



# Categorical proxy data fusion

- Alternatives
- Uncertainty analysis

*Li-Chun Zhang*

*Statistics Norway*

[lcz@ssb.no](mailto:lcz@ssb.no)

*(University of Southampton from 1.9.2012)*



## An example: Education and election turnout data

	Turnout		
Education	No	Yes	Total
Low	(0, $0.104 = 182/1743$ )	(885)	1039
High	(28)	(676)	704
Total	182	1551	1743

	Y2		
Y1	0	1	Proportion
0	$L = \max(0, p+q-1)$ $U = \min(p, q)$		p
1			1-p
Proportion	q	1-q	

## Proxy data fusion (PDF) given { (Z1, Y1), (Z1, Y2) }

Proxy conditional independence assumption (CIA):  $\lambda_{ij}^h = \lambda_i^h \lambda_j^h$

Target joint distribution:  $\theta_{ij}$

Surrogate joint distribution by proxy CIA-based PDF:  $\tilde{\theta}_{ij} = \lambda_i^h \phi_h \lambda_j^h$

$$\begin{pmatrix} \phi_{11} & \cdots & \phi_{1J} \\ \vdots & \ddots & \vdots \\ \phi_{H1} & \cdots & \phi_{HJ} \end{pmatrix} \leftarrow \begin{pmatrix} \phi_1(Z_1) \\ \vdots \\ \phi_H(Z_1) \end{pmatrix} \Rightarrow \text{IPF/Raking} \Rightarrow \begin{pmatrix} \tilde{\theta}_{11} & \cdots & \tilde{\theta}_{1J} \\ \vdots & \ddots & \vdots \\ \tilde{\theta}_{H1} & \cdots & \tilde{\theta}_{HJ} \end{pmatrix} \leftarrow \begin{pmatrix} \theta_1(Y_1) \\ \vdots \\ \theta_H(Y_1) \end{pmatrix}$$

$$\uparrow \qquad \qquad \qquad \uparrow$$

$$(\theta_1(Y_2) \quad \cdots \quad \theta_J(Y_2)) \qquad \qquad (\theta_1(Y_2) \quad \cdots \quad \theta_J(Y_2))$$

Structure preserving estimation (SPREE):  $\alpha_{ij}^{Y_1 Y_2} = \alpha_{hj}^{Z_1 Y_2}$

Minimum data requirement: {  $Y_1, (Z_1, Y_2)$  }



## Distribution calibration (DC) and Proxy DC

$$\mathbf{F} = \begin{pmatrix} f_{11} & \cdots & f_{1H} \\ \vdots & \ddots & \vdots \\ f_{H1} & \cdots & f_{HH} \end{pmatrix} \leftarrow \begin{pmatrix} \phi_1(Z) \\ \vdots \\ \phi_H(Z) \end{pmatrix}$$

↑

$$(\theta_1(Y) \quad \cdots \quad \theta_H(Y))$$

Example:  $\mathbf{Z} = (31, 28, 41)^T$  and  $\mathbf{Y} = (30, 30, 40)^T$

$$\mathbf{f}_0 = \begin{pmatrix} 31 & 0 & 0 \\ 0 & 28 & 0 \\ 0 & 0 & 41 \end{pmatrix} \Rightarrow \tilde{\mathbf{F}} = \begin{pmatrix} 30 & 1 & 0 \\ 0 & 28 & 0 \\ 0 & 1 & 40 \end{pmatrix}$$

$$\mathbf{f} = \begin{pmatrix} f_{11} \\ \vdots \\ f_{1H} \\ \vdots \\ f_{H1} \\ \vdots \\ f_{HH} \end{pmatrix} \quad \& \quad \mathbf{b} = \begin{pmatrix} \phi_1(Z) \\ \vdots \\ \phi_H(Z) \\ \theta_1(Y) \\ \vdots \\ \theta_H(Y) \end{pmatrix} \Rightarrow \mathbf{Af} = \mathbf{b}$$

Lagrangian  $L = \Delta + \lambda^T (\mathbf{Af} - \mathbf{b})$

$$\Delta = \frac{1}{2}(\mathbf{f} - \mathbf{f}_0)^T(\mathbf{f} - \mathbf{f}_0)$$

$$\mathbf{f}_0 = \mathbf{f}(\text{Diag}\{\phi_h(Z); h=1, \dots, H\})$$

$$\Rightarrow \tilde{\mathbf{f}} = \mathbf{f}_0 + \mathbf{A}^T(\mathbf{AA}^T)^{-1}(\mathbf{b} - \mathbf{Af}_0)$$

Surrogate joint distribution by Proxy DC:  $\tilde{\theta}_{ij} = \xi_i^h \phi_h \lambda_j^h$  for  $\xi_i^h$  by DC  $\mathbf{Z}_1 \rightarrow Y_1$

Minimum data requirement:  $\{Y_1, (Z_1, Y_2)\}$



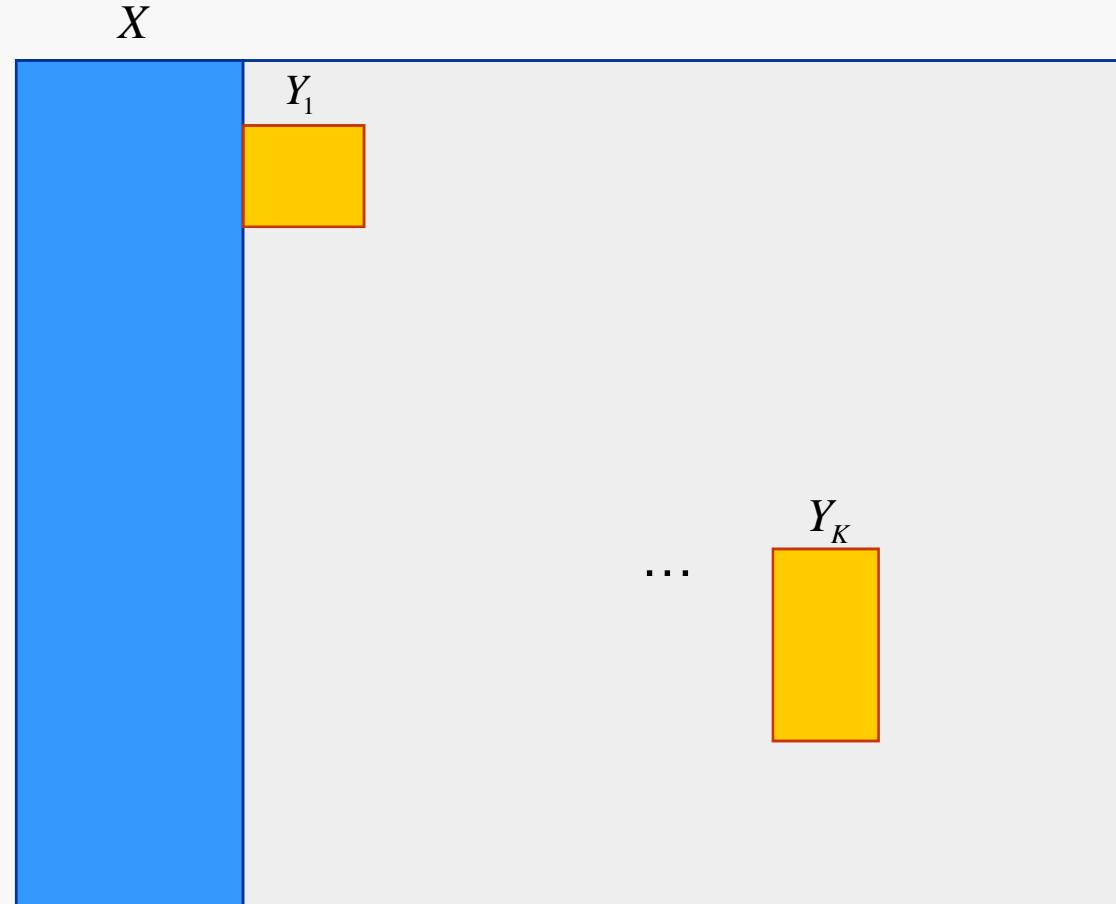
Table 1: Proxy data fusion for target  $(Y_1, Y_2)$  given proxy  $(Z_1, Z_2)$ .

Minimum Data	Assumptions for PDF
$\{(Z_1, Y_1, Z_2), (Z_1, Y_2, Z_2)\}$	CIA: $\lambda_{ij}^{hk} = \lambda_i^{hk} \lambda_j^{hk}$  SPREE: $\alpha_{(ijk)}^{Y_1 Y_2} = \alpha_{(hjk)}^{Z_1 Y_2}, \alpha_{(ijk)}^{Y_1 Y_2 Z_2} = \alpha_{(hjk)}^{Z_1 Y_2 Z_2}$  SPREE: $\alpha_{(hij)}^{Y_1 Y_2} = \alpha_{(hik)}^{Y_1 Z_2}, \alpha_{(hij)}^{Z_1 Y_1 Y_2} = \alpha_{(hik)}^{Z_1 Y_1 Z_2}$
$\{Y_1, (Z_1, Y_2, Z_2)\}$	SPREE: $\alpha_{(ijk)}^{Y_1 Z_2} = \alpha_{(hjk)}^{Z_1 Z_2}, \alpha_{(ijk)}^{Y_1 Y_2} = \alpha_{(hjk)}^{Z_1 Y_2}, \alpha_{(ijk)}^{Y_1 Y_2 Z_2} = \alpha_{(hjk)}^{Z_1 Y_2 Z_2}$  SPREE: $\alpha_{(ijk)}^{Y_2 Z_2} = \alpha_{(hjk)}^{Y_2 Z_2}, \alpha_{(ijk)}^{Y_1 Z_2} = \alpha_{(hjk)}^{Z_1 Z_2}, \alpha_{(ijk)}^{Y_1 Y_2} = \alpha_{(hjk)}^{Z_1 Y_2}, \alpha_{(ijk)}^{Y_1 Y_2 Z_2} = \alpha_{(hjk)}^{Z_1 Y_2 Z_2}$  CIA: $\lambda_{ik}^{hj} = \lambda_i^{hj} \lambda_k^{hj}$ and SPREE: $\alpha_{(ij)}^{Y_1 Y_2} = \alpha_{(hj)}^{Z_1 Y_2}$
$\{(Z_1, Y_1, Z_2), Y_2\}$	SPREE: $\alpha_{(hij)}^{Z_1 Y_2} = \alpha_{(hik)}^{Z_1 Z_2}, \alpha_{(hij)}^{Y_1 Y_2} = \alpha_{(hik)}^{Y_1 Z_2}, \alpha_{(hij)}^{Z_1 Y_1 Y_2} = \alpha_{(hik)}^{Z_1 Y_1 Z_2}$  SPREE: $\alpha_{(hij)}^{Z_1 Y_1} = \alpha_{(hik)}^{Z_1 Y_1}, \alpha_{(hij)}^{Z_1 Y_2} = \alpha_{(hik)}^{Z_1 Z_2}, \alpha_{(hij)}^{Y_1 Y_2} = \alpha_{(hik)}^{Y_1 Z_2}, \alpha_{(hij)}^{Z_1 Y_1 Y_2} = \alpha_{(hik)}^{Z_1 Y_1 Z_2}$  CIA: $\lambda_{hj}^{ik} = \lambda_h^{ik} \lambda_j^{ik}$
$\{Y_1, (Z_1, Z_2), Y_2\}$	SPREE: $\alpha_{(ij)}^{Y_1 Y_2} = \alpha_{(hk)}^{Z_1 Z_2}$  Separate Proxy DC: from $\{\phi_h\}$ to $\{\phi_i\}$ and from $\{\phi_k\}$ to $\{\phi_j\}$

# Integration of statistical data by Proxy Data Fusion

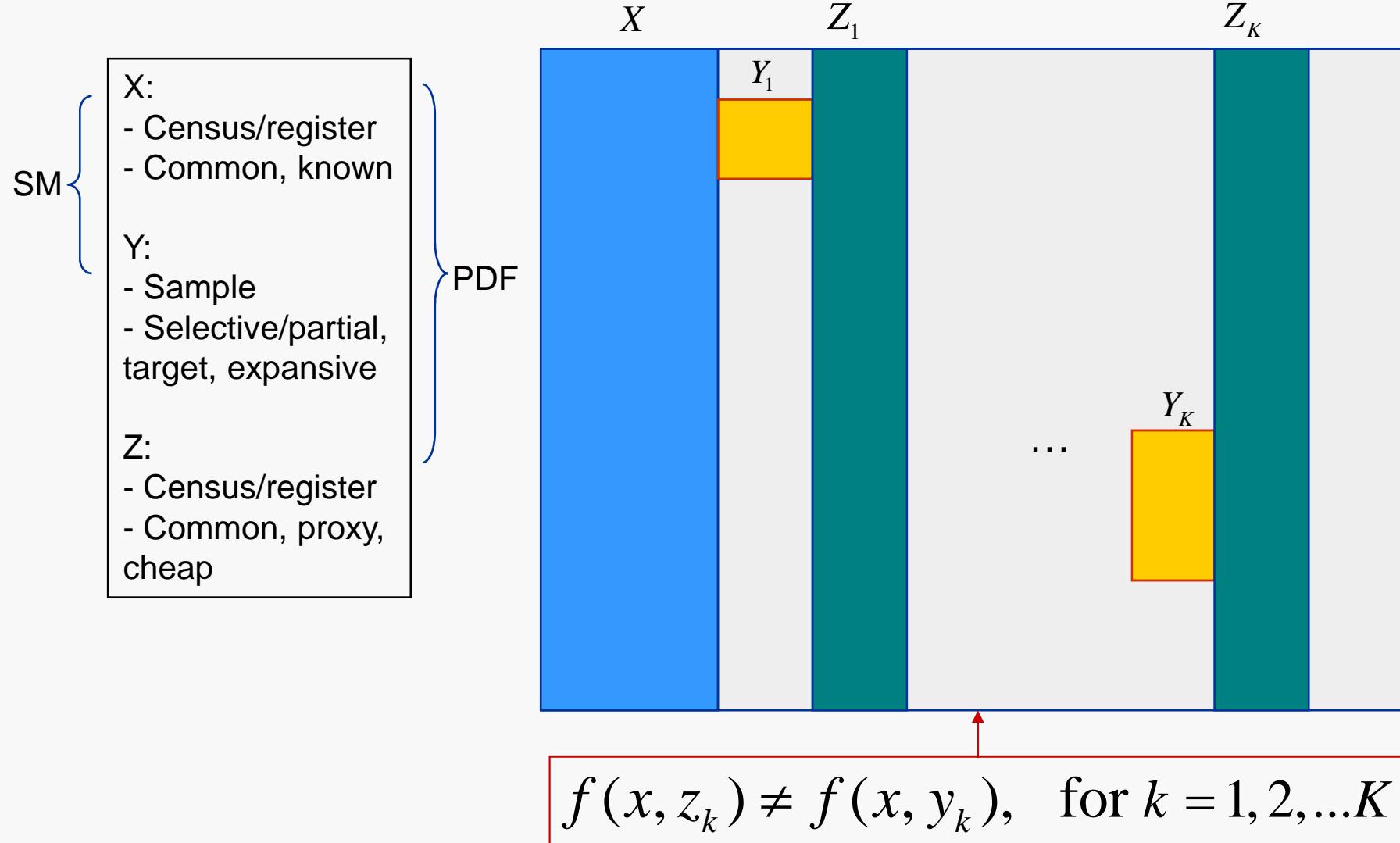
SM {

- X:
  - Census/register
  - Common, known
- Y:
  - Sample
  - Selective/partial, target, expansive



$$f(y_1, \dots, y_K | x) ? \Leftrightarrow f(x, y_1, \dots, y_K) ?$$

# Integration of statistical data by Proxy Data Fusion

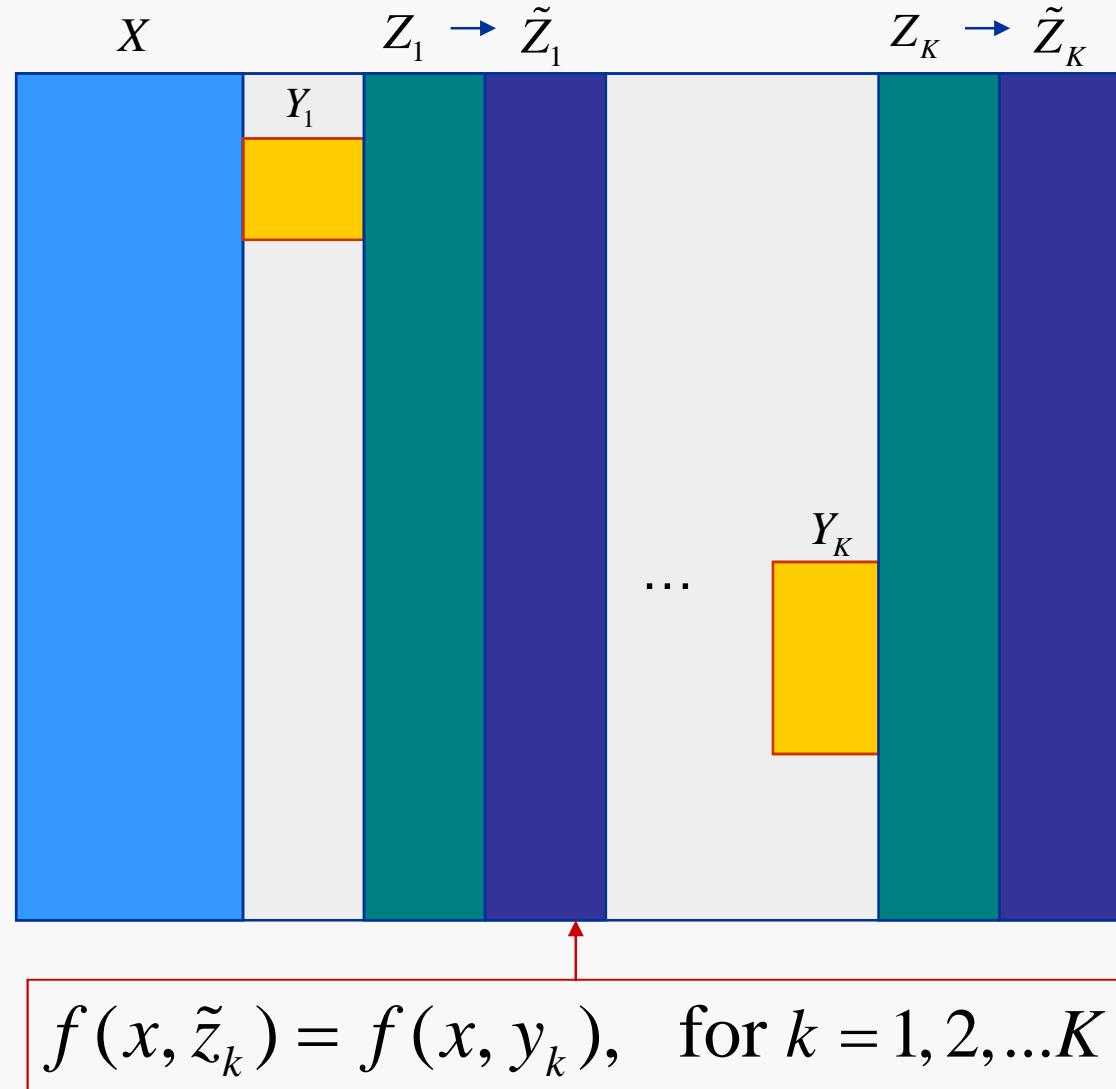


# Integration of statistical data by Proxy Data Fusion

SM {

- X:
  - Census/register
  - Common, known
- Y:
  - Sample
  - Selective/partial, target, expansive
- Z:
  - Census/register
  - Common, surrogate, cheap
- Z□:
  - Micro calibration
  - Valid (Zhang, 2012), complete

} PDF





## Uncertainty analysis: Relative efficiency of proxy data

$$\Delta = w^{ij} \Delta_{ij} \quad \text{where} \quad \Delta_{ij} = U_{ij} - L_{ij} \quad (7)$$

and  $w^{ij}$  is a (non-negative) weight function that sums to unity, i.e.  $w^{ij} \mathbf{1}_{ij} = 1$ , such as  $w^{ij} = \phi_i \phi_j$ .

$$\Delta_{ij}^{hk} = U_{ij}^{hk} - L_{ij}^{hk}$$

where

$$\max(0, \lambda_i^{hk} + \lambda_j^{hk} - 1) = L_{ij}^{hk} \leq \lambda_{ij}^{hk} = P_{ij|hk}^{Y_1 Y_2 | Z_1 Z_2} \leq U_{ij}^{hk} = \min(\lambda_i^{hk}, \lambda_j^{hk})$$

Let  $E_{Z_1 Z_2}$  denote expectation over the distribution of  $(Z_1, Z_2)$ . We have

$$\bar{L}_{ij} \stackrel{\text{def}}{=} L_{ij}^{hk} \phi_{hk} = E_{Z_1 Z_2}(L_{ij}^{hk}) \leq \theta_{ij} = E_{Z_1 Z_2}(\lambda_{ij}^{hk}) \leq E_{Z_1 Z_2}(U_{ij}^{hk}) = U_{ij}^{hk} \phi_{hk} \stackrel{\text{def}}{=} \bar{U}_{ij}$$

Meanwhile,  $\phi_i = E_{Z_1 Z_2}(\lambda_i^{hk})$  and  $\phi_j = E_{Z_1 Z_2}(\lambda_j^{hk})$ . It follows from the Jensen's inequality that

$$L_{ij} \leq \bar{L}_{ij} \quad \text{and} \quad U_{ij} \geq \bar{U}_{ij} \quad \implies \quad \bar{\Delta}_{ij} = \bar{U}_{ij} - \bar{L}_{ij} \leq \Delta_{ij}$$

The uncertainty at  $\theta_{ij}$  can thus be reduced in the presence of proxy data. So is the overall uncertainty, i.e.  $\bar{\Delta} = w^{ij} \bar{\Delta}_{ij} \leq w^{ij} \Delta_{ij} = \Delta$ . A general measure of the *relative efficiency* of the extra information due to the proxy data may be given as

$$\gamma(Z_1, Z_2; w) = \bar{\Delta}(w)/\Delta(w) = (w^{ij} \bar{\Delta}_{ij})/(w^{ij} \Delta_{ij}) \quad \text{where} \quad w^{ij} = \phi_i \phi_j \quad (8)$$



## Uncertainty analysis: Uncertainty upper bound of (proxy) data fusion

Take the setting  $\{(Z_1, Y_1, Z_2), (Z_1, Y_2, Z_2)\}$ . PDF necessarily stipulates a particular value  $\tilde{\lambda}_{ij}^{hk}$  between  $L_{ij}^{hk}$  and  $U_{ij}^{hk}$ , such as  $\tilde{\lambda}_{ij}^{hk} = \lambda_i^{hk} \lambda_j^{hk}$  under the proxy CIA (3). Its *maximum* difference to the true but unidentifiable  $\lambda_{ij}^{hk}$  is given by

$$\Lambda_{ij}^{hk} = \max(\tilde{\lambda}_{ij}^{hk} - L_{ij}^{hk}, U_{ij}^{hk} - \tilde{\lambda}_{ij}^{hk}) = \Delta_{ij}^{hk}/2 + \varepsilon_{ij}^{hk} \geq \Delta_{ij}^{hk}/2$$

where  $\varepsilon_{ij}^{hk} = |\tilde{\lambda}_{ij}^{hk} - \mu_{ij}^{hk}|$  and  $\mu_{ij}^{hk} = (L_{ij}^{hk} + U_{ij}^{hk})/2$ . This gives us a point-wise uncertainty upper bound for  $\tilde{\lambda}_{ij}^{hk}$ . It has a minimum value  $\Delta_{ij}^{hk}/2$ , regardless of the choice of PDF. Next, an uncertainty bound of the corresponding  $\theta_{ij}(\tilde{\lambda}) = \tilde{\lambda}_{ij}^{hk} \phi_{hk}$  can be given as

$$\bar{\Lambda}_{ij} = \Lambda_{ij}^{hk} \phi_{hk} = \bar{\Delta}_{ij}/2 + \bar{\varepsilon}_{ij} \geq \bar{\Delta}_{ij}/2$$

for  $\bar{\varepsilon}_{ij} = \varepsilon_{ij}^{hk} \phi_{hk}$ , and an overall uncertainty bound of PDF as

$$\bar{\Lambda} = w^{ij} \bar{\Lambda}_{ij} = \bar{\Delta}/2 + \bar{\varepsilon} \geq \bar{\Delta}/2 \quad \text{where } w^{ij} = \phi_i \phi_j \quad \text{and} \quad \bar{\varepsilon} = w^{ij} \bar{\varepsilon}_{ij}$$

It is now possible to define the *minimum uncertainty bound (MUB)* PDF as  $\theta_{ij}(\eta)$ , where

$$\{\eta_{ij}^{hk}\} = \arg \min_{\{\tilde{\lambda}_{ij}^{hk}\} \in \Omega} \bar{\varepsilon}(\tilde{\lambda}; w) \quad \text{and} \quad \bar{\varepsilon}(\tilde{\lambda}; w) = w^{ij} \varepsilon_{ij}^{hk}(\tilde{\lambda}) \phi_{hk} = w^{ij} |\tilde{\lambda}_{ij}^{hk} - \mu_{ij}^{hk}| \phi_{hk}$$



## The binary case

Table 2: Uncertainty analysis of binary data fusion. Notations:  $(L, U)$  for lower and upper bounds of cell probability,  $\tilde{\theta} = \theta(\xi)$  for stipulated cell probability and  $\Lambda$  for its uncertainty upper bound. Conditions:  $\theta_i \leq \theta_j$  and  $\theta_i + \theta_j \leq 1$ .

Cell	$L$	$U$	$U - L$	$\tilde{\theta}$	$\tilde{\theta} - L$	$U - \tilde{\theta}$	$\Lambda$
$(i, j)$	0	$\theta_i$	$\theta_i$	$\xi$	$\xi$	$\theta_i - \xi$	$\theta_i/2 +  \xi - \theta_i/2 $
$(i, j^c)$	0	$\theta_i$	$\theta_i$	$\theta_i - \xi$	$\theta_i - \xi$	$\xi$	$\theta_i/2 +  \xi - \theta_i/2 $
$(i^c, j)$	$\theta_j - \theta_i$	$\theta_j$	$\theta_i$	$\theta_j - \xi$	$\theta_i - \xi$	$\xi$	$\theta_i/2 +  \xi - \theta_i/2 $
$(i^c, j^c)$	$1 - (\theta_i + \theta_j)$	$1 - \theta_j$	$\theta_i$	$\xi + 1 - (\theta_i + \theta_j)$	$\xi$	$\theta_i - \xi$	$\theta_i/2 +  \xi - \theta_i/2 $

$$\Delta = \min(\theta_i, \theta_{i^c}, \theta_j, \theta_{j^c}) = \theta_i \wedge \theta_{i^c} \wedge \theta_j \wedge \theta_{j^c} \quad (9)$$

The same simplification of analysis applies to the setting  $\{(Z_1, Y_1, Z_2), (Z_1, Y_2, Z_2)\}$ . We have

$$\Delta^{hk} \equiv \Delta_{ij}^{hk} = U_{ij}^{hk} - L_{ij}^{hk} = \lambda_i^{hk} \wedge \lambda_{i^c}^{hk} \wedge \lambda_j^{hk} \wedge \lambda_{j^c}^{hk} \quad \text{and} \quad \bar{\Delta} = \Delta^{hk} \phi_{hk} \quad (11)$$

where  $\Delta^{hk}$  is the same for all  $(i, j)$ , and  $\bar{\Delta}$  does not depend on the weight function  $w^{ij}$ .



## Conditions for absolute reduction of uncertainty space $\Leftrightarrow$ Relative efficiency $< 1$

Table 3: Uncertainty analysis of binary PDF given  $\{(Z_1, Y_1), (Z_1, Y_2)\}$ .

$P(Y_1 Z_1)$		$Z_1$		$P(Y_2 Z_1)$		$Z_1$	
$Y_1$		Low	High	$Y_2$		Low	High
Low		$\alpha$	$1 - \alpha$	Yes		$\beta_L$	$\beta_H$
High		$1 - \alpha$	$\alpha$		No	$1 - \beta_L$	$1 - \beta_H$
Conditions		$\Delta^L$		$\Delta^H$		RE ( $\gamma$ )	
$1 - \alpha < \alpha < \beta_L < \beta_H$		$1 - \beta_L$		$1 - \beta_H$		1	
$1 - \alpha < \beta_L < \alpha < \beta_H$		$1 - \alpha$		$1 - \beta_H$		< 1	
$\beta_L < 1 - \alpha < \alpha < \beta_H$		$\beta_L$		$1 - \beta_H$		< 1	
$\beta_L < 1 - \alpha < \beta_H < \alpha$		$\beta_L$		$1 - \alpha$		< 1	
$\beta_L < \beta_H < 1 - \alpha < \alpha$		$\beta_L$		$\beta_H$		1	

Two conditions are sufficient and necessary for  $\gamma < 1$ , or absolute reduction of uncertainty:

$$\beta_L < \alpha \quad \text{and} \quad 1 - \alpha < \beta_H \quad (12)$$



Table 4: Education and election turnout data.

		Y <sub>2</sub>		Z <sub>2</sub>		
		No	Yes	No	Yes	
Y <sub>1</sub>		(154)	(885)	210	920	
Low		(154)	(885)	210	920	
High		(28)	(676)	44	569	
		182	1551	254	1489	
$Z_1 = \text{Low}$		$Z_1 = \text{High}$				
Y <sub>2</sub>		Y <sub>2</sub>				
Y <sub>1</sub>	No	Yes	Y <sub>1</sub>	No	Yes	
Low	(149)	(854)	1003	(5)	(31)	36
High	(9)	(118)	127	(19)	(558)	577
	158	972		24	589	
$Z_2 = \text{No}$		$Z_2 = \text{Yes}$				
Y <sub>2</sub>		Y <sub>2</sub>				
Y <sub>1</sub>	No	Yes	Y <sub>1</sub>	No	Yes	
Low	(140)	(61)	201	(14)	(824)	838
High	(26)	(27)	53	(2)	(649)	651
	166	88		16	1473	
$(Z_1, Z_2) = (\text{Low}, \text{No})$		$(Z_1, Z_2) = (\text{Low}, \text{Yes})$				
Y <sub>2</sub>		Y <sub>2</sub>				
Y <sub>1</sub>	No	Yes	Y <sub>1</sub>	No	Yes	
Low	(136)	(59)	195	(13)	(795)	808
High	(8)	(7)	15	(1)	(111)	112
	144	66		14	906	
$(Z_1, Z_2) = (\text{High}, \text{No})$		$(Z_1, Z_2) = (\text{High}, \text{Yes})$				
Y <sub>2</sub>		Y <sub>2</sub>				
Y <sub>1</sub>	No	Yes	Y <sub>1</sub>	No	Yes	
Low	(4)	(2)	6	(1)	(29)	30
High	(18)	(20)	38	(1)	(538)	539
	22	22		2	567	

$\Delta = \frac{182}{1743} \propto 182$	RE = 1
$\bar{\Delta} \propto (127 + 24) = 151$	RE = 0.83
$\bar{\Delta} \propto (53 + 16) = 69$	RE = 0.38
$\bar{\Delta} \propto (15 + 14 + 6 + 2) = 37$	RE = 0.20



Table 5: Estimated lower and upper bounds for education and election turnout.

	Without $(Z_1, Z_2)$	With $Z_1$	With $Z_2$	With $(Z_1, Z_2)$
$P[(Y_1, Y_2) = (\text{Low}, \text{No})]$	(0, 0.104)	(0.018, 0.104)	(0.065, 0.105)	(0.074, 0.095)

Estimates: True 0.0884 Proxy-CIA 0.0852 MUB 0.0845