Formal Privacy Guarantees and Analytical Validity of OnTheMap Public-Use Data

John M. Abowd, Fredrik Andersson, Matthew Graham, Lars Vilhuber, and Jeremy Wu

> IPAM Workshop Statistical and Learning-Theoretic Challenges in Data Privacy February 2010

Acknowledgements and Disclaimer

- This presentation has benefited enormously from collaboration with Johannes Gehrke, Daniel Kifer, Ashwin Machanavajjhala, and our colleagues at the Census Bureau.
- Research partially supported by the National Science Foundation (grants SES-0339191, SES-0427889 and CNS-0627680) and the Census Bureau.
- All statistical materials in this presentation have been reviewed for disclosure avoidance.
- The opinions are those of the author and not the National Science Foundation nor the Census Bureau.

Outline

- Motivation: formal privacy models and statistical disclosure limitation
- Detailed example of differential privacy concepts
- The basic OnTheMap application
- The statistical structure of the OnTheMap data
- Applying probabilistic differential privacy to OnTheMap
- The trade-off between analytical validity and confidentiality protection

Formal Privacy Models and Statistical Disclosure Limitation

- Formal privacy protection methods are based on open algorithms with provable properties
- The standard in privacy-preserving datamining is based on cryptography:
 - Only the private key (password, encryption key) is confidential; all algorithms and parameters are public
 - Attacker (= user) can have massive amounts of prior information

The Cryptographic Critique of SDL

- Standard SDL techniques fail because:
 - They and do not have provably protective properties when the attacker (= user) is allowed full access to the algorithm
 - They depend upon the realized data and not the algorithm
- Many standard SDL techniques are viewed as very risky when the cryptographic critique is applied

Point of Common Ground

Federal Committee on Statistical Methodology working paper
 22 offers the desirable disclosure avoidance property:

Disclosure relates to inappropriate attribution of information to a data subject, whether an individual or an organization. Disclosure occurs when a data subject is identified from a released file (identity disclosure), sensitive information about a data subject is revealed through the released file (attribute disclosure), or the released data make it possible to determine the value of some characteristic of an individual more accurately than otherwise would have been possible (inferential disclosure). (page 4)

 Evfimievski, Gehrke and Srikant (2003), Dwork (2006) show that disclosure avoidance in this sense is impossible to achieve in general.

Focus on Synthetic Data and Randomized Sanitizers

- The SDL technique known as synthetic data most closely resembles the cryptographic data protection techniques
- The cryptographic techniques are known as privacy-preserving datamining, randomized sanitizers, differential privacy, and e-privacy.

Definition of Synthetic Data

 $X \equiv \text{confidential data}$ $\Pr[\widetilde{X} | X] \equiv \text{PPD of } \widetilde{X} \text{ given } X$ Release data are samples of \widetilde{X}

- Synthetic data are created by estimating the posterior predictive distribution (PPD) of the release data given the confidential data; then sampling release data from the PPD conditioning on the actual confidential values.
- The PPD is a parameter-free forecasting model for new values of the complete data matrix that conditions on all values of the underlying confidential data.

Connection to Randomized Sanitizers

 $X \equiv \text{confidential data}$

 $U \equiv$ random noise

 $\operatorname{San}(X,U):(X,U)\to \widetilde{X}$ $\Pr[\widetilde{X}|X] \equiv \text{probability of } \widetilde{X} \text{ given } X$

- A randomized sanitizer creates a conditional probability distribution for the release data given the confidential data
- The randomness in a sanitizer is induced by the properties of the distribution of U
- The PPD is just a particular randomized sanitizer

E-Differential Privacy

Definition (ε - Differential Privacy): Let A be a randomized algorithm, let S be the set of all possible outputs of the algorithm, and let $\varepsilon > 0$. The algorithm A satisfies ε - differential privacy if for all pairs of data sets (D_1, D_2) that differ in exactly one row,

$$\forall S \in \mathsf{S}, \frac{P(\mathsf{A}(D_1)) = S}{P(\mathsf{A}(D_2)) = S} \le e^{\varepsilon} \text{ or } \ln \left| \frac{P(\mathsf{A}(D_1)) = S}{P(\mathsf{A}(D_2)) = S} \right| < \varepsilon$$

 Differential privacy (Dwork, and many coauthors) is difficult to maintain in sparse applications when geographically near blocks have very different posterior probabilities

Disclosure Set

Definition (Disclosure Set) : Let *D* be a table and D be the set of tables that differ from *D* in at most one row. Let A be a randomized algorithm and S be the space of outputs of the algorithm A. The disclosure set of *D*, denoted $\text{Disc}(D,\varepsilon)$, is

$$\left\{ S \in \mathsf{S} \middle| \exists X_1, X_2 \in \mathsf{D}(D), |X_1 \setminus X_2| = 1 \land \left| \ln \frac{P(\mathsf{A}(X_1) = S)}{P(\mathsf{A}(X_2) = S)} \right| > \varepsilon \right\}.$$

• This set describes the outcomes where differential privacy fails

Probabilistic Differential Privacy

- Definition (Probabilistic Differential Privacy): Let A be a randomized algorithm and S be the space of outputs of A. Let $\varepsilon > 0$ and $0 < \delta < 1$ be constants. Then A satisfies (ε, δ) probabilistic differential privacy (or (ε, δ) pdp) if for all tables *D*, $P(A(D) \in \text{Disc}(D, \varepsilon) \le \delta)$.
 - PDP allows us to control the probability that differential privacy fails
 - The analytical validity of sparse applications can be controlled with PDP because the restrictions on the prior used in the synthesizer are reasonable for use with sparse tables

Disclosure Limitation Definitions

$$X = x^{(1)}$$
 and $X = x^{(2)}$

 $\tilde{X} = \tilde{x}$, realization of the synthesizer

- Consider two confidential data matrices that differ in only a single row, $x^{(1)}$ and $x^{(2)}$
- Use the PPD to evaluate the probability of a particular release data set given the two different confidential data sets

Synthetic Data Can Leak Information about a Single Entity

$$\Pr[\tilde{X} = \tilde{x} | X = x^{(1)}] \neq \Pr[\tilde{X} = \tilde{x} | X = x^{(2)}]$$

- Changing a single row of the confidential data matrix changes the PPD or the random sanitizer
- The PPD or the random sanitizer define the transition probabilities from the confidential data to the release data
- True for all SDL procedures that infuse noise

Connection Between Synthetic Data and Differential Privacy

$$\frac{\Pr[X = x^{(1)} | \tilde{X} = \tilde{x}]}{\Pr[X = x^{(2)} | \tilde{X} = \tilde{x}]} = \frac{\Pr[\tilde{X} = \tilde{x} | X = x^{(1)}]}{\Pr[X = x^{(1)}]}$$

The posterior odds ratio for the gain in information about a single row of X is equal to the differential privacy from the randomized sanitizer that creates release data by sampling from the specified conditional distribution.

Connection Between Differential Privacy and Inferential Disclosure

$$\frac{\Pr[X = x^{(1)} | \tilde{X} = \tilde{x}]}{\Pr[X = x^{(2)} | \tilde{X} = \tilde{x}]} = \frac{\Pr[\tilde{X} = \tilde{x} | X = x^{(1)}]}{\Pr[X = x^{(1)}]}$$

The posterior odds ratio for the gain in information about a single row of X is the Dalenius (1977) definition of an inferential disclosure. Bounding the differential privacy therefore bounds the inferential disclosure.

Taking Account of Formal Privacy Models

- A variety of papers in the cryptographic data privacy literature (Dwork, Nissim and their many collaborators, Gehrke and his collaborators, and others) show that the confidentiality protection afforded by synthetic data or a randomized sanitizer depends upon properties of the transition probabilities that relate the confidential data to the release data.
- Exact data releases are not safe. Not surprising since

$$\Pr\left[\widetilde{X} \mid X\right] = I$$

implies that the sanitizer leaves the confidential data unchanged .

- Off-diagonal elements that are zero imply infinite differential privacy: exact disclosure in some cases with probability 1.
- For a full explanation of the relation between the transition matrix and differential privacy measures see Abowd and Vilhuber (2008).

Relationship to Post-randomization

• Post-randomization (Kooiman et al. 1997) focuses on the diagonal elements of

 $\Pr[\widetilde{X}|X]$

- When off-diagonal elements of this transition matrix are zero, infinite differential privacy usually results
- Swapping, shuffling, stratified sampling, and most noise-infusion methods result in off-diagonal elements that are zero

A DETAILED EXAMPLE: SYNTHETIC DATA

The Multinomial-Dirichlet Model

- The data matrix X consists of categorical variables that can by summarized by a contingency table with k categories.
- *n_i* are counts.
- π_i are probabilities

 $\mathbf{n} = (n_1, \ldots, n_k), n = \sum n_i$ $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_k), \alpha_0 = \sum \alpha_i$ $\boldsymbol{\pi} = (\pi_1, \ldots, \pi_k)$ $\mathbf{n} \sim \mathbf{M}(n, \boldsymbol{\pi})$ $\pi \sim D(\alpha)$, a priori $\pi \sim D(\alpha + n)$, a posteriori $\mathbf{m} = (m_1, \ldots, m_k), m = \sum m_i$ $\mathbf{m} \sim \mathbf{M}(m, \boldsymbol{\pi})$ 20

The Multinomial-Dirichlet Synthesizer

$$\Pr[\mathbf{m}|\mathbf{n}] = \mathrm{E}_{\boldsymbol{\pi}|\mathbf{n}}[\mathrm{M}(m,\boldsymbol{\pi})]$$

- The synthetic data are samples from the synthesizer, and can be summarized by their counts, m
- Since all the random variables are discrete, the synthesizer can be expressed as a simple transition probability matrix

		m_{1}	0	1	2	3	4	5
		m_2	5	4	3	2	1	0
n ₁	n ₂							
0	5		0.647228	0.294194	0.053490	0.004863	0.000221	0.000004
1	4		0.237305	0.395508	0.263672	0.087891	0.014648	0.000977
2	3		0.067544	0.241227	0.344610	0.246150	0.087911	0.012559
3	2		0.012559	0.087911	0.246150	0.344610	0.241227	0.067544
4	1		0.000977	0.014648	0.087891	0.263672	0.395508	0.237305
5	0		0.000004	0.000221	0.004863	0.053490	0.294194	0.647228

- $\alpha_i = \frac{1}{2}; \alpha_0 = 1$
- *n* = *m* = 5
- The table displays the transition probabilities that map n into m

$$\frac{\mathcal{E}\text{-Differential Privacy}}{\ln \frac{\Pr[m|n^{(1)}]}{\Pr[m|n^{(2)}]}} < \varepsilon$$

- The two confidential data matrices, n⁽¹⁾ and n⁽²⁾ differ by changing exactly one entity's data
- Bounding by ɛ the log inferential disclosure odds ratio in the M-D synthesizer amounts to controlling the probabilities in Pr[m|n] appropriately

				<i>m</i> ₁	0	1	2	3	4	5
				m_2	5	4	3	2	1	0
<mark>n⁽¹⁾1</mark>	n ⁽¹⁾ 2	n ⁽²⁾ 1	n ⁽²⁾ 2							
0	5	1	4		1.003353	0.29593	1.595212	2.894495	4.193778	5.493061
1	4	2	3		1.256572	0.494432	0.267708	1.029848	1.791988	2.554128
2	3	3	2		1.682361	1.009417	0.336472	0.336472	1.009417	1.682361
3	2	4	1		2.554128	1.791988	1.029848	0.267708	0.494432	1.256572
4	1	5	0		5.493061	4,193778	2.894495	1.595212	0.29593	1.003353

- The table shows all of the differential privacy ratios for the example problem
- The ε-differential privacy of this synthesizer is the maximum element in this table, 5.493061
- The differential privacy limit is attained when the synthesizer delivers (0,5) and the underlying data are either (5,0) or (4,1) (or (0,5) with original data (1,4) or (5,0))
- If I release (5,0) and you know 4 people are in category 2, then the odds are 243:1 (= exp(5.493061)) that the unknown person is in category 1

Probabilistic Differential Privacy

- This definition of differential privacy allows the *E*-differential privacy limit to fail with probability δ (Machanavajjhala *et al.* 2008)
- To compute the PDP, the joint distribution of **m** and **n** must be examined for outcomes with differential privacy that exceed the limit to ensure that they occur with total probability less than δ

		<i>m</i> ₁	0	1	2	3	4	5
		m_2	5	4	3	2	1	0
n_1	n ₂							
0	5		0.020226	0.009194	0.001672	0.000152	6.91E-06	1.26E-07
1	4		0.037079	0.061798	0.041199	0.013733	0.002289	0.000153
2	3		0.021107	0.075383	0.107691	0.076922	0.027472	0.003925
3	2		0.003925	0.027472	0.076922	0.107691	0.075383	0.021107
4	1		0.000153	0.002289	0.013733	0.041199	0.061798	0.037079
5	0		1.26E-07	6.91E-06	0.000152	0.001672	0.009194	0.020226

- The table is Pr[m,n], where the marginal Pr[n] is based on the prior D(α)
- If we want to have *E*-differential privacy of 2, then the synthesizer fails in the highlighted cells
- With prior D(α), probabilistic differential privacy has $\mathcal{E} = 2$ and $\delta = 0.000623$, which is just the sum of the highlighted cells

A DETAILED EXAMPLE: RANDOM SANITIZER

Laplace Sanitizer

- Dwork *et al.* (2006) show that ε-differential privacy can be achieved in the Multinomial model with a sanitizer using independent double exponential noise (Laplace noise) with mean zero and variance 2/ε
- Note that in our application the total *n* is released without noise

$$\mathbf{n} \sim \mathbf{M}(n, \boldsymbol{\pi})$$
$$u \sim i.i.d \operatorname{Lap}\left(0, \frac{2}{\varepsilon}\right)$$

		m ₁	0	1	2	3	4	5
		m ₂	5	4	3	2	1	0
n ₁	n ₂							
0	5		0.816060	0.159046	0.021525	0.002913	0.000394	0.000062
1	4		0.183940	0.632121	0.159046	0.021525	0.002913	0.000456
2	3		0.024894	0.159046	0.632121	0.159046	0.021525	0.003369
3	2		0.003369	0.021525	0.159046	0.632121	0.159046	0.024894
4	1		0.000456	0.002913	0.021525	0.159046	0.632121	0.183940
5	0		0.000062	0.000394	0.002913	0.021525	0.159046	0.816060

- *k* = 2
- *n* = *m* = 5
- ε = 2
- The table displays the transition probabilities that map n into m
- Note that the diagonals are larger than the M-D model and the extreme outcomes have greater probability

				m_{1}	0	1	2	3	4	5
				m_2	5	4	3	2	1	0
<mark>n⁽¹⁾1</mark>	n ⁽¹⁾ 2	n ⁽²⁾ 1	n ⁽²⁾ 2							
0	5	1	4		1.489880	1.379885	2.000000	2.000000	2.000000	2.000000
1	4	2	3		2.000000	1.379885	1.379885	2.000000	2.000000	2.000000
2	3	3	2		2.000000	2.000000	1.379885	1.379885	2.000000	2.000000
3	2	4	1		2.000000	2.000000	2.000000	1.379885	1.379885	2.000000
4	1	5	0		2.000000	2.000000	2.000000	2.000000	1.379885	1.489880

• The table confirms that the transition matrix on the previous page has $\mathcal{E} = 2$

Challenges and Applications

- Realistic problems are all very sparse
- Probabilistic differential privacy can solve the sparseness problem
 - But, it requires coarsening and domain shrinking to deliver acceptable analytical validity.
- The Laplace synthesizer can solve the sparseness problem by adaptive histogram coarsening
 - But the user cannot directly control the coarsening hence analytical validity for some hypotheses is low
- OnTheMap uses probabilistic differential privacy

A REAL APPLICATION: US CENSUS BUREAU'S ONTHEMAP



Sausalito, CA with high wages are employed

Home | Census 2000 | Subjects A to Z | FAQs | Search | Data Tools | Catalog | Quality | Privacy Policy | Policies | Source: U.S.Census Bureau, Center for Economic Studies | e-mail: <u>ces.local.employm</u>ent.dvnamics@

Brisbane

185

22910/2 8981

🔮 Internet

The OnTheMap Data Structure

- Set of linked data tables with a relational database schema
- Main tables (micro-data)
 - Job: [Person_ID, Employer_ID]
 - Residence: [Person_ID, Origin_Block, ...]
 - Workplace: [Employer_ID, Destination_Block, ...]
 - Geo-code: [Block, Tract, Latitude, Longitude, …]

Detailed Geo-spatial Data in OTM

- Workplace and residence geographies are defined using Census blocks
- Statistical analysis to estimate the PPD is based on Census tract-to-tract relations
- There are 8.2 million blocks and 65,000 tracts in the U.S.
- Every workplace block with positive employment has its own synthesizer

Dirichlet-Multinomial Synthesizer

- *I* origins
- Model each destination d separately for each demographic segment (age, earnings, industry)
- Sample data X tabulated into n
- Synthetic data tabulated into m
- Usually m = n, but not in the OTM application

 $\mathbf{n} = (n_1, \ldots, n_I), n = \sum n_i$ $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n), \alpha_0 = \sum \alpha_i = |\boldsymbol{\alpha}|$ $\boldsymbol{\pi} = (\pi_1, \ldots, \pi_I)$ $\mathbf{n} \sim \mathbf{M}(\boldsymbol{\pi}, n)$ $\pi \sim D(\alpha)$, a priori $\pi \sim D(\alpha + \mathbf{n})$, a posteriori $\mathbf{m} = (m_1, \ldots, m_I), m = \sum m_i$ $\mathbf{m} \sim \mathbf{M}(\boldsymbol{\pi}, m)$, a posteriori

Synthetic Data Model

 Likelihood of place of residence (index *i*) conditional on place of work (index *j*) and characteristics (index k):

$$p(n_{ijk} \mid \pi_{i|jk}) \propto \prod_{i=1}^{I} \pi_{i|jk}^{n_{ijk}}$$

- The resulting posterior for π is Dirichlet with parameter \mathbf{n}_{jk} + α_{jk} for each unique workplace and characteristic combination (age, earnings, industry).
- Synthesize residence counts by sampling from the posterior predictive distributions conditional on already protected (and published) destination employment counts, m_{ik}

Search Algorithm Implements PDP

- We rely on the concept of (ε, δ) -probabilistic differential privacy, where the search algorithm guarantees ε -differential privacy with 1- δ confidence (Machanavajjhala *et al.* (2008)).
- Search algorithm finds the minimum prior sample size to guarantee ε -differential privacy developed in with failure probability δ .
- This minimum prior sample size is then apportioned over points of support in the prior.
- The privacy-preserving algorithm implemented in *OnTheMap* guarantees \mathcal{E} -differential privacy protection of 8.99 with 99.99999% confidence (δ = 0.000001).

Measures to Improve Validity

- Coarsening of the outcome domain
 - Reducing the number of support points in the domain of the prior
- Editing the prior domain
 - Eliminating the most unlikely commute patterns (from prior and likelihood)
- Use of informative priors
 - Impose likely shape based on published data subject to minimum prior sample size that ensures (ε, δ)-PDP
- Pruning the prior
 - Randomly eliminating a fraction support points with no likelihood support.
 - Pruning comes with a penalty in terms of privacy protection

Refinement: Coarsening the Domain

- Blocks are collected into larger geographic areas- SuperPUMAs, PUMAs, Tracts
- Reduces the dimensionality of the domain of each destination's synthesizer
- Theorem 5.1 in Machanavajjhala et al. shows that \mathcal{E} -differential privacy, and (\mathcal{E}, δ)probabilistic differential privacy both survive coarsening with unchanged parameters

Coarsening Steps

- If origin block very far away from destination block (distance > 90th percentile of CTTP commute distribution) coarsened to Super-PUMA (400,000 population in Census 2000)
- Else if origin block far away from destination block (distance > 50th percentile of CTTP commute distribution) coarsened to PUMA (100,000 population in Census 2000)
- Else if origin block close to destination block (distance < 50th percentile of CTTP commute distribution) coarsened to Census Tract (4,000 population on average).
- Idea: "marginal differences in commute distances between candidate locations have less predictive power in allocating workers the farther away the locations are"

Effects of Coarsening

- Coarsening in formal privacy models is effectively the same as coarsening in traditional methods
- After coarsening, an entity (in this case a block) is chosen randomly to represent the coarsened unit (one block per SuperPUMA, PUMA, or tract, as appropriate)
- This ensures that the transition matrix has no zero elements at the block level
- Ratios of the elements of this transition matrix determine the differential privacy

Refinement: Editing the Prior Domain

- For each work tract:
 - if point in domain has zero probability in prior data then do:
 - eliminate point with p=0.98 if distance > 500 miles
 - eliminate point with p=0.9 if distance > 200 miles
 - eliminate point with p=0.5 if distance > 100 miles
 - do not eliminate if distance < 100 miles
 - else retain point
- Note: contribution of any likelihood data in eliminated points also eliminated

Fraction of Points in the Prior Domain with Positive Counts in Census Transportation Planning Package Data

	State A		State B		State C	
Distance (in miles)	Mean	SD	Mean	SD	Mean	SD
- low-10	0.47	0.37	0.40	0.32	0.92	0.18
- 10-25	0.30	0.26	0.19	0.19	0.63	0.29
- 25-100	0.01	0.13	0.09	0.10	0.15	0.16
- 100-500	0.01	0.03	0.01	0.02	0.02	0.04
- 500-high	0.00	0.01	0.00	0.01	0.00	0.00
All	0.18	0.28	0.14	0.23	0.34	0.40

Fraction of Points in the Domain with Positive Counts in CTPP after Eliminating Extremely Unlikely Commute Patterns

	Large State		Mediur	n State	Small State	
distance (in miles)	Mean	SD	Mean	SD	Mean	SD
- low-10	0.47	0.37	0.40	0.32	0.92	0.18
- 10-25	0.30	0.26	0.19	0.19	0.63	0.29
- 25-100	0.13	0.13	0.09	0.10	0.15	0.16
- 100-500	0.06	0.09	0.03	0.06	0.08	0.14
- 500-high	0.07	0.13	0.06	0.12	0.03	0.08
All	0.21	0.27	0.15	0.23	0.36	0.39

Fraction of likelihood data eliminated by eliminating unlikely commute patterns is about 3-7% depending on state and year

Support Points in Prior Domain (before pruning)

	Large State (A)			Medium State (B)			Small State (C)		
Support points:	Mean	Min	Max	Mean	Min	Max	Mean	Min	Max
Total	1,005	583	2,067	1,027	619	1,560	672	602	818
	-	-	By level	of coars	sening	-	-	-	-
- Super-PUMA	526	519	538	526	518	539	537	535	539
- PUMA	39	9	73	47	7	79	10	4	19
- Census Tract	438	32	1,506	453	72	998	125	56	272
	В	y distan	ce (in n	niles) bet	tween o	centroids	5		
- low-10	265	1	878	188	1	438	15	1	49
- 10-25	127	8	794	195	13	612	16	1	60
- 25-100	85	23	289	121	45	296	54	15	169
- 100-500	139	119	206	181	151	238	80	29	233
- 500-high	389	361	412	343	300	373	508	486	519

Refinement: Informative Priors

- In year 2002: Public-use CTTP data
- In year 2003-2008: Public-use previous year OnTheMap data (not posterior)
- α = max[min_alpha, f(prior density)] minimum prior sample size is the larger of the PDP value (min_alpha) or the informative prior value
- Priors unique to each employment tract
- Not strictly Bayesian because the posterior is not published, and published data are required for prior by PDP

Refinement: Domain Pruning

- Domain may still have too many blocks for good analytical validity
- Algorithm 2 prunes the domain for a given destination *j*:
 - Keep all origins in the likelihood support (confidential data)
 - For all other origins, add to domain with probability f_i;
 (generates min_p below)
 - From Machanavajjhala et al. 2008:

Theorem 5.2 (summary) : Applying the domain shrinking algorithm 2 changes the (ε, δ) - pdp to

$$\mathcal{E}' = \mathcal{E} + \max_{i \in \{i \mid n_i = 0\}} (\ln(1/f_i)) + \max_{i \in \{i \mid n_i = 0\}} [\alpha_i] \ln 2$$

Effects of Domain Pruning

- Domain pruning leaves all of the support points that appear in the likelihood function in the posterior
- Domain pruning removes some of the prior support points that have no likelihood
- Domain pruning improves analytical validity, but because it depends upon the confidential data, it increases the effective differential privacy limit

Final Privacy Settings for OnTheMap V3

- Unadjusted $\mathcal{E} = 4.6$
- Probability of failure δ = 0.000001
- Minimum retention probability min_p= -0.025
- Adjusted $\mathcal{E} = 8.9$
- Kullback-Leibler and Integrated Mean Squared Error loss functions used to set parameters of prior
- Multinomial-Dirichlet Posterior sampled for every workplace block in the U.S. (about 1.4 million)

Analytical Validity Measures

- The divergence between posterior and likelihood for a population is measured by the Kullback-Leibler Divergence index (KL) and the Integrated Mean Square Error (IMSE) over a 29 point grid defined by the cross product of:
 - 8 commute distance categories (in miles: 0, (0-1), [1-4), [4-10), [10-25),
 [25-100), [100,500), [500+]
 - 5 commute direction categories (NW, NE, SW, SE, "N/A")
- $D_{KL} = 0$ if identical; $D_{KL} = \infty$ if no overlap

$$D_{KL}(P \parallel L) = \sum_{i} L(i) \ln \frac{L(i)}{P(i)}$$





Summary: Varying *E*

- Figures show the population-weighted *D_{KL}* for all and small (1 to 9) workforce populations for *ε* = 2, 4, 4.6, 10 and 25
- Overall, D_{KL} close to zero for values of $\mathcal{E} > 4$
- Significant gains in analytical validity for small populations as we increase \mathcal{E} further to 4.6
- The marginal improvements in analytical validity from even higher values of *E* hard to justify in terms the costs in privacy protection loss



Kullback-Leibler Divergence and Prune-adjusted Epsilon by Minimum Retention Probability in Prune Function: All Populations



Kullback-Leibler Divergence and Prune-adjusted Epsilon by Minimum Retention Probability in Prune Function: Small Populations

Summary: Varying *min_p*

- Figures show the population-weighted DκL for all and small (1 to 9) workforce populations and ε for min_p = 0.1, 0.05, 0.025 and 0.001
- Large gains in analytical validity as min_p is decreased from 0.1 to 0.05 for all populations and further large gains for small populations as min_p is decreased to 0.025
- The marginal improvements in analytical validity from even lower values of min_p; hard to justify in terms the costs in privacy protection loss

Summary: Varying δ

- We evaluate δ = 0.001, 0.0001, 0.00001 and 0.000001
- Only very marginal improvements in analytical validity as we decrease confidence from 1 in a million to 1 in a 100
- No reason to consider values of δ > 0.00001

Posterior, Likelihood and Prior Mass across Commute Ranges for All and for Small Populations

Large State A											
		All		Sn	nall (min-1	.0)					
Distance	Post.	Lik.	Prior	Post.	Lik.	Prior					
0	0.07	0.07	0.01	0.30	0.32	0.18					
(0-1)	0.15	0.15	0.03	0.13	0.16	0.03					
[1-4)	0.23	0.23	0.07	0.25	0.27	0.17					
[4-10)	0.26	0.26	0.24	0.28	0.27	0.31					
[10-25)	0.28	0.28	0.39	0.21	0.22	0.17					
[25-100)	0.14	0.13	0.19	0.18	0.16	0.31					
[100-500)	0.03	0.03	0.07	0.04	0.03	0.11					
[500-high]	0.02	0.02	0.05	0.01	0.00	0.08					

Overall Summary

- Synthetic data as an privacy protection algorithm is a promising alternative to traditional disclosure avoidance methods, especially when data representation is sparse
- Hard to quantify degree of disclosure protection synthetic data methods may leak more information than intended
- OnTheMap version 3 demonstrates the successful implementation of formal privacy guarantees based on the concept of probabilistic *E*-differential privacy
- To achieve acceptable analytical validity results with privacy guarantees requires experimentation

References

- Theorems refer to A. Machanavajjhala, D. Kifer, J. Abowd, J. Gehrke, and L. Vilhuber. Privacy: From theory to practice on the map, ICDE, 2008.
- Federal Committee on Statistical Methodology, "Report on Statistical Disclosure Limitation Methodology," Working paper 22 (revised December 2005).
- Evfimievski, A., J. Gehrke, and R. Srikant. "Limiting privacy breaches in privacy-preserving data mining," PODS 2003.
- Kooiman, P., Willenborg, L.C.R.J. and Gouweleeuw, J.M., "PRAM: a method for disclosure limitation of microdata," Research paper no. 9705, Statistics Netherlands, (1997).
- Dalenius, T. "Towards a methodology for statistical disclosure control," *Statistik Tidskrift* (Statistical Review) (1977): 429-44.
- Chawla, S., C. Dwork F. McSherry, A. Smith, and H. Wee, "Towards privacy in public databases," in Proceedings of the 2nd Theory of Cryptography Conference (2005).
- Abowd, J. and L. Vilhuber, "How Protective are Synthetic Data," in J. Domingo-Ferrer and Y. Saygun, eds., Privacy in Statistical Databases, 2008" (Berlin: Springer-Verlag, 2008), pp. 239-246.
- Dwork, C "Differential Privacy," 33rd International Colloquium on Automata, Languages, and Programming—ICALP (2006): Part II, 1-12.