

# Reliable analysis of categorical data under epistemic imprecision

CFE-CMStatistics 2015

Julia Plass\*, Thomas Augustin\*,  
Marco Cattaneo\*\*, Georg Schollmeyer\*

\*Department of Statistics, Ludwigs-Maximilians University and

\*\*Department of Mathematics, University of Hull

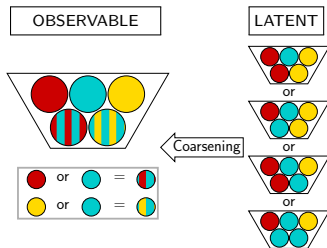


13<sup>th</sup> of December 2015

## Coarse data

*Data are not observed in the resolution originally intended*  
(epistemic vs. ontic interpretation)

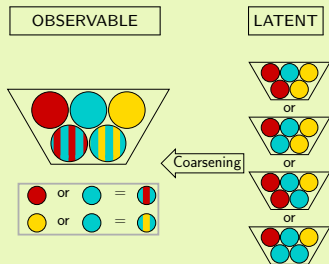
- Here: coarse data “=” data under epistemic imprecision
- Imprecise observation of something precise:



- 1 Where do data under epistemic imprecision typically arise?
- 2 How to deal with data under epistemic imprecision?
- 3 How to incorporate auxiliary information?
- 4 Are there possibilities to test on point identifying assumptions (coarsening at random, subgroup independence)?

## Epistemic imprecision:

*"Imprecise observation of something precise"*



⇒ Truth is hidden due to the underlying coarsening mechanism

## Examples:

- Matched data sets with partially overlapping variables
- Coarsening as anonymization technique
- Missing data as special case

**Here:** PASS-data

$\mathcal{Y}$ : income,  $\mathcal{X}$ : UBIII

$\Omega_{\mathcal{Y}} = \{<, \geq, na\}$

$\Omega_{\mathcal{X}} = \{0 \text{ (no)}, 1 \text{ (yes)}\}$

- Still common to **enforce precise results**
- Variety of **set-valued approaches**
  - via random sets  
(e.g. Nguyen, 2006, An Introduction to Random Sets)
  - using Bayesian approaches  
(e.g. de Cooman, Zaffalon, 2004, Artif. Intell.)
  - via likelihood-based belief function (Denœux, 2014, IJAR)
  - via profile likelihood  
(Cattaneo, Wiencierz, 2012, IJAR)

**Here:** Likelihood-based approach  
influenced by methodology of partial identification  
(Manski, 2003, Partial Identification of Probability Distributions)  
coarse categorical data

OBSERVABLE

coarse data

$\mathcal{Y}$

$$p_{\mathcal{Y}} = P(\mathcal{Y}_i = \mathcal{Y})$$

LATENT

latent variable

$Y$

$$\pi_{ij} = P(Y_i = j)$$

OBSERVABLE

coarse data  
 $\mathcal{Y}$

$$p_{\mathcal{Y}} = P(\mathcal{Y}_i = \mathcal{Y})$$

Main goal:

Maximum-Likelihood estimation of

Observation model  $Q$

error-freeness

$$q_{\mathcal{Y}|y} = P(\mathcal{Y} = \mathcal{Y} | Y = y)$$

LATENT

latent variable  
 $Y$

$$\pi_{ij} = P(Y_i = j)$$

$$\boldsymbol{\gamma} = (\mathbf{q}_{\mathcal{Y}|y}^T, \boldsymbol{\pi}_y^T)^T$$

# Basic problem (regression case)

OBSERVABLE

coarse data

$\mathcal{Y}$

$$p_{x_{\mathcal{Y}}} = P(\mathcal{Y} = \mathcal{Y} | X = x)$$

Main goal:

Maximum-Likelihood estimation of

$$\gamma = (\mathbf{q}_{\mathcal{Y}|\mathbf{x}}^T, \boldsymbol{\pi}_{\mathcal{Y}}^T)^T$$

LATENT

latent variable

$Y$

for  $j=1, \dots, K-1$

$$\pi_{ij} = P(Y_i = j | \mathbf{x}_i)$$

$$= \frac{\exp(\beta_{j0} + \mathbf{x}_i^T \boldsymbol{\beta}_j)}{1 + \sum_{s=1}^{K-1} \exp(\beta_{s0} + \mathbf{x}_i^T \boldsymbol{\beta}_s)}$$

for reference category  $K$

$$\pi_{iK} = \frac{1}{1 + \sum_{s=1}^{K-1} \exp(\beta_{s0} + \mathbf{x}_i^T \boldsymbol{\beta}_s)}$$

(multinomial logit model)

Observation model  $Q$

error-freeness

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



## OBSERVABLE

Use random-set perspective and determine ML estimator

$$\hat{p}_{\mathcal{Y}} = \hat{P}(\mathcal{Y} = \mathcal{y})$$

$$\rightarrow \hat{p}_{\mathcal{Y}} = \frac{n_{\mathcal{Y}}}{n}$$

Use the **connection** between  $\mathbf{p}$  and  $\boldsymbol{\gamma}$

$$\Phi(\boldsymbol{\gamma}) = \mathbf{p}$$

and the **invariance of the likelihood** under parameter transformations:

$$\hat{\Gamma} = \{\boldsymbol{\gamma} \mid \Phi(\boldsymbol{\gamma}) = \hat{\mathbf{p}}\}$$

## LATENT

$$\boldsymbol{\gamma} = (\mathbf{q}_{\mathcal{Y}|y}^T, \boldsymbol{\pi}_y^T)^T$$

$$\hat{\pi}_y \in \left[ \frac{n_{\{y\}}}{n}, \frac{\sum_{\mathcal{Y} \ni y} n_{\mathcal{Y}}}{n} \right]$$

$$\hat{q}_{\mathcal{Y}|y} \in \left[ 0, \frac{n_{\mathcal{Y}}}{n_{\{y\}} + n_{\mathcal{Y}}} \right]$$

## OBSERVABLE

Use random-set perspective  
and determine ML estimator

$$\hat{p}_y = \hat{P}(\mathcal{Y} = y)$$

$$\rightarrow \hat{p}_y = \frac{n_{\cdot y}}{n}$$

Use the **connection**  
**between**  $p$  and  $\gamma$

$$\Phi(\gamma) = p$$

and the **invariance of**  
**the likelihood** under  
parameter transformations:

$$\hat{\Gamma} = \{\gamma \mid \Phi(\gamma) = \hat{p}\}$$

## LATENT

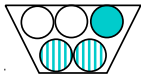
$$\gamma = (\mathbf{q}_{y|y}^T, \boldsymbol{\pi}_y^T)^T$$

$$\hat{\pi}_y \in \left[ \frac{n_{\{y\}}}{n}, \frac{\sum_{\mathcal{V} \ni y} n_{\mathcal{V}}}{n} \right]$$

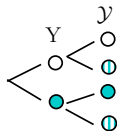
$$\hat{q}_{\mathcal{V}|y} \in \left[ 0, \frac{n_{\mathcal{V}}}{n_{\{y\}} + n_{\mathcal{V}}} \right]$$

OBSERVATION ( $\mathcal{Y}$ )

$p$

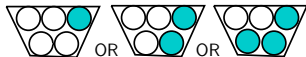


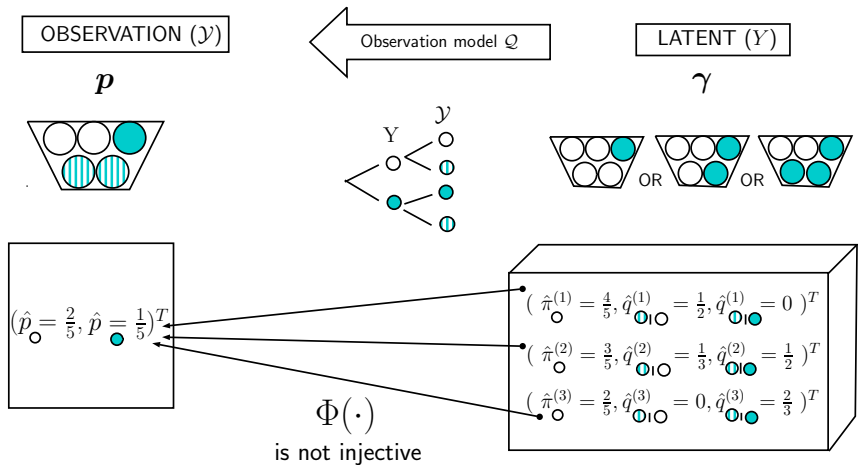
← Observation model  $\mathcal{Q}$

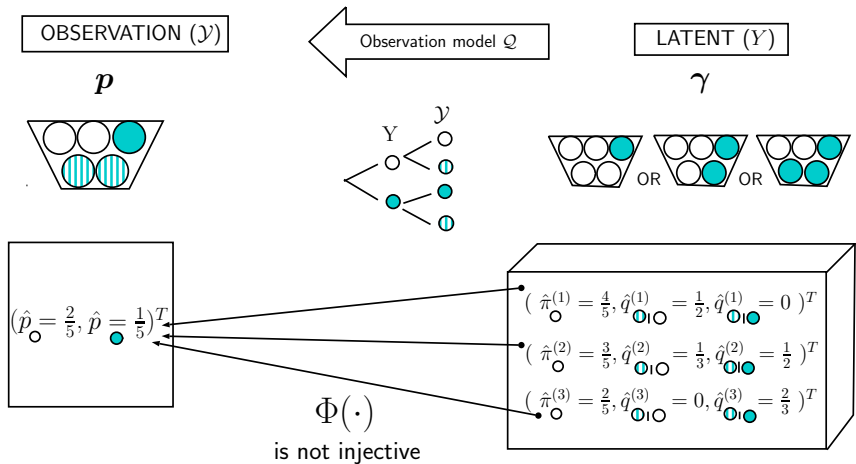


LATENT ( $\mathcal{Y}$ )

$\gamma$







$$\begin{pmatrix} \hat{p}_{\circ} \\ \hat{p}_{\bullet} \end{pmatrix} = \Phi \begin{pmatrix} \hat{\pi}_{\circ} \\ \hat{q}_{\circ|\circ} \\ \hat{q}_{\circ|\bullet} \end{pmatrix} = \begin{pmatrix} \hat{\pi}_{\circ} \cdot (1 - \hat{q}_{\circ|\circ}) \\ (1 - \hat{\pi}_{\circ}) \cdot (1 - \hat{q}_{\circ|\bullet}) \end{pmatrix}$$

## OBSERVABLE

Use random-set perspective  
and determine ML estimator

$$\hat{p}_{\mathcal{Y}} = \hat{P}(\mathcal{Y} = \mathcal{y})$$

$$\rightarrow \hat{p}_{\mathcal{Y}} = \frac{n_{\mathcal{Y}}}{n}$$

Use the **connection**  
**between**  $\mathbf{p}$  and  $\boldsymbol{\gamma}$

$$\Phi(\boldsymbol{\gamma}) = \mathbf{p}$$

and the **invariance of**  
**the likelihood** under  
parameter transformations:

$$\hat{\Gamma} = \{\boldsymbol{\gamma} \mid \Phi(\boldsymbol{\gamma}) = \hat{\mathbf{p}}\}$$

## LATENT

$$\boldsymbol{\gamma} = (\mathbf{q}_{\mathcal{Y}|\mathbf{y}}^T, \boldsymbol{\pi}_{\mathbf{y}}^T)^T$$

$$\hat{\pi}_{\mathbf{y}} \in \left[ \frac{n_{\{\mathbf{y}\}}}{n}, \frac{\sum_{\mathcal{Y} \ni \mathbf{y}} n_{\mathcal{Y}}}{n} \right]$$

$$\hat{q}_{\mathcal{Y}|\mathbf{y}} \in \left[ 0, \frac{n_{\mathcal{Y}}}{n_{\{\mathbf{y}\}} + n_{\mathcal{Y}}} \right]$$

## OBSERVABLE

Use random-set perspective  
and determine ML estimator

$$\hat{p}_{\mathcal{Y}} = \hat{P}(\mathcal{Y} = \mathcal{y})$$

$$\rightarrow \hat{p}_{\mathcal{Y}} = \frac{n_{\mathcal{Y}}}{n}$$

Use the **connection**  
**between**  $\mathbf{p}$  and  $\boldsymbol{\gamma}$

$$\Phi(\boldsymbol{\gamma}) = \mathbf{p}$$

and the **invariance of**  
**the likelihood** under  
parameter transformations:

$$\hat{\Gamma} = \{\boldsymbol{\gamma} \mid \Phi(\boldsymbol{\gamma}) = \hat{\mathbf{p}}\}$$

## LATENT

$$\boldsymbol{\gamma} = (\mathbf{q}_{\mathcal{Y}|y}^T, \boldsymbol{\pi}_y^T)^T$$

$$\hat{\pi}_y \in \left[ \frac{n_{\{y\}}}{n}, \frac{\sum_{\mathcal{Y} \ni y} n_{\mathcal{Y}}}{n} \right]$$

$$\hat{q}_{\mathcal{Y}|y} \in \left[ 0, \frac{n_{\mathcal{Y}}}{n_{\{y\}} + n_{\mathcal{Y}}} \right]$$

## OBSERVABLE

Use random-set perspective and determine ML estimator

$$\hat{p}_{\mathcal{Y}} = \hat{P}(\mathcal{Y} = \mathcal{y})$$

$$\rightarrow \hat{p}_{\mathcal{Y}} = \frac{n_{\mathcal{Y}}}{n}$$

## LATENT

Use the **connection** between  $\mathbf{p}$  and  $\boldsymbol{\gamma}$

$$\Phi(\boldsymbol{\gamma}) = \mathbf{p}$$

$$\boldsymbol{\gamma} = (\mathbf{q}_{\mathcal{Y}|\mathcal{Y}}^T, \boldsymbol{\pi}_{\mathcal{Y}}^T)^T$$

and the **invariance of the likelihood** under parameter transformations:

$$\hat{\Gamma} = \{\boldsymbol{\gamma} \mid \Phi(\boldsymbol{\gamma}) = \hat{\mathbf{p}}\}$$

$$\hat{\pi}_{\mathcal{Y}} \in \left[ \frac{n_{\{\mathcal{Y}\}}}{n}, \frac{\sum_{\mathcal{Y} \ni \mathcal{Y}} n_{\mathcal{Y}}}{n} \right]$$

$$\hat{q}_{\mathcal{Y}|\mathcal{Y}} \in \left[ 0, \frac{n_{\mathcal{Y}}}{n_{\{\mathcal{Y}\}} + n_{\mathcal{Y}}} \right]$$

### Illustration (PASS data, wave 1)

$$n_{<} = 238, \quad n_{\geq} = 835, \quad n_{\text{na}} = 338$$

$$\hat{\pi}_{<} \in \left[ \frac{238}{1411}, \frac{238+338}{1411} \right]$$



## OBSERVABLE

Use random-set perspective and determine ML estimator

$$\hat{p}_{x\mathcal{Y}} = \hat{P}(\mathcal{Y} = \mathcal{y} | X = x)$$

$$\rightarrow \hat{p}_{x\mathcal{Y}} = \frac{n_{x\mathcal{Y}}}{n_x}$$

Use the **connection** between  $\mathbf{p}$  and  $\boldsymbol{\gamma}$

$$\Phi(\boldsymbol{\gamma}) = \mathbf{p}$$

and the **invariance of the likelihood** under parameter transformations:

$$\hat{\Gamma} = \{\boldsymbol{\gamma} \mid \Phi(\boldsymbol{\gamma}) = \hat{\mathbf{p}}\}$$

## LATENT

$$\boldsymbol{\gamma} = (\mathbf{q}_{\mathcal{Y}|x\mathcal{Y}}^T, \boldsymbol{\pi}_{x\mathcal{Y}}^T)^T$$

$$\hat{\pi}_{x\mathcal{Y}} \in \left[ \frac{n_{x\{\mathcal{Y}\}}}{n_x}, \frac{\sum_{\mathcal{Y} \ni \mathcal{Y}} n_{x\mathcal{Y}}}{n_x} \right]$$

$$\hat{q}_{\mathcal{Y}|x\mathcal{Y}} \in \left[ 0, \frac{n_{x\mathcal{Y}}}{n_{x\{\mathcal{Y}\}} + n_{x\mathcal{Y}}} \right]$$



### Illustration (PASS data, wave 1)

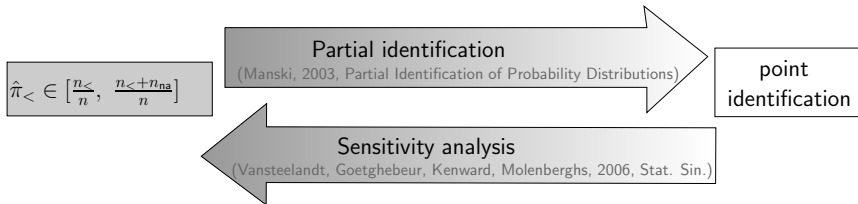
$$\hat{\pi}_{0<} \in [0.41, 0.64]$$

$$\hat{\pi}_{1<} \in [0.10, 0.34]$$

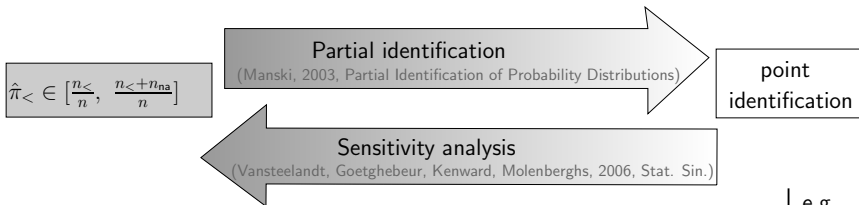
$$\hat{\beta}_{<0} \in [-0.37, 0.59]$$

$$\hat{\beta}_{<} \in [-1.83, -1.25]$$

# Reliable incorporation of auxiliary information



# Reliable incorporation of auxiliary information



↓ e.g.

$$R = \frac{q_{na \geq}}{q_{na <}}$$

(Nordheim, 1984,  
J. Am. Stat. Assoc.)

CAR:  $R=1$

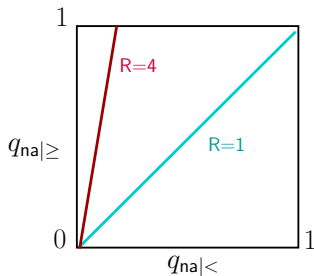
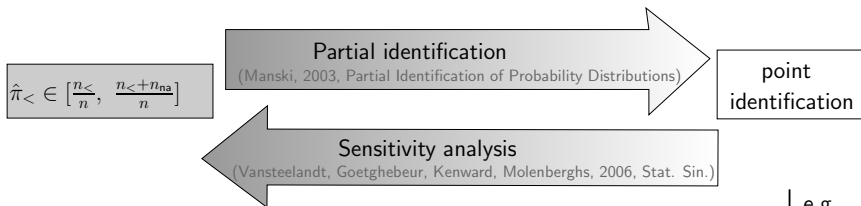
↓  $R=1$

$$\hat{\pi}_{<} = \frac{n_{<}}{n_{<} + n_{\geq}}$$

Auxiliary information:

Refined estimators:

# Reliable incorporation of auxiliary information



↓ e.g.

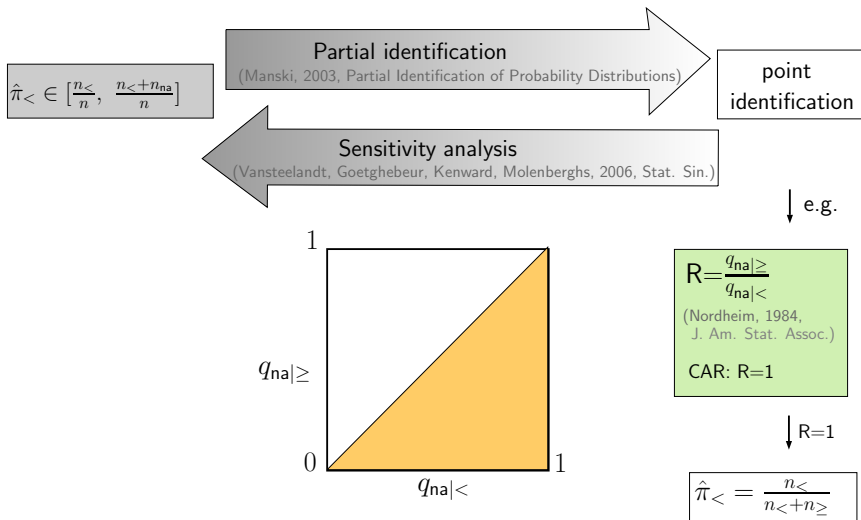
$$R = \frac{q_{na|>}}{q_{na|<}}$$

(Nordheim, 1984, J. Am. Stat. Assoc.)  
CAR:  $R=1$

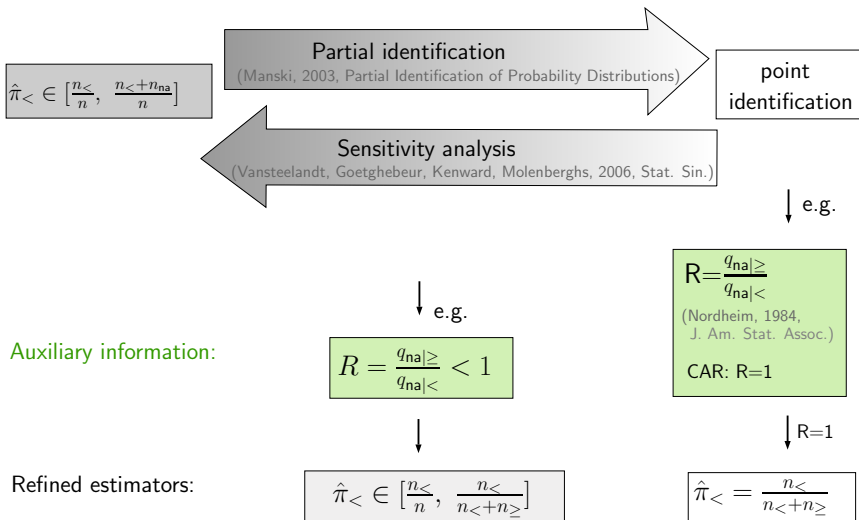
↓  $R=1$

$$\hat{\pi}_< = \frac{n_{<}}{n_{<} + n_{\geq}}$$

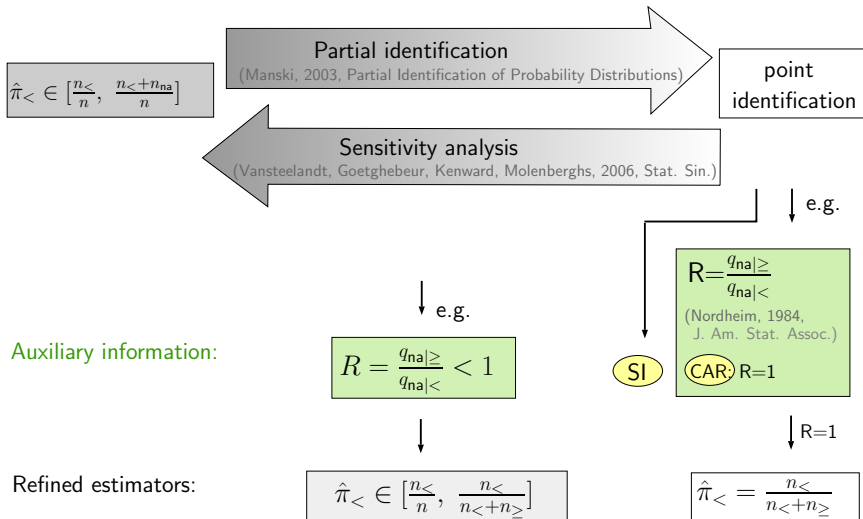
# Reliable incorporation of auxiliary information



# Reliable incorporation of auxiliary information



# Reliable incorporation of auxiliary information



# Coarsening at random & subgroup independence

## coarsening at random (CAR)

(Heitjan, Rubin, 1991, Ann. Stat.)

## subgroup independence (SI)

### Generally

For each fixed  $\mathcal{Y}$ ,  
 $q_{\mathcal{Y}|y}$  takes the same values for  
all  $y$  that are consistent with  $\mathcal{Y}$

Coarsening does not  
depend on the value  
of the covariate

### Example

$$q_{na|<} = q_{na|\geq}$$

$$q_{na|0<} = q_{na|1<} \\ \text{and}$$

$$q_{na|0\geq} = q_{na|1\geq}$$



# Coarsening at random & subgroup independence

|   | <  | ≥   | na | total |
|---|----|-----|----|-------|
| 0 | 38 | 385 | 95 | 518   |
| 1 | 36 | 42  | 9  | 87    |

Table: PASS data, wave 5

$$\begin{aligned}\hat{\pi}_{0<} &\in [0.07; 0.26] & \hat{\pi}_{1<} &\in [0.41; 0.52] \\ \hat{q}_{na|0<} &\in [0; 0.71] & \hat{q}_{na|1<} &\in [0; 0.2] \\ \hat{q}_{na|0\geq} &\in [0; 0.20] & \hat{q}_{na|1\geq} &\in [0; 0.18]\end{aligned}$$

$X=0$

$$\begin{aligned}\hat{\pi}_{0<} &= 0.07 \\ \hat{q}_{na|0<} &= 0 \\ \hat{q}_{na|0\geq} &= 0.20\end{aligned}$$

...

$$\begin{aligned}\hat{\pi}_{0<}^{(CAR)} &= 0.09 \\ \hat{q}_{na|0<}^{(CAR)} &= 0.18 \\ \hat{q}_{na|0\geq}^{(CAR)} &= 0.18\end{aligned}$$

...

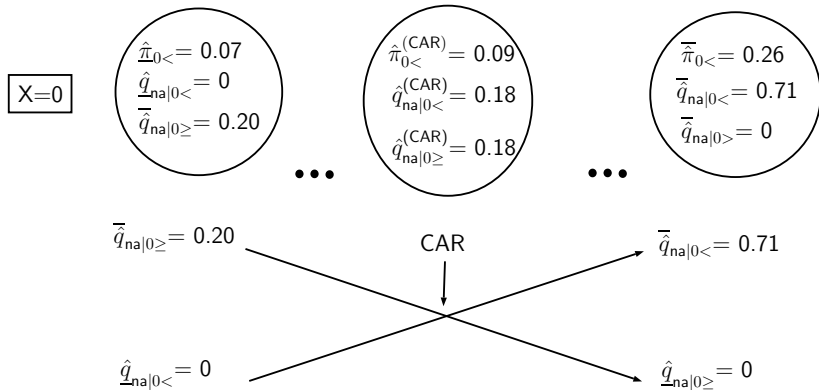
$$\begin{aligned}\hat{\pi}_{0<} &= 0.26 \\ \hat{q}_{na|0<} &= 0.71 \\ \hat{q}_{na|0>} &= 0\end{aligned}$$

# Coarsening at random & subgroup independence

|   | <  | ≥   | na | total |
|---|----|-----|----|-------|
| 0 | 38 | 385 | 95 | 518   |
| 1 | 36 | 42  | 9  | 87    |

Table: PASS data, wave 5

$$\begin{aligned} \hat{\pi}_{0<} &\in [0.07; 0.26] & \hat{\pi}_{1<} &\in [0.41; 0.52] \\ \hat{q}_{na|0<} &\in [0; 0.71] & \hat{q}_{na|1<} &\in [0; 0.2] \\ \hat{q}_{na|0\geq} &\in [0; 0.20] & \hat{q}_{na|1\geq} &\in [0; 0.18] \end{aligned}$$

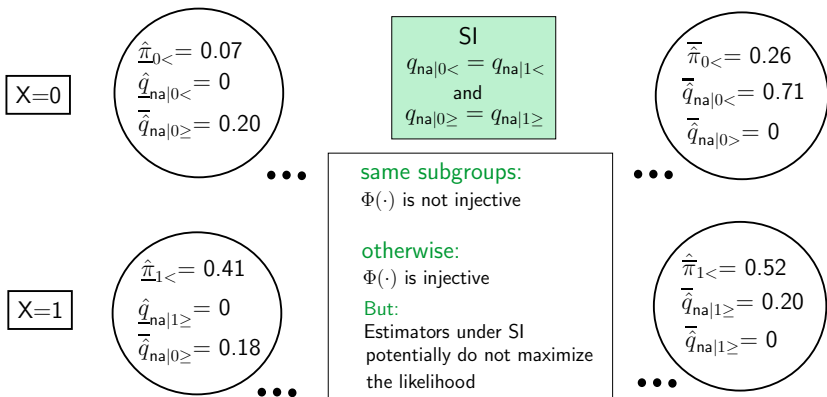


# Coarsening at random & subgroup independence

|   | <  | ≥   | na | total |
|---|----|-----|----|-------|
| 0 | 38 | 385 | 95 | 518   |
| 1 | 36 | 42  | 9  | 87    |

Table: PASS data, wave 5

$$\begin{aligned} \hat{\pi}_{0<} &\in [0.07; 0.26] & \hat{\pi}_{1<} &\in [0.41; 0.52] \\ \hat{q}_{na|0<} &\in [0; 0.71] & \hat{q}_{na|1<} &\in [0; 0.20] \\ \hat{q}_{na|0\geq} &\in [0; 0.20] & \hat{q}_{na|1\geq} &\in [0; 0.18] \end{aligned}$$



# Coarsening at random & subgroup independence

|   | <  | ≥   | na | total |
|---|----|-----|----|-------|
| 0 | 38 | 385 | 95 | 518   |
| 1 | 36 | 42  | 9  | 87    |

Table: PASS data, wave 5

$$\begin{aligned}\hat{\pi}_{0<} &\in [0.07; 0.26] & \hat{\pi}_{1<} &\in [0.41; 0.52] \\ \hat{q}_{na|0<} &\in [0; 0.71] & \hat{q}_{na|1<} &\in [0; 0.20] \\ \hat{q}_{na|0\geq} &\in [0; 0.20] & \hat{q}_{na|1\geq} &\in [0; 0.18]\end{aligned}$$

CAR

SI

One can never  
reject CAR

$$\begin{aligned}\hat{\pi}_{1<}^{(SI)} &= \frac{n_{1<} n_{1\geq} n_{0<} - n_{1<} n_{0\geq} n_{1<}}{n_{1<} n_{0<} n_{1\geq} - n_{0\geq} n_{1<} n_{1<}} \\ &= 0.39 \quad (\text{cf. } \hat{\pi}_{1<} \in [0.41; 0.52])\end{aligned}$$

⇒ Construction of hypothesis test with

$$H_0 : q_{na|0<} = q_{na|1<} \quad \& \quad q_{na|0\geq} = q_{na|1\geq}$$







$$H_1 : \text{no restrictions on the coarsening parameters}$$

- Via the observation model  $Q$  maximum-likelihood estimators referring to the latent variable may be obtained for both cases
  - ... the homogeneous case
  - ... the case with categorical covariates
- Proper inclusion of auxiliary information via further restrictions on  $Q$

## Next steps:

- Likelihood-based hypothesis tests and uncertainty regions
- Comparison to Bayesian approaches
- Applying the observation model to coarse ordinal data

# References

-  Couso, Dubois, Sánchez.  
*Random Sets and Random Fuzzy Sets as Ill-Perceived Random Variables*, Springer, 2014.
-  Heitjan, Rubin.  
*Ignorability and Coarse Data*, *Annals of Statistics*, 1991.
-  Manski.  
*Partial Identification of Probability Distributions*, Springer, 2003.
-  E. Nordheim.  
*Inference from nonrandomly missing categorical data: An example from a genetic study on Turner's syndrome*, *J. Am. Stat. Assoc.*, 1984.
-  Plass, Augustin, Cattaneo, Schollmeyer.  
*Statistical modelling under epistemic data imprecision: some results on estimating multinomial distributions and logistic regression for coarse categorical data*, *ISIPTA*, 2015.
-  Vansteelandt, Goetghebeur, Kenward, Molenberghs.  
*Ignorance and uncertainty regions as inferential tools in a sensitivity analysis*, *Stat. Sin.*, 2006.