

Distinction between Epistemic and Ontic Interpretation (Couso, Dubois, Sánchez, 2014, Springer)

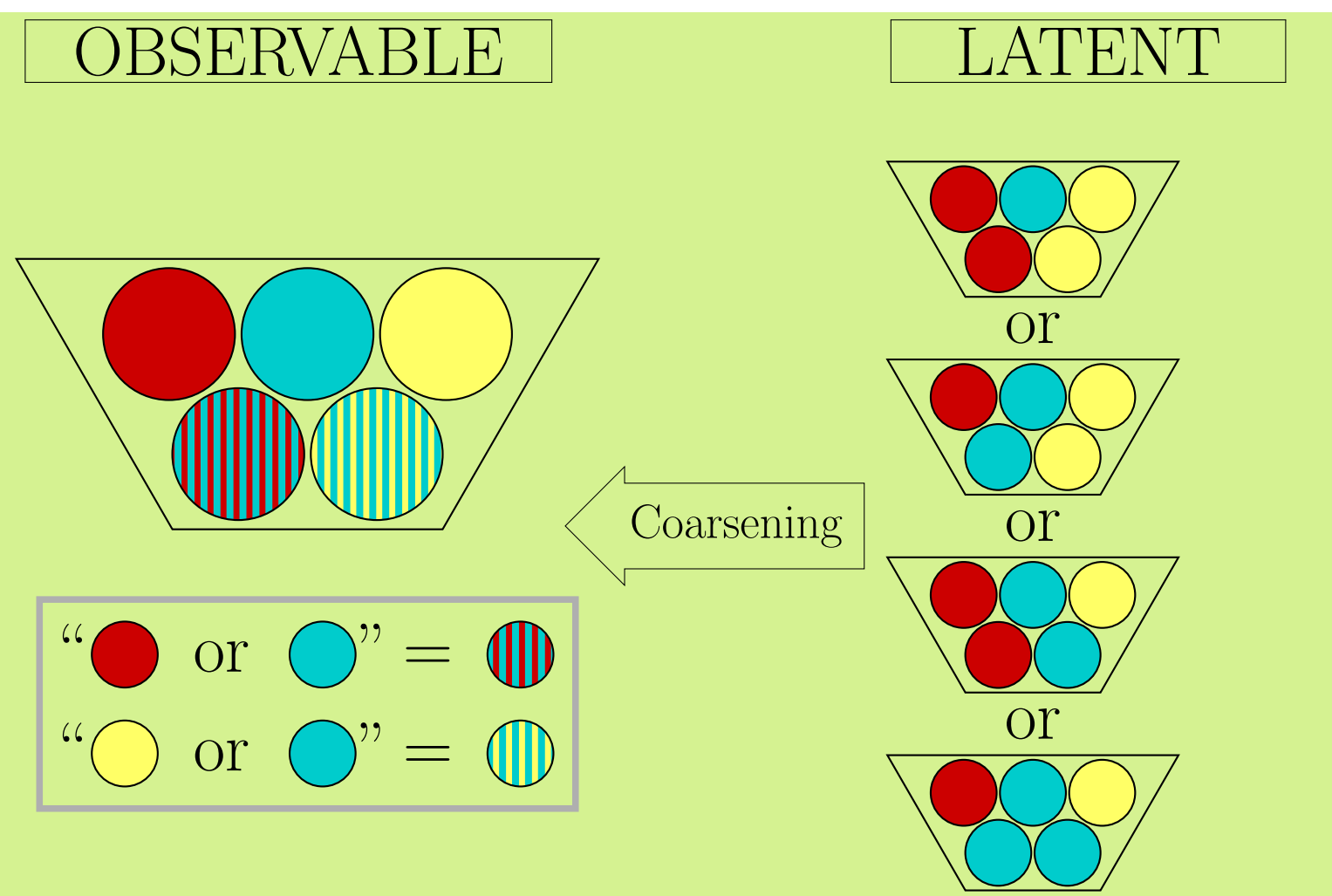
Ontic poster in session on Wednesday

Epistemic data imprecision:

- Imprecise observation of something precise
- Actually precise values may only be observed in a coarse form, due to an underlying coarsening mechanism

Examples:

- Missing data as a special case
- Coarsening deliberately applied as an anonymization technique
- Matched data sets with not completely identical categories

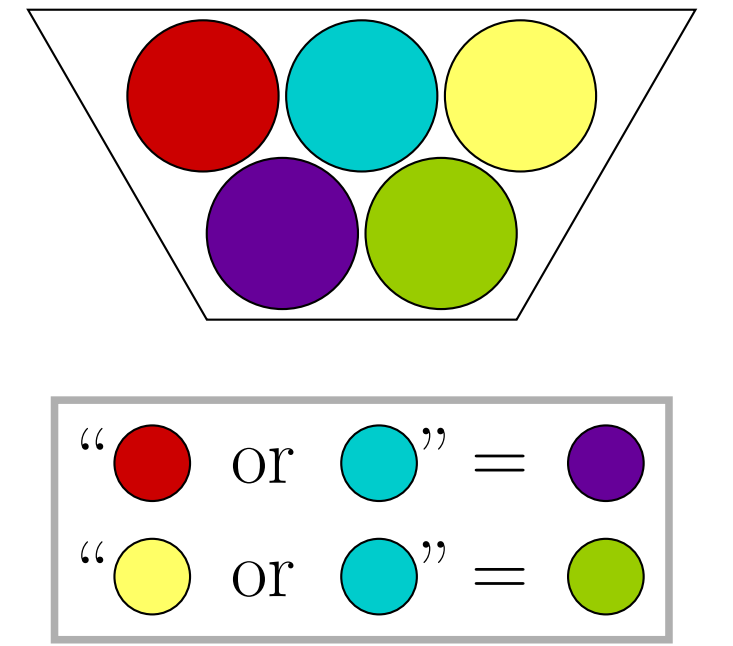


Ontic data imprecision:

- Precise observation of something imprecise
- Truth is represented by coarse observations

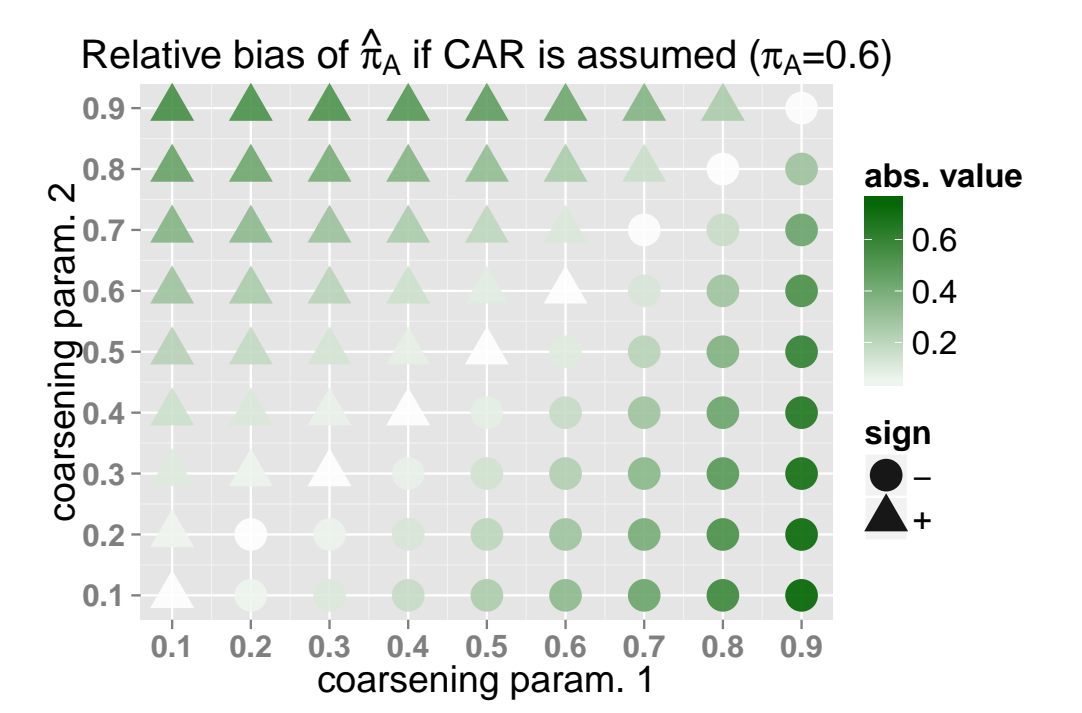
Example:

Answers of indecisive respondents



Already existing approaches

- Still common to impose strict assumptions and thus **to enforce precise results** \Rightarrow Problem: results may be substantially biased (cf. Figure; π_A is parameter of interest, CAR is only satisfied if coarsening parameter 1 = coarsening parameter 2)
- There is a variety of different **set-valued approaches** aiming at a proper reflection of the available information
 - using a Bayesian point of view (e.g. de Cooman, Zaffalon, Artif. Intell)
 - via likelihood-based belief function (Denoeux, 2014, IJAR)
 - via the profile likelihood (e.g. Cattaneo, Wiencierz, 2012, IJAR)
 - **Here:** Likelihood-based approach, strongly influenced by the methodology of partial identification, coarse categorical data only



Basic idea

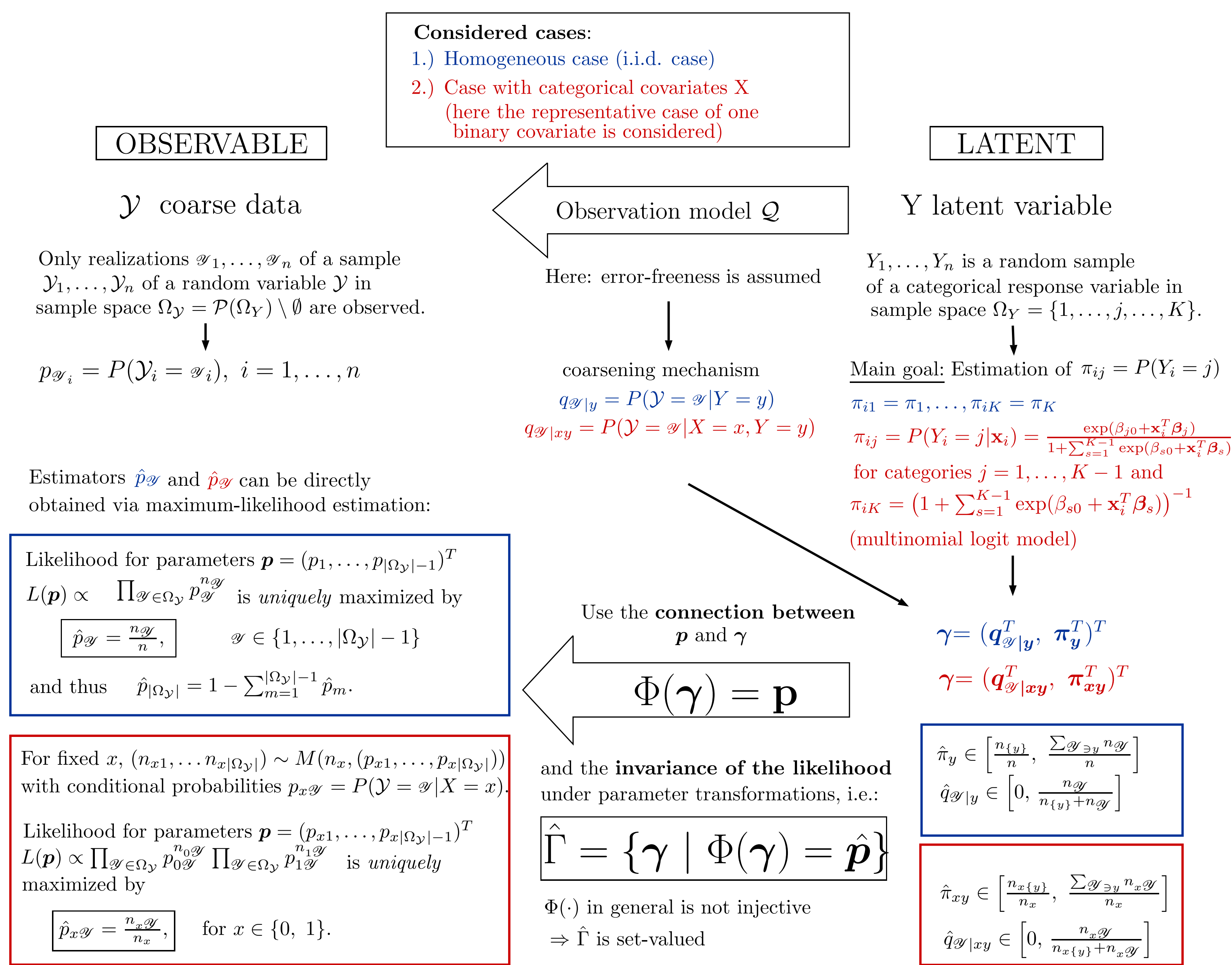


Illustration 1:

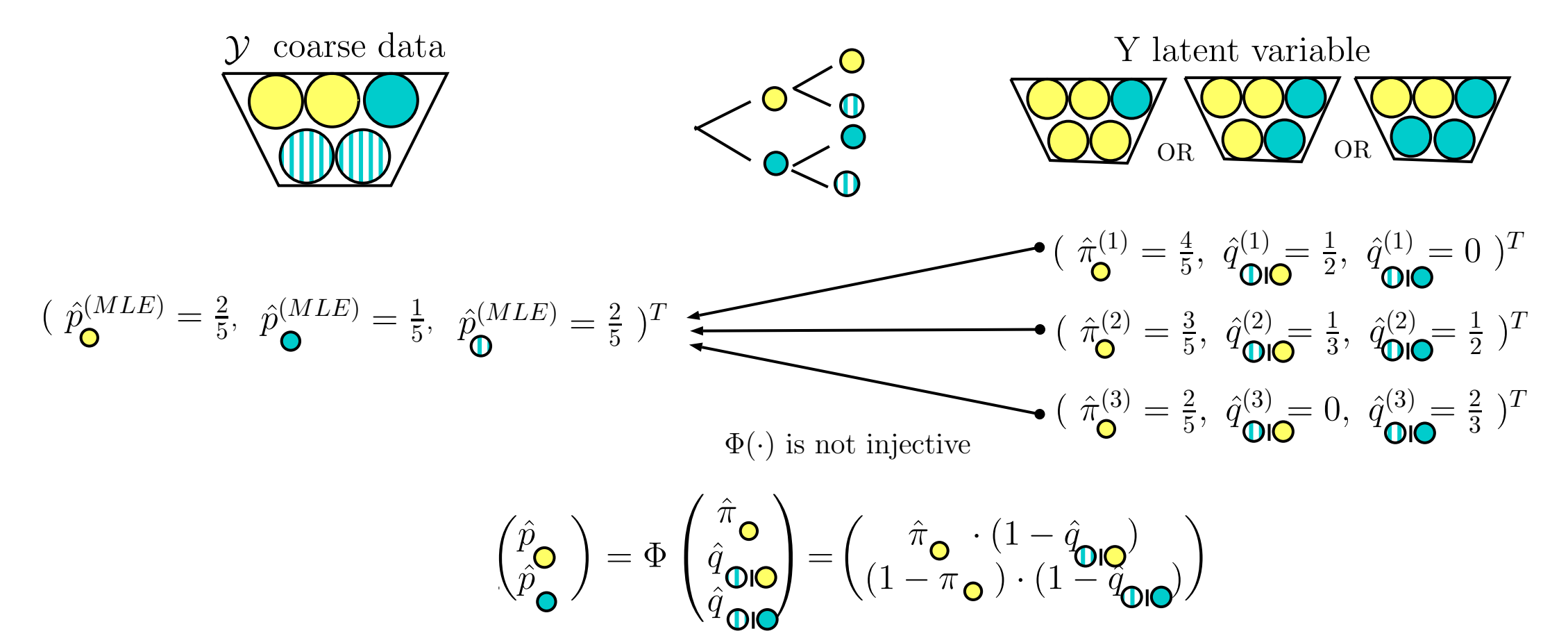


Illustration 2: PASS-Data

(German panel study "Labour market and Social Security"; Trappmann, Gundert, Wenzig, Gebhardt, 2010, Schmollers Jahrbuch)

- Here: $\Omega_Y = \{< 1000\text{€}(<), \geq 1000\text{€}(\geq)\}$
- Some respondents give no suitable answer ("na"; i.e. coarse answer "< 1000€ or \geq 1000€"): $\Omega_Y = \{<, \geq, \text{na}\}$

i.i.d. case
 $n_{<} = 238, n_{\geq} = 835, n_{\text{na}} = 338 \Rightarrow \hat{\pi}_{<} \in \left[\frac{238}{1411}, \frac{238+338}{1411} \right]$

Categorical covariate Unemployment Benefit II (UBII)

		income			
		<	\geq	na	total
UBII	yes (0)	130	114	75	319
	no (1)	108	108	263	1092
	total	238	835	338	1411

$\hat{\pi}_{0<} \in \left[\frac{130}{319}, \frac{130+75}{319} \right]$
 $\hat{\pi}_{1<} \in \left[\frac{108}{1092}, \frac{108+263}{1092} \right]$ or
 $\hat{\beta}_{<0} \in [-0.37, 0.59]$
 $\hat{\beta}_{<1} \in [-1.83, -1.25]$

Reliable Incorporation of Auxiliary Information

Starting from point-identifying assumptions, we use sensitivity parameters to allow inclusion of partial knowledge.

Assumption	Coarsening at random (CAR) and its generalization $R = \frac{1-q_{na >}}{1-q_{na <}}$	Subgroup independent coarsening (SIC) and its generalizations $R_1 = \frac{1-q_{na 0<}}{1-q_{na 1<}}$ & $R_2 = \frac{1-q_{na 0\geq}}{1-q_{na 1\geq}}$
Illustration	<p>CAR $q_{na <} = q_{na >}$ i.e. probability of "na" does not depend on true income category</p> <p>Assumptions about exact value of R e.g. $R=1$ $R=4$, where $R=1$ corresponds to CAR.</p> <p>Rough evaluation of R low income group has a higher tendency to report in a precise way e.g. $R \leq 1$</p>	<p>SIC $q_{na 0<} = q_{na 1<}$ and $q_{na 0\geq} = q_{na 1\geq}$ Reporting "na" does not depend on the receipt of the UBII</p> <p>Assumption about the exact value of R_1 and R_2 Knowledge about the relative magnitude of precise observations in both subgroups allows more flexible inclusion of information $R_1 = R_2 = 1$ represents SIC</p> <p>Rough evaluation of R_1 and R_2 Partial knowledge about the relative magnitude</p>
Estimators	$\hat{\pi}_{<} = \frac{n_{<}}{n_{<} + n_{\geq}}$ $\hat{\pi}_{<} = \frac{n_{<4}}{n_{>} + n_{<4}}$ $\hat{\pi}_{<} = \left[\frac{n_{<}}{n}, \frac{n_{<}}{n_{<} + n_{\geq}} \right]$	<p>Exemplary for SIC</p> $\hat{\pi}_{0<} = \frac{n_{0<}}{n_0} = \frac{n_{1>}n_0 - n_1n_{0>}}{n_0n_0 + n_1n_0 - n_0n_1 - n_1n_{0>}}$ $\hat{\pi}_{1<} = \frac{n_{1<}}{n_1} = \frac{n_{1>}n_1 - n_1n_{1>}}{n_1n_0 + n_1n_1 - n_0n_1 - n_1n_{1>}}$

Summary and Outlook

- Via the observation model \mathcal{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 restrictions on \mathcal{Q}
- Next steps:**
- Likelihood-based hypothesis tests and uncertainty regions for coarse categorical data
 - Inclusion of auxiliary information by sets of priors
 - Consideration of other "deficiency" processes
 - Extension to metric covariates?