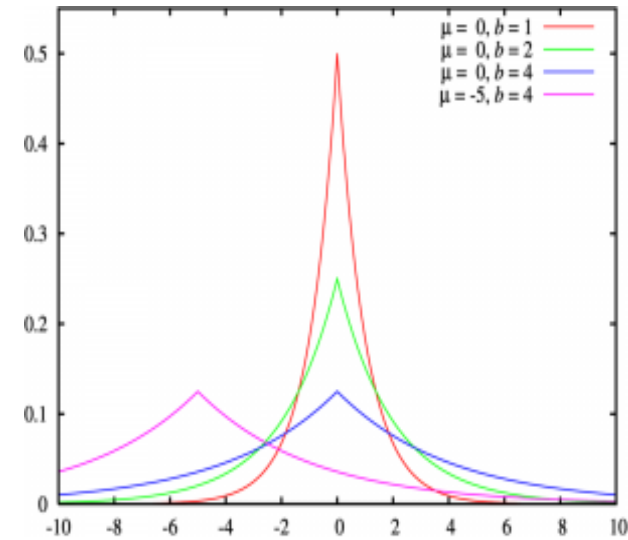
Calibrating noise: what scale of noise b is large enough to ensure privacy on a query q ?

$q(D)+Z$ and Z sampled from $\text{Lap}(0,b)$

Amount of noise depends on ϵ and sensitivity of q denoted Δq

$$\Delta q = \max_{D,D'} |q(D) - q(D')|$$

where D, D' any neighbouring databases



Example:

query	Δq
count	1
max(age)	120
avg(age)	120/n

Theorem: setting scale (b) of Laplace noise to $\Delta q/\epsilon$ ensures ϵ -differential privacy

Mechanisms in Differential Privacy

Non-interactive Mechanism

Data custodian produces a 'safe' object, such as a synthetic database or collection of summary statistics

After this *release* all post-perturbative analyses are safe (no privacy budget spent after the original object)

Interactive Mechanisms

Data analyst sends queries (functions applied to a database) adaptively, deciding which query to pose next based on observed responses to previous queries

Accuracy will deteriorate with the number of questions asked, and providing accurate answers to all possible questions will be infeasible

Differential Privacy in the SDL Tool-kit at Statistical Agencies

Non-interactive mechanisms as agencies unable to monitor queries

DP useful when perturbative methods are needed with stricter privacy guarantees such as outputs disseminated via the internet where agencies relinquish control of the releases

Examples: flexible table builder, synthetic data, and multiple data products released from survey microdata

Agencies should still maintain 'safe access' as an SDL approach via RDCs for 'trusted' users

Differential Privacy vs. SDL

No distinction between key variables and sensitive variables, types of disclosure risks, sample or population or prior intruder knowledge

Designed for output perturbation and in this case a sum/average is disclosive and needs to be protected (same as disclosure by differencing)

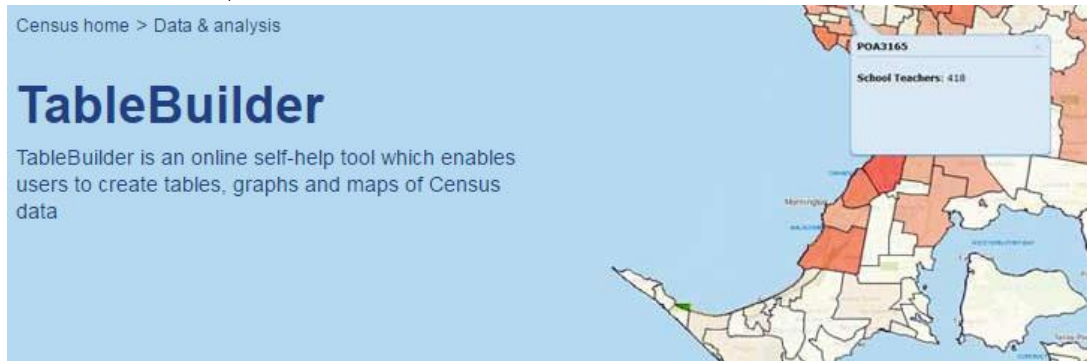
Zeros need to be perturbed

Perturbation mechanism not hidden and can be used to correct statistical analysis

Dissemination Strategies for Open Data

Online Flexible Table Builder

- Increasing demands for online dissemination and open access of census tables (ABS, USA, EU)



- Web-based platform (drop down lists) with restrictions:
 - number of dimensions, population thresholds, no sparse tables
- SDL on-the-fly: pre-tabular (hypercubes, swapping) and/or post-tabular methods (noise addition, rounding)
- Perturbation matrix $p_{ij} = P(\text{perturb cell to } j | \text{original cell is } i)$
- Change (or do not change) value according to p_{ij} and random draw

Online Flexible Table Builder

- Other principles in SDL:

Perturbations unbiased, bounded, maximal entropy, non-negative and zeros not perturbed

Microdata keys for same perturbations on same cells across tables (Fraser and Wooton 2005)

Additivity - probability perturbation matrix with property of 'invariance' (ensures margins in expectations) and IPF (Shlomo and Young 2008)

- Differential Privacy (DP) for flexible table builders (Rinott, O'Keefe, Shlomo and Skinner 2018)

Differential Privacy Algorithm

- List space $a = (a_1, \dots, a_k)$, eg. internal cells and margins (overlapping individuals) in a non-interactive mechanism
- Consider $M(\cdot)$ such that $M(a) = b = (b_1, \dots, b_k)$ where $p(b_k | a_k)$ set of conditional probabilities and cells perturbed independently (assume perturbed list has same structure as original list) and in our case, M is discrete
- Definition: $M(\cdot)$ satisfies **(ϵ, δ) -differential privacy** if for all neighbouring lists a, a' differing by one individual:

$$P(M(a) = b) \leq e^\epsilon P(M(a') = b) + \delta$$

and this is true for all potential lists and all possible outcomes

Differential Privacy Algorithm

- **Exponential Mechanism** (McSherry and Talwar 2007) - perturbation mechanism selected that produces values with high utility

- Define two loss functions:

$$l_1 = \sum_{i=1}^K |a_k - b_k| \quad (\text{motivated by discretized Laplace})$$

$$l_2 = \sum_{i=1}^K (a_k - b_k)^2 \quad (\text{motivated by discretized Normal})$$

Then define $u_i = -l_i, \quad i = 1, 2$

Perturb with probability proportional to $e^{\left(\frac{\epsilon}{2}\right)u/\Delta u}$

$$\Delta u = \max_{b \in B} \max_{a \sim a' \in A} |u(a, b) - u(a', b)|$$

No Margins: $\Delta u = 1$
t-way table margins: $\Delta u = 2^t - 1$

Bound the perturbations $|a_k - b_k| \leq m, \forall k$ leads to (ϵ, δ) - DP

Exponential Mechanism

Original Value	Range for $\epsilon=1.5$ and $\delta=0.00002$					Range for $\epsilon=0.5$ and $\delta=0.008$				
	± 0	± 1	± 2	± 3	± 4	± 0	± 1	± 2	± 3	± 4
	Laplace m=7					Laplace m=7				
0	0.82	0.96	0.99	1	1	0.63	0.78	0.87	0.93	0.96
1	0.64	0.96	0.99	1	1	0.25	0.78	0.87	0.93	0.96
2	0.64	0.92	0.99	1	1	0.25	0.55	0.87	0.93	0.96
3	0.64	0.92	0.98	1	1	0.25	0.55	0.74	0.93	0.96
4	0.64	0.92	0.98	1	1	0.25	0.55	0.74	0.85	0.96
≥ 5	0.64	0.92	0.98	1	1	0.25	0.55	0.74	0.85	0.92
	Normal m=12					Normal m=10				
0	0.57	0.70	0.81	0.89	0.94	0.54	0.63	0.71	0.78	0.84
1	0.14	0.70	0.81	0.89	0.94	0.09	0.62	0.71	0.78	0.84
2	0.14	0.40	0.81	0.89	0.94	0.09	0.26	0.71	0.78	0.84
3	0.14	0.40	0.62	0.89	0.94	0.09	0.26	0.42	0.78	0.84
4	0.14	0.40	0.62	0.78	0.94	0.09	0.26	0.42	0.57	0.84
≥ 5	0.14	0.40	0.62	0.78	0.88	0.09	0.26	0.42	0.57	0.69

*Negative values to 0

Exponential Mechanism

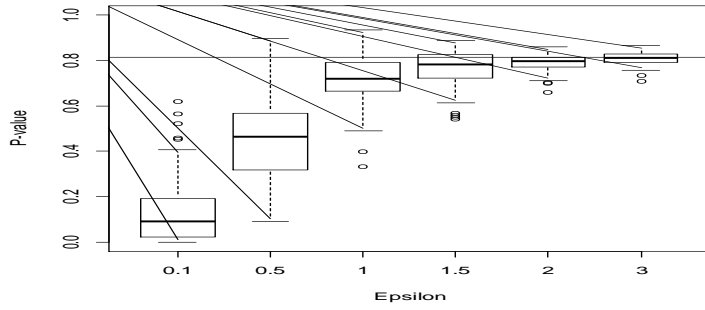
Implications:

- DP leads to negative values, setting to zero still ensures DP but biased perturbations
- All (non-structural) zeroes must be perturbed
- If list-space has internal cells only $\Delta u = 1$, margins summed from internal cells DP but low utility
- In a t -way table all margins, $\Delta u = 2^t - 1$ (not including total) much larger perturbations implying smaller utility
- Margins can be perturbed (with appropriate sensitivity) and prorated to ensure additivity (post-processing does not violate DP)

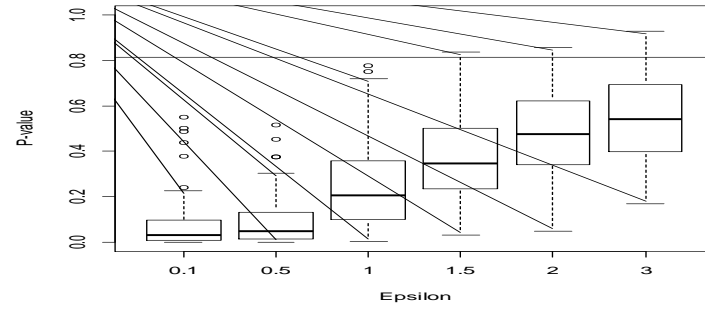
Parameters of Differential Privacy not secret and can be used to adjust statistical analysis

Generated independent table, N=10000, K=100 (average cell size=100)

Laplace Perturbations

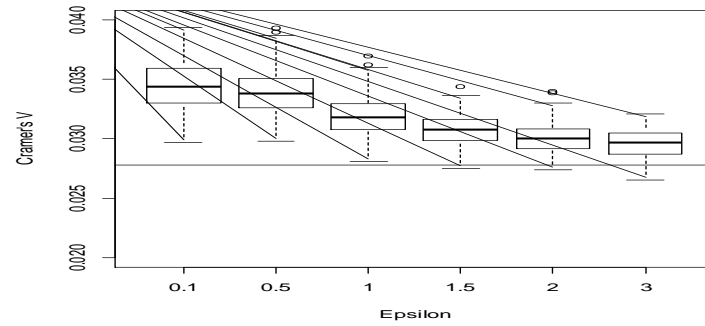
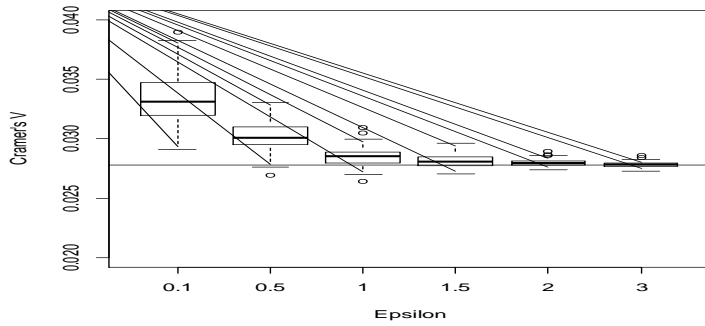


Normal Perturbations

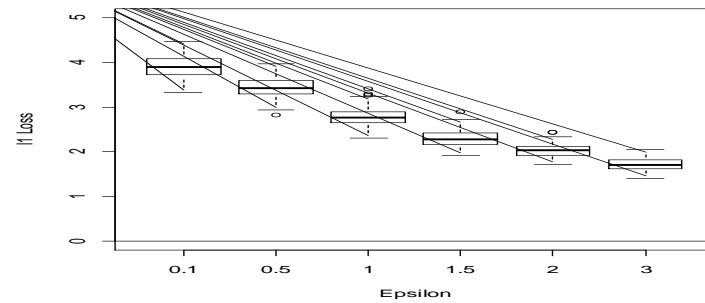
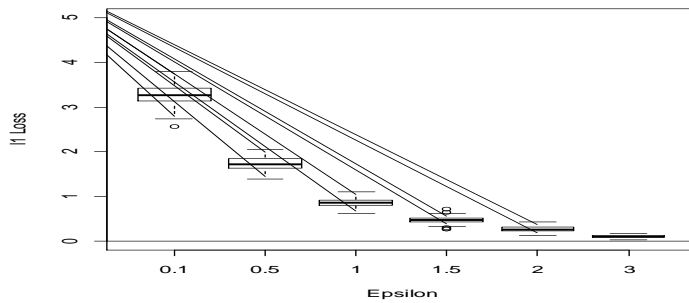


P-Value

Cramer's V

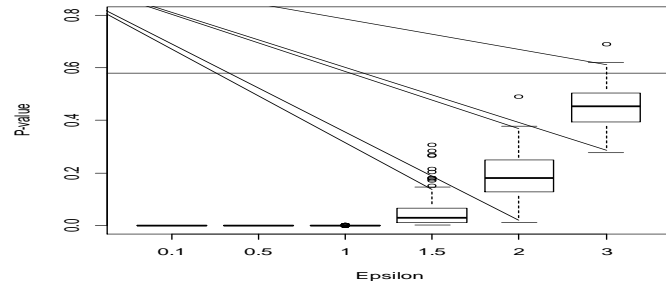


l_1 Loss Function

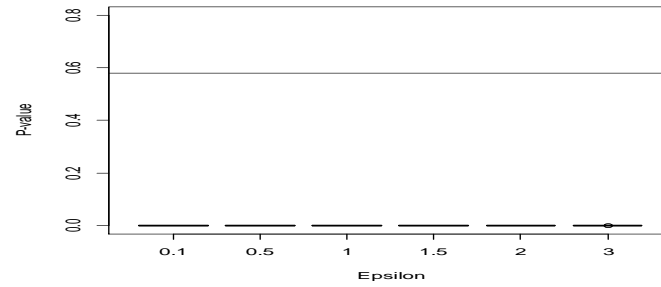


Generated independent table, N=10000, K=1000 (average cell size=10)

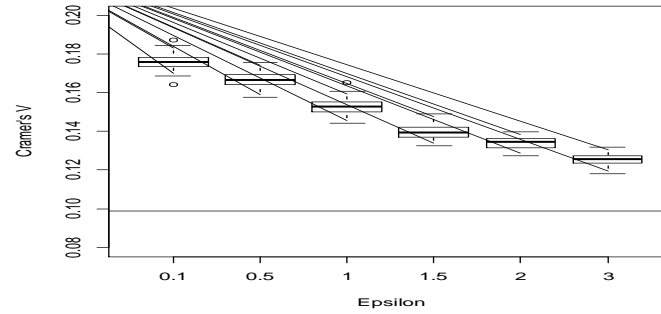
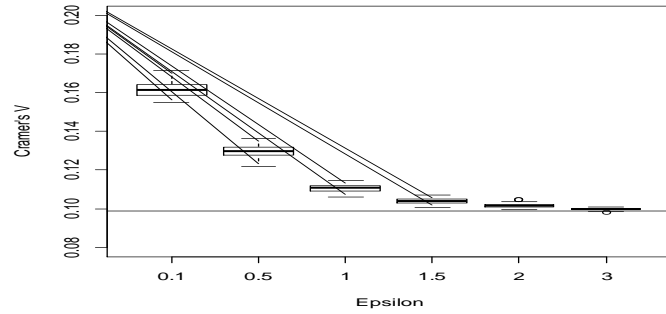
Laplace Perturbations



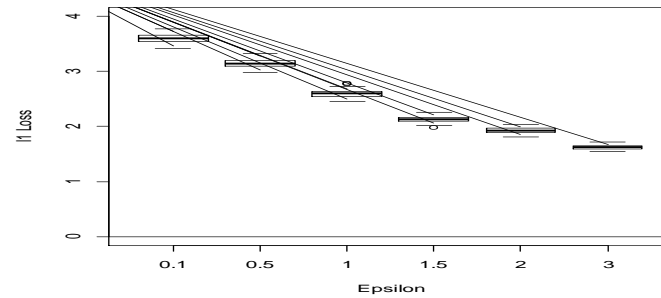
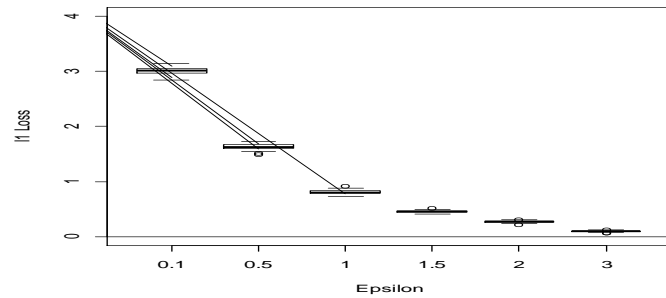
Normal Perturbations



Cramer's V



l_1 Loss Function

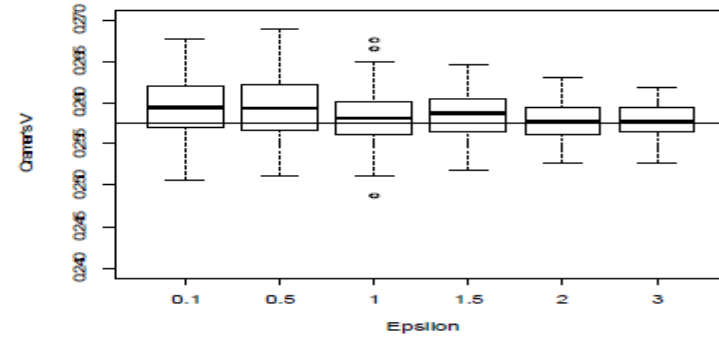
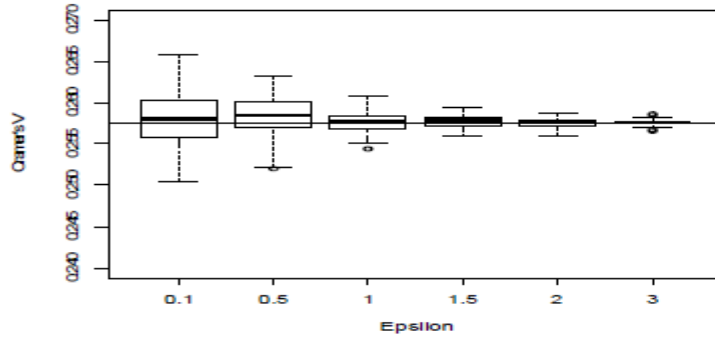


Real (dependent) Table from UK Census Data

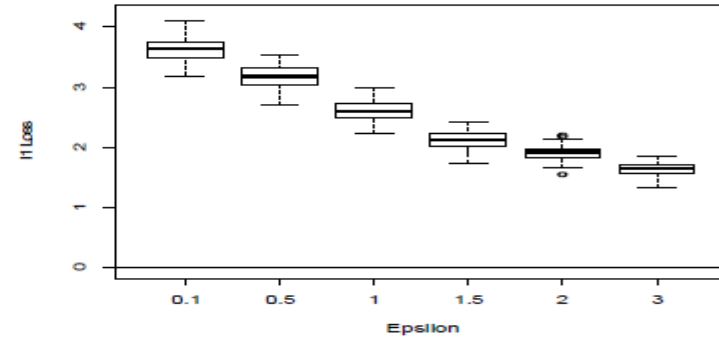
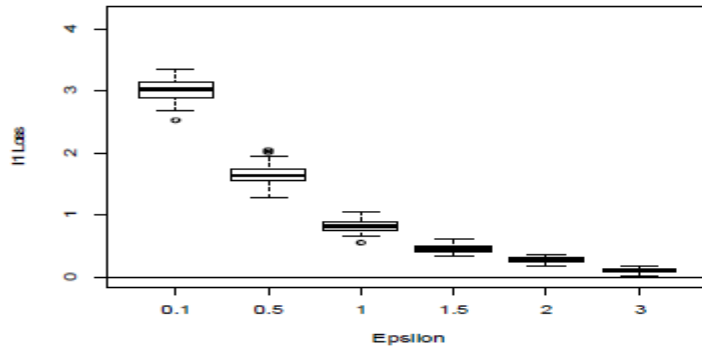
Laplace Perturbations

Cramer's V

Normal Perturbations



l_1 Loss Function



l_2 Loss Function

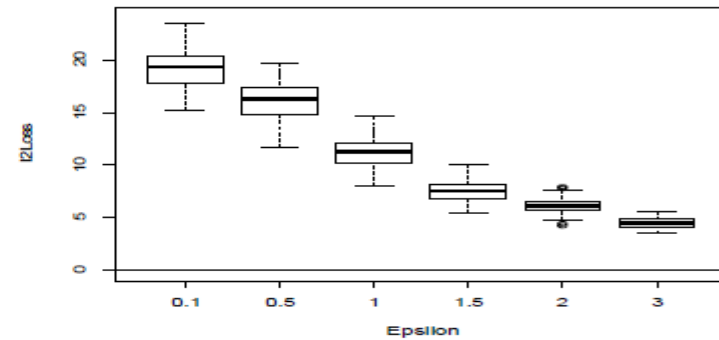
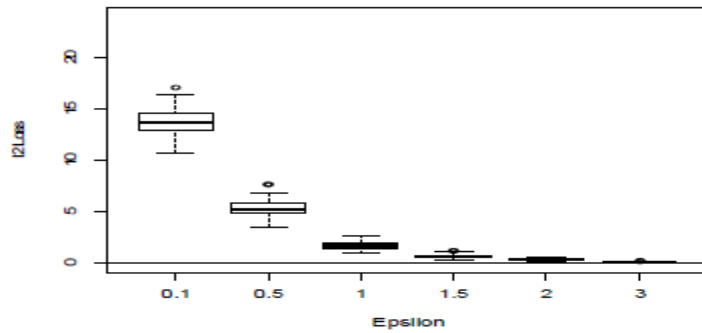


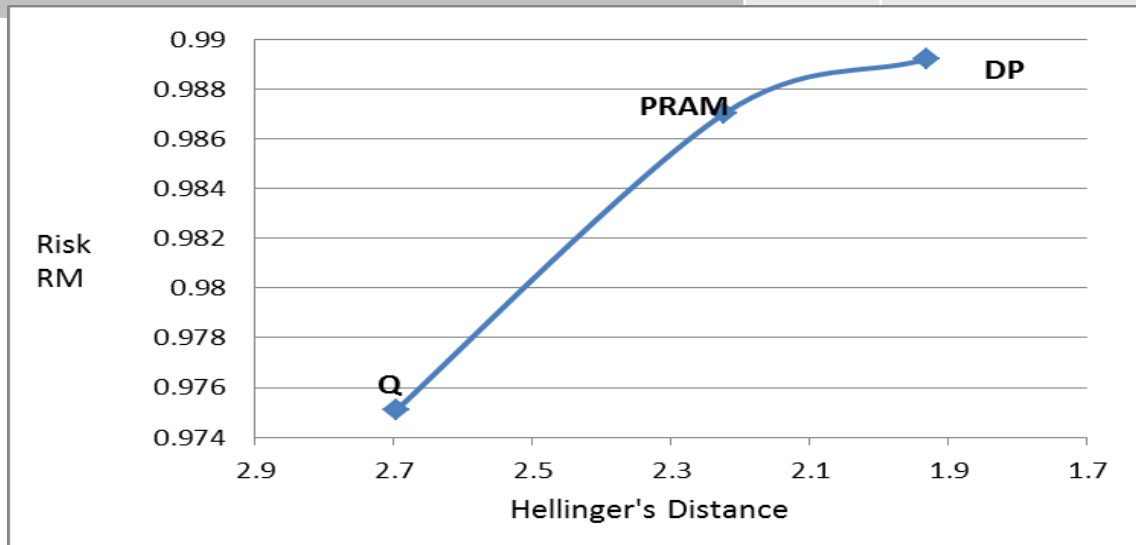
Table Builder for Survey Weighted Counts

- Inferential disclosure from multiple data products: tables from original data and public-use files (Shlomo, Krenzke and Li 2019)
 - $l_1 = \sum_{i=1}^K |a_k - b_k| a_k, b_k$ weighted survey counts and Δu is maximum weight
 - For survey weights with little variation ($CV < 10\%$) consider replacing weights by the average weight so that the exponential mechanism (for internal cells) reduces to $e^{-\frac{\epsilon}{2} \bar{w} \sum |a_k - b_k| / \bar{w}}$
 - Perturb sample counts, eg. add/subtract p from the sample count, and adjust weighted count by $p\bar{w}$

Table Builder for Survey Weighted Counts

Example: Compare DP with Post-randomisation and Drop/Add-up-to-q (Li and Krenzke 2016) on generated tables 7×7 and $\text{cap} \pm 7$

Risk Measures		DP	PRAM	Q
DP parameters when original sample count=1	ϵ	2	2	0.01
	δ	0.00000063	0.1192	0.333
Percent Cells Perturbed		23.8	27.7	67.1
1-Proportion of conditional entropy (RM) (Antal, et al., 2014)		0.9891	0.9870	0.9751



Y-axis: $RM = 1 - \left(\frac{H(a|b)}{H(a)} \right)$

X-axis: $HD(a, b)$ (reverse order)

Microdata

- Sampling (and other non-perturbative SDL methods) are not DP
 - ‘slippage’ parameter δ depends on $\tau_1 = \sum_k I(f_k = 1, F_k = 1)$ which is very small in social surveys
- Stochastic perturbation can be made DP if every record has a non-zero chance of being perturbed (Shlomo and Skinner 2012)

Synthetic Data

- Fit models from original data, eg. posterior predictive distributions
Can be implemented on parts of data where a mixture is obtained of real and synthetic data
- Draw and release several samples to account for the uncertainty and obtain 'proper' variance estimates (Reiter 2005)

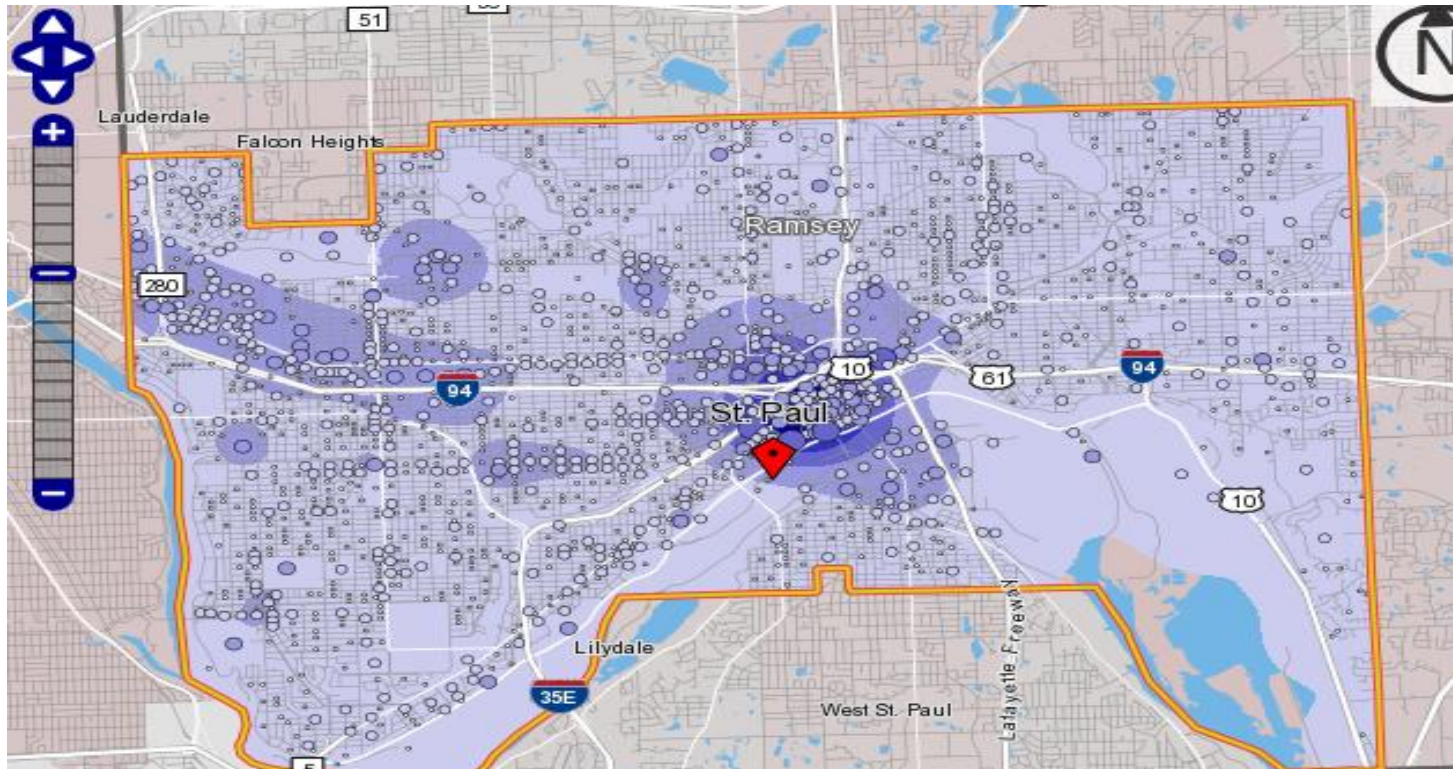
Techniques for synthetic magnitude tables

(CTA – controlled tabular adjustment) (Dandekar and Cox 2002)

- In practice, difficult to capture all conditional relationships between variables and within sub-populations
 - If models of interest are sub-models of the synthesis model, then the analysis of (multiple) synthetic samples should give valid inferences

Synthetic Data

- On the Map (Abowd and Vilhuber 2008) – one of the first applications of DP synthetic data
 - Simple trajectory where you live where you work
 - Dirichlet-Multinomial where a DP smoothing parameter was added to the Dirchlet hyperparameters



Differential Privacy for Synthetic Data

- Synthetic data

Ongoing Research:

- Bayesian Modeling with differentially private priors
- Current work on adding noise to estimating equations and also looking at ridge regression to regularize linear regression by adding a constraint to likelihood function: use in Sequential Regression modeling (Ragunathan et al. 2001)
- Reproducing microdata from differentially private counts
- Remote Analysis Servers

The Unlinkable Data Challenge: Advancing Methods in Differential Privacy

📌 Data Science, Government, Non-Profit & Social Impact, Technology

Propose a mechanism to enable the protection of personally identifiable information while maintaining a dataset's utility for analysis. [Read Overview...](#)

FOLLOW



STAGE

Submission Deadline



\$50,000

Research Data Centres

- Secure environment for trusted users, eg. Virtual Microdata Lab (VML)
- Minimise risk of disclosure:
 - No removal of data, no printers, no link to internet
 - All outputs checked manually by staff
 - Training course for understanding security rules



Remote Access

- Access to data through remote connection to secure server, typically at Universities and Research Institutes
- Carry out analysis as if on personal PC and view results on screen
- Outputs dropped in a folder to be manually checked and emailed back to researchers

Remote Analysis

- Initial research in developing platforms for remote analysis or allowing researchers to submit code
- Aim to protect outputs without the need for human intervention

Example (O'keefe and Shlomo 2012):

- 338 Sugar Canes Farm Data from a 1982 survey of sugar cane industry in Queensland, Australia: Region (4 categories) and 5 continuous variables: Area, Harvest, Receipts, Costs, Profits (=Receipts-Costs)
- Confidentializing Input:
 - 5 outliers removed resulting in 333 farms
 - Area (identifying variable) coarsened to 9 categories
 - Remaining continuous variables perturbed with multivariate random Gaussian noise within quintiles of receipts

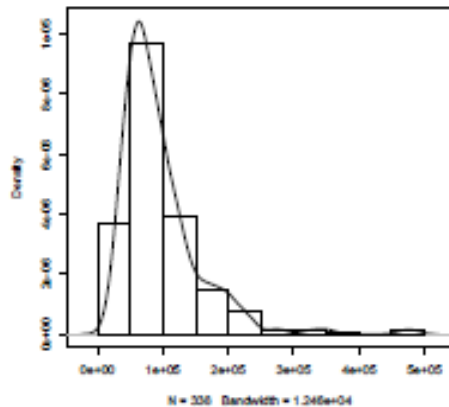
Remote Analysis

Receipts:

Original

	Receipts
Minimum	11703
1st Quartile	57607
Median	80391
Mean	95965
3rd Quartile	117062
Maximum	484346
Standard Deviation	61609.105256

(a) Summary Statistics

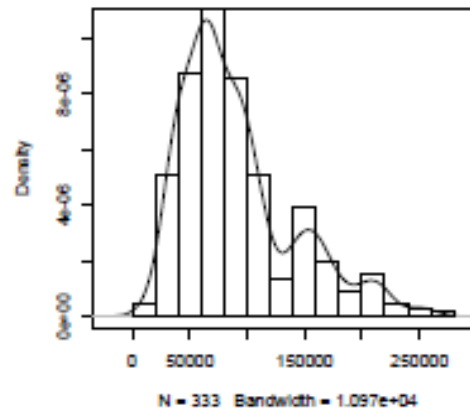


(c) Histogram and Density

Input

	Receipts
No. observations	333
Minimum	11140
1st Quartile	57473
Median	77144
Mean	90643
3rd Quartile	109637
Maximum	260098
Standard Deviation	49214.06

(a) Summary Statistics

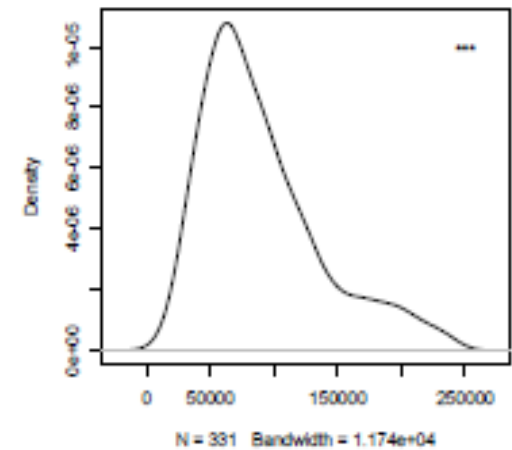


(c) Histogram and Density

Output

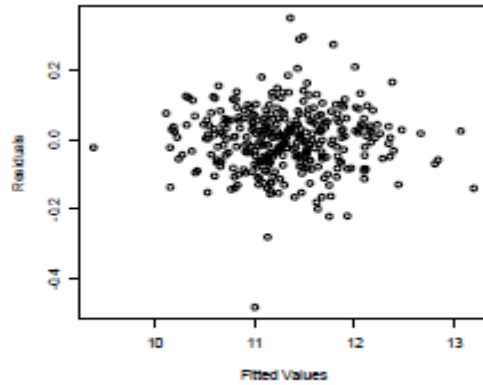
	Receipts
1st Quartile	57600
Median	80400
Mean	96000
3rd Quartile	117100
Standard Deviation	61600

(a) Confidentialised Summary Statistics

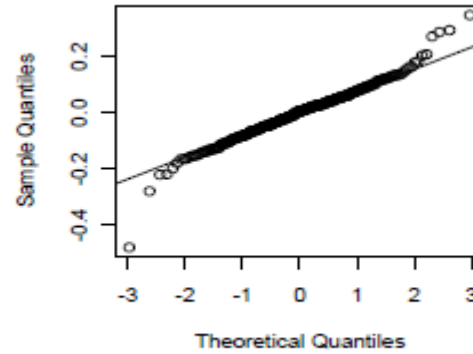


(c) Confidentialised Density Estimate

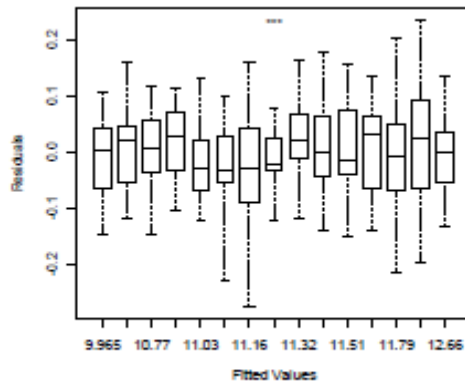
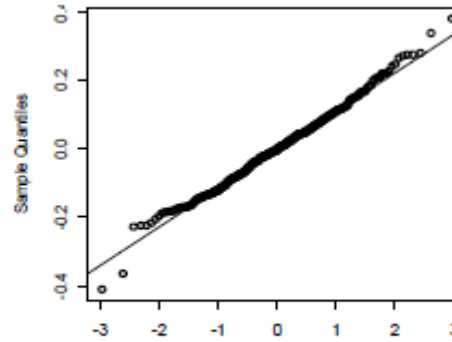
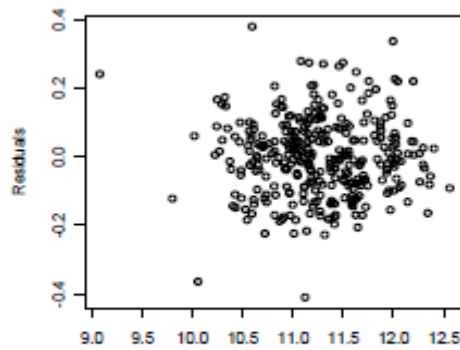
Remote Analysis



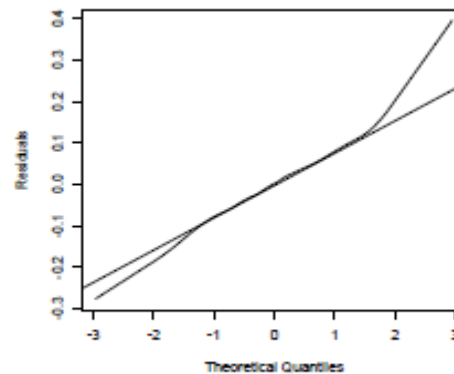
(a) Residuals vs Fitted Values



(b) Normal Q-Q Plot of Residuals



(a) Residuals by fitted values



(b) Normal QQ Plot of Residuals

Differential Privacy and Remote Analysis

- Remote analysis a natural extension to table builders
- All statistics, eg. means, medians, etc. can have DP noise added
- Regression models will typically use robust regression but use a functional DP mechanism, eg. adding noise to estimating equations, ridge regression
- Research on differentially private graphs, scatterplots

Remote analysis servers were once a topic of research
but lately been ignored
This research should be revived

Challenges and Discussion

Differential Privacy with formal privacy guarantees may provide solutions for SDL

Allows statistical agencies to consider new ways of disseminating open data via the internet

It provides a formal 'by-design' privacy guarantee against inferential disclosure

Combined with other SDL approaches of coarsening, subsampling, variable suppression etc. impacts on the privacy budget

Further research is needed to set these privacy budgets

Additive noise perturbation of DP can provide more utility than other additive SDL noise perturbations

Agencies should release parameters of the perturbation and DP parameters are not secret and can be used to adjust analyses

References

- Abowd, J.M. and Vilhuber, L., (2008). How Protective Are Synthetic Data? In *PSD'2008 Privacy in Statistical Databases*, (Eds. J.Domingo-Ferrer and Y. Saygin), Springer LNCS 5262, 239-246.
- Antal, L., Shlomo, N. and Elliot, M. (2014). Measuring Disclosure Risk with Entropy in Population Based Frequency Tables. In *Privacy in Statistical Databases 2014*, (Ed. J. Domingo-Ferrer), Springer LNCS 8744, pp. 62-78.
- Dandekar, R.A. and Cox L. H. (2002). Synthetic Tabular Data: An Alternative to Complementary Cell Suppression. *Manuscript, Energy Information Administration*, U. S. Department of Energy.
- Dwork, C. and Roth, A. (2014). The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science* 9, 211-407.
- Fraser, B. and Wooton, J. (2005). A Proposed Method for Confidentialising Tabular Output to Protect Against Differencing. *Joint UNECE/Eurostat work session on statistical data confidentiality*, Geneva, 9-11 November.
- Li, J., and Krenzke, T. (2016). Confidentiality approaches for real-time systems generating aggregated results. *Proceedings of the Survey Research Methods Section of the American Statistical Association*.
- McSherry, F. and Talwar, K. (2007). Mechanism Design via Differential Privacy. In *Foundations of Computer Science, 2007, FOCS'07, 48th Annual IEEE Symposium on* 94-103. IEEE, New York.
- O'Keefe, C.M. and Shlomo, N. (2012). Comparison of Remote Analysis with Statistical Disclosure Control for Protecting the Confidentiality of Business Data. *Transactions on Data Privacy*, Vol. 5, Issue 2, 403-432.
- Raghunathan T.E., Lepkowski J.M., van Hoewyk J., Solenbeger P. (2001). A multivariate technique for multiply imputing missing values using a sequence of regression models. *Survey Methodology*, Vol. 27, 85-95.
- Reiter, J.P. (2005), Releasing Multiply Imputed, Synthetic Public-Use Microdata: An Illustration and Empirical Study. *Journal of the Royal Statistical Society, A*, Vol.168, No.1, 185-205.
- Rinott, Y., O'Keefe, C., Shlomo, N., and Skinner, C. (2018). Confidentiality and Differential Privacy in the Dissemination of Frequency Tables. *Statistical Sciences*, Vol. 33, No. 3, 358-385.
- Shlomo, N., Krenzke, T. and Li, J. (2019) Confidentiality Protection Approaches for Survey Weighted Frequency Tables. Submitted.
- Shlomo, N. and Skinner. C.J. (2012). Privacy Protection from Sampling and Perturbation in Survey Microdata. *Journal of Privacy and Confidentiality*, Vol. 4, Issue 1.
- Shlomo, N. and Young, C. (2008). Invariant Post-tabular Protection of Census Frequency Counts. In *PSD'2008 Privacy in Statistical Databases*, (Eds. J.Domingo-Ferrer and Y. Saygin), Springer LNCS 5262, 77-89.

Thank you for your attention