

Categorical regression analysis for coarse data

Julia Plass*, Marco Cattaneo**, Paul Fink*,
Georg Schollmeyer*, Thomas Augustin*

*Department of Statistics, Ludwig-Maximilians University and

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



02nd of December 2016

Coarse data:

Data are not observed in the resolution originally intended in the subject matter context

Categorical regression analysis:

- Modelling the (not necessarily causal) relation between some covariates X (input variables) and a dependent categorical variable Y (output variable)
- Here considering ...
 - ... Y is partly only observed in a coarse(ned) way (\mathcal{Y})
 - ... precisely observed covariates

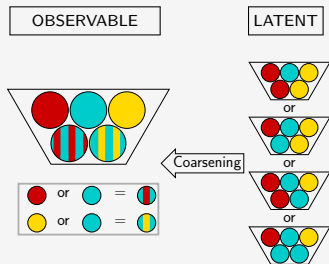
- 1.) Stressing the distinction between **ontic** and **epistemic** data imprecision
- 2.) “Disambiguation” strategy
- 3.) Incorporation of **coarsening assumptions**
 - error freeness
 - superset assumption
 - coarsening at random
 - subgroup independence

COMPARISON 1:

Distinction between epistemic
and ontic data imprecision

Epistemic imprecision:

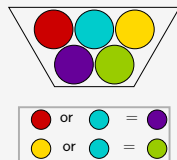
"Imprecise observation of something precise"



⇒ Truth is hidden due to the underlying coarsening mechanism

Ontic imprecision:

"Precise observation of something imprecise"

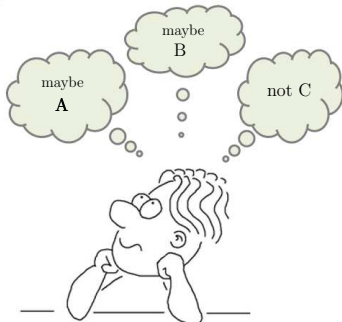


⇒ Truth is represented by coarse observation

Example of data under ontic imprecision

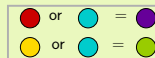
Which party are you considering to elect?

A B C Don't know



Ontic imprecision:

"Precise observation of something imprecise"

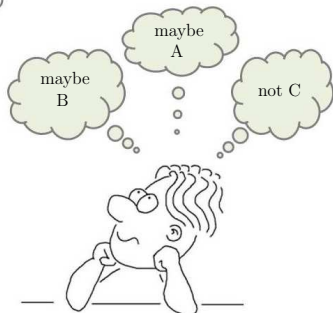


⇒ Truth is represented by coarse observation

Example of data under ontic imprecision

Which party are you considering to elect?

A B C Don't know



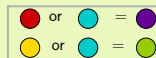
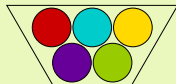
↓

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^*
- Multi-label classification

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

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

- 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

- 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

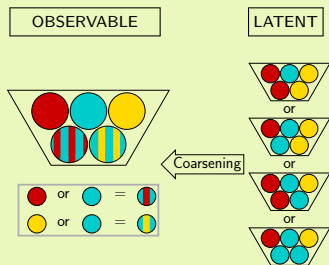
Now: **Epistemic** data imprecision

COMPARISON 2: Disambiguation strategy

A first comparison of the disambiguation strategy

Epistemic imprecision:

"Imprecise observation of something precise"



⇒ Truth is hidden due to the underlying coarsening mechanism

You

- Machine Learning
- Simultaneous model identification and data disambiguation
- Generalized loss function

We

- Survey statistics
- First: information, then: inference
- Likelihood approach

Extension principle:

- Consider all models that are compatible with the observations
- All models are assessed as equally plausible

Basic idea:

- Accounting for model assumptions
- “Model induction and data disambiguation go hand in hand”
- Instead of “ambiguation” of the learning algorithm (extension principle), “ambiguation” of the loss functions

⇒ Disambiguation strategy: The most plausible precise value is the one that minimizes the generalized loss function

LATENT

$$\pi_{xy} := \\ P(Y_i = y | X_i = x)$$

(error-freeness)

Observation model

$$q_{\mathcal{Y}|xy} := \\ P(\mathcal{Y}_i = y | X_i = x, Y_i = y)$$

OBSERVABLE

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

LATENT

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

OBSERVABLE

$$\mathbf{p}_{xy} := P(\mathcal{Y}_i = \mathbf{y} | X_i = x)$$

LATENT

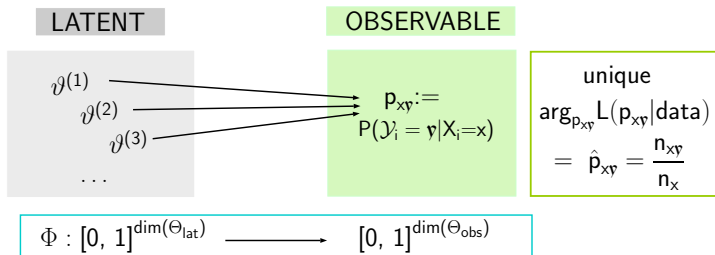
$$\vartheta = (\pi_{xy}^T, \mathbf{q}_{y|xy}^T)^T$$

OBSERVABLE

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

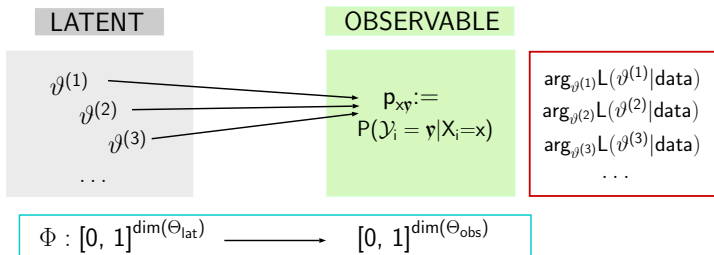
$$\begin{aligned} & \text{unique} \\ & \arg_{p_{xy}} L(p_{xy} | \text{data}) \\ & = \hat{p}_{xy} = \frac{n_{xy}}{n_x} \end{aligned}$$

1.) Determine MLE of observed variable distribution



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

$$\mathbf{p}_{\mathbf{x}\mathbf{y}} = \sum_{\mathbf{y} \in \mathcal{Y}} (\pi_{\mathbf{xy}} \cdot \mathbf{q}_{\mathbf{y}|\mathbf{xy}}).$$

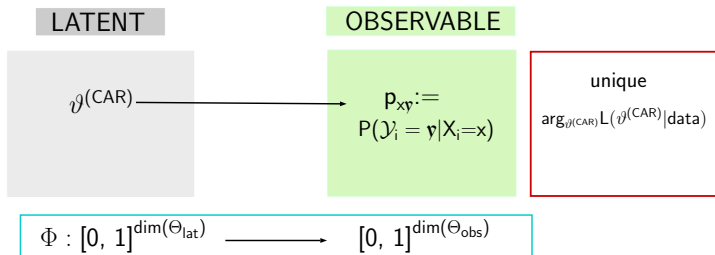


- 1.) Determine MLE of observed variable distribution
- 2.) Use connection between both worlds

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

- 3.) Use invariance of the likelihood

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



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

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

- 3.) Use invariance of the likelihood

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

Estimation of regression coefficients

You

- Model assumptions via the specification of the loss function
- Learning the data and the model simultaneously

We

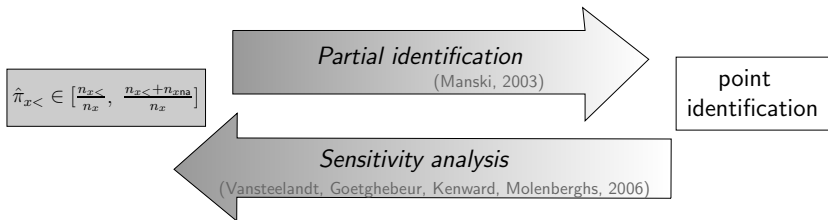
- Including model assumptions via the response function
- No learning of the disambiguation process
- Only external assumptions about the coarsening behaviour

COMPARISON 3:

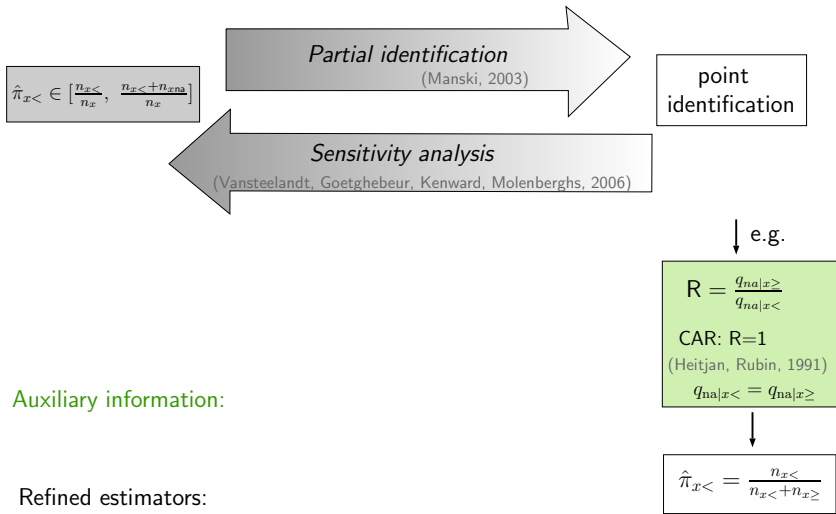
Incorporation of coarsening assumptions

Error freeness, superset assumption, CAR, SI

Reliable incorporation of auxiliary information



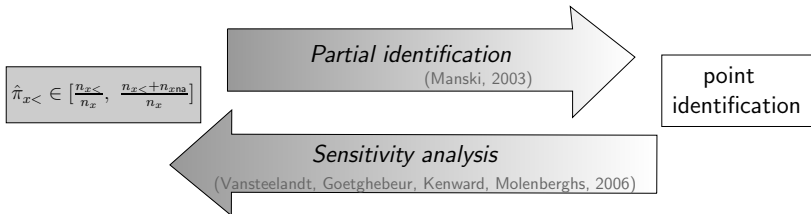
Reliable incorporation of auxiliary information



Auxiliary information:

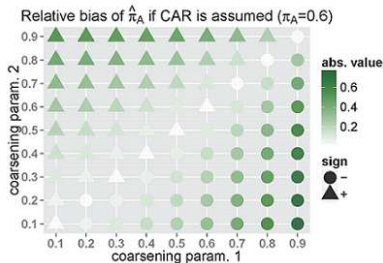
Refined estimators:

Reliable incorporation of auxiliary information

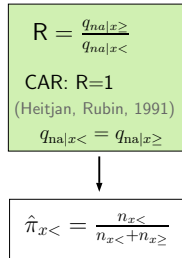


Auxiliary information:

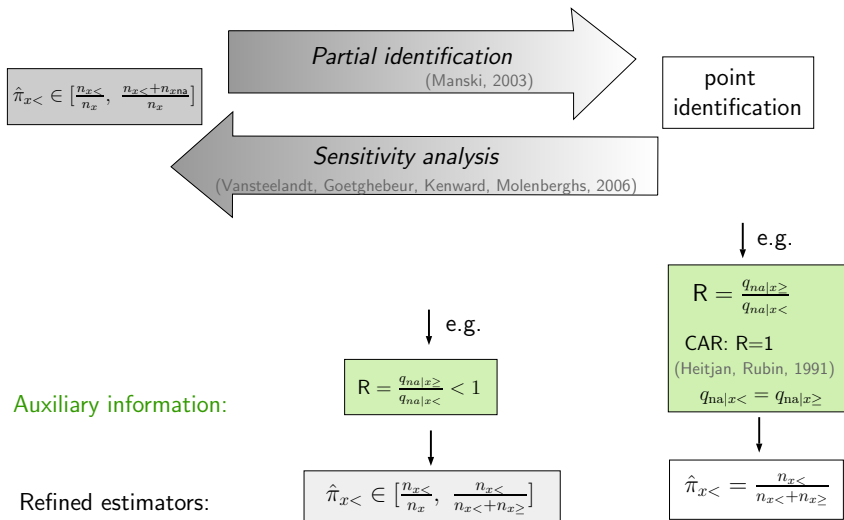
Refined estimators:



e.g.



Reliable incorporation of auxiliary information



Comparison: Incorporation of assumptions

You

- Superset assumption
- Model assumptions

We