# Coarse categorical data
# under ontologic and epistemic uncertainty

Julia Plaß

Ludwigs-Maximilians-University (LMU)

08th of September 2014

## Outline

1. Introduction to the problem

2. Data under ontologic uncertainty
   - Motivation
   - Basic idea of analysis

3. Data under epistemic uncertainty
   - Motivation
   - Analysis of two different situations

4. Summary

# Epistemic vs. ontic/ontologic uncertainty (I. Couso, D. Dubois, 2014)

▶ Ontologic uncertainty

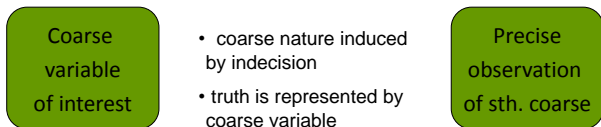<table>
<tr><td>Coarse<br>variable<br>of interest</td><td>• coarse nature induced<br>by indecision<br><br>• truth is represented by<br>coarse variable</td><td>Precise<br>observation<br>of sth. coarse</td></tr>
</table>

# Epistemic vs. ontic/ontologic uncertainty (I. Couso, D. Dubois, 2014)

▶ Ontologic uncertainty

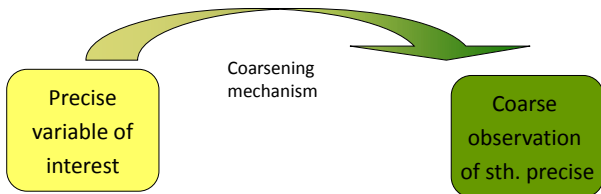| Coarse variable of interest | • coarse nature induced by indecision <br> • truth is represented by coarse variable | Precise observation of sth. coarse |
|---|---|---|

▶ Epistemic uncertainty

| Precise variable of interest | Coarsening mechanism | Coarse observation of sth. precise |
|---|---|---|

# Why should data under ontologic uncertainty be collected?

**Example:**

***Which party will you give your vote?***

☐ *A*  ☐ *B*  ☐ *C*  ☐ *Don't know*

# Why should data under ontologic uncertainty be collected?

**Example:**

***Which party will you give your vote?***

☐ *A* ☐ *B* ☐ *C* ☐ *Don't know*

# Why should data under ontologic uncertainty be collected?

**Example:**

*Which party will you give your vote?*

☐ *A*  ☐ *B*  ☐ *C*  ☐ *Don't know*

# Why should data under ontologic uncertainty be collected?

# Why should data under ontologic uncertainty be collected?



**Example:**

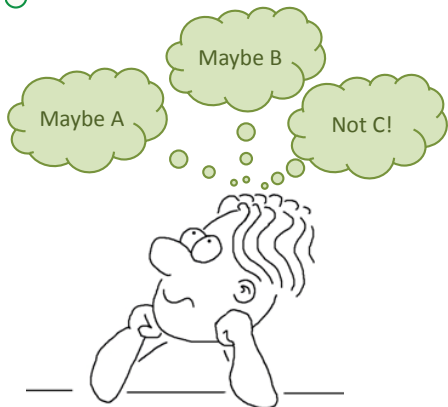*Which party will you give your vote?*

☐ A  ☐ B  ☐ C  ☒ *Don't know*
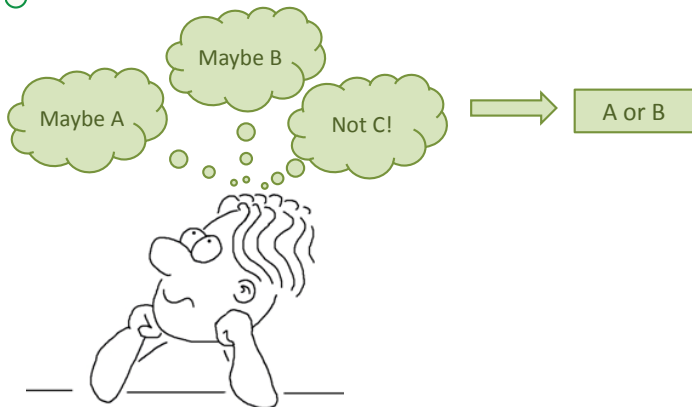
Maybe B

Maybe A

Not C!

⟹ A or B

# Why should data under ontologic uncertainty be collected?

**Example:**

### Which party will you give your vote?

☒ *A*        ☒ *B*        ☐ *C*        ☐ *Don't know*



Dealing with ontologic uncertainty:

→ Allow multiple anwers

## Basic idea (illustrated by GLES 2013)

| certainty | vote | assesCD | assesSPD | assesGREEN | assesLEFT | ... |
|---|---|---|---|---|---|---|
| very certain | SPD | -1 | 5 | 2 | 1 | ... |
| certain | CD | 4 | 3 | 3 | 1 | ... |
| not that certain | GREEN | 3 | 4 | 4 | -1 | ... |
| not certain at all | CD | -3 | 2 | 2 | 2 | ... |

Exemplary
extraction of
the dataset

# Basic idea (illustrated by GLES 2013)

| certainty | vote | assesCD | assesSPD | assesGREEN | assesLEFT | ... |
|-----------|------|---------|----------|------------|-----------|-----|
| very certain | SPD | -1 | 5 | 2 | 1 | ... |
| certain | CD | 4 | 3 | 3 | 1 | ... |
| not that certain | GREEN | 3 | 4 | 4 | -1 | ... |
| not certain at all | CD | -3 | 2 | 2 | 2 | ... |

Exemplary extraction of the dataset

**Analysis:**     *Prediction*

- Classical analysis - Vote of respondents who are certain:

  Prediction for "CD": $\frac{\#\text{respondents who will vote for "CD" with certainty}}{\#\text{respondents who are certain}} = 0.46$

# Basic idea (illustrated by GLES 2013)

| certainty | vote | assesCD | assesSPD | assesGREEN | assesLEFT | ... |
|-----------|------|---------|----------|------------|-----------|-----|
| very certain | SPD | -1 | 5 | 2 | 1 | ... |
| certain | CD | 4 | 3 | 3 | 1 | ... |
| not that certain | GREEN | 3 | 4 | 4 | -1 | ... |
| not certain at all | CD | -3 | 2 | 2 | 2 | ... |

Exemplary extraction of the dataset

**Analysis:**  *Prediction*

- Classical analysis - Vote of respondents who are certain:
  Prediction for "CD": $\frac{\#\text{respondents who will vote for "CD" with certainty}}{\#\text{respondents who are certain}} = 0.46$

- Dealing with ontologic uncertainty:

| CD | SPD | GREEN | LEFT | OTHER |
|----|-----|-------|------|-------|
| 519 | 287 | 105 | 101 | 62 |
| **SPD-CD** | **GREEN-SPD** | **CD-OTHER** | **LEFT-SPD** | **GREEN-LEFT** |
| 34 | 34 | 24 | 14 | 13 |
| **GREEN-SPD-CD** | **SPD-CD-OTHER** | **LEFT-GREEN-SPD** | **SPD-OTHER** | rare comb. |
| 13 | 13 | 13 | 12 | 77 |

$$\hat{Bel}(CD) = \frac{519}{1321} = 0.39$$

$$\hat{Pl}(CD) = \frac{519 + 34 + 24 + 13 + 13}{1321} = 0.45$$

$\Rightarrow$  Prediction for "CD": [0.39, 0.45]

# Basic idea (illustrated by GLES 2013)

| certainty | vote | assesCD | assesSPD | assesGREEN | assesLEFT | ... |
|-----------|------|---------|----------|------------|-----------|-----|
| very certain | SPD | -1 | 5 | 2 | 1 | ... |
| certain | CD | 4 | 3 | 3 | 1 | ... |
| not that certain | GREEN | 3 | 4 | 4 | -1 | ... |
| not certain at all | CD | -3 | 2 | 2 | 2 | ... |

Exemplary extraction of the dataset

**Analysis:**         *Regression*

- For reasons of explanation interpretation of selected $\beta$ estimators only

- Reference category: "SPD"

- Classical analysis - Vote of respondents who are certain:

|  | Intercept | sexFEM | InfoNEWSPAPER | InfoRADIO | InfoWEB |
|----|-----------|--------|---------------|-----------|---------|
| **CD** | 0.2914 | 0.2456 | -0.0037 | -0.1487 | -0.6774 |

# Model under ontologic uncertainty

**Data under ontologic uncertainty:**

- $Y_i$: categorical random variable of nominal scale of measurement with $Y_i \subseteq \underbrace{\{a, b, ...\}}_{\text{precise categories}}$

- $m = |\mathcal{P}(\Omega) \setminus \emptyset|$: number of categories of $Y_i$

**Model under ontologic uncertainty:**

$\Rightarrow$ classical multinomial logit model with different number of categories:

The probability of occurence for category $r = 1, 2, 3, ..., m-1$ can be calculated by

$$P(Y_i = r | \boldsymbol{x}_i) = \frac{\exp(\boldsymbol{x}_i^{\boldsymbol{T}} \beta_r)}{1 + \sum_{s=1}^{m-1} \exp(\boldsymbol{x}_i^{\boldsymbol{T}} \beta_s)}$$

and for category $m$ by

$$P(Y_i = m | \boldsymbol{x}_i) = \frac{1}{1 + \sum_{s=1}^{m-1} \exp(\boldsymbol{x}_i^{\boldsymbol{T}} \beta_s)}$$

# Basic idea (illustrated by GLES 2013)

| certainty | vote | assesCD | assesSPD | assesGREEN | assesLEFT | ... |
|-----------|------|---------|----------|------------|-----------|-----|
| very certain | SPD | -1 | 5 | 2 | 1 | ... |
| certain | CD | 4 | 3 | 3 | 1 | ... |
| not that certain | GREEN | 3 | 4 | 4 | -1 | ... |
| not certain at all | CD | -3 | 2 | 2 | 2 | ... |

Exemplary extraction of the dataset

**Analysis:** *Regression*

- For reasons of explanation interpretation of selected $\beta$ estimators only
- Reference category: "SPD"

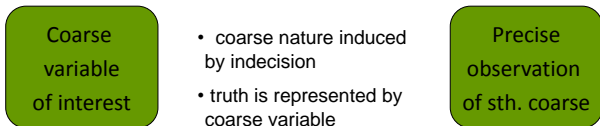- Classical analysis - Vote of respondents who are certain:

| | Intercept | sexFEM | InfoNEWSPAPER | InfoRADIO | InfoWEB |
|---|-----------|--------|---------------|-----------|---------|
| **CD** | 0.2914 | 0.2456 | -0.0037 | -0.1487 | -0.6774 |

- Dealing with ontologic uncertainty:

| | Intercept | sexFEM | InfoNEWSPAPER | InfoRADIO | InfoWEB |
|---|-----------|--------|---------------|-----------|---------|
| **CD** | 0.4706 | 0.3135 | -0.0784 | 0.1856 | -0.5025 |
| **SPD-CD** | -2.2007 | 0.4034 | -1.2556 | 0.9476 | 0.1886 |
| **GREEN-SPD** | -2.4432 | 0.9393 | -1.3053 | 0.2300 | 0.1844 |

# Basic idea (illustrated by GLES 2013)

| certainty | vote | assesCD | assesSPD | assesGREEN | assesLEFT | ... |
|-----------|------|---------|----------|------------|-----------|-----|
| very certain | SPD | -1 | 5 | 2 | 1 | ... |
| certain | CD | 4 | 3 | 3 | 1 | ... |
| not that certain | GREEN | 3 | 4 | 4 | -1 | ... |
| not certain at all | CD | -3 | 2 | 2 | 2 | ... |

Exemplary extraction of the dataset

**Analysis:** *Regression*

- For reasons of explanation interpretation of selected $\beta$ estimators only
- Reference category: "SPD"

- Classical analysis - Vote of respondents who are certain:

  |    | Intercept | sexFEM | InfoNEWSPAPER | InfoRADIO | InfoWEB |
  |----|-----------|--------|---------------|-----------|---------|
  | CD | 0.2914 | 0.2456 | -0.0037 | -0.1487 | -0.6774 |

- Dealing with ontologic uncertainty:

  |    | Intercept | sexFEM | InfoNEWSPAPER | InfoRADIO | InfoWEB |
  |----|-----------|--------|---------------|-----------|---------|
  | CD | 0.4706 | 0.3135 | -0.0784 | 0.1856 | -0.5025 |
  | SPD-CD | -2.2007 | 0.4034 | -1.2556 | 0.9476 | 0.1886 |
  | GREEN-SPD | -2.4432 | 0.9393 | -1.3053 | 0.2300 | 0.1844 |

# Epistemic vs. ontologic uncertainty

▶ Ontologic uncertainty

Coarse variable of interest

- coarse nature induced by indecision
- truth is represented by coarse variable

Precise observation of sth. coarse

▶ Epistemic uncertainty

Precise variable of interest

Coarsening mechanism

Coarse observation of sth. precise

# When do data under epistemic uncertainty occur?

**Reasons for coarse categorical data:**

- Guarantee of anonymization, prevention of refusals

  **Example:**
  "Which kind of party did you elect?"
  ☐ rather left    ☐ center    ☐ rather right

# When do data under epistemic uncertainty occur?

**Reasons for coarse categorical data:**

- Guarantee of anonymization, prevention of refusals

  **Example:**
  "Which kind of party did you elect?"
  □ rather left    □ center    □ rather right

- Different levels of reporting accuracy
  (lack of knowledge, vague question formulation)

  **Examples:**
  "Which car do you drive?"

# Addressed data situations



$$q_1 = P(\mathcal{Y} = AB | Y = A)$$
$$q_2 = P(\mathcal{Y} = AB | Y = B)$$

## Addressed data situations



Two different situations will be regarded:

### Case 1 - No covariates:

- IID-assumption

$$\pi_{iA} = \pi_A$$

- Constant coarsening mechanisms $q_1$ and $q_2$

# Addressed data situations



Two different situations will be regarded:

### Case 1 - No covariates:

- IID-assumption

$$\pi_{iA} = \pi_A$$

- Constant coarsening mechanisms $q_1$ and $q_2$

### Case 2 - Binary covariates, no intercept

- $$\pi_{iA|X_i=1} = \frac{\exp(\beta_A)}{1 + \exp(\beta_A)}$$

  $$\pi_{iA|X_i=0} = \frac{1}{2}$$

- Constant coarsening mechanisms $q_1$ and $q_2$

# Case 1 - No covariates + IID-assumption

**log-Likelihood under the iid assumption :**

$$
\begin{aligned}
l(\pi_A, q_1, q_2) &= \ln\Big( \prod_{i:\mathcal{Y}_i=A} \underbrace{P(\mathcal{Y}=A|Y=A)}_{(1-q_1)} \pi_{iA} \prod_{i:\mathcal{Y}_i=B} \underbrace{P(\mathcal{Y}=B|Y=B)}_{(1-q_2)}(1-\pi_{iA}) \\
&\qquad \prod_{i:\mathcal{Y}_i=AB} \underbrace{P(\mathcal{Y}=AB|Y=A)}_{q_1} \pi_{iA} + \underbrace{P(\mathcal{Y}=AB|Y=B)}_{q_2}(1-\pi_{iA}) \Big) \\
&\overset{iid}{=} n_A \cdot [\ln(1-q_1) + \ln(\pi_A)] + n_B \cdot [ln(1-q_2) + ln(1-\pi_A)] \\
&\qquad n_{AB} \cdot [q_1\pi_A + q_2(1-\pi_A))]
\end{aligned}
$$

**FOC:**

I.) $\dfrac{\partial}{\partial \pi_A} = \dfrac{n_{AB}}{q_1\pi_A + q_2(1-\pi_A)}(q_1 - q_2) + \dfrac{n_A}{\pi_A} - \dfrac{n_B}{1-\pi_A} \overset{!}{=} 0$

II.) $\dfrac{\partial}{\partial q_1} = \dfrac{n_{AB}}{q_1\pi_A + q_2(1-\pi_A)}\pi_A - \dfrac{n_A}{1-q_1} \overset{!}{=} 0$

III.) $\dfrac{\partial}{\partial q_2} = \dfrac{n_{AB}}{q_1\pi_A + q_2(1-\pi_A)}(1-\pi_A) - \dfrac{n_B}{1-q_2} \overset{!}{=} 0$

# Case 1 - No covariates + IID-assumption

**log-Likelihood under the iid assumption :**

$$l(\pi_A, q_1, q_2) = \ln\Big( \prod_{i:\mathcal{Y}_i=A} \underbrace{P(\mathcal{Y}=A|Y=A)}_{(1-q_1)}\pi_{iA} \prod_{i:\mathcal{Y}_i=B} \underbrace{P(\mathcal{Y}=B|Y=B)}_{(1-q_2)}(1-\pi_{iA})$$

$$\prod_{i:\mathcal{Y}_i=AB} \underbrace{P(\mathcal{Y}=AB|Y=A)}_{q_1}\pi_{iA} + \underbrace{P(\mathcal{Y}=AB|Y=B)}_{q_2}(1-\pi_{iA}) \Big)$$

$$\overset{iid}{=} n_A \cdot [\ln(1-q_1) + \ln(\pi_A)] + n_B \cdot [ln(1-q_2) + ln(1-\pi_A)]$$

$$n_{AB} \cdot [q_1\pi_A + q_2(1-\pi_A))]$$

**FOC:**

$$\text{I.) } \frac{\partial}{\partial \pi_A} = \frac{n_{AB}}{q_1\pi_A + q_2(1-\pi_A)}(q_1-q_2) + \frac{n_A}{\pi_A} - \frac{n_B}{1-\pi_A} \overset{!}{=} 0$$

$$\text{II.) } \frac{\partial}{\partial q_1} = \frac{n_{AB}}{q_1\pi_A + q_2(1-\pi_A)}\pi_A - \frac{n_A}{1-q_1} \overset{!}{=} 0$$

$$\text{III.) } \frac{\partial}{\partial q_2} = \frac{n_{AB}}{q_1\pi_A + q_2(1-\pi_A)}(1-\pi_A) - \frac{n_B}{1-q_2} \overset{!}{=} 0$$

$\Longrightarrow$

Neccessary and sufficient condition for estimators $(\hat\pi_A, \hat q_1, \hat q_2)$

$$\frac{n_{AB}}{n} = \hat q_1 \cdot \hat\pi_A + \hat q_2 \cdot (1-\hat\pi_A)$$

# Case 2 - Binary covariates, no intercept

**Involved assumptions:**

- *No intercept $\beta_0$*

  $\Rightarrow \pi_{iA|X_i=1} = \frac{\exp(\beta_A)}{1+\exp(\beta_A)}$ and $\pi_{iA|X_i=0} = \frac{1}{2}$

# Case 2 - Binary covariates, no intercept

**Involved assumptions:**

- *No intercept $\beta_0$*

  $\Rightarrow \pi_{iA|X_i=1} = \frac{\exp(\beta_A)}{1+\exp(\beta_A)}$ and $\pi_{iA|X_i=0} = \frac{1}{2}$

- *Constant coarsening mechanisms $q_1$ and $q_2$:*

  **Example:** "Do you regularly steal candy ($\bowtie$) out of your mother's candy box?"
  Asked: girls ($g$) and boys ($b$)

| X | Y | $\mathcal{Y}$ |
|---|---|---|
| $g$ | $\bowtie$ | $\bowtie$ or $\boxtimes$ |
| $g$ | $\boxtimes$ | $\boxtimes$ |
| $g$ | $\bowtie$ | $\bowtie$ |
| $b$ | $\bowtie$ | $\bowtie$ or $\boxtimes$ |
| $b$ | $\bowtie$ | $\bowtie$ |
| $b$ | $\boxtimes$ | $\boxtimes$ |
| $b$ | $\boxtimes$ | $\boxtimes$ |

$P(\mathcal{Y} = \boxtimes \text{ or } \boxtimes | Y = \bowtie, X = g) = P(\mathcal{Y} = \boxtimes \text{ or } \boxtimes | Y = \bowtie, X = b)$

$P(\mathcal{Y} = \boxtimes \text{ or } \boxtimes | Y = \boxtimes, X = g) = P(\mathcal{Y} = \boxtimes \text{ or } \boxtimes | Y = \boxtimes, X = b)$
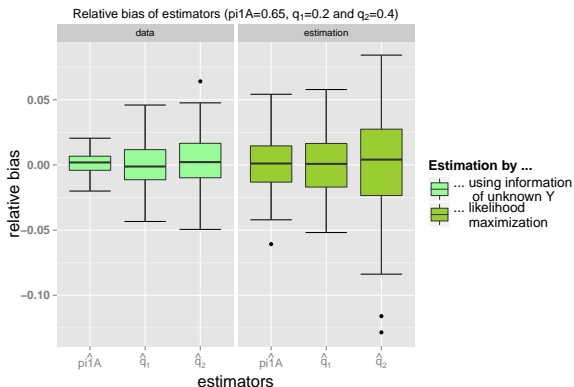
$$\Downarrow$$

the coarsening mechanisms do not
depend on subgroup $x$

# Case 2 - Binary covariates, no intercept

### log-Likelihood involving binary covariate $X$:

$$
\begin{aligned}
l(\pi_{0A}, \pi_{1A}, q_1, q_2) &= n_{1A} \cdot [\ln(1 - q_1) + \ln(\pi_{1A})] + n_{0A} \cdot [\ln(1 - q_1) + \ln(\pi_{0A})] \\
&+ n_{1B} \cdot [\ln(1 - q_2) + \ln(1 - \pi_{1A})] + n_{0B} \cdot [\ln(1 - q_2) + \ln(1 - \pi_{0A})] \\
&+ n_{1AB} \cdot [q_1 \pi_{1A} + q_2(1 - \pi_{1A})] + n_{0AB} \cdot [q_1 \pi_{0A} + q_2(1 - \pi_{0A})]
\end{aligned}
$$

with $\pi_{01} = \frac{1}{2}$



Relative bias of estimators (pi1A=0.65, $q_1$=0.2 and $q_2$=0.4)

Estimation by ...
- ... using information of unknown Y
- ... likelihood maximization

## Case 2 - Binary covariates, no intercept

**Questions / Discussion suggestions:**

- Is this result reasonable? (Remember: The coarsening mechanism does not depend on the values of $X$)

## Case 2 - Binary covariates, no intercept

**Questions / Discussion suggestions:**

- Is this result reasonable? (Remember: The coarsening mechanism does not depend on the values of $X$)

- Is it possible to derive those estimators by means of equations like in the *iid* case, but now for each subgroup of $X$?

$$
\begin{aligned}
\frac{n_{1AB}}{n_1} &= \hat{q}_1 \cdot \hat{\pi}_{1A} + \hat{q}_2 \cdot (1 - \hat{\pi}_{1A}) \\
\frac{n_{0AB}}{n_0} &= \hat{q}_1 \cdot \hat{\pi}_{0A} + \hat{q}_2 \cdot (1 - \hat{\pi}_{0A}) \\
\frac{n_{1A}}{n_1} &= (1 - \hat{q}_1) \cdot \hat{\pi}_{1A} \\
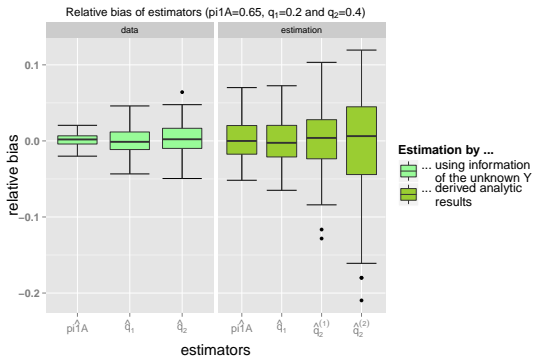\frac{n_{0A}}{n_0} &= (1 - \hat{q}_1) \cdot \hat{\pi}_{0A}
\end{aligned}
$$

# Case 2 - Binary covariates, no intercept

**Questions / Discussion suggestions:**

- Is this result reasonable? (Remember: The coarsening mechanism does not depend on the values of $X$)

- Is it possible to derive those estimators by means of equations like in the *iid* case, but now for each subgroup of $X$?

$$
\begin{aligned}
\frac{n_{1AB}}{n_1} &= \hat{q}_1 \cdot \hat{\pi}_{1A} + \hat{q}_2 \cdot (1 - \hat{\pi}_{1A}) \\
\frac{n_{0AB}}{n_0} &= \hat{q}_1 \cdot \hat{\pi}_{0A} + \hat{q}_2 \cdot (1 - \hat{\pi}_{0A}) \\
\frac{n_{1A}}{n_1} &= (1 - \hat{q}_1) \cdot \hat{\pi}_{1A} \\
\frac{n_{0A}}{n_0} &= (1 - \hat{q}_1) \cdot \hat{\pi}_{0A}
\end{aligned}
\qquad\Longrightarrow
$$

Resulting estimators

$$
\begin{aligned}
\hat{\pi}_{1A} &= \frac{n_{1A} \cdot n_0}{2 \cdot n_1 \cdot n_{0A}} \\
\hat{q}_1 &= 1 - 2 \cdot \frac{n_{0A}}{n_0}
\end{aligned}
$$

and $\hat{q}_2^{(1)}$ and $\hat{q}_2^{(2)}$

# Case 2 - Binary covariates, no intercept



Relative bias of estimators (pi1A=0.65, q₁=0.2 and q₂=0.4)

$$\hat{q}_2^{(1)} = \frac{2 \cdot n_{0AB} - n_0 + 2 \cdot n_{0A}}{n_0}$$

$$\hat{q}_2^{(2)} = \frac{\frac{n_{1AB}}{n_1} - \frac{n_{1A}(n_0 - 2n_{0A})}{2n_1 n_{0A}}}{\frac{2n_1 n_{0A} - n_{1A} n_0}{2n_1 n_{0A}}}$$

## Summary

- Important to distinguish between epistemic and ontologic uncertainty
- One can deal with ontologic uncertainty by redefining the sample space
- In case of ...
    - ... iid variables under epistemic uncertainty, a set of estimators results characterized by a special condition
    - ... being a binary covariate available, precise real valued point estimators seem to result

## References

📎 Couso, I. and Dubois, D.

*Statistical reasoning with set-valued information: Ontic vs. epistemic views, IJAR, 2014.*

📎 Heitjan, D. and Rubin, D.

*Ignorability and Coarse Data, Annals of Statistics, 1991.*

📎 Matheron, G.

*Random Sets and Integral Geometry, Wiley New York, 1975.*

📎 Plaß, J.

*Coarse categorical data under epistemic and ontologic uncertainty: Comparison and extension of some approaches, master's thesis, 2013*