



Cornell University
ILR School

Cornell University ILR School
DigitalCommons@ILR

Labor Dynamics Institute

Centers, Institutes, Programs

12-7-2016

Why Statistical Agencies Need to Take Privacy-loss Budgets Seriously, and What It Means When They Do

John M. Abowd
Cornell University, John.Abowd@cornell.edu

Follow this and additional works at: <https://digitalcommons.ilr.cornell.edu/ldi>

Thank you for downloading an article from DigitalCommons@ILR.

Support this valuable resource today!

This Article is brought to you for free and open access by the Centers, Institutes, Programs at DigitalCommons@ILR. It has been accepted for inclusion in Labor Dynamics Institute by an authorized administrator of DigitalCommons@ILR. For more information, please contact catherwood-dig@cornell.edu.

If you have a disability and are having trouble accessing information on this website or need materials in an alternate format, contact web-accessibility@cornell.edu for assistance.

Why Statistical Agencies Need to Take Privacy-loss Budgets Seriously, and What It Means When They Do

Comments

To appear on fcsm.sites.usa.gov, as presented to the 2016 FCSM Statistical Policy Seminar.

Abowd acknowledges partial support through NSF Grants [1131848](#) (NCRN) and [1012593](#) (TC:Large), support by the Alfred P. Sloan Foundation, and from the U.S. Census Bureau.

At the time this talk was given, Abowd was Associate Director for Research and Methodology and Chief Scientist, U.S. Census Bureau. The opinions expressed in this talk are his own.

Why Statistical Agencies Need to Take Privacy-loss Budgets Seriously, and What It Means When They Do

John M. Abowd
Associate Director for Research and Methodology and Chief Scientist,
U.S. Census Bureau

2016 FCSM Statistical Policy Seminar
*The Future of Federal Statistics – Use of Multiple Data Sources, Anchored
in Fundamental Principles and Practices*
December 6-7, 2016

Acknowledgments and Disclaimer

- Parts of this talk were supported by the National Science Foundation, the Sloan Foundation, and the Census Bureau (before and after my appointment started)
- The opinions expressed in this talk are the my own

Outline

- The database reconstruction theorem, a.k.a. the fundamental law of information recovery
- What is a privacy-loss budget?
- How do you respect a privacy-loss budget?
- How do you prove that the rate of privacy loss in published data is consistent with the budget?
- What does it mean to prove that the released data are robust to all future attacks?

The Database Reconstruction Theorem

- Powerful result from Dinur and Nissim (2003) [[link](#)]
- *Too many statistics published too accurately from a confidential database exposes the entire database with certainty*
- How accurately is “too accurately”?
 - Cumulative noise must be of the order \sqrt{N}

Database Reconstruction II

- Led quickly to “differential privacy”:
 - Dwork, McSherry, Nissim, and Smith (2006) [[link](#)]
 - Dwork (2006) [[link](#)]
- Leading formal privacy model

Database Reconstruction III

- “The Fundamental Law of Information Recovery”
 - Dwork and Roth, 2014 [[link](#)]
 - Dwork, undated [[link](#)]
- Includes extensions found in
 - Dwork, McSherry and Talwar (2007) [[link](#)]
 - Muthukrishnan and Nikolov (2012) [[link](#)]
 - Kasiviswanathan, Rudelson and Smith (2013) [[link](#)]
 - Dwork, Smith, Steinke, Ullman, and Vadhan (2015) [[link](#)]

Historical Note

- The U.S. Census Bureau: first organization in the world to use a formally private confidentiality protection system in production
 - [OnTheMap](#) (residential side)
- Machanavajjhala, Kifer, Abowd, Gehrke, and Vilhuber (2008)
[\[link\]](#)

What is a Privacy-loss Budget?

- Not a dollar budget, but works the same way
- Constrains aggregate risk of partial database reconstruction given all published statistics
- Worst-case limit to the inferential disclosure of any identity or item
- In differential privacy, worst case is over all possible databases with the same schema for all individuals and items

Why Use Worst-case Protection?

- “Worst case” is “equal protection under the law”
 - Protects every person in the population the same way
 - Anyone who might have been selected for the census or survey, whether in the database or not
- “Average-case” protection does not
 - Can identify who is advantaged or disadvantaged *a priori*

Respecting a Privacy-loss Budget

- All released statistics can *never* permit a database reconstruction more accurate than the budget
- Protection into the indefinite future
- For differential privacy, guarantee is over all future attackers and any database with the same schema

Current Context

- Don't current confidentiality laws require data stewards to respect a privacy-loss budget, at least implicitly?
- Unclear
- Law are silent on limitations of what can be learned about the confidential data from the released statistics (database reconstruction)
- All data publication inherently involves some inferential disclosure risk; otherwise, it is useless
 - Dwork and Naor (2008) [[link](#)]: impossibility theorem
 - Kifer and Machanavajjhala (2011) [[link](#)]: no free lunch theorem

This Is Not a New Problem

- Ratio of the circumference of a circle to its diameter is constant
- Ancients didn't understand irrational numbers:
 - Babylonians: $\pi = 3 \frac{1}{8}$
 - Egyptians: $\pi = 4 \times (\frac{8}{9})^2$
 - Israelites: $\pi = 3$ [Talmud legislated value]
 - Hindu: $\pi = \frac{62,832}{20,000} = 3.1416$
 - Euclid: no rational number is exact for this problem
 - Archimedes: sequences can approximate π with increasing accuracy
- But legal documents continued to use crude approximations
- Takes time to process abstract ideas into practical laws
- Legal guidance on inferential disclosure limitation is important
- But must be constructed sensibly

Source: Beckman, Petr "A History of Pi" (1971) [\[link\]](#)

Example: Randomized Response

- Randomized response is provably privacy-loss protective
- Privacy loss bounded by the maximum Bayes factor

$$\max BF = \frac{\frac{Pr[SQ = Yes|A = Yes]}{Pr[SQ = No|A = Yes]}}{\frac{Pr[SQ = Yes]}{Pr[SQ = No]}} = \frac{Pr[A = Yes|SQ = Yes]}{Pr[A = Yes|SQ = No]} = \frac{(1/2) + (1 - 1/2)^{1/2}}{(1 - 1/2)^{1/2}} = 3$$

- Bound is the logarithm of the maximum Bayes factor
- If
 - Sensitive question asked with probability $1/2$
 - And innocuous question is “yes” with probability $1/2$
 - Then the maximum Bayes factor is 3, and $\ln 3 = 1.1$
- The privacy-loss expenditure (ϵ -differential privacy) is 1.1
- Sources: Warner (1965) [[link](#)] and Greenberg, Abdel-Latif, Simmons, and Horvitz (1969) [[link](#)]. SDL uses: Fienberg and Steele (1998) [[link](#)], Du and Zhan (2003) [[link](#)] and Erlingsson, Vasyl and Korolova (2014) [[link](#)].

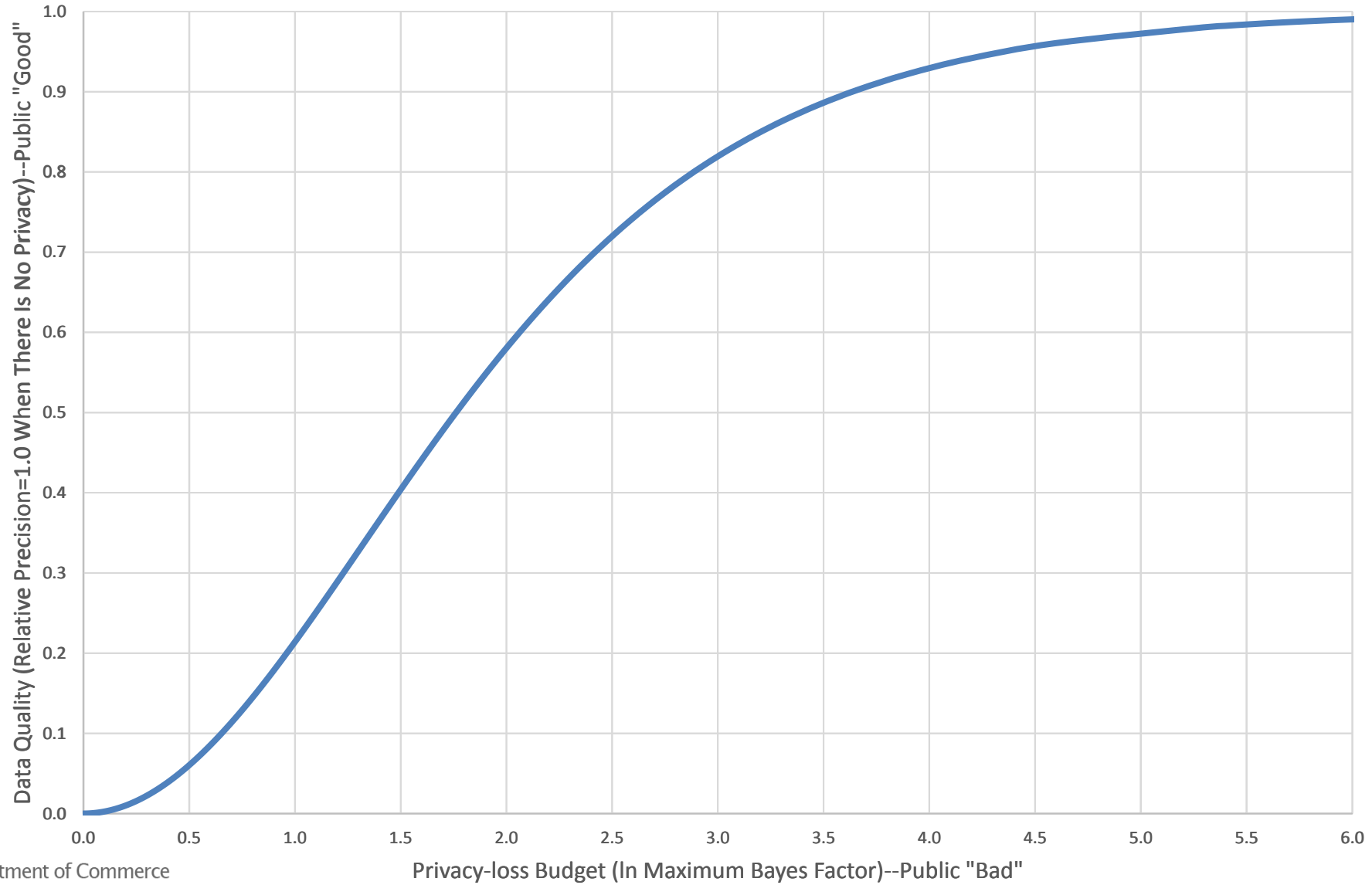
What Happens to Data Quality?

- Use relative sampling precision

$$\text{Rel. Precision} = \frac{\{Pr[\text{Ask Sensitive } Q]\}^2 \frac{n}{\theta(1-\theta)}}{\frac{n}{\theta(1-\theta)}} = \left\{\frac{1}{2}\right\}^2 = 0.25$$

- If
 - Privacy loss is $\ln 3$
 - Then, relative sampling precision is 25% of the most accurate estimator

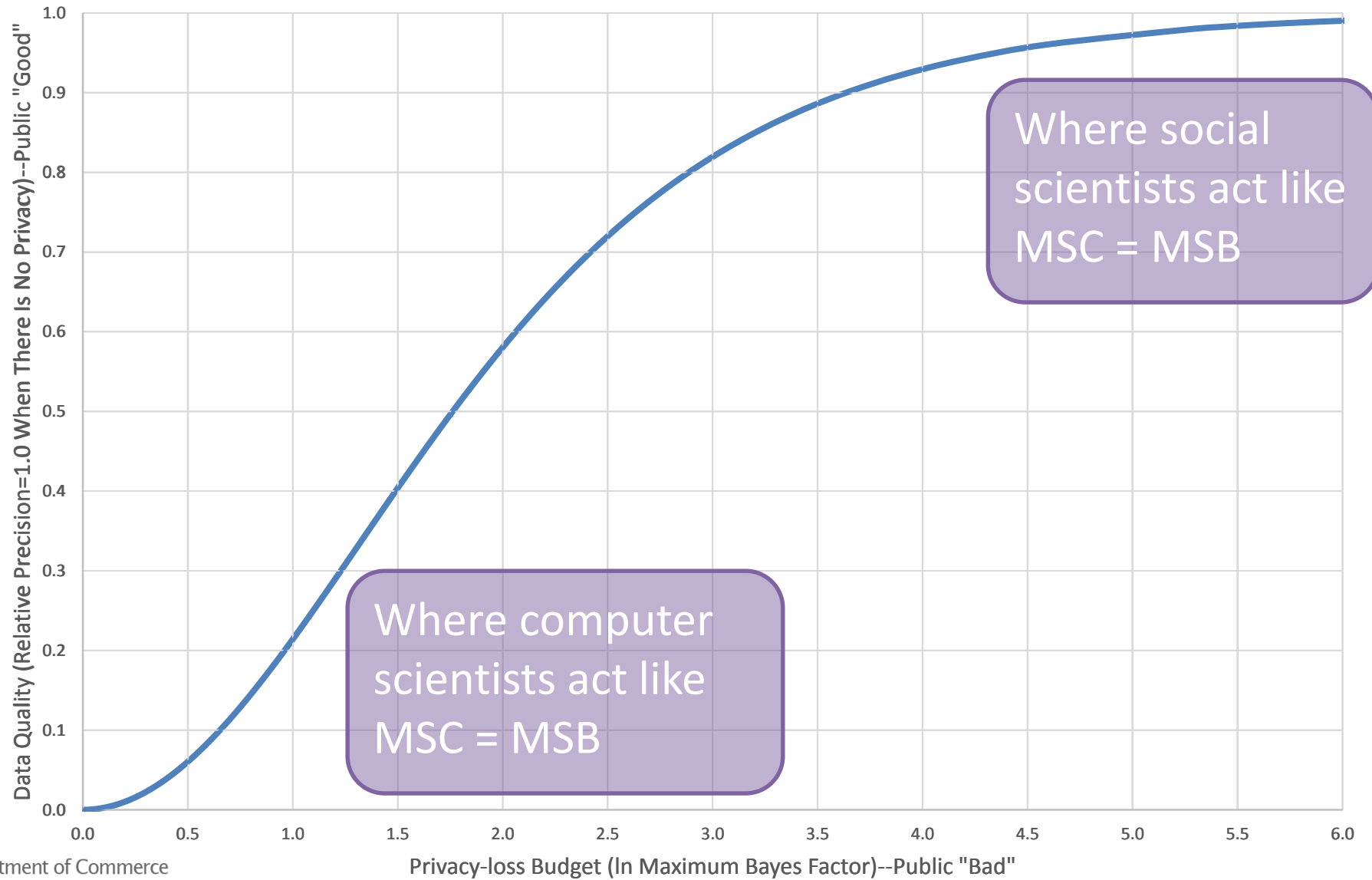
Production Possibilities Frontier/Risk-Utility/Receiver Operation Characteristics for Statistical Disclosure Limitation via Randomized Response



Disclosure Limitation is Technology

- The price of increasing data quality (public “good”) in terms of increased privacy loss (public “bad”) is the slope of the technology frontier:
 - Economics: [Production Possibilities Frontier \(Risk-Return in finance\)](#)
 - Forecasting models: [Receiver Operating Characteristics Curve](#)
 - Statistical Disclosure Limitation: [Risk-Utility Curve \(with risk on the x-axis\)](#)
- All exactly the same thing
- None able to select an optimal point

Production Possibilities Frontier/Risk-Utility/Receiver Operation Characteristics for Statistical Disclosure Limitation via Randomized Response



Some Examples

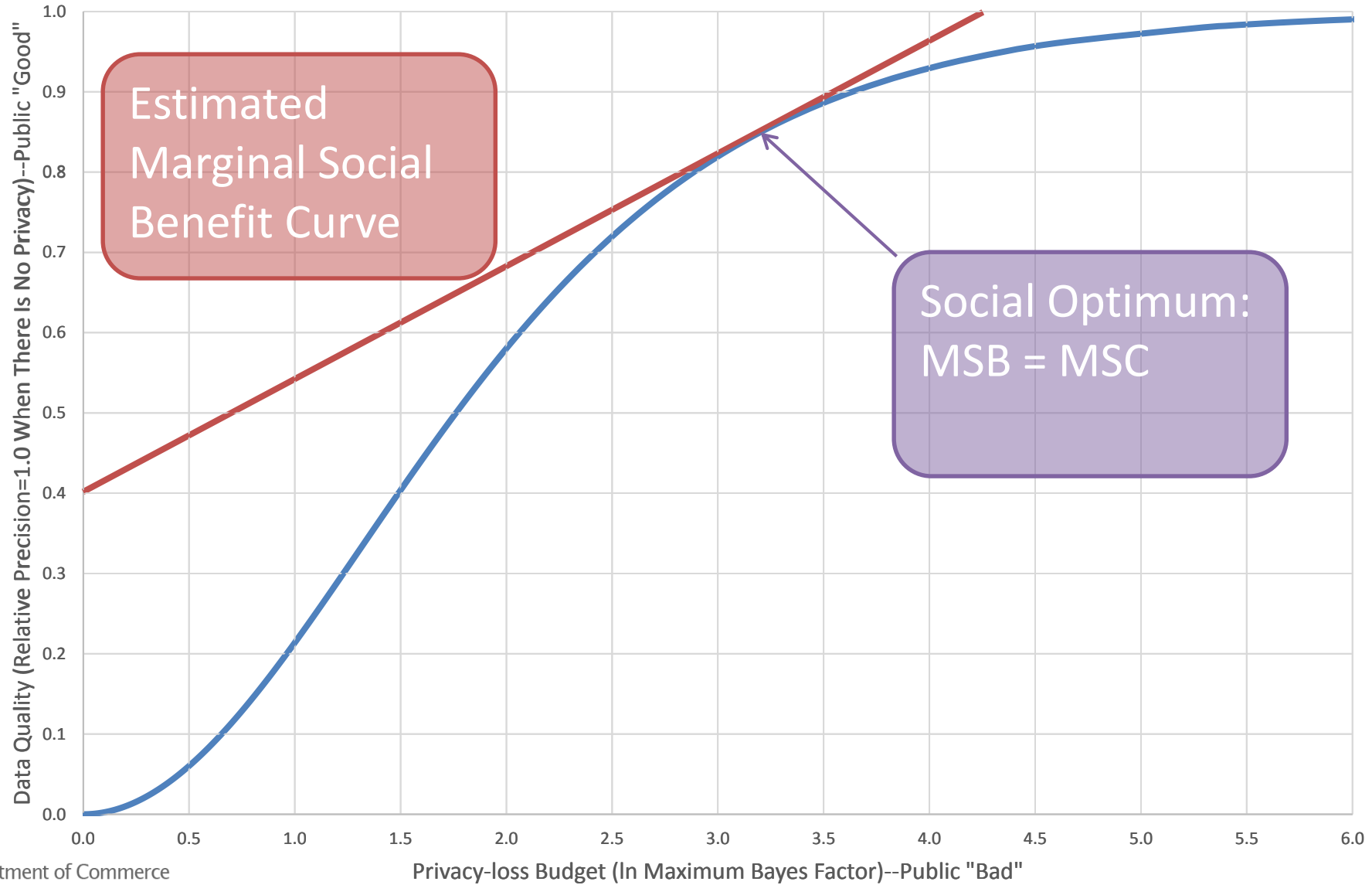
- Dwork (2008): “The parameter e in Definition 1 is public. The choice of e is essentially a social question and is beyond the scope of this paper.” [[link](#), p. 3]
- Dwork (2011): “The parameter e is public, and its selection is a social question. We tend to think of e as, say, 0.01, 0.1, or in some cases, $\ln 2$ or $\ln 3$.” [[link](#), p. 91]
- In OnTheMap, $e = 8.9$, was required to produce tract-level estimates with acceptable accuracy

How to Think about the Social Choice Problem

- The marginal social benefit is the sum of all citizens' willingness-to-pay for data quality with increased privacy loss
- Can be estimated from survey data
- The next slide shows how

See Abowd and Schmutte (2015) [[link](#)].

Production Possibilities Frontier/Risk-Utility/Receiver Operation Characteristics for Statistical Disclosure Limitation via Randomized Response



How to Prove That a Privacy-loss Budget Was Respected

- Must quantify the privacy-loss expenditure of each publication
- The collection of the algorithms taken altogether must satisfy the privacy-loss budget
- Requires methods that compose

How to Prove That the Algorithms are Resistant to All Future Attacks

- Information environment is changing much faster than before
- *It may no longer be reasonable to assert that a product is empirically safe given best-practice disclosure limitation prior to its release*
- Formal privacy models replace empirical assessment with designed protection
- Resistance to all future attacks is a property of the design

The Silver Lining

- American Statistical Association on p-values [[link](#)]
- Call for more nuanced use
- Data analysis conducted using privacy-preserving methods:
 - Control the false discovery rate
 - Reduce inferential errors due to multiple comparisons
 - Examples: Erlingsson, Vasyl and Korolova (2014) [[link](#)]; Dwork et al. (2015) [[link](#)]; Apple (2016) [[link](#)]

A Long Row to Hoe

- Concerted research and engineering effort needed to bring disclosure limitation into the 21st century
- Scientific integrity requires that we tackle this challenge
- First step is experimentation with the technologies known to work:
 - Synthetic data with validation using formally private synthesizers
 - Privacy-preserving data analysis via pre-specified query systems

Thank you.

john.maron.abowd@census.gov

References (in order of appearance)

- Dinur, Irit and Kobbi Nissim. 2003. Revealing information while preserving privacy. In *Proceedings of the twenty-second ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*(PODS '03). ACM, New York, NY, USA, 202-210. DOI: 10.1145/773153.773173.
- Dwork, Cynthia, Frank McSherry, Kobbi Nissim, and Adam Smith. 2006. in Halevi, S. & Rabin, T. (Eds.) *Calibrating Noise to Sensitivity in Private Data Analysis Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings, Springer Berlin Heidelberg*, 265-284, DOI: 10.1007/11681878_14.
- Dwork, Cynthia. 2006. Differential Privacy, *33rd International Colloquium on Automata, Languages and Programming, part II (ICALP 2006), Springer Verlag, 4052*, 1-12, ISBN: 3-540-35907-9.
- Dwork, Cynthia and Aaron Roth. 2014. *The Algorithmic Foundations of Differential Privacy*. Foundations and Trends in Theoretical Computer Science. Vol. 9, Nos. 3–4. 211–407, DOI: 10.1561/04000000042. [\[download\]](#)
- Dwork, Cynthia. Undated. The State of the Art. Slide presentation. [\[download\]](#)
- Dwork, Cynthia, Frank McSherry and Kunal Talwar. 2007. The price of privacy and the limits of LP decoding. In *Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*(STOC '07). ACM, New York, NY, USA, 85-94. DOI:10.1145/1250790.1250804.
- Muthukrishnan, S. and Aleksandar Nikolov. 2012. Optimal Private Halfspace Counting via Discrepancy, *CoRR (Computing Research Repository)*, abs/1203.5453. [\[download\]](#)
- Kasiviswanathan, Shiva Prasad, Mark Rudelson and Adam Smith. 2012. The Power of Linear Reconstruction Attacks, *CoRR (Computing Research Repository)*, abs/1210.2381.
- Dwork, Cynthia, Adam Smith, Thomas Steinke, Jonathan Ullman, and Salil Vadhan. 2015. “[Robust Traceability from Trace Amounts.](#)” In IEEE Symposium on Foundations of Computer Science (FOCS 2015). Berkeley, California, 10/18-20/2015.
- Machanavajjhala, Ashwin, Daniel Kifer, John M. Abowd , Johannes Gehrke, and Lars Vilhuber. 2008. Privacy: Theory Meets Practice on the Map, International Conference on Data Engineering (ICDE) 2008: 277-286, doi:10.1109/ICDE.2008.4497436.
- Dwork, Cynthia and Moni Naor. 2008. On the Difficulties of Disclosure Prevention in Statistical Databases or The Case for Differential Privacy, *Journal of Privacy and Confidentiality*: Vol. 2: Iss. 1, Article 8. Available at: <http://repository.cmu.edu/jpc/vol2/iss1/8>.
- Kifer, Daniel and Ashwin Machanavajjhala. 2011. No free lunch in data privacy. In *Proceedings of the 2011 ACM SIGMOD International Conference on Management of data* (SIGMOD '11). ACM, New York, NY, USA, 193-204. DOI:10.1145/1989323.1989345.
- Beckman, Petr. 1971. *A History of Pi*. Barnes and Noble, New York, NY. ISBN:0-88029-418-3.
- Warner, Stanley. L. 1965. Randomized Response: A Survey Technique for Eliminating Evasive Answer Bias *Journal of the American Statistical Association*, 60, 63-69.
- Greenberg, Bernard G., Abdel-Latif Abul-Ela, Walt R. Simmons, and Daniel G. Horvitz. 1969. The Unrelated Question Randomized Response Model: Theoretical Framework. *Journal of the American Statistical Association*, 64, 520-539.
- Fienberg, Stephen E. and Russel J. Steele. 1998. *Journal of Official Statistics*, Vol. 14, No. 4 (December): 485-502.
- Du, Wenliang and Zhijun Zhan. 2003. Using randomized response techniques for privacy-preserving data mining. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining* (KDD '03). ACM, New York, NY, USA, 505-510. DOI:10.1145/956750.956810.
- Erlingsson, Úlfar, Vasily Pihur and Aleksandra Korolova. 2014. RAPPOR: Randomized Aggregatable Privacy-Preserving Ordinal Response. In *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security* (CCS '14). ACM, New York, NY, USA, 1054-1067. DOI:10.1145/2660267.2660348.
- Dwork, Cynthia. 2008. Differential Privacy: A Survey of Results. In Agrawal, M., Du, D.; Duan, Z. and Li, A. (Eds.). *Theory and Applications of Models of Computation: 5th International Conference, TAMC 2008, Xi'an, China, April 25-29, 2008. Proceedings, Springer Berlin Heidelberg*, 1-19. [\[download\]](#)
- Dwork, Cynthia. 2011. A firm foundation for private data analysis. *Communications ACM* 54, 1 (January): 86-95. DOI:10.1145/1866739.1866758.
- Abowd, John M. and Ian M. Schmutte. 2015. Revisiting the economics of privacy: Population statistics and confidentiality protection as public goods. Labor Dynamics Institute, Cornell University, Labor Dynamics Institute, Cornell University. [\[download\]](#)
- Wasserstein, Ron L. and Nicole A. Lazar. 2016. The ASA's Statement on p-Values: Context, Process, and Purpose. *The American Statistician*, 70, 129-133. DOI: 10.1080/00031305.2016.1154108.
- Dwork, Cynthia, Feldman V., Hardt M., Pitassi T., Reingold O., Roth A. 2015. The reusable holdout: Preserving validity in adaptive data analysis. *Science*, 349, 6248 (August 7):636-8. DOI:10.1126/science.aaa9375.
- Apple, Inc. 2016. Apple previews iOS 10, the biggest iOS release ever. Press Release (June 13). URL=<http://www.apple.com/newsroom/2016/06/apple-previews-ios-10-biggest-ios-release-ever.html>.