Methods of Confidentiality Protection

John M. Abowd April 2005

© John M. Abowd 2005, all rights reserved

General Principles of Disclosure Review and Control

- Protection of the confidentiality of the underlying micro data
 - Avoiding identity disclosure
 - Avoiding attribute disclosure
- Traditional methods
 - Suppression
 - Coarsening
 - Adding noise (explicit or via swapping)
- New methods
 - Partially-synthetic data
 - Fully-synthetic data

Guidance on Traditional Methods

- Use the Federal Committee on Statistical Methodology working paper series
- In particular: "Statistical Policy Working Paper 22 - Report on Statistical Disclosure Limitation Methodology"
- <u>http://www.fcsm.gov/working-papers/wp22.html</u>
- <u>http://www.fcsm.gov/working-papers/spwp22.html</u>

Tabular Methods

- Tables summarize counts, magnitudes, ratios, differences
- Micro data are aggregated to cells

Coarsening

- Modifying the cell boundaries so that they are broader
- Example: basic Census geography
 - Block
 - Block Group
 - Tract
 - Political divisions
 - Public-use microsample areas

Primary Suppression

- Suppression rule: do not publish a data item if it is based on too few entities or if it is dominated by too few entities
- Example: n,k rules
 - Suppress the data item if it is base on fewer than n entities or if fewer than n entities account for k percent of the total

Complementary Suppression

- Suppress additional data items so that the primary suppressions cannot be recovered from other published items
- Example: if a small county is suppressed then either the margin across all counties or an additional county must be suppressed to prevent recovering the suppressed value from the published values

Noise Addition

- Adding noise to the published item or to the underlying micro data to disguise the true value
- Example: confidentiality edit of Decennial Census
- Example: QWIs

The Quarterly Workforce Indicator System

- Multiplicative noise infusion system
- Establishment level micro data are distorted according to a permanent distortion factor
- Distortion factor always moves the fuzzed item away from the actual item by a minimum and maximum percentage
- All release data are computed from the fuzzed items

Fuzz Factor Distribution

$$p(\delta_j) = \begin{cases} (b-\delta)/(b-a)^2, \ \delta \in [a,b] \\ (b+\delta-2)/(b-a)^2, \ \delta \in [2-b,2-a] \end{cases}$$

$$F(\delta_{j}) = \begin{cases} 0.5 + [(b-a)^{2} - (b-\delta)^{2}]/[2(b-a)^{2}], \delta \in [a,b] \\ [(\delta+b-2)^{2}]/[2(b-a)^{2}], \delta \in [2-b,2-a] \end{cases}$$

© John M. Abowd 2005, all rights reserved

Types of Data Distortion

- Counts: *B, E, M, F, A, S, H, R, FA, FS, W*1, *W*2, *W*3, *NA, NH, NR,* and *NS*
- Ratios: *ZW*2, *ZW*3, *ZWA*, *ZWS*, *ZNA*, *ZNH*, *ZNR*, and *ZNS*
- Differences: JF, JC, JD, FJF, FJC, FJD, DWA, DWFA, DWS, DWFS

Multiplicative Noise Addition

- B is beginning of quarter employment; E is end of period; E-bar is the average.
- ZW2 is end of quarter employee earnings, W2 is total payroll for end of quarter employees.
- JF is net job flows
- Asterisk indicates distorted values.

$$B_{jt}^{*} = \delta_{j} \times B_{jt}$$

$$ZW_{2jt}^{*} = \frac{W_{2jt}^{*}}{E_{jt}} = \frac{\delta_{j} \times W_{2jt}}{E_{jt}}$$

$$JF_{kt}^{*} = G_{kt} \times \overline{E}_{kt}^{*} = JF_{kt} \times \frac{\overline{E}_{kt}^{*}}{\overline{E}_{kt}}$$

$$Z\Delta WA_{kt}^{*} = \frac{\Delta WA_{kt}}{A_{kt}} \times \frac{A_{kt}^{*}}{A_{kt}}$$

© John M. Abowd 2005, all rights reserved

Weighting

- Each fuzzed micro data item is weighted by the QWI final weight before aggregation
- This means that all input data are real numbers (not integers)
- Final disclosure control formulas must reflect rounding of the counts

Interpreting the Algorithm

 Based on the public use version of the spreadsheet

Analytic Validity of Time Series Properties

Distribution of the Error in the First Order Serial Correlation Coefficient Due to Multiplicative Noise Distortion (r* - r)					
	Beginning of				
	Quarter			Full Quarter	Net Job
Quantile	Employment	Accessions	Separations	Employment	Flows
99%	0.07894	0.07153	0.06711	0.06644	0.01104
95%	0.04338	0.04253	0.04070	0.03465	0.00503
90%	0.02610	0.03043	0.02826	0.01972	0.00314
75%	0.00946	0.01387	0.01326	0.00718	0.00124
50%	-0.00043	0.00103	0.00004	-0.00003	0.00000
25%	-0.01026	-0.01271	-0.01179	-0.00641	-0.00096
10%	-0.02520	-0.03012	-0.02592	-0.01720	-0.00281
5%	-0.03695	-0.04100	-0.03569	-0.02806	-0.00471
1%	-0.06984	-0.06863	-0.06645	-0.06185	-0.01038

Aggregating Formulas

- Counts and magnitudes:
 - Add
- Ratios and differences
 - Multiply by released base
 - Aggregate numerator and denominator separately
 - Add
- Job creations and destructions
 - Handle like counts but understand that there is an inherent loss of information

Handling the Suppressions

- Type 2 suppressions
 Randomly insert 1 or 2
- Type 3 suppressions
 - Must be treated as missing data and estimated
- Type 4 suppressions
 - Use the aggregating formulas

LEHD Research Papers

- <u>Review the LEHD Disclosure Review</u>
 <u>memo</u>
- A non-confidential version was included with today's materials

Business and Establishment-Based Models

- <u>Review the RDC disclosure guide and</u>
 <u>associated programs</u>
- A non-confidential version was included with today's materials

Special Tabulations

- <u>Review the DRB memo on special</u> tabulations
- A non-confidential version was included with today's materials