

# Cautious statistical modelling for categorical data under epistemic and ontic data imprecision

Julia Plass, Supervision: Prof. Thomas Augustin

Department of Statistics, Ludwig-Maximilians University

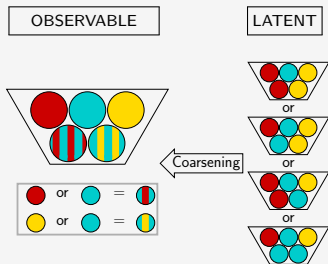


14<sup>th</sup> of September 2015

21. DStatG Nachwuchsworkshop, Hamburg

## Epistemic imprecision:

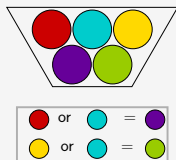
*"Imprecise observation of something precise"*



⇒ Truth is hidden due to the underlying coarsening mechanism

## Ontic imprecision:

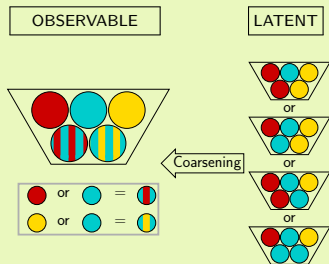
*"Precise observation of something imprecise"*



⇒ Truth is represented by coarse observation

## 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,  $X$ : UBI

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

$\Omega_X = \{0 \text{ (no)}, 1 \text{ (yes)}\}$

OBSERVABLE

coarse data  
 $\mathcal{Y}$

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

Observation model  $\mathcal{Q}$

error-freeness

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

LATENT

latent variable  
 $Y$

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

$$\begin{aligned} \pi_{ij} &= P(Y_i = j | \mathbf{x}_i) \\ &= \frac{\exp(\beta_{j0} + \mathbf{x}_i^T \beta_j)}{1 + \sum_{s=1}^{K-1} \exp(\beta_{s0} + \mathbf{x}_i^T \beta_s)} \end{aligned}$$

for reference category  $K$

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

(multinomial logit model)

Main goal:

Maximum-Likelihood estimation of  $\gamma = (\mathbf{q}_{\mathcal{Y}|xy}^T, \boldsymbol{\pi}_y^T)^T$

## OBSERVABLE

Use random-set perspective  
and determine ML estimator

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

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

## LATENT

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

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

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

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

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

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

OBSERVABLE

LATENT

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**  $p$  and  $\gamma$

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

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

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

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

## OBSERVABLE

Use random-set perspective  
and determine ML estimator

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

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

## LATENT

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

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

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

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

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

## OBSERVABLE

Use random-set perspective  
and determine ML estimator

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

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

## LATENT

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

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

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

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

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

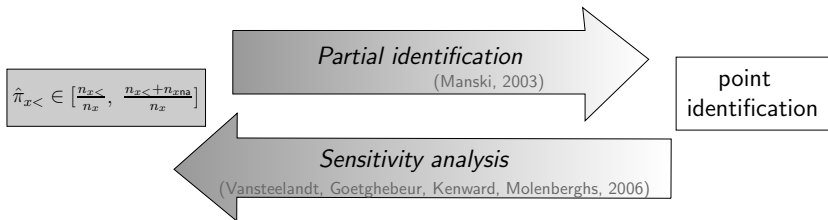


### Illustration (PASS data)

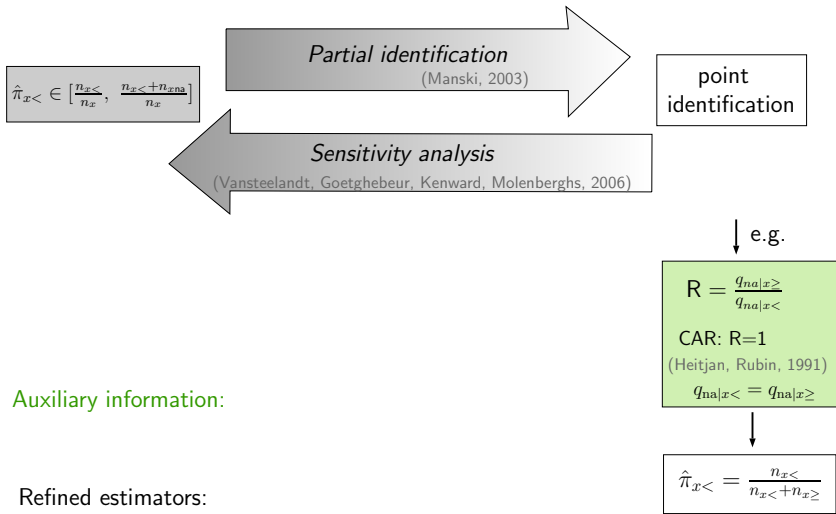
$$\hat{\pi}_{0<} \in [0.41, 0.64] \quad \hat{\pi}_{1<} \in [0.10, 0.34]$$
$$\hat{\beta}_{<} \in [-0.37, 0.59] \quad \hat{\beta}_{<} \in [-1.83, -1.25]$$



# Reliable incorporation of auxiliary information



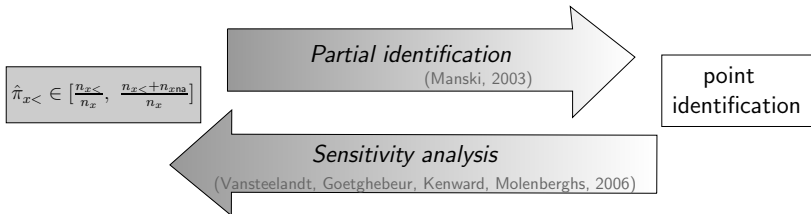
# Reliable incorporation of auxiliary information



Auxiliary information:

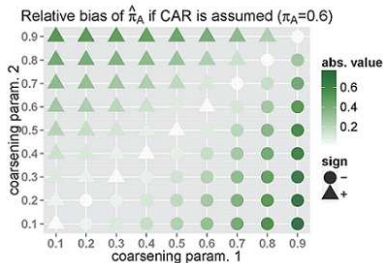
Refined estimators:

# Reliable incorporation of auxiliary information

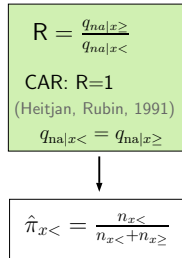


Auxiliary information:

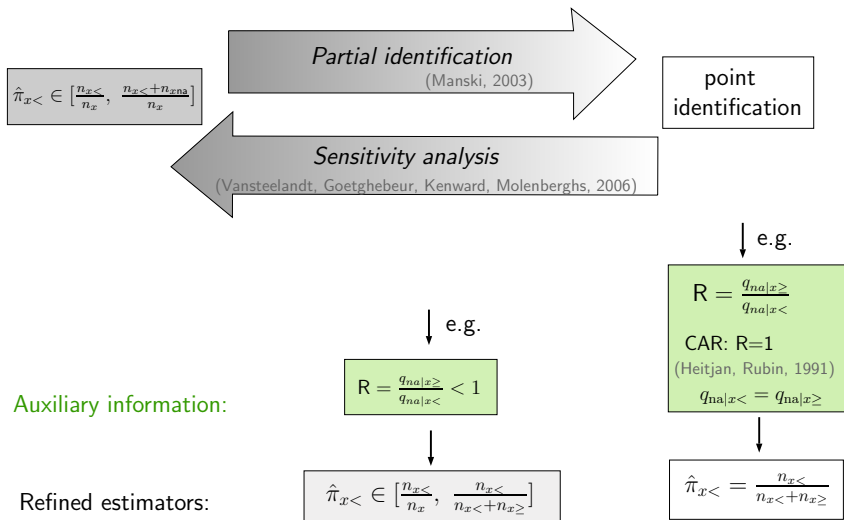
Refined estimators:



e.g.



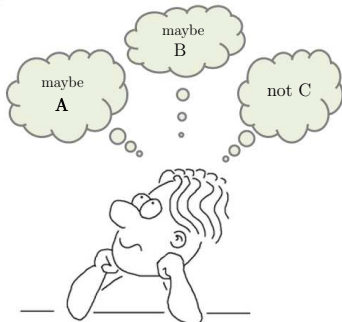
# Reliable incorporation of auxiliary information



# Example of data under ontic imprecision

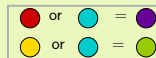
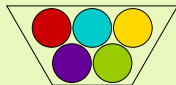
Which party are you considering to elect?

A    B    C    Don't know



## Ontic imprecision:

*"Precise observation of something imprecise"*

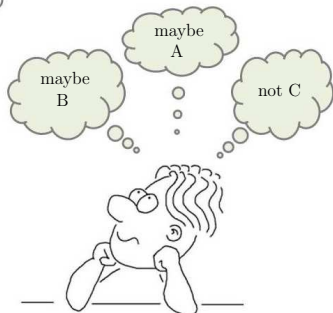


⇒ Truth is represented by coarse observation

# Example of data under ontic imprecision

Which party are you considering to elect?

A    B    C    Don't know

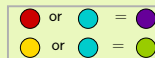
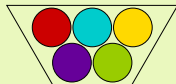


↓  
**A or B**

Allow for multiple answers!

## Ontic imprecision:

*"Precise observation of something imprecise"*



⇒ Truth is represented by coarse observation

## General analysis:

- Interpretation of coarse answers as ontic sets (random sets)  
(Couso, Dubois, Sánchez, 2014)
- Regard coarse answers like “A or B” as own categories
- Extension of state space  $S$  of  $Y$  to  $S^* = \mathcal{P}(S) \setminus \emptyset$  of  $Y^*$

## Example: Multinomial logistic regression

For each category  $s \in S^* = \{1, \dots, m-1\}$ ,  $m = |S^*|$ , probabilities of response  $Y^*$  given covariates  $\mathbf{x}_i$  are modelled by

$$P^*(Y_i^* = s | \mathbf{x}_i) = \frac{\exp(\tilde{\mathbf{x}}_i^T \boldsymbol{\beta}_s^*)}{1 + \sum_{r=1}^{m-1} \exp(\tilde{\mathbf{x}}_i^T \boldsymbol{\beta}_r^*)}$$

with  $\tilde{\mathbf{x}}_i^T = (1, \mathbf{x}_i^T)$  and for reference category  $m$  by

$$P^*(Y_i^* = m | \mathbf{x}_i) = \left(1 + \sum_{r=1}^{m-1} \exp(\tilde{\mathbf{x}}_i^T \boldsymbol{\beta}_r^*)\right)^{-1}.$$

# Illustration by the GLES'13 data

- $Y$ : first vote (reference category S)
- $X$ : religious denomination, most important information source

Coefficient	ontic		classical
	CD	G:S	CD
intercept	0.33	-1.41 **	-0.12
rel.christ	0.37 **	-0.25	0.52 ***
info.tv	-0.02	-0.32	0.25
info.np	-0.12	-1.69 **	0.13

⇒ Own categories for coarse categories

⇒ remarkable differences partly associated with a change in sign



## Illustration by the GLES'13 data

- Y: first vote (reference category S)
- X: religious denomination, most important information source

Coefficient	ontic		classical
	CD	G:S	CD
intercept	0.33	-1.41 **	-0.12
rel.christ	0.37 **	-0.25	0.52 ***
info.tv	-0.02	-0.32	0.25
info.np	-0.12	-1.69 **	0.13

⇒ Own categories for coarse categories

⇒ remarkable differences partly associated with a change in sign

## Illustration by the GLES'13 data

- $Y$ : first vote (reference category S)
- $X$ : religious denomination, most important information source

Coefficient	ontic		classical
	CD	G:S	CD
intercept	0.33	-1.41 **	-0.12
rel.christ	0.37 **	-0.25	0.52 ***
info.tv	-0.02	-0.32	0.25
info.np	-0.12	-1.69 **	0.13

⇒ Own categories for coarse categories

⇒ remarkable differences partly associated with a sign change

## EPISTEMIC







- Obtain MLE referring to the latent variable via the observation model  $\mathcal{Q}$
- Inclusion of auxiliary information via further restrictions on  $\mathcal{Q}$

## Next steps:

- Bayesian approach
- Likelihood-based hypothesis tests, uncertainty regions
- Other “deficiency” processes

## ONTIC

- Coarse categories as own categories  
⇒ Change in state space
- Statistical methods do not change, only interpretation
  
- Ontic imprecision in covariates
- Adaptation to ordinal scale

-  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.
-  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.
-  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.