

# Testing of coarsening mechanisms: Coarsening at random versus subgroup independence

SMPS 2016, Rome

Julia Plass\*, Marco Cattaneo\*\*,  
Georg Schollmeyer\*, Thomas Augustin\*

\*Department of Statistics, Ludwig-Maximilians University and

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



12<sup>th</sup> of September 2016

# What's the problem?

Common dealing with incomplete data: assumptions

⇒ Missing at random (MAR) / coarsening at random (CAR)

⇒ Frequently: assumptions only for pragmatic reasons

# What's the problem?

Common dealing with incomplete data: assumptions

- ⇒ Missing at random (MAR) / coarsening at random (CAR)
- ⇒ Frequently: assumptions only for pragmatic reasons



## Two types of uninformative coarsening:

The coarsening is independent of ...

- ... the true underlying value (CAR)
- ... the covariate's value (subgroup independence (SI))

## Two types of uninformative coarsening:

The coarsening is independent of ...

- ... the true underlying value (CAR)
- ... the covariate's value (subgroup independence (SI))

**Here:** categorical setting

- 1.) No assumptions about the coarsening
- 2.) CAR: intuitive insight why not testable
- 3.) SI: construction of Likelihood Ratio test

# Coarse data (categorical setting)

Epistemic interpretation of coarse data (Couso, Dubois, 2014):

LATENT

random sample

$$Y_1, \dots, Y_n$$
$$y_i \in \Omega_Y$$
$$i = 1, \dots, n$$

coarsening

$$y_i \in \mathfrak{y}_i$$

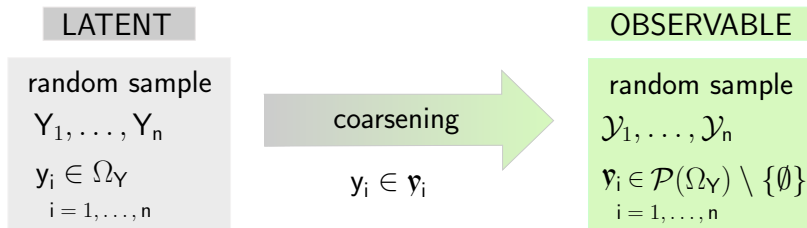
OBSERVABLE

random sample

$$\mathcal{Y}_1, \dots, \mathcal{Y}_n$$
$$\mathfrak{y}_i \in \mathcal{P}(\Omega_Y) \setminus \{\emptyset\}$$
$$i = 1, \dots, n$$

# Coarse data (categorical setting)

Epistemic interpretation of coarse data (Couso, Dubois, 2014):

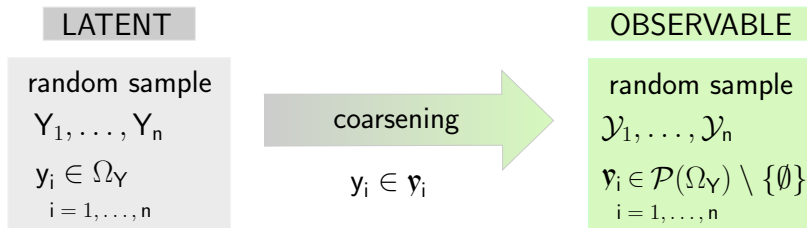


		$\mathcal{Y}$							sum
		{a}	{b}	{c}	{a,b}	{a,c}	{b,c}	{a,b,c}	
X	0	$n_{0\{a\}}$	$n_{0\{b\}}$	$n_{0\{c\}}$	$n_{0\{a,b\}}$	$n_{0\{a,c\}}$	$n_{0\{b,c\}}$	$n_{0\{a,b,c\}}$	$n_0$
	1	$n_{1\{a\}}$	$n_{1\{b\}}$	$n_{1\{c\}}$	$n_{1\{a,b\}}$	$n_{1\{a,c\}}$	$n_{1\{b,c\}}$	$n_{1\{a,b,c\}}$	$n_1$

Table: Contingency table for coarse data

# Coarse data (categorical setting)

Epistemic interpretation of coarse data (Couso, Dubois, 2014):



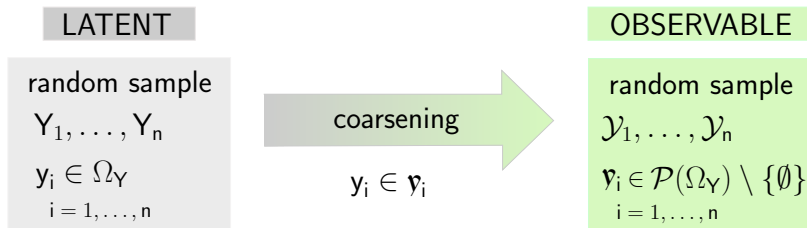
		$\mathcal{Y}$			sum
		{a}	{b}	{a,b}	
X	0	$n_{0\{a\}}$	$n_{0\{b\}}$	$n_{0\{a,b\}}$	$n_0$
	1	$n_{1\{a\}}$	$n_{1\{b\}}$	$n_{1\{a,b\}}$	$n_1$

Table: Contingency table for missing data



# Coarse data (categorical setting)

Epistemic interpretation of coarse data (Couso, Dubois, 2014):



- **UBII**: receipt of Unemployment Benefit II

- $\mathfrak{y}$ : categorical income

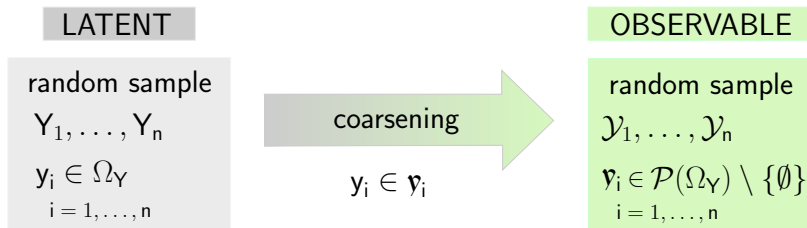
- $\{a\}$ :  $< 1000\text{€}$
- $\{b\}$ :  $\geq 1000\text{€}$
- $\{a,b\}$ :  $< \text{or } \geq$

		$\mathfrak{y}$			
		$\{a\}$	$\{b\}$	$\{a,b\}$	sum
UBII	0	38	385	95	518
	1	36	42	9	87

Table: PASS, w5 (Trappmann, 2010)

# Coarse data (categorical setting)

Epistemic interpretation of coarse data (Couso, Dubois, 2014):



		$\mathcal{Y}$		sum
		a	b	
UBII	0	38+95	385	518
	1	36+9	42	87

Table: potential true table

		$\mathcal{Y}$			sum
		{a}	{b}	{a,b}	
UBII	0	38	385	95	518
	1	36	42	9	87

Table: observed table

# Coarse data (categorical setting)

Epistemic interpretation of coarse data (Couso, Dubois, 2014):

LATENT

random sample

$$Y_1, \dots, Y_n$$

$$y_i \in \Omega_Y$$

$$i = 1, \dots, n$$

coarsening

$$y_i \in \mathfrak{y}_i$$

OBSERVABLE

random sample

$$\mathcal{Y}_1, \dots, \mathcal{Y}_n$$

$$\mathfrak{y}_i \in \mathcal{P}(\Omega_Y) \setminus \{\emptyset\}$$

$$i = 1, \dots, n$$

		$\mathcal{Y}$		sum
		a	b	
UBII	0	38	385+ <b>95</b>	518
	1	36	42+ <b>9</b>	87

Table: potential true table

		$\mathcal{Y}$			sum
		{a}	{b}	{a,b}	
UBII	0	38	385	<b>95</b>	518
	1	36	42	<b>9</b>	87

Table: observed table

# Coarse data (categorical setting)

Epistemic interpretation of coarse data (Couso, Dubois, 2014):

LATENT

random sample

$$Y_1, \dots, Y_n$$

$$y_i \in \Omega_Y$$

$$i = 1, \dots, n$$

coarsening

$$y_i \in \mathfrak{y}_i$$

OBSERVABLE

random sample

$$\mathcal{Y}_1, \dots, \mathcal{Y}_n$$

$$\mathfrak{y}_i \in \mathcal{P}(\Omega_Y) \setminus \{\emptyset\}$$

$$i = 1, \dots, n$$

		$\mathcal{Y}$		sum
		a	b	
UBII	0	38+20	385+75	518
	1	36+5	42+4	87

Table: potential true table

		$\mathcal{Y}$			sum
		{a}	{b}	{a,b}	
UBII	0	38	385	95	518
	1	36	42	9	87

Table: observed table

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

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

LATENT

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

OBSERVABLE

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

LATENT

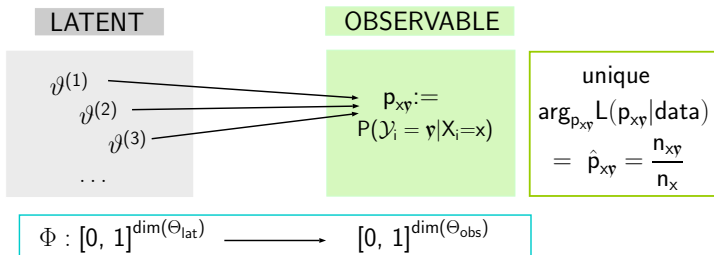
$$\vartheta = (\pi_{xy}^T, \mathbf{q}_{y|xy}^T)^T$$

OBSERVABLE

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

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

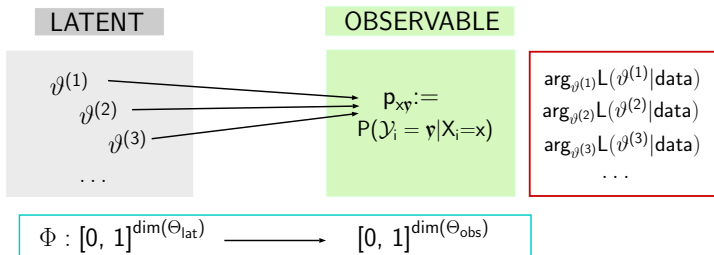
1.) Determine MLE of observed variable distribution



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

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



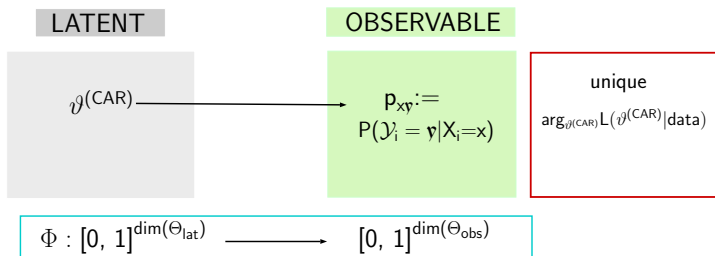


- 1.) Determine MLE of observed variable distribution
- 2.) Use connection between both worlds

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

- 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

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

- 3.) Use invariance of the likelihood

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

# Coarsening at random (CAR)

Definition of CAR (Heitjan, Rubin, 1991):

For each fixed  $\mathbf{y}$  and  $x$ ,  $q_{\mathbf{y}|xy}$  takes the same value  $\forall \mathbf{y} \in \mathcal{Y}$ .

# Coarsening at random (CAR)

Definition of CAR (Heitjan, Rubin, 1991):

For each fixed  $y$  and  $x$ ,  $q_{y|xy}$  takes the same value  $\forall y \in \mathcal{Y}$ .

Illustrated by the example:

- The probability of  $\{a,b\}$  is taken to be independent of the true income category in both subgroups split by UBI:

$$q_{\{a,b\}|0a} = q_{\{a,b\}|0b} \quad \text{and} \quad q_{\{a,b\}|1a} = q_{\{a,b\}|1b}$$



# Coarsening at random (CAR)

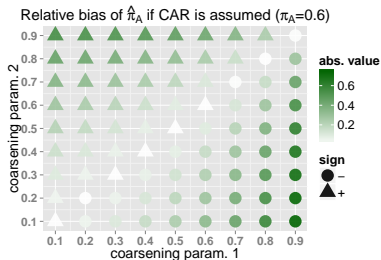
Definition of CAR (Heitjan, Rubin, 1991):

For each fixed  $\mathbf{y}$  and  $\mathbf{x}$ ,  $q_{\mathbf{y}|\mathbf{x}\mathbf{y}}$  takes the same value  $\forall \mathbf{y} \in \mathbf{y}$ .

Illustrated by the example:

- The probability of  $\{a,b\}$  is taken to be independent of the true income category in both subgroups split by UBII:

$$q_{\{a,b\}|0a} = q_{\{a,b\}|0b} \quad \text{and} \quad q_{\{a,b\}|1a} = q_{\{a,b\}|1b}$$



# Coarsening at random (CAR)

Definition of CAR (Heitjan, Rubin, 1991):

For each fixed  $y$  and  $x$ ,  $q_{y|xy}$  takes the same value  $\forall y \in \mathcal{Y}$ .

Illustrated by the example:

- The probability of  $\{a,b\}$  is taken to be independent of the true income category in both subgroups split by UBII:

$$q_{\{a,b\}|0a} = q_{\{a,b\}|0b} \quad \text{and} \quad q_{\{a,b\}|1a} = q_{\{a,b\}|1b}$$



- Resulting estimators:

$$\hat{\pi}_{xa}^{(CAR)} = \frac{n_{x\{a\}}}{n_{x\{a\}} + n_{x\{b\}}}, \quad \hat{q}_{\{a,b\}|xa}^{(CAR)} = \hat{q}_{\{a,b\}|xb}^{(CAR)} = \frac{n_{x\{a,b\}}}{n_x}$$

Estimators for subgroup 0:

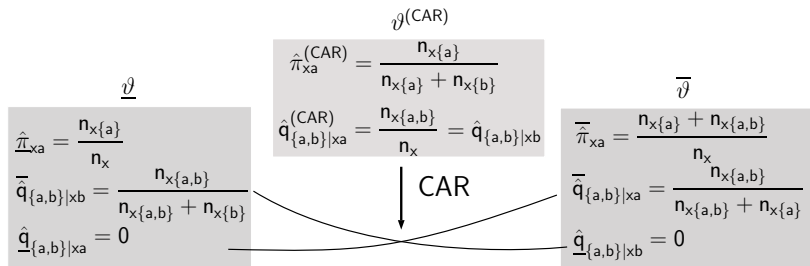
$$\hat{\pi}_{0a} \in [0.07, 0.26], \quad \hat{q}_{\{a,b\}|0a} \in [0, 0.71], \quad \hat{q}_{\{a,b\}|0b} \in [0, 0.20]$$

$$\hat{\pi}_{0a}^{(CAR)} = 0.09, \quad \hat{q}_{\{a,b\}|0a}^{(CAR)} = 0.18, \quad \hat{q}_{\{a,b\}|0b}^{(CAR)} = 0.18$$

Estimators for subgroup 0:

$$\hat{\pi}_{0a} \in [0.07, 0.26], \quad \hat{q}_{\{a,b\}|0a} \in [0, 0.71], \quad \hat{q}_{\{a,b\}|0b} \in [0, 0.20]$$

$$\hat{\pi}_{0a}^{(CAR)} = 0.09, \quad \hat{q}_{\{a,b\}|0a}^{(CAR)} = 0.18, \quad \hat{q}_{\{a,b\}|0b}^{(CAR)} = 0.18$$



$\Rightarrow$  CAR is generally not testable, unless further assumptions about the coarsening are justified



# Subgroup independence (SI)

**Subgroup independence** (Plass, Augustin, Cattaneo, Schollmeyer, 2015):

For each fixed  $\mathfrak{y}$  and  $y \in \mathfrak{y}$ ,  $q_{\mathfrak{y}|xy}$  takes the same value  $\forall x \in \Omega_X$ .

# Subgroup independence (SI)

**Subgroup independence** (Plass, Augustin, Cattaneo, Schollmeyer, 2015):

For each fixed  $\mathfrak{y}$  and  $y \in \mathfrak{y}$ ,  $q_{\mathfrak{y}|xy}$  takes the same value  $\forall x \in \Omega_X$ .

Illustrated by the example:

- The probability of  $\{a, b\}$  is taken to be **independent of the receipt of the UBI** given  $y$ :

$$q_{\{a,b\}|0a} = q_{\{a,b\}|1a} \quad \text{and} \quad q_{\{a,b\}|0b} = q_{\{a,b\}|1b}$$



# Subgroup independence (SI)

**Subgroup independence** (Plass, Augustin, Cattaneo, Schollmeyer, 2015):

For each fixed  $\mathfrak{y}$  and  $y \in \mathfrak{y}$ ,  $q_{\mathfrak{y}|xy}$  takes the same value  $\forall x \in \Omega_X$ .

Illustrated by the example:

- The probability of  $\{a, b\}$  is taken to be **independent of the receipt of the UBI** given  $y$ :

$$q_{\{a,b\}|0a} = q_{\{a,b\}|1a} \quad \text{and} \quad q_{\{a,b\}|0b} = q_{\{a,b\}|1b}$$



- Resulting estimators (if well-defined and inside  $[0,1]$ ):

$$\hat{\pi}_{xa}^{(SI)} = \frac{n_{x\{a\}}}{n_x} \frac{v}{w}, \quad \hat{q}_{\{a,b\}|xa}^{(SI)} = 1 - \frac{w}{v}, \quad \hat{q}_{\{a,b\}|xb}^{(SI)} = 1 - \frac{w}{z}$$

with  $v = n_{0\{a\}}n_{1\{b\}} - n_{0\{b\}}n_{1\{a\}}$ ,  $w = n_{0\{a\}}n_{1\{b\}} - n_{0\{b\}}n_{1\{a\}}$  and  $z = n_{0A}n_{1\{b\}} - n_{1A}n_{0\{b\}}$

Estimators for subgroup 0:

$$\hat{\pi}_{0a} \in [0.07, 0.26], \quad \hat{q}_{\{a,b\}|0a} \in [0, 0.71], \quad \hat{q}_{\{a,b\}|0b} \in [0, 0.198]$$

$$\hat{\pi}_{0a}^{(SI)} = 0.42, \quad \hat{q}_{\{a,b\}|0a}^{(SI)} = -0.04, \quad \hat{q}_{\{a,b\}|0b}^{(SI)} = 0.20$$

Estimators for subgroup 0:

$$\hat{\pi}_{0a} \in [0.07, 0.26], \quad \hat{q}_{\{a,b\}|0a} \in [0, 0.71], \quad \hat{q}_{\{a,b\}|0b} \in [0, 0.198]$$

$$\hat{\pi}_{0a}^{(SI)} = 0.42, \quad \hat{q}_{\{a,b\}|0a}^{(SI)} = -0.04, \quad \hat{q}_{\{a,b\}|0b}^{(SI)} = 0.20$$

- $\Rightarrow$  There are data situations that might hint to (partial) incompatibility with SI
- $\Rightarrow$  SI is testable in our setting

- Hypothesis:

$H_0$  :  $q_{\mathbf{y}|0\mathbf{y}} = q_{\mathbf{y}|1\mathbf{y}}$  for all  $y \in \Omega_Y = \{a,b\}$ ,  $\mathbf{y} \in \mathcal{P}(\Omega_Y) \setminus \emptyset$

$H_1$  :  $q_{\mathbf{y}|0\mathbf{y}} \neq q_{\mathbf{y}|1\mathbf{y}}$  for some  $y \in \Omega_Y = \{a,b\}$ ,  $\mathbf{y} \in \mathcal{P}(\Omega_Y) \setminus \emptyset$ .

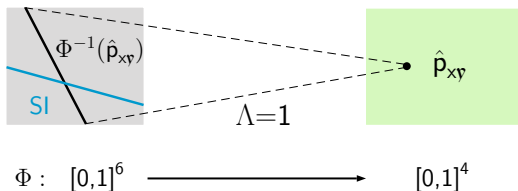
- Test based on test statistic:

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

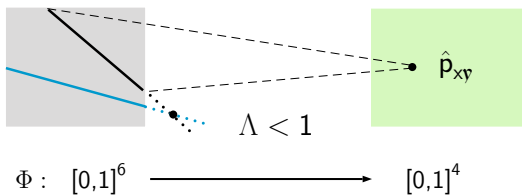
with  $\vartheta = (\pi_{0a}, \pi_{1a}, \mathbf{q}_{\{a,b\}|0a}, \mathbf{q}_{\{a,b\}|1a}, \mathbf{q}_{\{a,b\}|0b}, \mathbf{q}_{\{a,b\}|1b})^T$

# Testing SI: Sensitivity of $\Lambda$

- No evidence to reject SI



- Some evidence to reject SI



- Result of the data example:

$\Lambda \approx 0.93 \Rightarrow$  slight evidence against SI

# Studying more general settings...

Now: number of subgroups  $m$ ,  $|\Omega_Y| = k$

**CAR** is generally known to be ...

- ... point-identifying  
(Heitjan, Rubin, 1991)
- ... not testable  
(e.g. Jaeger, 2006)



# Studying more general settings...

Now: number of subgroups  $m$ ,  $|\Omega_Y| = k$

**CAR** is generally known to be ...

- ... point-identifying  
(Heitjan, Rubin, 1991)
- ... not testable  
(e.g. Jaeger, 2006)

**SI:** A determination of the number of degrees of freedom  $df$  is crucial:

$$df = \dim(\Theta_{\text{obs}}) - \dim(\Theta_{\text{lat, SI}}) = 2^{m-1}(2k-m) - (m+1)(k-1) - 1$$

- Point-identification and testability are only valid if sufficient subgroups are available inducing  $df \geq 0$
  - If  $df \geq 0$ , a LR-Test on SI may be constructed with ...
    - ... test statistic  $T = -2 \cdot \log(\Lambda)$  (Wilks, 1938)
    - ... asymptotic distribution of  $T$  under  $H_0$ 
      - $\frac{1}{2} \delta_0 + \frac{1}{2} \chi_1^2$  if  $df=0$
      - $\chi_{df}^2$  if  $df > 0$ ,
- where  $\delta_0$  is the Dirac distribution at 0 (Chernoff, 1954)

# Summary: CAR versus SI





	CAR	SI
Point-identifying?	always	in specific settings
Testability	generally impossible	possible in specific settings

## Construction of a hypothesis test:

- $H_0$ : SI,  $H_1$ : no SI
- Test statistic based on the Likelihood Ratio

**Next step:** Generalized version of LR-Test on SI

# References

-  Couso, Dubois.  
*Statistical Reasoning with Set-Valued Information: Ontic vs. Epistemic Views, IJAR, 2014.*
-  Chernoff.  
*On the distribution of the likelihood ratio, Ann. Stat. Math., 1954.*
-  Heitjan, Rubin.  
*Ignorability and Coarse Data, Annals of Statistics, 1991.*
-  Jaeger.  
*On testing the missing at random assumption, ECML, 2006.*
-  Manski.  
*Partial Identification of Probability Distributions, Springer, 2003.*
-  Plass, Augustin, Cattaneo, Schollmeyer.  
*Statistical modelling under epistemic data imprecision, ISIPTA, 2015.*
-  Trappmann, Gundert, Wenzig, Gebhardt.  
*PASS: a household panel survey for research on unemployment and poverty, Schmollers Jahrbuch, 2010.*
-  Wilks.  
*The large-sample distribution of the likelihood ratio for testing composite hypotheses, Ann. Stat., 1938.*