

Reliable likelihood inference under coarse data and the (non-)testability of coarsening assumptions

Leibniz-Institut für Bildungsverläufe
Arbeitsbereich: Methoden der Survey-Statistik

Julia Plass*, Marco Cattaneo**, Thomas Augustin*
Georg Schollmeyer*, Christian Heumann*

*Department of Statistics, Ludwig-Maximilians University and

**School of Mathematics and Physical Sciences, University of Hull



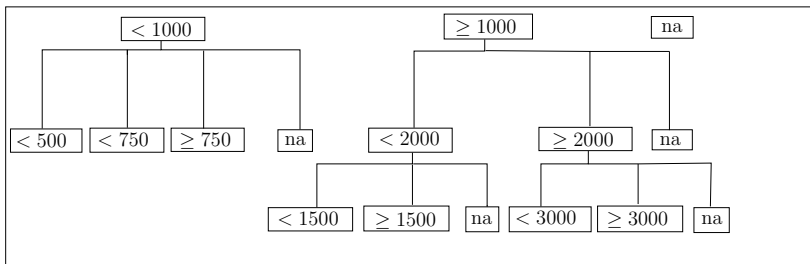
&



Bamberg, 29th of January 2018

Example for coarse categorical data

- **coarse data**: Data that can not be observed in the required accuracy originally intended in the substance matter context
- German Panel Study “Labour Market and Social Security” conducted by the IAB (**PASS** study)
- sensitive income question precisely asked, follow-up questions are directed to nonresponders:



- missing data as a special case

What's the problem?

Common dealing with incomplete data: assumptions

- ⇒ Missing at random (MAR) / coarsening at random (CAR)
- ⇒ Frequently: assumptions only for pragmatic reasons

What's the problem?

Common dealing with incomplete data: assumptions

⇒ Missing at random (MAR) / coarsening at random (CAR)

⇒ Frequently: assumptions only for pragmatic reasons



What's the problem?

Common dealing with incomplete data: assumptions

⇒ Missing at random (MAR) / coarsening at random (CAR)

⇒ Frequently: assumptions only for pragmatic reasons



Here:

- 1.) **Reliable likelihood inference**
- 2.) Testability of coarsening assumptions
- 3.) Reliable regression estimators

LATENT

i.i.d. random sample
 $(x_1, y_1), \dots, (x_n, y_n)$ of
categ. variables (X, Y)
with $\Omega_X \times \Omega_Y$

coarsening

$$y_i \in \mathcal{Y}, i=1, \dots, n$$

OBSERVABLE

i.i.d. random sample
 $(x_1, y_1), \dots, (x_n, y_n)$ of
categ. variables (X, \mathcal{Y})
with $\Omega_X \times \Omega_Y$,
 $\Omega_Y = \mathcal{P}(\Omega_Y) \setminus \{\emptyset\}$

Coarse data (categorical setting)

LATENT

i.i.d. random sample
 $(x_1, y_1), \dots, (x_n, y_n)$ of
categ. variables (X, Y)
with $\Omega_X \times \Omega_Y$

coarsening

$$y_i \in \mathcal{Y}, i=1, \dots, n$$

OBSERVABLE

i.i.d. random sample
 $(x_1, y_1), \dots, (x_n, y_n)$ of
categ. variables (X, \mathcal{Y})
with $\Omega_X \times \Omega_Y$,
 $\Omega_Y = \mathcal{P}(\Omega_Y) \setminus \{\emptyset\}$

		\mathcal{Y}							
		{a}	{b}	{c}	{a,b}	{a,c}	{b,c}	{a,b,c}	sum
X	0	$n_{0\{a\}}$	$n_{0\{b\}}$	$n_{0\{c\}}$	$n_{0\{a,b\}}$	$n_{0\{a,c\}}$	$n_{0\{b,c\}}$	$n_{0\{a,b,c\}}$	n_0
	1	$n_{1\{a\}}$	$n_{1\{b\}}$	$n_{1\{c\}}$	$n_{1\{a,b\}}$	$n_{1\{a,c\}}$	$n_{1\{b,c\}}$	$n_{1\{a,b,c\}}$	n_1

Contingency table for coarse data

Coarse data (categorical setting)

LATENT

i.i.d. random sample
 $(x_1, y_1), \dots, (x_n, y_n)$ of
categ. variables (X, Y)
with $\Omega_X \times \Omega_Y$

coarsening

$$y_i \in \mathcal{Y}_i, i=1, \dots, n$$

OBSERVABLE

i.i.d. random sample
 $(x_1, \mathcal{Y}_1), \dots, (x_n, \mathcal{Y}_n)$ of
categ. variables (X, \mathcal{Y})
with $\Omega_X \times \Omega_{\mathcal{Y}}$,
 $\Omega_{\mathcal{Y}} = \mathcal{P}(\Omega_Y) \setminus \{\emptyset\}$

		\mathcal{Y}			sum
		{a}	{b}	{a,b}	
X	0	$n_{0\{a\}}$	$n_{0\{b\}}$	$n_{0\{a,b\}}$	n_0
	1	$n_{1\{a\}}$	$n_{1\{b\}}$	$n_{1\{a,b\}}$	n_1

Contingency table for missing data

Coarse data (categorical setting)

LATENT

i.i.d. random sample
 $(x_1, y_1), \dots, (x_n, y_n)$ of
categ. variables (X, Y)
with $\Omega_X \times \Omega_Y$

coarsening

$$y_i \in \mathcal{Y}, i=1, \dots, n$$

OBSERVABLE

i.i.d. random sample
 $(x_1, y_1), \dots, (x_n, y_n)$ of
categ. variables (X, \mathcal{Y})
with $\Omega_X \times \Omega_{\mathcal{Y}}$,
 $\Omega_{\mathcal{Y}} = \mathcal{P}(\Omega_Y) \setminus \{\emptyset\}$

- **UBII**: receipt of Unemployment Benefit II

- \mathcal{Y} : categorical income

- $\{<\}$: $< 1000\text{€}$
- $\{\geq\}$: $\geq 1000\text{€}$
- $\{<, \geq\}$: $< \text{ or } \geq$

		\mathcal{Y}			
		$\{<\}$	$\{\geq\}$	$\{<, \geq\}$	sum
UBII	0	38	385	95	518
	1	36	42	9	87

PASS, w5 (Trappmann, 2010)

Coarse data (categorical setting)

LATENT

i.i.d. random sample
 $(x_1, y_1), \dots, (x_n, y_n)$ of
categ. variables (X, Y)
with $\Omega_X \times \Omega_Y$

coarsening

$$y_i \in \mathcal{Y}, i=1, \dots, n$$

OBSERVABLE

i.i.d. random sample
 $(x_1, y_1), \dots, (x_n, y_n)$ of
categ. variables (X, \mathcal{Y})
with $\Omega_X \times \Omega_Y$,
 $\Omega_Y = \mathcal{P}(\Omega_Y) \setminus \{\emptyset\}$

		\mathcal{Y}		
		$<$	\geq	sum
UBII	0	38+ 95	385	518
	1	36+ 9	42	87

potential true table

		\mathcal{Y}			
		$\{<\}$	$\{\geq\}$	$\{<, \geq\}$	sum
UBII	0	38	385	95	518
	1	36	42	9	87

observed table

Coarse data (categorical setting)

LATENT

i.i.d. random sample
 $(x_1, y_1), \dots, (x_n, y_n)$ of
categ. variables (X, Y)
with $\Omega_X \times \Omega_Y$

coarsening

$$y_i \in \mathcal{Y}, i=1, \dots, n$$

OBSERVABLE

i.i.d. random sample
 $(x_1, y_1), \dots, (x_n, y_n)$ of
categ. variables (X, \mathcal{Y})
with $\Omega_X \times \Omega_Y$,
 $\Omega_Y = \mathcal{P}(\Omega_Y) \setminus \{\emptyset\}$

		\mathcal{Y}		
		$<$	\geq	sum
UBII	0	38	385+ 95	518
	1	36	42+ 9	87

potential true table

		\mathcal{Y}			
		$\{<\}$	$\{\geq\}$	$\{<, \geq\}$	sum
UBII	0	38	385	95	518
	1	36	42	9	87

observed table

Coarse data (categorical setting)

LATENT

i.i.d. random sample
 $(x_1, y_1), \dots, (x_n, y_n)$ of
categ. variables (X, Y)
with $\Omega_X \times \Omega_Y$

coarsening

$$y_i \in \mathcal{Y}, i=1, \dots, n$$

OBSERVABLE

i.i.d. random sample
 $(x_1, y_1), \dots, (x_n, y_n)$ of
categ. variables (X, \mathcal{Y})
with $\Omega_X \times \Omega_Y$,
 $\Omega_Y = \mathcal{P}(\Omega_Y) \setminus \{\emptyset\}$

		\mathcal{Y}		
		$<$	\geq	sum
UBII	0	38+ 20	385+ 75	518
	1	36+ 5	42+ 4	87

potential true table

		\mathcal{Y}			
		$\{<\}$	$\{\geq\}$	$\{<, \geq\}$	sum
UBII	0	38	385	95	518
	1	36	42	9	87

observed table

LATENT

$$\pi_{xy} := \\ P(Y=y|X=x)$$

(error-freeness)

Observation model

$$q_{y|x,y} := \\ P(\mathcal{Y} = y|X=x, Y=y)$$

OBSERVABLE

$$p_{xy} := \\ P(\mathcal{Y} = y|X=x)$$

LATENT

$$\vartheta = (\pi_{xy}^T, \mathbf{q}_{y|xy}^T)^T$$

OBSERVABLE

$$p_{xy} := \\ P(\mathcal{Y} = y | \mathbf{X} = \mathbf{x})$$

LATENT

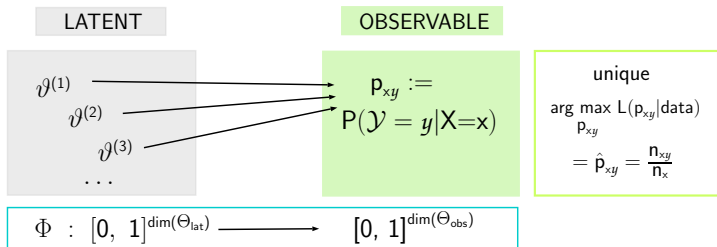
$$\vartheta = (\pi_{xy}^T, \mathbf{q}_{y|xy}^T)^T$$

OBSERVABLE

$$p_{xy} := \\ P(\mathcal{Y} = y | X=x)$$

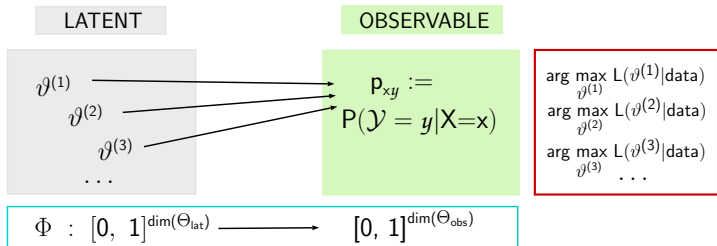
$$\begin{aligned} & \text{unique} \\ & \arg \max_{p_{xy}} L(p_{xy} | \text{data}) \\ & = \hat{p}_{xy} = \frac{n_{xy}}{n_x} \end{aligned}$$

1.) Determine MLE of observed variable distribution



- 1.) Determine MLE of observed variable distribution
- 2.) Use connection between both worlds

$$p_{xy} = \sum_{y \in \mathcal{Y}} (\pi_{xy} \cdot q_{y|x})$$

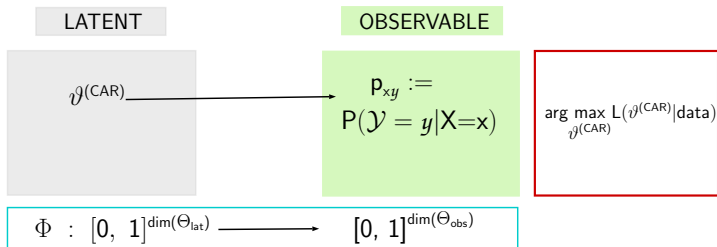


- 1.) Determine MLE of observed variable distribution
- 2.) Use connection between both worlds

$$p_{xy} = \sum_{y \in \mathcal{Y}} (\pi_{xy} \cdot q_{y|xy}).$$

- 3.) Use invariance of the likelihood

$$\hat{\pi}_{xy} \in \left[\frac{n_{x\{y\}}}{n_x}, \frac{\sum_{y \ni y} n_{xy}}{n_x} \right], \quad \hat{q}_{y|xy} \in \left[0, \frac{n_{xy}}{n_{x\{y\}} + n_{xy}} \right].$$



- 1.) Determine MLE of observed variable distribution
- 2.) Use connection between both worlds

$$p_{xy} = \sum_{y \in \mathcal{Y}} (\pi_{xy} \cdot q_{y|xy}).$$

- 3.) Use invariance of the likelihood

$$\hat{\pi}_{xy} \in \left[\frac{n_{x\{y\}}}{n_x}, \frac{\sum_{y \ni y} n_{xy}}{n_x} \right], \quad \hat{q}_{y|xy} \in \left[0, \frac{n_{xy}}{n_{x\{y\}} + n_{xy}} \right].$$

Coarsening at random (CAR)

Definition of CAR (Heitjan, Rubin, 1991):

For each fixed y and x , $q_{y|xy}$ takes the same value $\forall y \in y$.

Coarsening at random (CAR)

Definition of CAR (Heitjan, Rubin, 1991):

For each fixed y and x , $q_{y|xy}$ takes the same value $\forall y \in y$.

Illustrated by the example:

- The probability of $\{<, \geq\}$ is taken to be independent of the true income category in both subgroups split by UBI:

$$q_{\{<,\geq\}|0<} = q_{\{<,\geq\}|0\geq} \quad \text{and} \quad q_{\{<,\geq\}|1<} = q_{\{<,\geq\}|1\geq}$$



Coarsening at random (CAR)

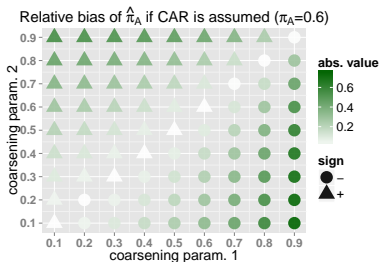
Definition of CAR (Heitjan, Rubin, 1991):

For each fixed y and x , $q_{y|x}$ takes the same value $\forall y \in y$.

Illustrated by the example:

- The probability of $\{<, \geq\}$ is taken to be independent of the true income category in both subgroups split by UBII:

$$q_{\{<,\geq\}|0<} = q_{\{<,\geq\}|0\geq} \quad \text{and} \quad q_{\{<,\geq\}|1<} = q_{\{<,\geq\}|1\geq}$$



Coarsening at random (CAR)

Definition of CAR (Heitjan, Rubin, 1991):

For each fixed y and x , $q_{y|xy}$ takes the same value $\forall y \in y$.

Illustrated by the example:

- The probability of $\{<, \geq\}$ is taken to be independent of the true income category in both subgroups split by UBI:

$$q_{\{<,\geq\}|0<} = q_{\{<,\geq\}|0\geq} \quad \text{and} \quad q_{\{<,\geq\}|1<} = q_{\{<,\geq\}|1\geq}$$



- Resulting estimators:

$$\hat{\pi}_{x<}^{(CAR)} = \frac{n_{x\{<\}}}{n_{x\{<\}} + n_{x\{\geq\}}}, \quad \hat{q}_{\{<,\geq\}|x<}^{(CAR)} = \hat{q}_{\{<,\geq\}|x\geq}^{(CAR)} = \frac{n_{x\{<,\geq\}}}{n_x}$$

Summary: Reliable likelihood inference

Estimators for subgroup 0 (PASS data example) under ...

- ... no assumptions

$$\hat{\pi}_{0<} \in [0.07, 0.26], \quad \hat{q}_{\{<,\geq\}|0<} \in [0, 0.71], \quad \hat{q}_{\{<,\geq\}|0\geq} \in [0, 0.20]$$

- ... CAR

$$\hat{\pi}_{0<}^{(\text{CAR})} = 0.09, \quad \hat{q}_{\{<,\geq\}|0<}^{(\text{CAR})} = 0.18, \quad \hat{q}_{\{<,\geq\}|0\geq}^{(\text{CAR})} = 0.18$$

- all tenable auxiliary information formalized via (Nordheim, 1984)

$$R_{0,<,\geq,\{<,\geq\}} = \frac{q_{\{<,\geq\}|0<}}{q_{\{<,\geq\}|0\geq}} \in [0, 1] :$$

$$\hat{\pi}_{0<} \in [0.07, 0.09], \quad \hat{q}_{\{<,\geq\}|0<} \in [0, 0.18], \quad \hat{q}_{\{<,\geq\}|0\geq} \in [0, 0.18]$$

What's the problem?

Common dealing with incomplete data: assumptions

⇒ Missing at random (MAR) / coarsening at random (CAR)

⇒ Frequently: assumptions only for pragmatic reasons



Here:

- 1.) Reliable likelihood inference
- 2.) **Testability of coarsening assumptions**
- 3.) Reliable regression estimators

Subgroup independence (SI)

Subgroup independence (Plass, Augustin, Cattaneo, Schollmeyer, 2015):

For each fixed y and $y \in \mathcal{Y}$, $q_{y|xy}$ takes the same value $\forall x \in \Omega_X$.

Subgroup independence (SI)

Subgroup independence (Plass, Augustin, Cattaneo, Schollmeyer, 2015):

For each fixed y and $y \in \mathcal{Y}$, $q_{y|xy}$ takes the same value $\forall x \in \Omega_X$.

Illustrated by the example:

- The probability of $\{<, \geq\}$ is taken to be **independent of the receipt of the UBI** given y :

$$q_{\{<,\geq\}|\mathbf{0}<} = q_{\{<,\geq\}|\mathbf{1}<} \quad \text{and} \quad q_{\{<,\geq\}|\mathbf{0}\geq} = q_{\{<,\geq\}|\mathbf{1}\geq}$$




Subgroup independence (SI)

Subgroup independence (Plass, Augustin, Cattaneo, Schollmeyer, 2015):

For each fixed y and $y \in \mathcal{Y}$, $q_{y|xy}$ takes the same value $\forall x \in \Omega_X$.

Illustrated by the example:

- The probability of $\{<, \geq\}$ is taken to be **independent of the receipt of the UBI** given y :

$$q_{\{<,\geq\}|\mathbf{0}<} = q_{\{<,\geq\}|\mathbf{1}<} \quad \text{and} \quad q_{\{<,\geq\}|\mathbf{0}\geq} = q_{\{<,\geq\}|\mathbf{1}\geq}$$


- Resulting estimators (if well-defined and inside $[0,1]$):

$$\hat{\pi}_{x<}^{(SI)} = \frac{n_{x\{<\}} v}{n_x w}, \quad \hat{q}_{\{<,\geq\}|\mathbf{x}<}^{(SI)} = 1 - \frac{w}{v}, \quad \hat{q}_{\{<,\geq\}|\mathbf{x}\geq}^{(SI)} = 1 - \frac{w}{z}$$

with $v = n_{0\{>= \}} n_{1\{>= \}} - n_{0\{>= \}} n_{1\{< \}}$, $w = n_{0\{< \}} n_{1\{>= \}} - n_{0\{>= \}} n_{1\{< \}}$
and $z = n_{0\{< \}} n_{1\{< \}} - n_{1\{< \}} n_{0\{< \}}$

Estimators for subgroup 0:

$$\hat{\pi}_{0<} \in [0.07, 0.26], \quad \hat{q}_{\{<,\geq\}|0<} \in [0, 0.71], \quad \hat{q}_{\{<,\geq\}|0\geq} \in [0, 0.198]$$

$$\hat{\pi}_{0<}^{(SI)} = 0.42, \quad \hat{q}_{\{<,\geq\}|0<}^{(SI)} = -0.04, \quad \hat{q}_{\{<,\geq\}|0\geq}^{(SI)} = 0.20$$

Estimators for subgroup 0:

$$\hat{\pi}_{0<} \in [0.07, 0.26], \quad \hat{q}_{\{<,\geq\}|0<} \in [0, 0.71], \quad \hat{q}_{\{<,\geq\}|0\geq} \in [0, 0.198]$$

$$\hat{\pi}_{0<}^{(SI)} = 0.42, \quad \hat{q}_{\{<,\geq\}|0<}^{(SI)} = -0.04, \quad \hat{q}_{\{<,\geq\}|0\geq}^{(SI)} = 0.20$$

- \Rightarrow There are data situations that might hint to (partial) incompatibility with SI
- \Rightarrow SI is testable in our setting

- Number of **degrees of freedom** under SI:

$$\begin{aligned}df^{\text{SI}} &= \dim(\Theta_{\text{obs}}) - \dim(\Theta_{\text{lat, SI}}) \\ &= k \cdot (2^m - 2) - [k \cdot (m - 1) + m \cdot (2^{m-1} - 1)]\end{aligned}$$

with k as the number of subgroups and $m = |\Omega_Y|$

- Point-identification and testability are only valid if **sufficient subgroups** are available inducing $df \geq 0$

m	2	3,4,5	6,7	8,9	...	20	...	50	...
minimum k	2	3	4	5	...	11	...	26	...

- Hypotheses:

H_0 : $q_{y|xy} = q_{y|x'y}$ for all $y \in \Omega_Y$, $x, x' \in \Omega_X$, $y \in \Omega_Y$,

H_1 : $q_{y|xy} \neq q_{y|x'y}$ for some $y \in \Omega_Y$, $x, x' \in \Omega_X$, $y \in \Omega_Y$.

- Hypotheses:

H_0 : $q_{y|xy} = q_{y|x'y}$ for all $y \in \Omega_Y$, $x, x' \in \Omega_X$, $y \in \Omega_Y$,

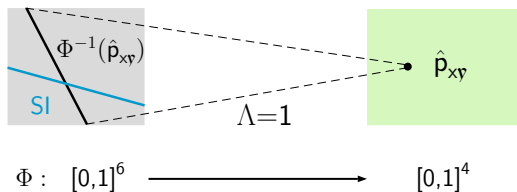
H_1 : $q_{y|xy} \neq q_{y|x'y}$ for some $y \in \Omega_Y$, $x, x' \in \Omega_X$, $y \in \Omega_Y$.

- Test based on **test statistic** (Wilks, 1938)

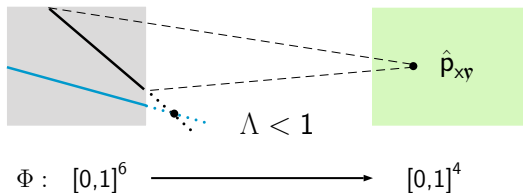
$$T = -2 \cdot \ln(\Lambda(y_1, \dots, y_n, x_1, \dots, x_n)) \quad \text{with}$$

$$\Lambda(y_1, \dots, y_n, x_1, \dots, x_n) = \frac{\sup_{H_0} L(\vartheta | y_1, \dots, y_n, x_1, \dots, x_n)}{\sup_{H_0 \cup H_1} L(\vartheta | y_1, \dots, y_n, x_1, \dots, x_n)} .$$

- No evidence to reject SI



- Some evidence to reject SI



- Asymptotic distribution of T under H_0 :

- $\frac{1}{2}\delta_0 + \frac{1}{2}\chi_1^2$ if $df^{SI} = 0$
- $\chi_{df^{SI}}^2$ if $df^{SI} > 0$,

where δ_0 is the Dirac distribution at 0 (Chernoff, 1954)

- Critical values for test decision:

- $\chi_{1,1-2\cdot\alpha}^2$ if $df^{SI}=0$
- $\chi_{df^{SI},1-\alpha}^2$ if $df^{SI} > 0$,

- Asymptotic distribution of T under H_0 :

- $\frac{1}{2}\delta_0 + \frac{1}{2}\chi_1^2$ if $df^{SI} = 0$
- $\chi_{df^{SI}}^2$ if $df^{SI} > 0$,

where δ_0 is the Dirac distribution at 0 (Chernoff, 1954)

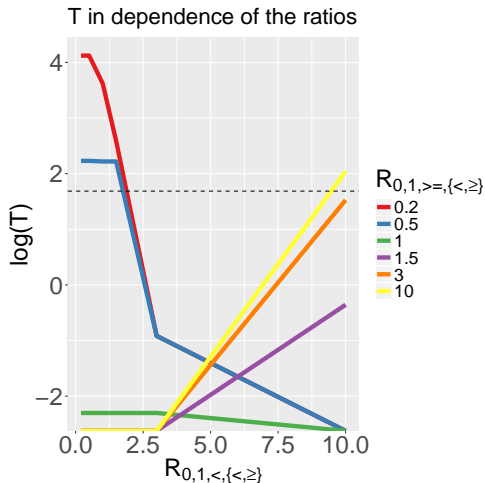
- Critical values for test decision:

- $\chi_{1,1-2\cdot\alpha}^2$ if $df^{SI} = 0$
- $\chi_{df^{SI},1-\alpha}^2$ if $df^{SI} > 0$,

- In data example (with $\alpha = 0.01$):

$$T = 0.14 < 5.4 = \chi_{1,1-2\cdot 0.01}^2 \Rightarrow H_0 \text{ can not be rejected}$$

LR-Test for generalized SI



with ratios

$$R_{0,1,<,>} = \frac{q_{\{<,\ge\}}|0<}{q_{\{<,\ge\}}|1<}$$

and

$$R_{0,1,\ge,>} = \frac{q_{\{<,\ge\}}|0\ge}{q_{\{<,\ge\}}|1\ge}$$

Summary: CAR versus SI

	CAR	SI
Point-identifying?	always	in specific settings
Testability	generally impossible	possible in specific settings

Construction of a hypothesis test:

- $H_0: (g)SI, H_1: \text{no } (g)SI$
- Test statistic based on the Likelihood Ratio

- global log-likelihood

$$l(\Phi^{-1}(p_{00<}, p_{00\geq}, \dots)) = l(\pi_{00<}, q_{na|00<}, q_{na|00\geq}, \dots)$$

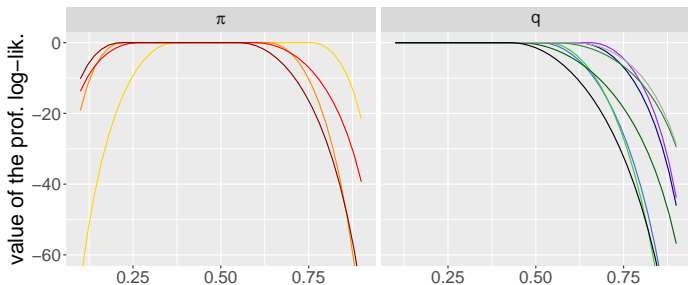
(Relative) profile log-likelihood

- global log-likelihood

$$l(\Phi^{-1}(p_{00<}, p_{00\geq}, \dots)) = l(\pi_{00<}, q_{na|00<}, q_{na|00\geq}, \dots)$$

- profile log-likelihood with nuisance param. $\xi = \vartheta \setminus \pi_{00<}$ or $\xi = \vartheta \setminus q_{na|00<}$

$$l(\pi_{00<}) = \max_{\xi} l(\pi_{00<}, \xi) \quad \text{or} \quad l(q_{y|xy}) = \max_{\xi} l(q_{na|00<}, \xi)$$



What's the problem?

Common dealing with incomplete data: assumptions

⇒ Missing at random (MAR) / coarsening at random (CAR)

⇒ Frequently: assumptions only for pragmatic reasons



Here:

- 1.) Reliable likelihood inference
- 2.) Testability of coarsening assumptions
- 3.) **Reliable regression estimators**

Logit model with linear predictor

$$\eta_{\mathbf{x}} = \beta_0 + \beta_1 \cdot x_{\text{Abitur}} + \beta_2 \cdot x_{\text{age}} + \beta_{12} \cdot x_{\text{Abitur}} x_{\text{age}} :$$

Response function h

$$\pi_{\mathbf{x}<} = h(\eta_{\mathbf{x}}) = \frac{\exp(\eta_{\mathbf{x}})}{1 + \exp(\eta_{\mathbf{x}})}$$

Link function $g = h^{-1}$

$$\text{and } \eta_{\mathbf{x}} = g(\pi_{\mathbf{x}<}) = \ln\left(\frac{\pi_{\mathbf{x}<}}{1 - \pi_{\mathbf{x}<}}\right).$$

Reliable regression estimators (Direct method)

Logit model with linear predictor

$$\eta_{\mathbf{x}} = \beta_0 + \beta_1 \cdot x_{\text{Abitur}} + \beta_2 \cdot x_{\text{age}} + \beta_{12} \cdot x_{\text{Abitur}} x_{\text{age}} :$$

Response function h

$$\pi_{\mathbf{x}<} = h(\eta_{\mathbf{x}}) = \frac{\exp(\eta_{\mathbf{x}})}{1 + \exp(\eta_{\mathbf{x}})}$$

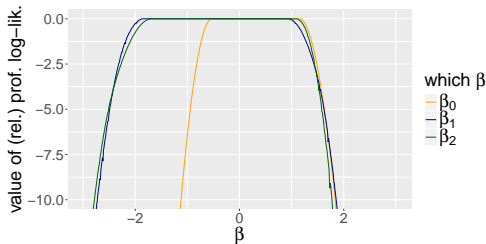
Link function $g = h^{-1}$

$$\text{and } \eta_{\mathbf{x}} = g(\pi_{\mathbf{x}<}) = \ln\left(\frac{\pi_{\mathbf{x}<}}{1 - \pi_{\mathbf{x}<}}\right).$$

⇒ global likelihood

$$l(\beta_0, \beta_1, \dots, q_{na|00<}, \dots)$$

⇒ Consider (rel.)
profile log-lik. for
each β :



No parametric assumption on the regression model:

- Saturated model, all interactions included
- $\ln\left(\frac{\pi_{x<}}{1-\pi_{x<}}\right) = \beta_0 + \beta_1 \cdot x_{Abitur} + \beta_2 \cdot x_{age} + \beta_{12}x_{Abitur}x_{age}$
- $\pi_{x<}$ and regression coefficients basically represent same information

Parametric assumption on the regression model:

- (some) interactions are equal to 0 (e.g. $\beta_{12} = 0$)
- $\ln\left(\frac{\pi_{x<}^*}{1-\pi_{x<}^*}\right) = \beta_0 + \beta_1 \cdot x_{Abitur} + \beta_2 \cdot x_{age}$
- $\pi_{x<}$ and regression coefficients do not represent same information

Aim:

We want to study the impact of the parametric assumption on the regression model in the coarse data problem

Two-step method:

- 1.) Estimate the bounds $\hat{\pi}_{\mathbf{x}<}$ and $\overline{\hat{\pi}}_{\mathbf{x}<}$ (no param. assump.)
- 2.) Use these bounds to obtain the reliable regression estimators

No parametric assumption:

Function h is bijective \Rightarrow transform $\hat{\pi}_{\mathbf{x}<}$ and $\overline{\hat{\pi}}_{\mathbf{x}<}$ via link function

Parametric assumption: solve an optimization problem

$\beta_1 \rightarrow \min/\max$ given

$$\hat{\pi}_{00<} \leq \frac{\exp(\beta_0)}{1 + \exp(\beta_0)} \leq \bar{\pi}_{00<},$$

$$\hat{\pi}_{10<} \leq \frac{\exp(\beta_0 + \beta_1)}{1 + \exp(\beta_0 + \beta_1)} \leq \bar{\pi}_{10<},$$

$$\hat{\pi}_{01<} \leq \frac{\exp(\beta_0 + \beta_2)}{1 + \exp(\beta_0 + \beta_2)} \leq \bar{\pi}_{01<},$$

$$\hat{\pi}_{11<} \leq \frac{\exp(\beta_0 + \beta_1 + \beta_2)}{1 + \exp(\beta_0 + \beta_1 + \beta_2)} \leq \bar{\pi}_{11<}.$$

Different types of results

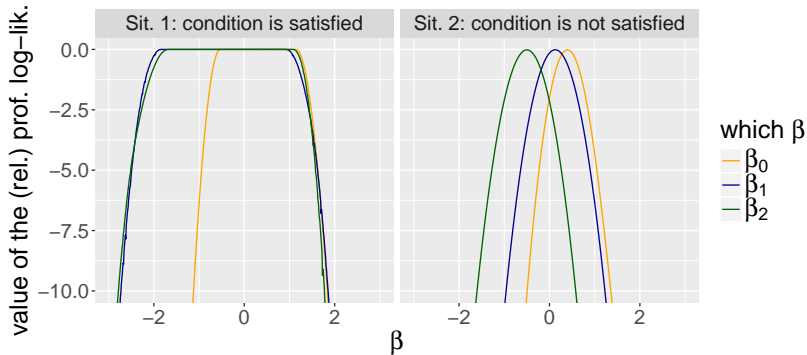
1.) There is a solution.

- Regression estimators obtainable that produce $\hat{\pi}_{\mathbf{x}<}$ and $\bar{\pi}_{\mathbf{x}<}$
 \Rightarrow **No impact on coarsening** (same $\hat{q}_{\{<,\geq\}|\mathbf{x}<}$ and $\bar{q}_{\{<,\geq\}|\mathbf{x}<}$)
- The resulting regression estimators can only represent tighter bounds of the estimated latent variable distribution
 \Rightarrow **Tighter bounds** compared to $\hat{q}_{\{<,\geq\}|\mathbf{x}<}$ and $\bar{q}_{\{<,\geq\}|\mathbf{x}<}$

2.) There is no solution.

Impact on the compatible coarsening scenarios?

Relative profile log-likelihood for the regression coefficients:



Results and conclusion

Reliable estimation and estimation under CAR (MAR)

- No parametric assumption on the regression model

$\hat{\beta}_0 \in [-0.53, 1.15]$	$\hat{\beta}_1 \in [-2.16, 0.92]$	$\hat{\beta}_2 \in [-2.42, 1.08]$	$\hat{\beta}_{12} \in [-2.76, 3.64]$
$\hat{\beta}_0 = 0.43$	$\hat{\beta}_1 = -0.85$	$\hat{\beta}_2 = -0.94$	$\hat{\beta}_{12} = 0.63$

- Parametric assumption on the regression model

$\hat{\beta}_0 \in [-0.53, 1.15]$	$\hat{\beta}_1 \in [-1.84, 0.92]$	$\hat{\beta}_2 \in [-1.68, 1.08]$
$\hat{\beta}_0 = 0.35$	$\hat{\beta}_1 = 0.05$	$\hat{\beta}_2 = 0.00$

Conclusions:

- Two methods for a reliable regression estimators have been studied
- Principally different impact on the estimated coarsening parameters
- Compatibility with observed data can be “repaired” in the presence of coarse data

Some important references



Couso, Dubois.

Statistical Reasoning with Set-Valued Information: Ontic vs. Epistemic Views, IJAR, 2014.



Chernoff.

On the distribution of the likelihood ratio, Ann. Stat. Math., 1954.



Heitjan, Rubin.

Ignorability and Coarse Data, Annals of Statistics, 1991.



Jaeger.

On testing the missing at random assumption, ECML, 2006.



Manski.

Partial Identification of Probability Distributions, Springer, 2003.



Nordheim.

Inference from nonrandomly missing categorical data: an example from a genetic study on Turner's syndrome, JASA, 1984.








Trappmann, Gundert, Wenzig, Gebhardt.

PASS: a household panel survey for research on unemployment and poverty, Schmollers Jahrbuch, 2010.



Wilks.

The large-sample distribution of the likelihood ratio for testing composite hypotheses, Ann. Stat., 1938.

-  Plass, Augustin, Cattaneo, Schollmeyer.
Statistical modelling under epistemic data imprecision, ISIPTA, 2015.
-  Plass, Cattaneo, Schollmeyer, Augustin.
On the testability of coarsening assumptions: A hypothesis test for subgroup independence, IJAR, 2017.
-  Plass, Cattaneo, Schollmeyer, Augustin.
Towards a reliable categorical regression analysis for nonrandomly coarsened observations: An analysis with German labour market data, under Review.
-  Plass, Omar, Augustin.
Towards a cautious modelling of missing data in Small Area Estimation, ISIPTA, 2017.
-  Plass, Fink, Schöning, Augustin.
Statistical modelling in surveys without neglecting The Undecided: Multinomial logistic regression models and imprecise classification trees under ontic data imprecision, ISIPTA, 2015.