

Reliable incorporation of auxiliary information in regression analysis with a coarse categorical response variable

Sommerklausur 2016, Holzhausen

Julia Plass, Thomas Augustin,
Marco Cattaneo, Georg Schollmeyer

10th of June 2016

Motivation of the problem

- Different kinds of uncertainty
 - ⇒ due to sampling (confidence intervals, . . .)
 - ⇒ due to incompleteness (is often forgotten!)

- Different kinds of uncertainty
 - ⇒ due to sampling (confidence intervals, . . .)
 - ⇒ due to incompleteness (is often forgotten!)
- **Missingness mechanism can not be tested**
 - **But:** common procedures (e.g. imputation) require strong (point-identifying) assumptions on the missingness ⇒ bias
 - **Why not?:** Proper reflection of the available information on the missingness
(e.g. Manki (2016), Kenward, Goetghebeur, Molenberghs (2001), Nguyen (2006))

- Different kinds of uncertainty
 - ⇒ due to sampling (confidence intervals, . . .)
 - ⇒ due to incompleteness (is often forgotten!)
- **Missingness mechanism can not be tested**
 - **But:** common procedures (e.g. imputation) require strong (point-identifying) assumptions on the missingness ⇒ bias
 - **Why not?:** Proper reflection of the available information on the missingness
(e.g. Manki (2016), Kenward, Goetghebeur, Molenberghs (2001), Nguyen (2006))

⇒ Here: Focus on the tradeoff between information and reliability in the context of **coarse categorical data**

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

Coarse data:

- Data are not observed in the resolution originally intended
- Here: imprecise observation of something precise
(epistemic interpretation, cf. Couso, Dubois, Sánchez, 2014)

Examples:

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

Here:

- response variable is coarse only
- categorical response variable and categorical covariate(s)

The PASS data example

- German Panel Study “Labour market and social security” (IAB)
- X : self-employed, $\Omega_X = \{0 \text{ (no)}, 1 \text{ (yes)}\}$
- \mathcal{Y} : observed income, Y : latent true income, where $\Omega_Y \subseteq \mathcal{P}(\Omega_Y) \setminus \emptyset$

The PASS data example

- German Panel Study “Labour market and social security” (IAB)
- X : self-employed, $\Omega_X = \{0 \text{ (no)}, 1 \text{ (yes)}\}$
- \mathcal{Y} : observed income, Y : latent true income, where $\Omega_Y \subseteq \mathcal{P}(\Omega_Y) \setminus \emptyset$

Example 1: ordinal scaled Y

$$\Omega_Y = \{< 500, \geq 500, < 2000, \geq 2000\}$$

$$\Omega_y = \{\{< 500\}, \{\geq 500\}, \{< 2000\}, \{\geq 2000\}, \\ \{< 500, \geq 500\}^{<1000}, \{< 2000, \geq 2000\}^{\geq 1000}, \text{na}\}$$

The PASS data example

- German Panel Study “Labour market and social security” (IAB)
- X : self-employed, $\Omega_X = \{0 \text{ (no)}, 1 \text{ (yes)}\}$
- \mathcal{Y} : observed income, Y : latent true income, where $\Omega_Y \subseteq \mathcal{P}(\Omega_Y) \setminus \emptyset$

Example 1: ordinal scaled Y

$$\Omega_Y = \{< 500, \geq 500, < 2000, \geq 2000\}$$

$$\Omega_{\mathcal{Y}} = \{\{< 500\}, \{\geq 500\}, \{< 2000\}, \{\geq 2000\}, \\ \{< 500, \geq 500\}^{<1000}, \{< 2000, \geq 2000\}^{\geq 1000}, \text{na}\}$$

Example 2: nominal scaled Y (reduces to missing case)

$$\Omega_Y = \{<, \geq\}, \text{ (abbr. for: } < 1000, \geq 1000)$$

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

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}

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

Observation model Q

error-freeness

$$q_{\mathcal{Y}|y} =$$

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

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

$$q_{\mathcal{Y}|y} =$$

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

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

OBSERVABLE

coarse data
 \mathcal{Y}

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

Main goal:

Maximum-Likelihood estimation of

Observation model Q

error-freeness

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

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

LATENT

latent variable

Y (nominal scaled)

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

$$\pi_{iK} = 1 - \pi_{i1} - \dots - \pi_{iK-1}$$

(multinomial logit model)

OBSERVABLE

LATENT

coarse data
 \mathcal{Y}

latent variable

Y (ordinal scaled)

Observation model Q

error-freeness

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

π_{ij} with $P(Y_i \leq j | \mathbf{x}_i)$

$$= \frac{\exp(\theta_j + \mathbf{x}_i^T \boldsymbol{\beta})}{1 + \exp(\theta_j + \mathbf{x}_i^T \boldsymbol{\beta})}$$

for reference category K

$$\pi_{iK} = 1 - \pi_{i1} - \dots - \pi_{iK-1}$$

(cumulative logit model)

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

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

Main goal:

Maximum-Likelihood estimation of

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

OBSERVABLE

LATENT

Determine ML estimator

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

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

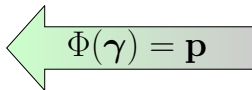
OBSERVABLE

LATENT

Determine ML estimator

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

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


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

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

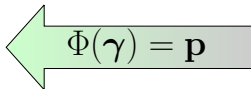
OBSERVABLE

LATENT

Determine ML estimator

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

Use the **connection**
between p and γ


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

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

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

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

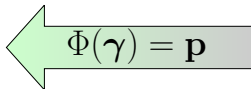
OBSERVABLE

LATENT

Determine ML estimator

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

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


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

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

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

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

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

$$\begin{aligned}\hat{\pi}_{x<} &= \frac{n_{x\{<\}}}{n_x} \\ \hat{q}_{na|x<} &= 0 \\ \hat{q}_{na|x\geq} &= \frac{n_{xna}}{n_{xna} + n_{x\{\geq\}}}\end{aligned}$$

• • •

$$\begin{aligned}\hat{\pi}_{x<} &= \frac{n_{x\{<\}} + n_{xna}}{n_x} \\ \hat{q}_{na|x<} &= \frac{n_{xna}}{n_{xna} + n_{x\{<\}} \\ \hat{q}_{na|x\geq} &= 0\end{aligned}$$

$$\hat{\pi}_{x<} = \frac{n_{x\{<\}}}{n_x}$$

$$\hat{q}_{na|x<} = 0$$

$$\hat{q}_{na|x\geq} = \frac{n_{xna}}{n_{xna} + n_{x\{\geq\}}}$$

• • •

$$\hat{\pi}_{x<} = \frac{n_{x\{<\}} + n_{xna}}{n_x}$$

$$\hat{q}_{na|x<} = \frac{n_{xna}}{n_{xna} + n_{x\{<\}}}$$

$$\hat{q}_{na|x\geq} = 0$$

- Results for subgroup $x = 0$ (not self-employed)

$$\hat{\pi}_{0<} \in [0.16, 0.36], \quad \hat{q}_{na|0\geq} \in [0, 0.23], \quad \hat{q}_{na|0<} \in [0, 0.55]$$

$$\hat{\pi}_{x<} = \frac{n_{x\{<\}}}{n_x}$$

$$\hat{q}_{na|x<} = 0$$

$$\hat{q}_{na|x\geq} = \frac{n_{xna}}{n_{xna} + n_{x\{\geq\}}}$$

• • •

$$\hat{\pi}_{x<} = \frac{n_{x\{<\}} + n_{xna}}{n_x}$$

$$\hat{q}_{na|x<} = \frac{n_{xna}}{n_{xna} + n_{x\{<\}}}$$

$$\hat{q}_{na|x\geq} = 0$$

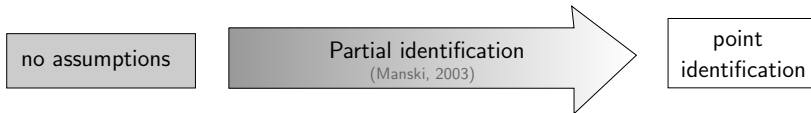
- Results for subgroup $x = 0$ (not self-employed)

$$\hat{\pi}_{0<} \in [0.16, 0.36], \quad \hat{q}_{na|0\geq} \in [0, 0.23], \quad \hat{q}_{na|0<} \in [0, 0.55]$$

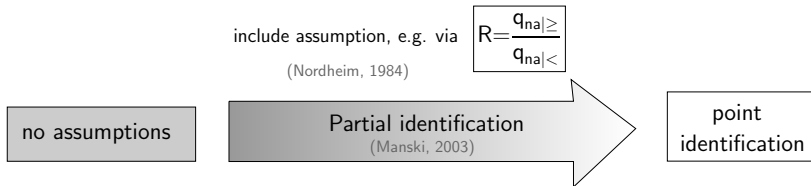
- Results in terms of the regression coefficients:

$$\hat{\beta}_{<0} \in [-1.66, -0.59] \text{ (Intercept)}, \quad \hat{\beta}_{<} \in [0.61, 1.15]$$

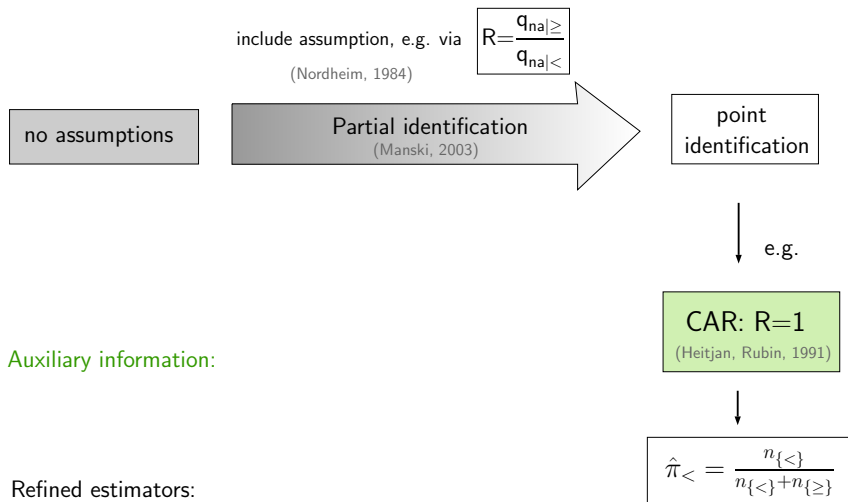
Reliable incorporation of auxiliary information



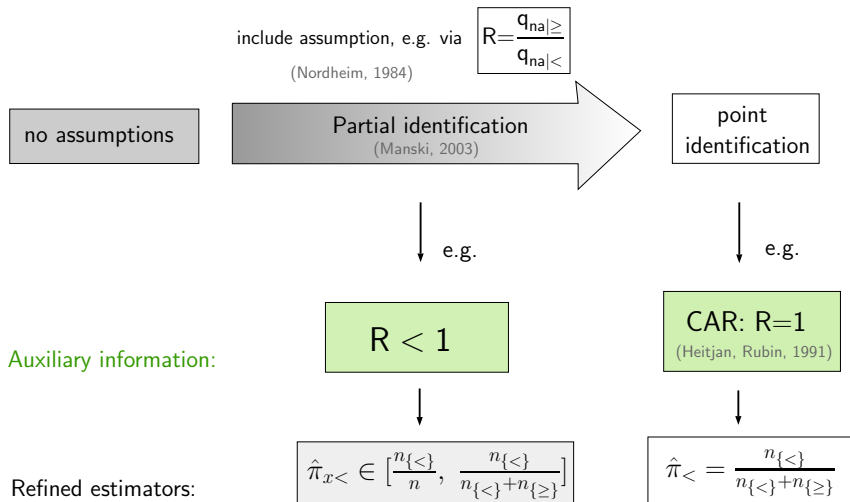
Reliable incorporation of auxiliary information



Reliable incorporation of auxiliary information



Reliable incorporation of auxiliary information



Coarsening at random & subgroup independence

coarsening at random (CAR)

(Heitjan, Rubin, 1991)

Generally

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

Example

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

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

subgroup independence (SI)

Coarsening does not
depend on the value
of the covariate

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

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

- Data situations that might hint to incompatibility with SI
- Hypotheses:

$$H_0 : q_{\mathfrak{y}|xy} = q_{\mathfrak{y}|x'y} \text{ for all } \mathfrak{y} \in \Omega_{\mathcal{Y}}, x, x' \in \Omega_{\mathcal{X}}, y \in \Omega_{\mathcal{Y}},$$

$$H_1 : q_{\mathfrak{y}|xy} \neq q_{\mathfrak{y}|x'y} \text{ for some } \mathfrak{y} \in \Omega_{\mathcal{Y}}, x, x' \in \Omega_{\mathcal{X}}, y \in \Omega_{\mathcal{Y}}.$$

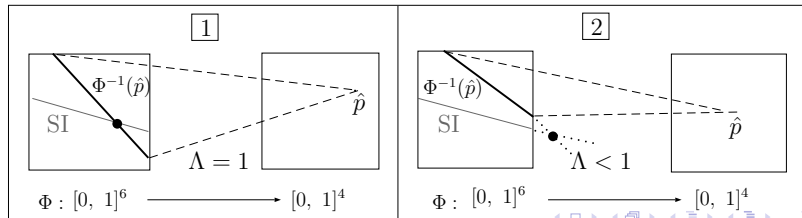
- Data situations that might hint to incompatibility with SI
- Hypotheses:

$$H_0 : q_{\mathbf{y}|xy} = q_{\mathbf{y}|x'y} \text{ for all } \mathbf{y} \in \Omega_{\mathcal{Y}}, x, x' \in \Omega_{\mathcal{X}}, y \in \Omega_{\mathcal{Y}},$$

$$H_1 : q_{\mathbf{y}|xy} \neq q_{\mathbf{y}|x'y} \text{ for some } \mathbf{y} \in \Omega_{\mathcal{Y}}, x, x' \in \Omega_{\mathcal{X}}, y \in \Omega_{\mathcal{Y}}.$$

- Test statistic based on the Likelihood ratio Λ may be constructed:

$$\Lambda(\mathbf{y}_1, \dots, \mathbf{y}_n, x_1, \dots, x_n) = \frac{\sup_{H_0} L(\vartheta || \mathbf{y}_1, \dots, \mathbf{y}_n, x_1, \dots, x_n)}{\sup_{H_0 \cup H_1} L(\vartheta || \mathbf{y}_1, \dots, \mathbf{y}_n, x_1, \dots, x_n)},$$










- We studied a cautious regression approach with categorical precisely observed covariates in case of a
 - ... coarse response of nominal scale (multinomial logit model)
 - ... coarse response of ordinal scale (cummulative model)
- Include auxiliary information via restrictions on the coarsening
- Support information via hypothesis tests on the coarsening

Next steps:

- Decision rule for test on SI
- Comparison to traditional approaches
- Application of the approach to other problems (e.g. propensity score matching)

References

-  Couso, Dubois, Sánchez.
Random Sets and Random Fuzzy Sets as Ill-Perceived Random Variables, Springer, 2014.
-  Heitjan, Rubin.
Ignorability and Coarse Data, *Ann. Stat.*, 1991.
-  Kenward, Goetghebeur, Molenberghs.
Sensitivity analysis for incomplete categorical data, *Stat. Modelling*, 2001
-  Manski.
Credible interval estimates for official statistics with survey nonresponse, *J Econometrics* , 2016.
-  Nordheim.
Inference from nonrandomly missing categorical data: An example from a genetic study on Turner's syndrome, *J. Am. Stat. Assoc.*, 1984.
-  Nguyen.
An introduction to random sets, CRC press, 2006
-  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.
-  Plass, Cattaneo, Schollmeyer, Augustin.
Testing of coarsening mechanisms: Coarsening at random versus subgroup independence, *SMPS*, 2016
-  Vansteelandt, Goetghebeur, Kenward, Molenberghs.
Ignorance and uncertainty regions as inferential tools in a sensitivity analysis, *Stat. Sin.*, 2006.