



Cornell University
ILR School

Cornell University ILR School
DigitalCommons@ILR

Labor Dynamics Institute

Centers, Institutes, Programs

1-17-2013

Presentation: Revisiting the Economics of Privacy: Population Statistics and Privacy as Public Goods

John Abowd

Cornell University, John.Abowd@cornell.edu

Follow this and additional works at: <http://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 hlmdigital@cornell.edu.

Presentation: Revisiting the Economics of Privacy: Population Statistics and Privacy as Public Goods

Abstract

Anonymization and data quality are intimately linked. Although this link has been properly acknowledged in the Computer Science and Statistical Disclosure Limitation literatures, economics offers a framework for formalizing the linkage and analyzing optimal decisions and equilibrium outcomes.

Keywords

Data Privacy, Confidentiality Protection, Economics

Comments

Data Linkage and Anonymisation Scoping Meeting, Issac Newton Institute for Mathematical Sciences, Cambridge, UK

Revisiting the Economics of Privacy: Population Statistics and Privacy as Public Goods

John M. Abowd
Cornell University
January 17, 2013

Acknowledgements and Disclaimer

- This research uses data from the Census Bureau's Longitudinal Employer-Household Dynamics (LEHD) Program, which was partially supported by the following grants: National Science Foundation (NSF) SES-9978093, SES-0339191 and ITR-0427889; National Institute on Aging AG018854; and grants from the Alfred P. Sloan Foundation
- I also acknowledge partial direct support by NSF grants CNS-0627680, SES-0820349, SES-0922005, SES-0922494, BCS-0941226, SES-1042181, TC-1012593, and SES 1131848; and by the Census Bureau
- All confidential data used for this presentation were reviewed using the Census Bureau's disclosure avoidance protocols
- The opinions expressed in this presentation are those of the author and neither the National Science Foundation nor the Census Bureau

Colleagues and Collaborators

Fredrik Andersson, Matthew Armstrong, Sasan Bakhtiari, Patti Becker, Gary Benedetto, Melissa Bjelland, Chet Bowie, Holly Brown, Evan Buntrock, Hyowook Chiang, Stephen Ciccarella, Cynthia Clark, Rob Creecy, Lisa Dragoset, Chuncui Fan, John Fattaleh, Colleen Flannery, Lucia Foster, Matthew Freedman, Monica Garcia-Perez, Johannes Gehrke, Nancy Gordon, Kaj Gittings, Matthew Graham, Robert Groves, Owen Haaga, Hermann Habermann, John Haltiwanger, Heath Hayward, Tomeka Hill, Henry Hyatt, Emily Isenberg, Ron Jarmin, Dan Kifer, C. Louis Kincannon, Shawn Klimek, Fredrick Knickerbocker, Mark Kutzbach, Walter Kydd, Julia Lane, Paul Lengermann, Tao Li, Cindy Ma, Ashwin Machanavajjhala, Erika McEntarfer, Kevin McKinney, Thomas Mesenbourg, Jeronimo Mulato, Nicole Nestoriak, Camille Norwood, Ron Prevost, Kenneth Prewitt, George Putnam, Kalyani Raghunathan, Uma Radhakrishnan, Arnie Reznek, Bryan Ricchetti, Jerry Reiter, Marc Roemer, Kristin Sandusky, Ian Schmutte, Matthew Schneider, Rob Sienkiewicz, Liliana Sousa, Bryce Stephens, Martha Stinson, Michael Strain, Stephen Tibbets, Lars Vilhuber, J. Preston Waite, Chip Walker, Doug Webber, Dan Weinberg, Bill Winkler, Simon Woodcock, Jeremy Wu, Laura Zayatz, Chen Zhao, Nellie Zhao, Lingwen Zheng, and Chaoling Zheng

Italics = earned Ph.D. while interning at LEHD

Overview

- Anonymization and data quality are intimately linked
- Although this link has been properly acknowledged in the CS and SDL literatures, economics offers a framework for formalizing the linkage and analyzing optimal decisions and equilibrium outcomes

Technology

$$f(I, \phi; b) = 0 \text{ from } \left\{ (I, \phi) \mid \max_{\tilde{\phi} \leq \phi} I(\tilde{\phi}) \text{ s.t. } I(0) \leq b \right\}$$

where

I is a measure of the information in a data publication

ϕ is a measure of the privacy/confidentiality protection

b is the total resource limit

Technology of Anonymization

- Privacy (CS)/confidentiality (SDL) controls on data publication can be described formally as a production possibility frontier
- A PPF measures the maximum attainable data quality when the privacy controls are parameterized as ϕ , ($-\varepsilon$ from the differential privacy viewpoint)
- This is related to risk-utility curves in statistics but the formalization is more demanding

Preferences

$$v_i(y_i; I, \phi, p) = \max_x u_i(x, I, \phi) \text{ s.t. } x^T p \leq y_i$$

$$i = 1, \dots, N$$

where

u_i is consumer i 's direct utility function

v_i is consumer i 's indirect utility function

y_i is consumer i 's income

I, ϕ are the public goods (data information and privacy)

x_i is the chosen private good bundle

p is the vector of private good prices

Public Goods and Private Goods

- My formulation of the problem makes both the data publication (I) and the privacy associated with the publication (ϕ) public goods.
- No privileged access to the data (think: public-use tables or series)
- Equal protection of all consumer/citizens

Samuelson (1954) Equilibrium

$$SWF : \sum_{i=1}^n v_i(y_i, I, \phi, p)$$

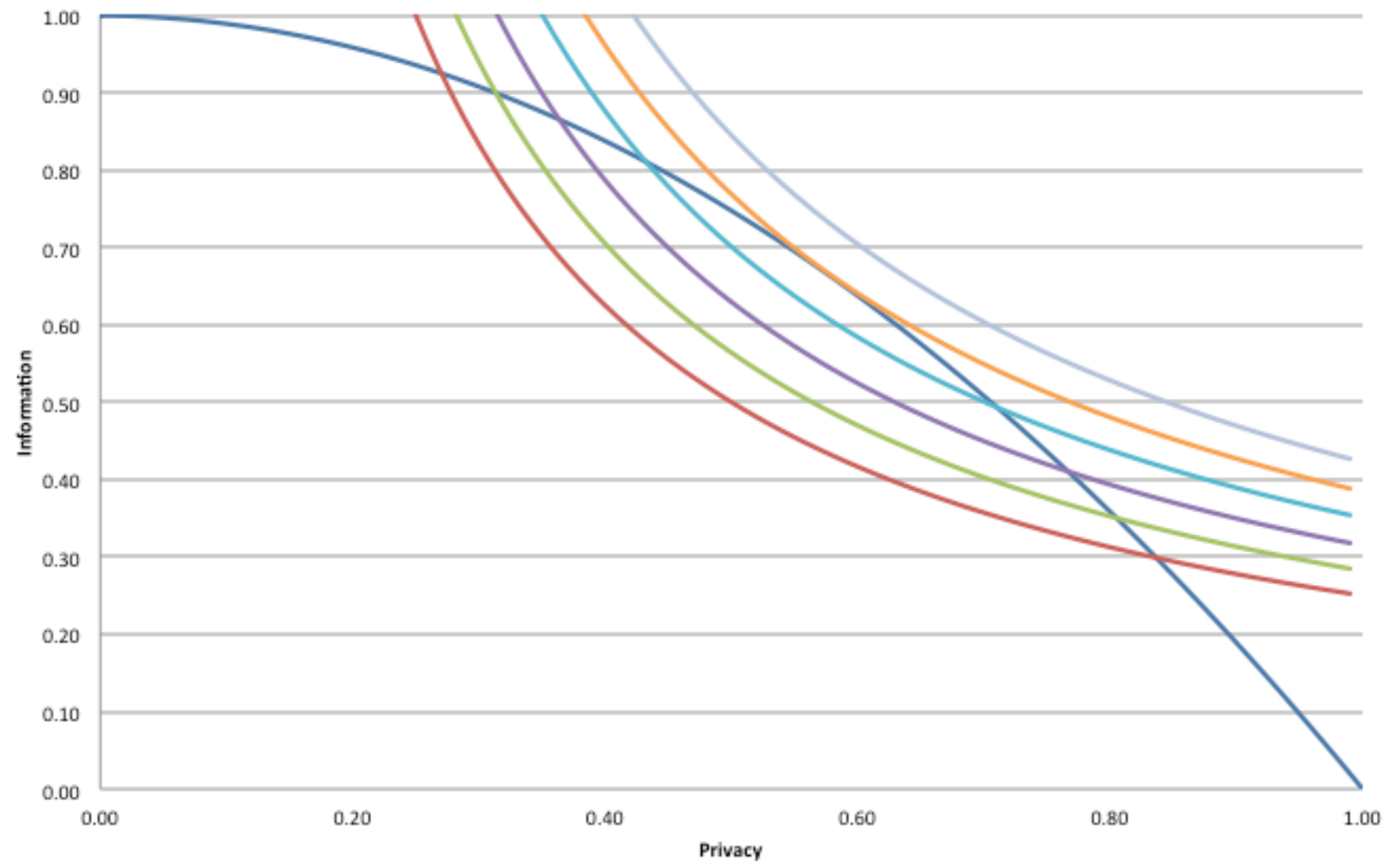
$$PPF : f(I, \phi, b) = 0$$

Optimal production of the public goods (I^0, ϕ^0)

maximize SWF subject to PPF .

$$\frac{\frac{\partial f(I^0, \phi^0, b)}{\partial \phi}}{\frac{\partial f(I^0, \phi^0, b)}{\partial I}} = \frac{\frac{\partial}{\partial \phi} \sum_{i=1}^N v_i(y_i, I^0, \phi^0, p)}{\frac{\partial}{\partial I} \sum_{i=1}^N v_i(y_i, I^0, \phi^0, p)}$$

General Equilibrium



Implications of Public Good Model

- With zero collection costs (PPF depends only on the privacy technology), always conduct a census (or, use all the administrative records)
- Straightforward to relax this, but not helpful
- Set the marginal rate of transformation (slope of the PPF) equal to the ratio of the sums of the marginal utilities of the consumers (not the marginal rate of substitution as with a private good)
- Private provision of I fails; it is undersupplied, privacy is oversupplied

Special Case: Separable Utility

$$\frac{\frac{\partial f(I^0, \phi^0, b)}{\partial \phi}}{\frac{\partial f(I^0, \phi^0, b)}{\partial I}} = \frac{\frac{\partial \sum_{i=1}^N v_i(y_i, I^0, \phi^0, p)}{\partial \phi}}{\frac{\partial \sum_{i=1}^N v_i(y_i, I^0, \phi^0, p)}{\partial I}} = \frac{\sum_{i=1}^N \frac{\partial v_i(I^0, \phi^0, p)}{\partial \phi}}{\sum_{i=1}^N \frac{\partial v_i(I^0, \phi^0, p)}{\partial I}} = \frac{\overline{\frac{\partial v}{\partial \phi}}}{\overline{\frac{\partial v}{\partial I}}}$$

- The optimal choice of data information and privacy depends upon the ratio of average marginal utilities.
- Optimal choice caters to the average consumer (not an extreme consumer)

Special Case: Separable, Identical Utility

$$\frac{\frac{\partial f(I^0, \phi^0, b)}{\partial \phi}}{\frac{\partial f(I^0, \phi^0, b)}{\partial I}} = \frac{\frac{\partial \sum_{i=1}^N v_i(y_i, I^0, \phi^0, p)}{\partial \phi}}{\frac{\partial \sum_{i=1}^N v_i(y_i, I^0, \phi^0, p)}{\partial I}} = \frac{\sum_{i=1}^N \frac{\partial v}{\partial \phi}(I^0, \phi^0, p)}{\sum_{i=1}^N \frac{\partial v}{\partial I}(I^0, \phi^0, p)} = \frac{\frac{\partial v}{\partial \phi}(I^0, \phi^0, p)}{\frac{\partial v}{\partial I}(I^0, \phi^0, p)}$$

- The optimal choice can be determined by the representative consumer even though all consumers are not identical, so there is still demand for the information

Special Case: Non-separable Quadratic Utility I

$$\frac{\frac{\partial f(I^0, \phi^0, b)}{\partial \phi}}{\frac{\partial f(I^0, \phi^0, b)}{\partial I}} = \frac{\frac{\partial}{\partial \phi} \sum_{i=1}^N v_i(y_i, I^0, \phi^0, p)}{\frac{\partial}{\partial I} \sum_{i=1}^N v_i(y_i, I^0, \phi^0, p)} = \frac{\sum_{i=1}^N \delta_i y_i \phi^0}{\sum_{i=1}^N \eta_i y_i I^0} = \frac{\bar{y}_\delta \phi^0}{\bar{y}_\eta I^0}$$

- The optimal choice depends on the ratio of weighted means of income, weighted by privacy preferences in the numerator and by information preferences in the denominator

Special Case: Non-separable Quadratic Utility II

$$\frac{\frac{\partial f(I^0, \phi^0, b)}{\partial \phi}}{\frac{\partial f(I^0, \phi^0, b)}{\partial I}} = \frac{\frac{\partial \sum_{i=1}^N v_i(y_i - \bar{y}, I^0, \phi^0)}{\partial \phi}}{\frac{\partial \sum_{i=1}^N v_i(y_i - \bar{y}, I^0, \phi^0)}{\partial I}} = \frac{\sum_{i=1}^N [\delta_i (y_i - \bar{y}) \phi^0]}{\sum_{i=1}^N [\eta_i (y_i - \bar{y}) I^0]} = \frac{\text{Cov}[\delta_i, y_i] \phi^0}{\text{Cov}[\eta_i, y_i] I^0}$$

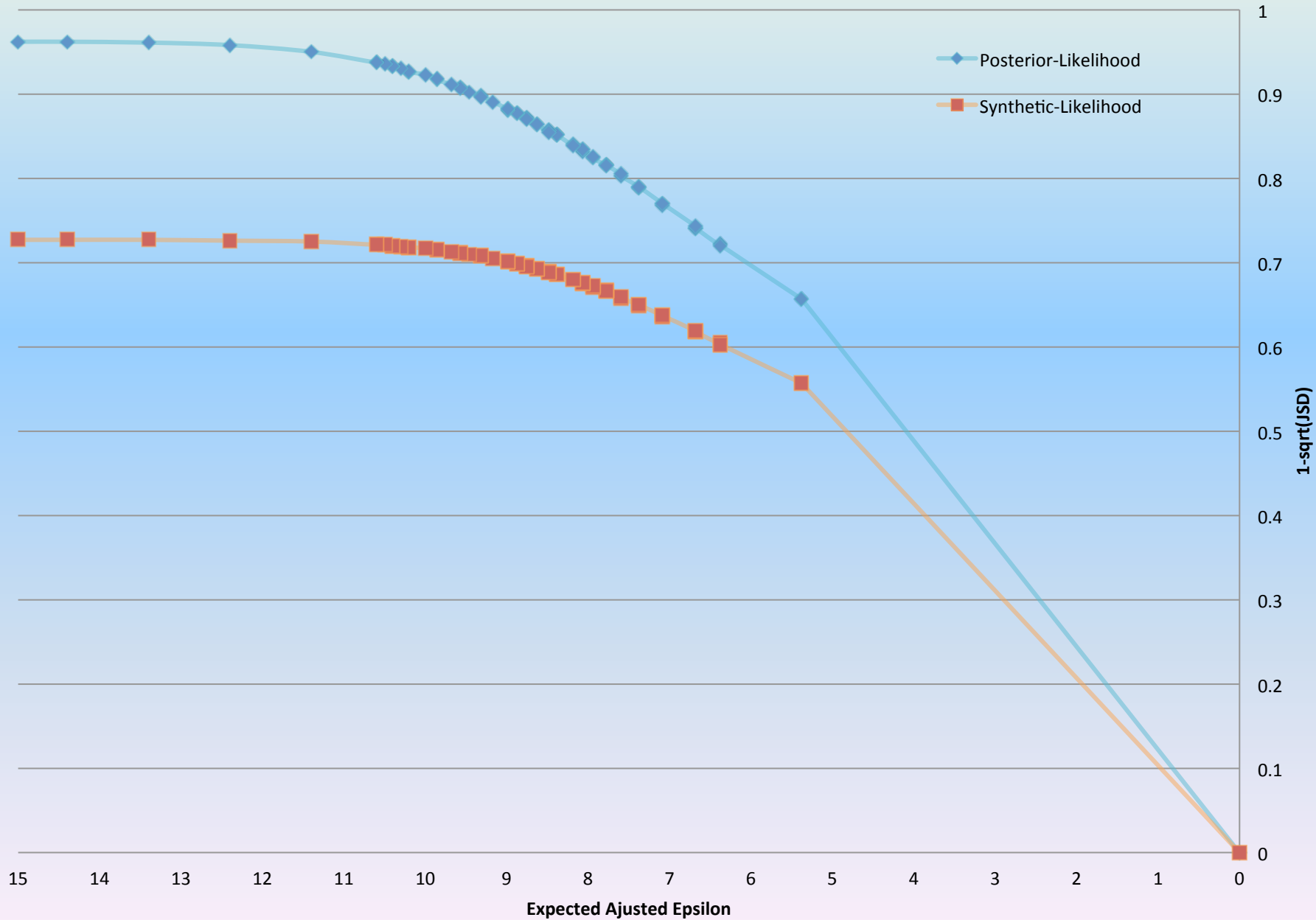
- The optimal choice depends on the ratio of covariances of preferences towards privacy (numerator) and information (denominator) with income

Example 1 PPF

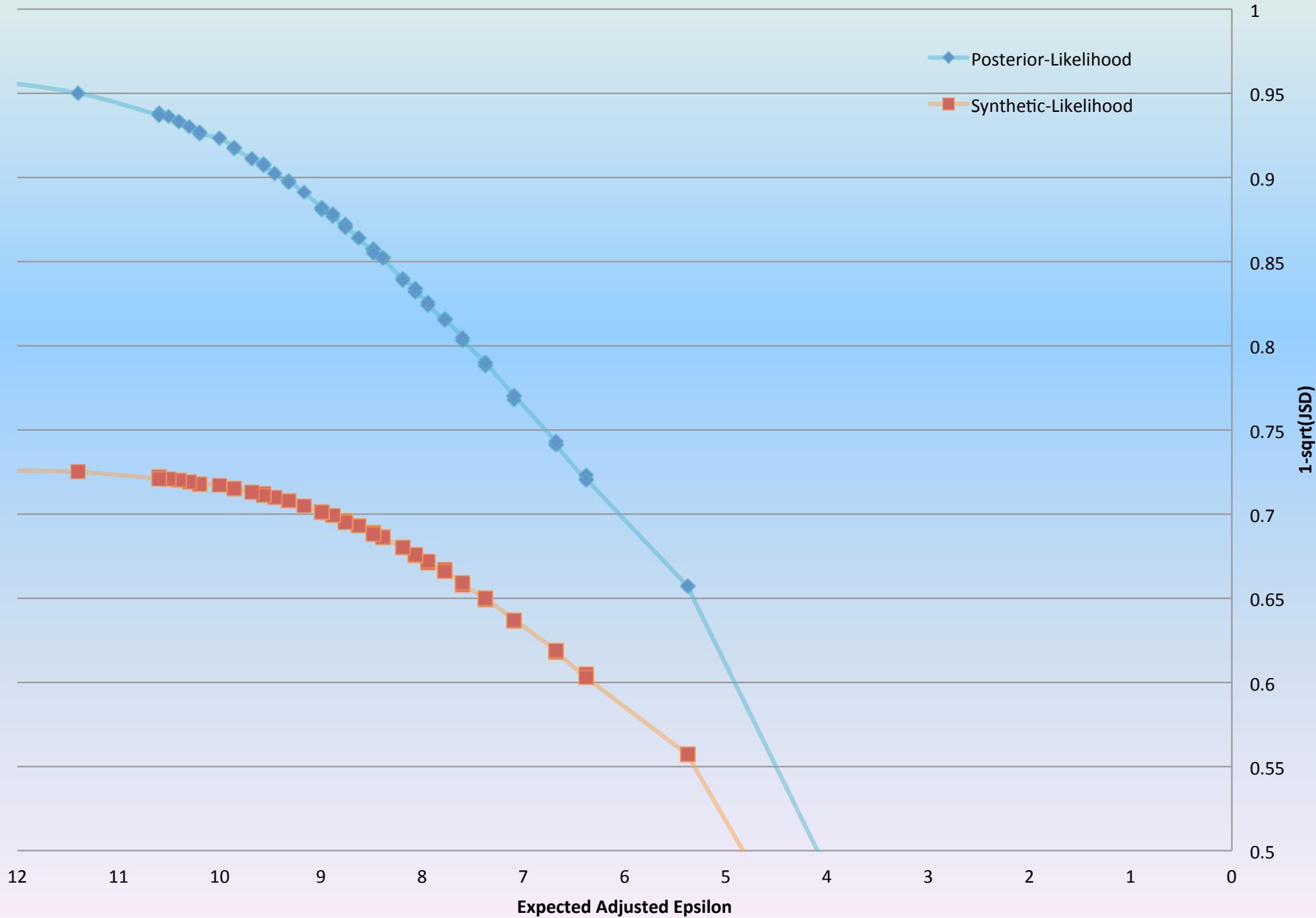
$$I_{JSD}(\phi) = 1 - \sqrt{0.5 \sum_k \pi_k \log_2 \frac{\pi_k}{0.5 \hat{\pi}_k(\phi) + 0.5 \pi_k} + 0.5 \sum_k \hat{\pi}_k(\phi) \log_2 \frac{\hat{\pi}_k(\phi)}{0.5 \hat{\pi}_k(\phi) + 0.5 \pi_k}}$$

- Based on the Jensen-Shannon distance between the true probabilities over a grid $k = 1, \dots, K$ (π_k) and the probability in each cell after protection ($\hat{\pi}_k(\phi)$)
- Note that the total information from a census of N individuals is normalized to 1, this would change if the size of the population changes, general form is $b(N)$

PPF: LODES Quality Measured by Jensen-Shannon Distance



PPF: LODES Quality Measured by Jensen-Shannon Divergence (Zoomed)

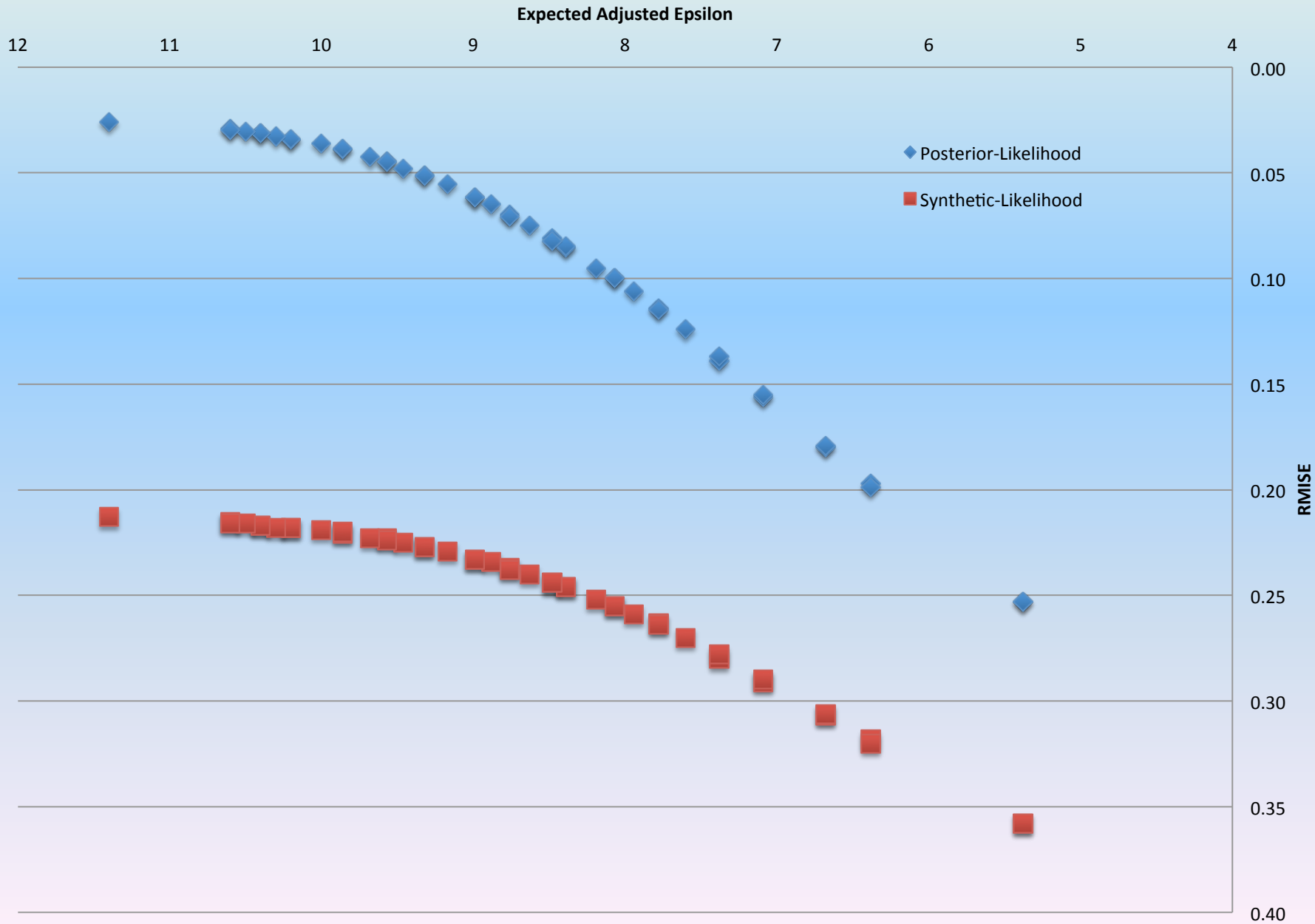


Example 2 PPF

$$I_{RMISE}(\phi) = -\sqrt{\sum_k \left(\hat{\pi}_k(\phi) - \pi_k \right)^2}$$

- Based on the root mean integrated squared error from the same census of N individuals published with privacy ϕ

PPF: LODES Quality Root Mean Integrated Squared Error



Some Implications for Optimal Data Publication/Privacy

- The OnTheMap application at the U.S. Census Bureau published with $\phi = 6.0$ attaining data quality of $I = 0.7$
- Using the separable quadratic utility model (specification II) above, this implies that the Bureau considered the ratio of preference covariances to be 0.002, which means that it assumed preferences for information were much more correlated with income than were preferences for privacy.

Alternative Specification

$$SWF : \min_i v_i(y_i, I, \phi, p)$$

$$PPF : f(I, \phi, b) = 0$$

- Known as the Rawlsian social welfare function
- Conjecture: differential privacy with $\varepsilon = -\phi$ chosen for the correct marginal individual (the one whose utility is the minimum at the optimum) is the global optimum privacy

Wrapping Up

- I have tried to pose an old problem (public good provision) in a manner that might incite mathematicians to consider models of optimal data production and protection
- This work would build on the existing CS and SDL protection methods by explicitly examining how the protection technology interacts with the data quality measure (PPF), and how preferences interact with the publication choices (SWF)