

# The Formal Privacy Research Agenda for Complex Survey Statistics

John M. Abowd

Chief Scientist and Associate Director for Research and Methodology

U.S. Census Bureau

Joint Statistical Meetings 2020

Invited Session 147: Private Data for the Public Good: Formal Privacy in Survey Organizations

Tuesday, August 4, 2020 10:00am

*The views expressed in this talk are my own and not those of the U.S. Census Bureau.*

Shape  
your future  
START HERE >

United States<sup>®</sup>  
Census  
2020

# Privacy Protection out of the Shadows

- As I have argued since 2015, traditional statistical disclosure limitation is broken and must be fixed
- I am unaware of a single argument against using formal privacy methods that does not apply *a fortiori* to traditional SDL methods as well
- Privacy practices for many current statistical products depend upon obfuscation
- 2020 Census Disclosure Avoidance System is the most transparent view into Census Bureau privacy practices ever
- That system and all future systems depend on constant feedback and interaction with the user community
- Those feedback mechanisms must be built into the overall statistical design

# Lessons learned to date

- It is far easier to implement modern disclosure limitation when the alternative is full suppression: insist on modern methods for new products
- Sampling does not imply privacy protection by itself, you have to bound the information leakage first
- Tacking formal privacy onto the current design of official statistical products usually fails
- Formal privacy requires careful definition of the estimand and modification of the traditional estimators, *but so do all statistical disclosure limitation methods if they are to be used reliably*
- Understanding how to incorporate SDL into the statistical workload and how to evaluate the resulting estimators is the first order problem
- My thoughts are largely based on Abowd and Schmutte (*Brookings Papers on Economic Activity*, Spring 2015, <https://www.brookings.edu/wp-content/uploads/2015/03/AbowdText.pdf> ) which gives a full Bayesian analysis of the problem in the online appendix here: <https://digitalcommons.ilr.cornell.edu/ldi/24/>

# Complete data likelihood function

$$\mathcal{L}_\theta (\theta_p, \theta_D | Y, R) = p_Y (Y | \theta_p) p_{R|Y} (R | Y, \theta_D) = p_{YR} (Y, R | \theta_p, \theta_D)$$

$Y$  = complete data matrix  $N \times K$  (index:  $i, j$ )

$\theta_p$  = process parameters

$R$  = inclusion matrix  $N \times K$  ( $r_{ij} = 1$  if  $y_{ij}$  is included in the sample design; 0, otherwise)

$\theta_D$  = design parameters

# Observed data likelihood function

$$\begin{aligned}\mathcal{L}_{\theta}^{(obs)}(\theta_p, \theta_D | Y^{(obs)}, R) &= p_{Y^{(obs)}R}(Y^{(obs)}, R | \theta_p, \theta_D) \\ &= \int p_{YR}(Y, R | \theta_p, \theta_D) dY^{(mis)}\end{aligned}$$

$Y^{(obs)}$  = observed data elements  $N \times K$  (index:  $i, j$ )

$Y^{(mis)}$  = missing data elements  $N \times K$  (index:  $i, j$ )

These correspond to  $r_{ij} = 1$  and  $r_{ij} = 0$ , respectively

# Inference and estimation without SDL

$$\begin{aligned} p_{\theta_p \theta_D | Y^{(obs)} R}(\theta_p, \theta_D | Y^{(obs)}, R) &\propto p_{\theta_D | \theta_p}(\theta_D | \theta_p) p_{\theta_p}(\theta_p) p_{Y^{(obs)} R}(Y^{(obs)}, R | \theta_p, \theta_D) \\ p_{\theta_P | Y^{(obs)} R}(\theta_p | Y^{(obs)}, R) &= \int p_{\theta | Y^{(obs)} R}(\theta_p, \theta_D | Y^{(obs)}, R) d\theta_D \\ &\propto \int \int p_Y(Y | \theta_p) p_{R|Y}(R | Y, \theta_D) p_{\theta_D | \theta_p}(\theta_D | \theta_p) p_{\theta_p}(\theta_p) dY^{(mis)} d\theta_D \end{aligned} \quad (\text{A.5})$$

The data inclusion model is *ignorable* if

$$p_{\theta_P | Y^{(obs)} R}(\theta_p | Y^{(obs)}, R) \equiv p_{\theta_P | Y^{(obs)}}(\theta_p | Y^{(obs)}). \quad (\text{A.6})$$

Ignorability here covers ignorable sampling and/or missing data.

# Published data likelihood function

$$\begin{aligned}\mathcal{L}_{\theta}^{(pub)}(\theta_p, \theta_D, \theta_S | Z, R) &= \int p_{Z|YR}(Z | Y, R, \theta_S) p_{YR}(Y, R | \theta_p, \theta_D) dY \\ &= \int p_{Z|YR}(Z | Y, R, \theta_S) p_{R|Y}(R | Y, \theta_D) p_Y(Y | \theta_p) dY\end{aligned}$$

$Z$  = published data matrix  $N \times K$  (index:  $i, j$ )

$\theta_S$  = SDL parameters

# Inference and estimation including SDL

$$\begin{aligned} p_{\theta|ZR}(\theta_p, \theta_D, \theta_S | Z, R) &\propto \int p_{Z|YR}(Z | Y, R, \theta_S) p_{YR}(Y, R | \theta_p, \theta_D) p_{\theta}(\theta) dY \\ &= p_{\theta}(\theta) \mathcal{L}_{\theta}^{(pub)}(\theta_p, \theta_D, \theta_S | Z, R), \end{aligned}$$

$$p_{\theta_P|ZR}(\theta_p | Z, R) = \int \int p_{\theta|ZR}(\theta_p, \theta_D, \theta_S | Z, R) d\theta_D d\theta_S$$

$$p_{\theta_P|ZR}(\theta_p | Z, R) = \int p_{\theta_P|Y^{(obs)}R}(\theta_p | Y^{(obs)}, R) p_{Y^{(obs)}|ZR}(Y^{(obs)} | Z, R) dY^{(obs)}$$

The inference or estimation should be conditioned on  $Z$  and  $R$ , which are not the same as the confidential observed data  $Y^{(obs)}$ .



# Inference and estimation including SDL-II

We define *ignorable statistical disclosure limitation* as

$$p_{\theta_P|Y^{(obs)}R}(\theta_p | Y^{(obs)} = Z, R) \equiv p_{\theta_P|ZR}(\theta_p | Z, R)$$

for all  $Y^{(obs)}$ ,  $Z$ , and  $R$ .

If the model possesses both ignorable inclusion and ignorable SDL then

$$p_{\theta_P|Y^{(obs)}}(\theta_p | Y^{(obs)} = Z) \equiv p_{\theta_P|Z}(\theta_p | Z)$$

# SDL-aware inference and estimation

$$\begin{aligned} p_{\theta_P|ZR}(\theta_p | Z, R) &= p_{\theta_P|Z}(\theta_p | Z) \\ &= \int p_{\theta_P|Y^{(obs)}}(\theta_p | Y^{(obs)}) p_{Y^{(obs)}|Z}(Y^{(obs)} | Z) dY^{(obs)} \end{aligned} \quad (\text{A.12})$$

$$p_{\theta_p\theta_D|Y^{(obs)}R}(\theta_p, \theta_D | Y^{(obs)}, R) = p_{\theta_p|Y^{(obs)}}(\theta_p | Y^{(obs)}) \quad (\text{A.13})$$

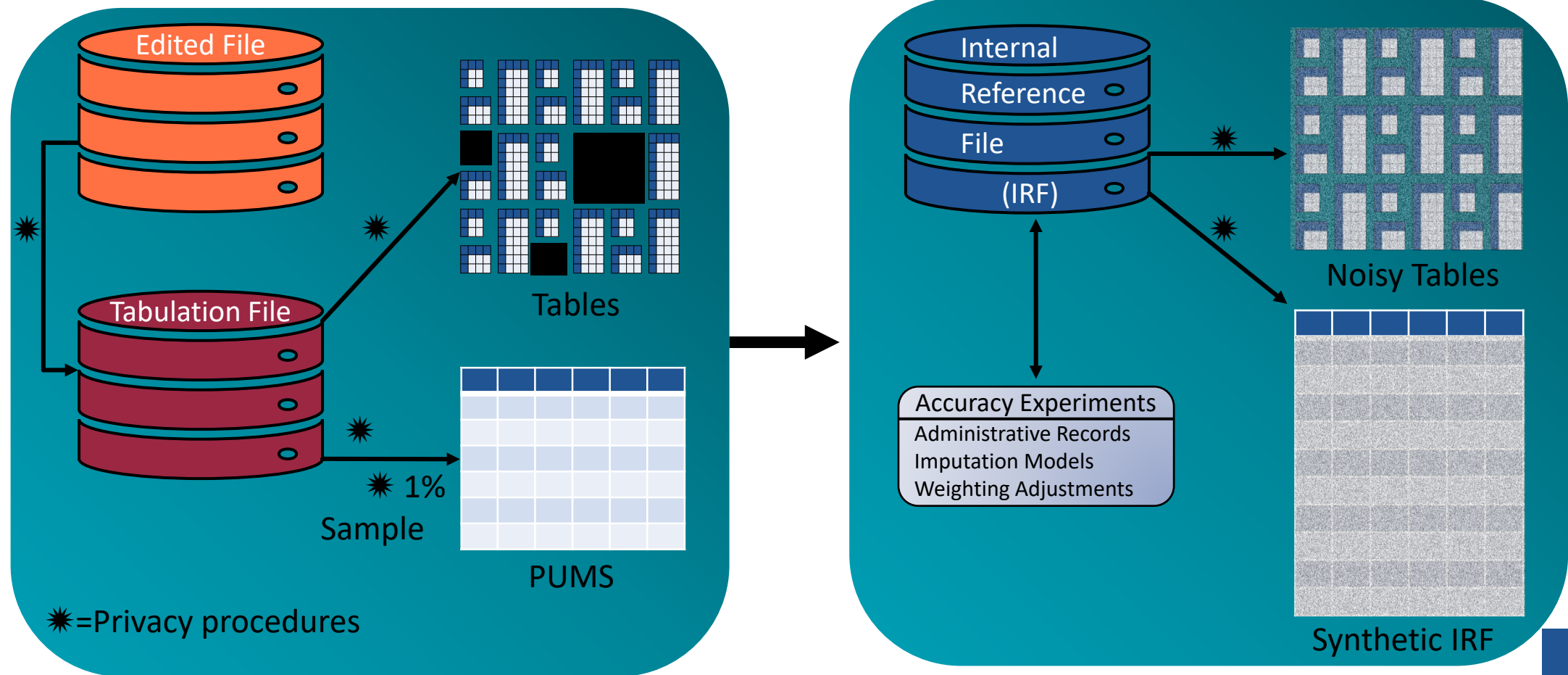
$$p_{\theta_S|ZR\theta_p\theta_D}(\theta_S | Z, R, \theta_p, \theta_D) = p_{\theta_S|Z\theta_p}(\theta_S | Z, \theta_p) \quad (\text{A.14})$$

and

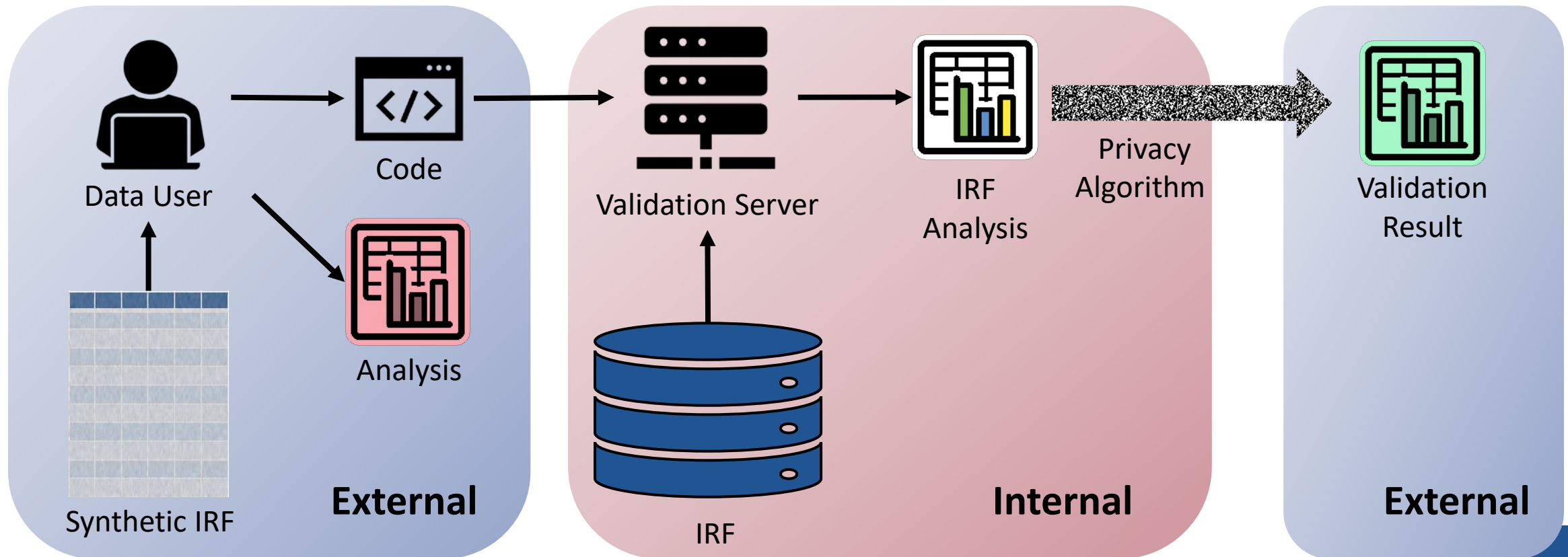
$$p_{Y^{(obs)}|ZR\theta_p\theta_D\theta_S}(Y^{(obs)} | Z, R, \theta_p, \theta_D, \theta_S) = p_{Y^{(obs)}|Z\theta_p\theta_S}(Y^{(obs)} | Z, \theta_p, \theta_S). \quad (\text{A.15})$$

These MCMC equations implement non-ignorable SDL, assuming an ignorable, known inclusion mechanism (sampling probabilities are public).

# Fewer privacy procedures allow for simpler and more adaptive workflows



# Data users will gain the ability to validate modeled output





Thank you

[John.Maron.Abowd@census.gov](mailto:John.Maron.Abowd@census.gov)