

# **Imprecise Measurement Error Models and Partial Identification - Towards a Unified Approach for Non-Idealized Data**

**Thomas Augustin**

Department of Statistics  
Ludwig-Maximilians University Munich

# Table of Contents

## 1. Background and Sketch of the Arguments

## 2. Measurement Error Correction based on Precise Error Models

2.1 Measurement Error Modelling

2.2 Unbiased Estimating Equations and Corrected Score Functions for Classical Measurement Error (in the Cox Model)

2.3 Extended Corrected Score Functions – A Unified View at Measurement Error and Censoring

2.4 Corrected Score Functions for Berkson Models

(2.5) (Unconditionally Corrected Score Functions and Rounding) (Felderer)

## 3. Overcoming the Dogma of Ideal Precision in Deficiency Models

3.1 Credal Deficiency Model as Imprecise Measurement Error Models

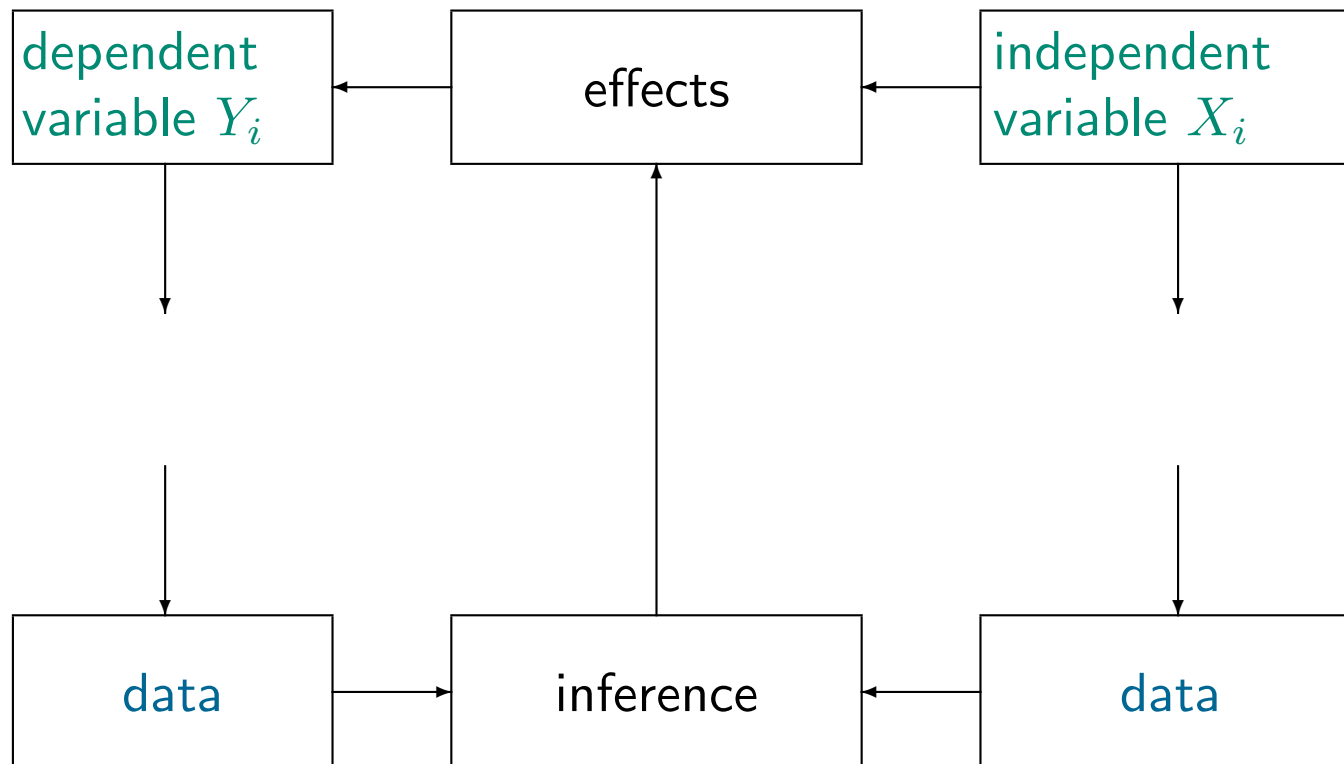
2.2 Credal Consistency of Set-Valued Estimators

2.3 Minimal and Complete Sets of Unbiased Estimating Functions

# 1. Background and Sketch of the Arguments

# 1. Background and Sketch of the Arguments

Applied Statistics: Learning from data by sophisticated models  
Complex relationships between variables

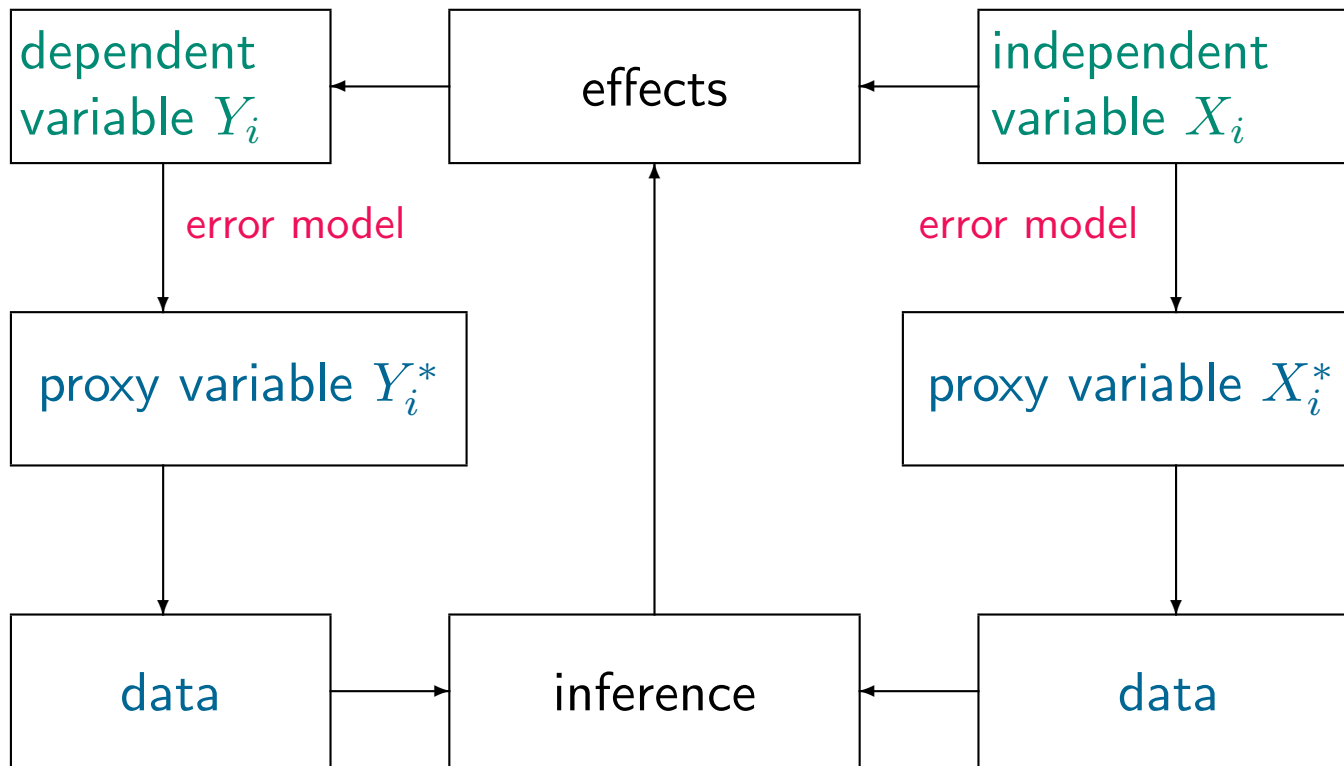


Often the relationship between **variables** and **data** is complex, too:

- \* Often **variables of interest (gold standard)** are not ascertainable.
- \* Only **proxy variables** (surrogates) are available instead.

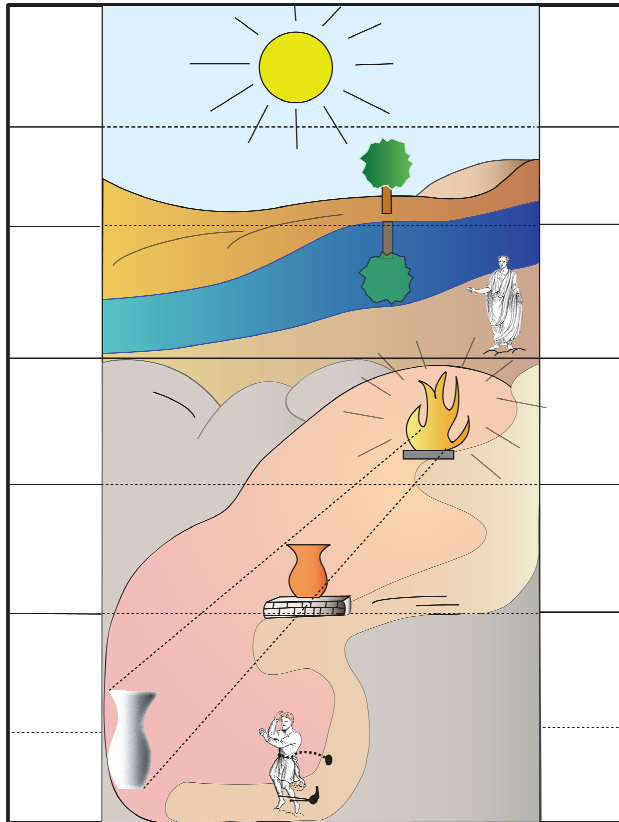
# Typical examples: Measurement Error

- Error-prone measurements of true quantities
  - \* error in technical devices
  - \* indirect measurement
  - \* response effects
  - \* use of aggregated quantities, averaged values, imputation, rough estimates etc.
  - \* anonymization of data by deliberate contamination
- Operationalization of complex constructs; latent variables
  - \* long term quantities: permanent income,
  - \* importance of a patent
  - \* extent of motivation, degree of customer satisfaction
  - \* severeness of undernutrition



# Cave Allegory

[http://commons.wikimedia.org/wiki/File:Allegory\\_of\\_the\\_Cave\\_blank.png](http://commons.wikimedia.org/wiki/File:Allegory_of_the_Cave_blank.png)





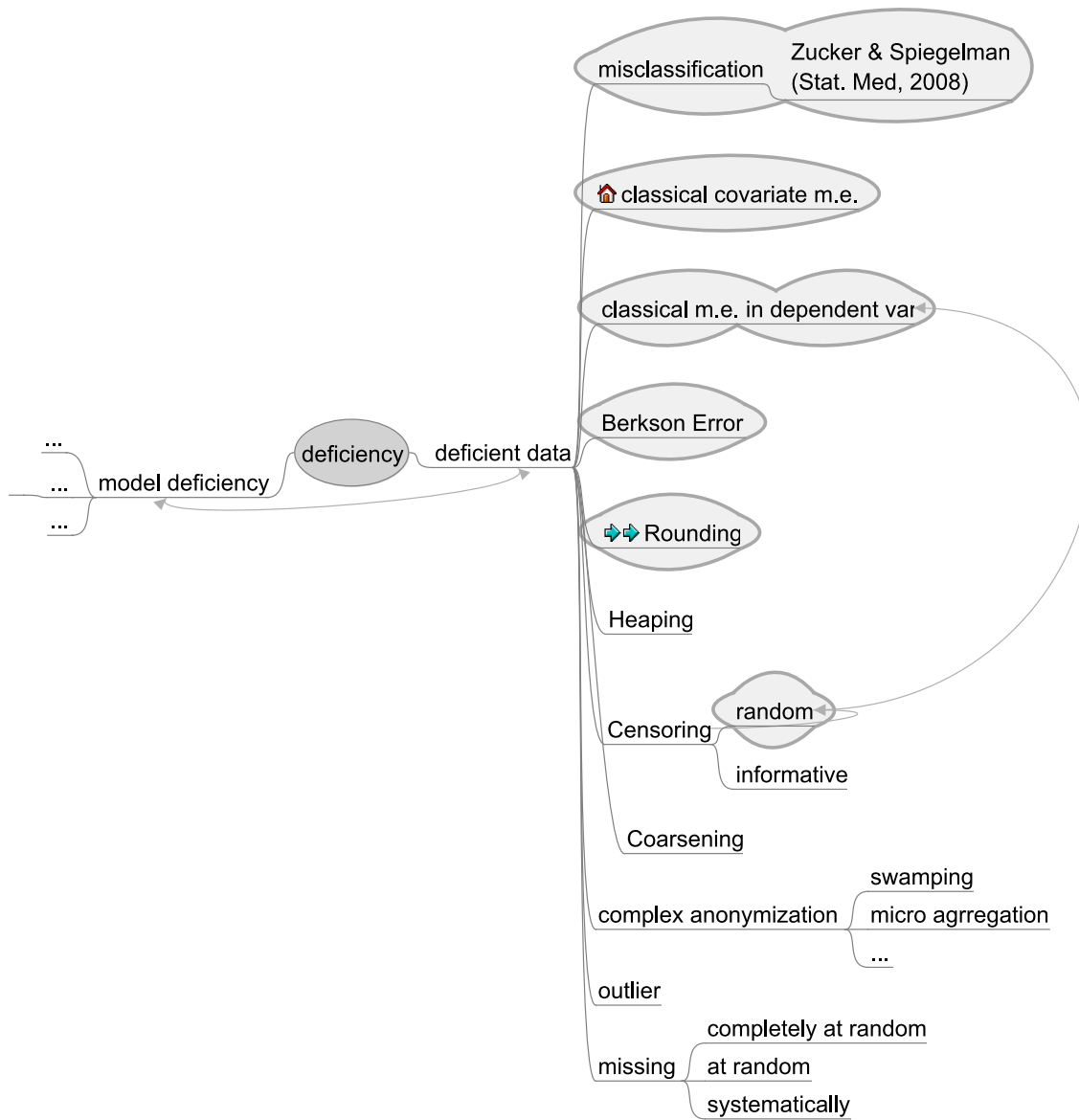
# Notation

We have to distinguish between true (correctly measured) variable and its (possible incorrect) measurement, i.e. between the **gold standard** and the corresponding **surrogate**.

\* - Notation (here)

$X, Z$  : (unobservable) variable, gold standard

$X^*, Z^*$ : corresponding possibly incorrect measurements analogously:  $Y, Y^*$  and  $T, T^*$



# Partial Identification and Beyond

## Manski's Law of Decreasing Credibility

### Reliability !? Credibility ?

"The credibility of inference decreases with the strength of the assumptions maintained." (Manski (2003, p. 1))

**Identifying Assumptions** Very strong assumptions needed to ensure identifiability = precise solution

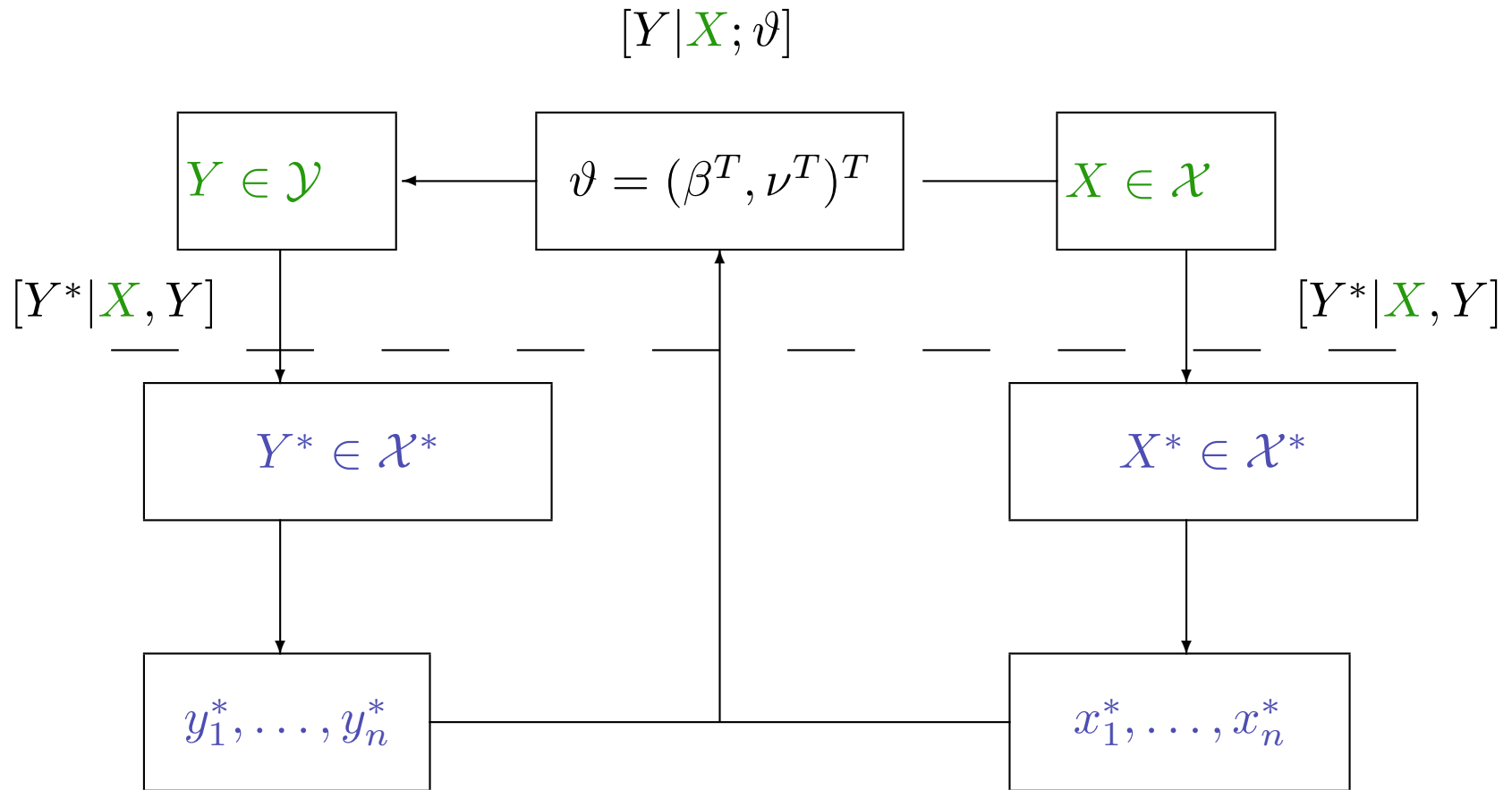
- Measurement error model completely known
  - type of error, in particular assumptions on (conditional) independence
  - type of error distribution
  - moments of error distribution
- validation studies often not available

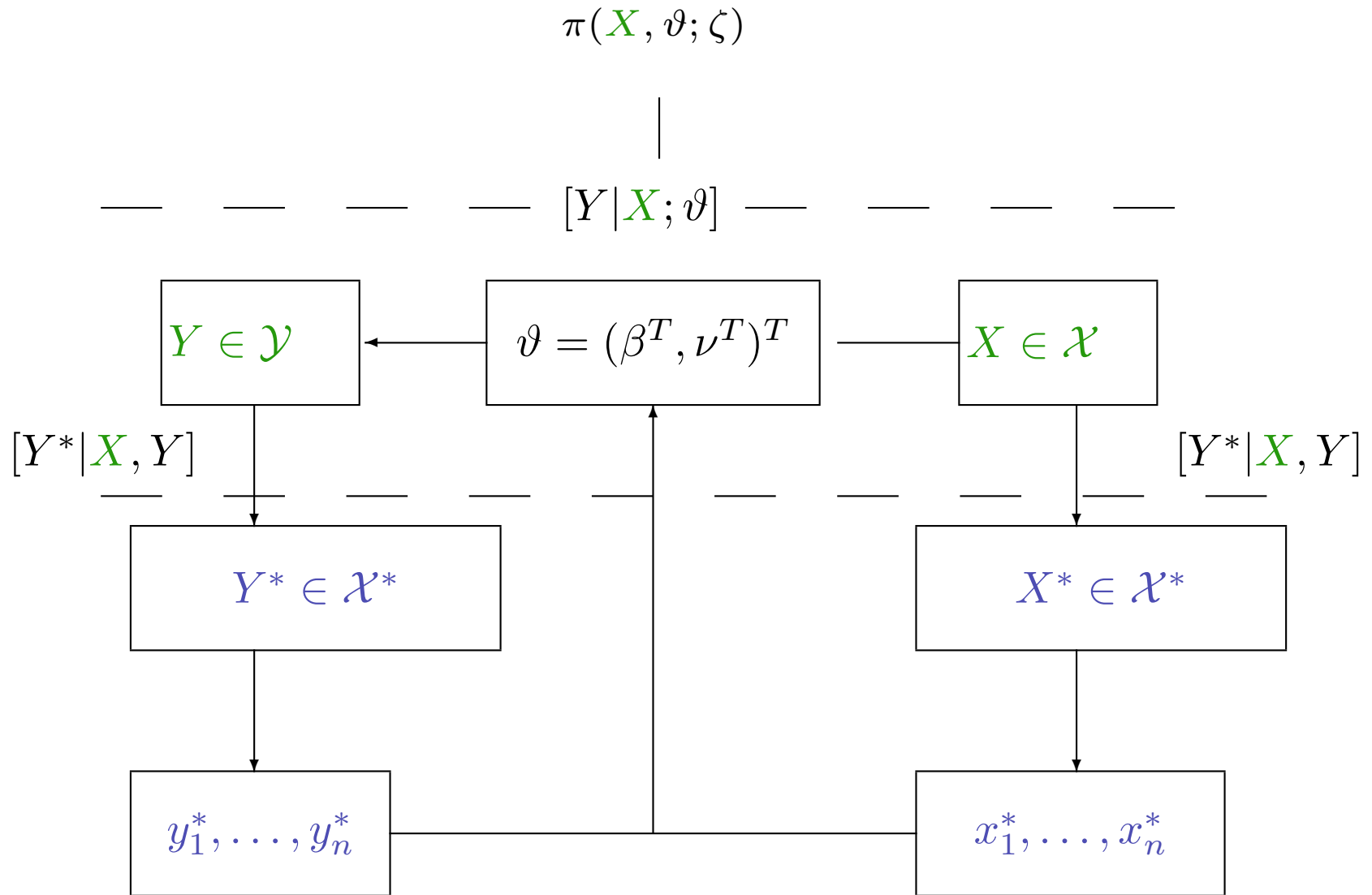
# Reliable Inference Instead of Overprecision!!

- Make more „realistic“ assumptions and let the data speak for themselves!
- Consider the *set* of *all* models that are compatible with the data (and then add successively additional assumptions, if desirable)
- The results may be imprecise, but are more reliable for sure
- **The extent of imprecision is related to the data quality!**
- As a welcome by-product: clarification of the implication of certain assumptions
- parallel development (missing data; transfer to measurement error context!)
  - \* econometrics: *partial identification*: e.g., Manski (2003, Springer)
  - \* biometrics: *systematic sensitivity analysis*: e.g., Vansteelandt, Goetghebeur, Kenward, Molenberghs (2006, Stat. Sinica)
- current developments, e.g.,
  - \* Cheng, Small (2006, JRSSB)
  - \* Henmi, Copas, Eguchi (2007, Biometrics)
  - \* Stoye (2009, Econometrica)
- Kleyer (2009, MSc., LMU); Kunz, Augustin, Küchenhoff (2010, TR)

# Credal Estimation

- Natural idea: sets of traditional models  $\longrightarrow$  sets of traditional estimators
- Construct estimators  $\hat{\Theta} \subseteq \mathbb{R}^p$ , set-valued, typically interval-valued, point estimators appropriately reflecting the ambiguity (non-stochastic uncertainty, ignorance) in the credal set  $\mathcal{P}$ .
- $\hat{\Theta}$  small if, and only if (!),  $\mathcal{P}$  "small"
  - \* Usual point estimator as the border case of precise probabilistic information
  - \* Connection to Manski's (2003) *identification regions* and Vansteelandt, Goetghebeur, Kenward & Molenberghs (Stat Sinica, 2006) *ignorance regions*.
- Construction of unbiased sets of estimating functions
- Credal consistency





# Credal Deficiency Models

Different types of deficiency can be expressed

- Measurement error problems
- Misclassification
- If  $\mathcal{Y}^* \subseteq \mathcal{P}(\mathcal{Y}) \times \{0, 1\}$  : coarsening, rounding, censoring, missing data
- Outliers

*Credal set*: convex set of traditional probability distributions

$$\begin{aligned} [Y|X, \vartheta] &\in \mathcal{P}_{Y|X, \vartheta} \\ [Y^*|X, Y] &\in \mathcal{P}_{Y^*|X, Y} \\ [X^*|X, Y] &\in \mathcal{P}_{X^*|X, Y} \end{aligned}$$

$$\mathcal{P} := \text{conv} \left( \mathcal{P}_{Y|X, \vartheta} \otimes \mathcal{P}_{Y^*|X, Y} \otimes \mathcal{P}_{X^*|X, Y} \right)$$



## **2. Measurement Error Correction based on Precise Error Models**

### **2.1. Measurement Error Modelling**

“Measure what is measurable,  
and make measurable what is not so.”  
(Galilei Galileo)

# Typical examples

- Error-prone measurements of true quantities
  - \* error in technical devices
  - \* indirect measurement
  - \* response effects
  - \* use of aggregated quantities, averaged values, imputation, rough estimates etc.
  - \* anonymization of data by deliberate contamination
- Operationalization of complex constructs; latent variables
  - \* long term quantities: permanent income,
  - \* importance of a patent
  - \* extent of motivation, degree of customer satisfaction
  - \* severeness of undernutrition

## Notation again

We have to distinguish between true (correctly measured) variable and its (possible incorrect) measurement, i.e. between the **gold standard** and the corresponding **surrogate**.

\* - Notation (here)

$X, Z$  : (unobservable) variable, gold standard

$X^*, Z^*$ : corresponding possibly incorrect measurements analogously:  $Y, Y^*$  and  $T, T^*$

# Sources of measurement error

- Induced by an instrument (laboratory value, blood pressure)
- Induced by the study participants (medical doctors or patients; interviewers and respondents)
- Surrogate variables, e.g. "Job exposure matrix: typical, instead of individual, exposure", "economic wealth of a district instead of individual income"
- Measurement error induced by definition, e.g. "long term mean of daily fat intake", "average income"
- Operationalization of complex constructs ("quality of life", "consumer satisfaction")

## Note:

- 'Measurement error' and 'misclassification' are not just a matter of sloppiness.
- Latent variables are eo ipso not exactly measurable.
- “Almost all economic variables are measured with error. [...] Unfortunately, the statistical consequences of errors in explanatory variables are severe.”  
(Davidson and Mackinnon (1993),  
Estimation and Inference in Econometrics.)
- In nonlinear models, the later statement **may apply (!?)** to the dependent variable, too. (Dependence on the DGP: Torelli & Trivellato (1993, J. Econometrics))

# The triple whammy effect of measurement error

Carroll, Ruppert, Stefanski, Crainiceanu (2006, Chap.H.)

- bias
  - masking of features
  - loss of power
- **classical error: "attenuation"**

Results

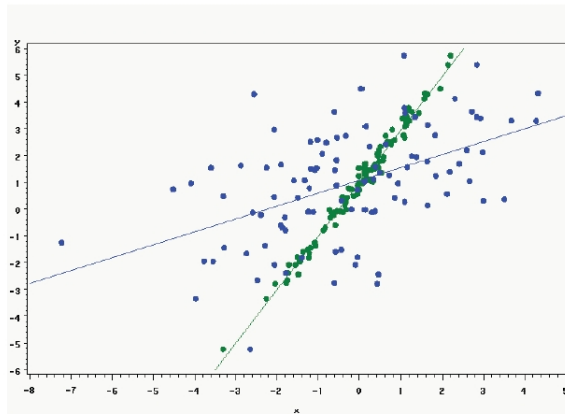
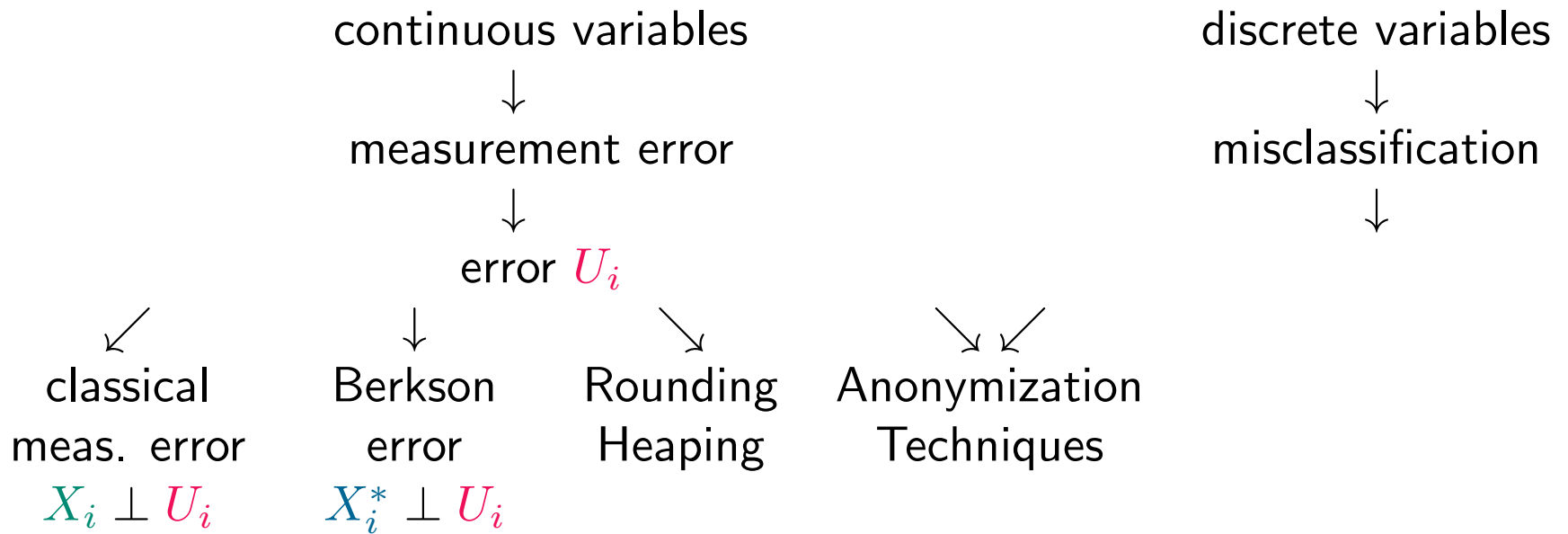


Figure 1: Effect of additive measurement error on linear regression

# Terminology





## Typical Examples: Berkson Error

$$X_i^* = X_i + U_i \text{ but } X_i^* \perp U_i$$

- Experimental design:  $X^*$  target value,  $X$  truly absorbed value
- Aggregated data:
  - \* aggregation for confidentiality, e.g., average income on "block level"
  - \*  $X^*$  mean exposure,  $X$  true individual exposure, JEM
- Omitted variable bias
- Remainder term of regression calibration

## Aggregated data

- $X_1, \dots, X_n$  i.i.d.  $\sim (\mu, \sigma^2)$

- $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \quad \text{plim}_{n \rightarrow \infty} \bar{X} = \mu$

- $X_i \approx (\bar{X}, \sigma^2)$

$$\iff X_i = \bar{X} + \delta \quad \text{with } \delta \sim (0; \sigma^2 (1 - \frac{1}{n})) \quad \delta \perp \bar{X}$$

REFERENCES ARE MISSING!!

## Omitted variable ???

- Linear predictor with  $p = 2$

$$\eta_i = \beta_0 + \beta_1 X_i + \beta_2 Z_i, \quad Z_i \perp X_i$$

- Omitting  $Z_i$  means to work with

$$\eta^* = \beta_0^* + \beta_1^* X_i^*$$

where

$$X_i^* = X_i + \underbrace{\frac{\beta_2}{\beta_1} Z_i}_{\delta_i}$$

- WHAT HAPPENS FOR  $p > 2$ , HOW IS  $Z_i$  SHARED BETWEEN  $X_{i1}, X_{i2}$ ?

# Fehlklassifikation

$Y^*$  beobachtete,  $Y$  wahre Klassenzugehörigkeit

- $Y^* = 1$  Test positiv,  $Y = 1$  krank,
- Betrugsverdachtsfall, zu Schulungsmaßnahme angemeldet, geschätzte Entwicklung des Auftragseingang

	$Y = 1$	$Y = 0$
$Y^* = 1$	$\mathbb{P}(Y^* = 1 Y = 1)$ sensitivity ☺	$\mathbb{P}(Y^* = 1 Y = 0)$ false positive ☹
$Y^* = 0$	$\mathbb{P}(Y^* = 0 Y = 1)$ false negative ☹	$\mathbb{P}(Y^* = 0 Y = 0)$ specificity ☺

# Fehlklassifikationsbias

$$p^* = p \cdot \text{sens} + (1 - p) \cdot (1 - \text{spec}) = p(\text{sens} + \text{spec} - 1) + (1 - \text{spec})$$

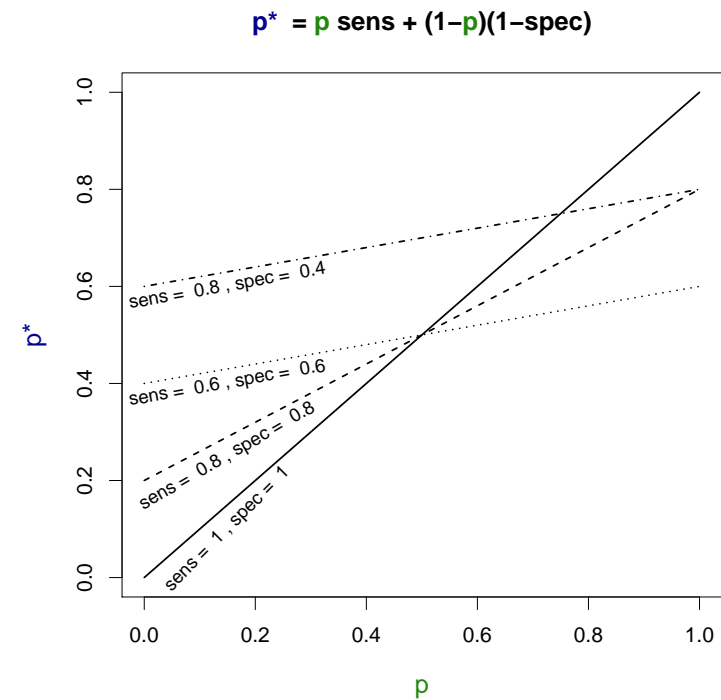
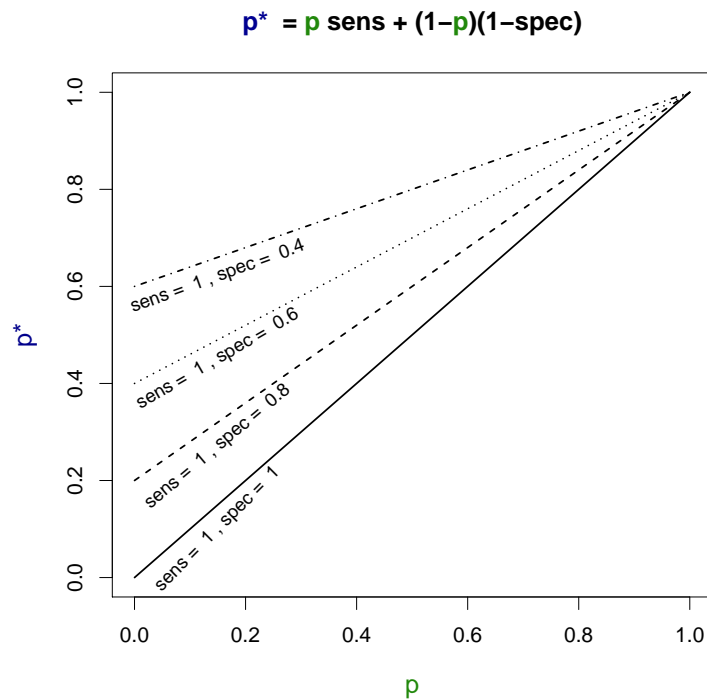


Figure 1: Veranschaulichung Fehlklassifikationsbias: Abweichung von der Winkelhalbierenden

## Typical Examples: Rounding and Heaping

- duration data are commonly collected in a retrospective way
- strong memory effects when time spans have to be remembered, e.g., *Skinner & Humphreys (1999, Lifetime Data Analysis)*
- *Holt et al (1991, Biemer et al(eds.))*: age at menarche
- heaping in episode / spell-based designs: *Torelli & Trivellato (1993, J. Econometrics)*:  
concentration of values of unemployment duration at multiples of six (“identification problem”: heaping versus effect of different levels of compensation)  
strong dependence of the bias on the DGP

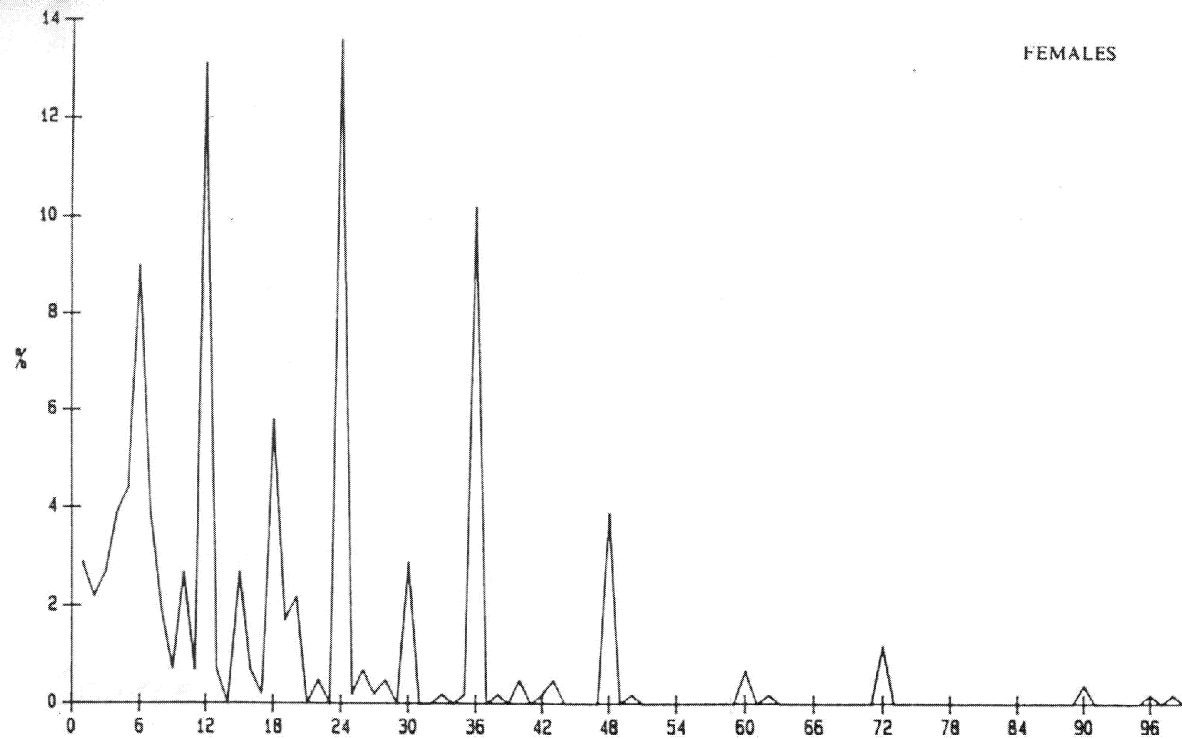


Fig. 2: Percentage distribution of unemployed individuals aged 14–29 years by reported unemployment duration (in months) at initial survey, Italian LFS, matched data for Lombardy, 1986.I–II: males,  $N = 267$  and females,  $N = 411$ .

- heaping in calendar-based designs (German socio-economic panel SOEP)
  - \* distorted values for entry and leave of state of unemployment
  - \* bias analysis under simplified assumptions: Augustin & Wolff (2004, Stat.Papers)
  - \* simulation study (with data constellation based on the SOEP): Wolff & Augustin (2003, ASTA); Jürgens (2007, JRSS A)



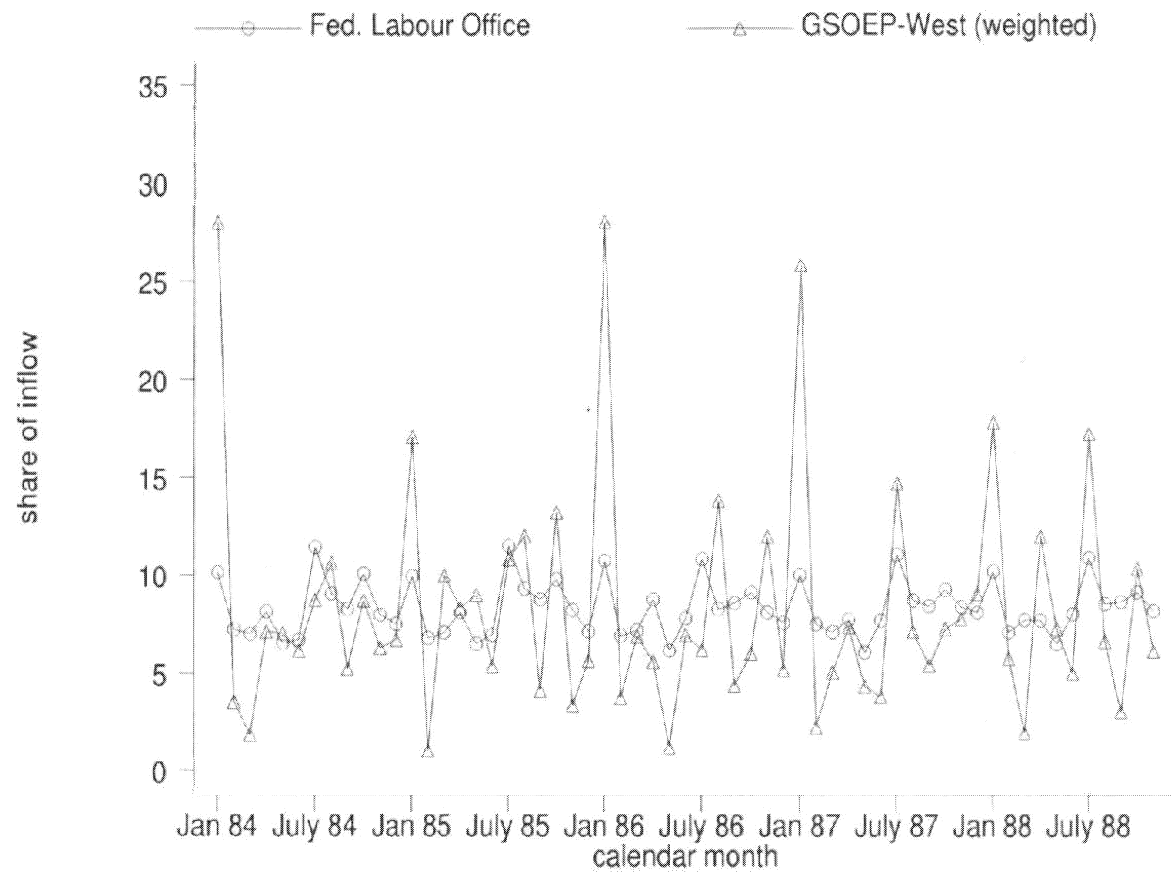


FIGURE 1. Proportion of the annual inflow into registered unemployment in each calendar month – Women, West-Germany.

## Anonymisation Techniques

- Recent trend in official statistics: public use files (Statistisches Bundesamt, 2005, Statistik und Wissenschaft, Bd. 4) and big economic research institutes (IAB, Nuremberg)
- Error mechanism known!
- Distortion by classical measurement error or misclassification
- Often other techniques, e.g. micro-aggregation (Schmid, Schneeweiß, Küchenhoff (2007) Stat. Neerl.; Schmid, (2007, Diss. LMU))
- Growing importance in biometrics
- Discussion on "privacy preserving data mining"

# 2.2 Unbiased Estimating Equations and Corrected Score Functions for Classical Measurement Error (in the Cox Model)

# Corrected Score Functions (and the Cox Model under Covariate Measurement Error)

- Frame the problem in terms of *unbiased estimating functions (score functions)* for the parameter  $\vartheta$

$$s^{\mathbf{X}}(\mathbf{Y}; \mathbf{X}; \vartheta) \quad \text{such that} \quad \mathbb{E}_{\vartheta_0}(s^{\mathbf{X}}(\mathbf{Y}; \mathbf{X}; \vartheta_0)) = 0$$

at the true parameter value  $\vartheta_0$

- \* (Huber: (1981, Wiley) M-estimators; Godambe (1991, Oxford UP); Rieder (1994, Springer): Robust Asymptotic Statistics; Wansbeek & Meijer (2000, Elsevier): GMM)
- \* Under regularity conditions still consistency and asymptotic normality.

- For the moment classical **covariate** measurement error only

$$X_i^* = X_i + U_i, \quad X_i \perp U_i.$$

- Note that typically, even if  $\mathbb{E}(X^*) = \mathbb{E}(X)$   
then  $\mathbb{E}((X^*)^r) \neq \mathbb{E}(X^r), \quad r > 1.$

- Therefore *naive estimation* by simply replacing  $\mathbf{X}$  with  $\mathbf{X}^*$ , leads in general to

$$|\mathbb{E}_{\vartheta_0} (s^{\mathbf{X}}(\mathbf{Y}; \mathbf{X}^*; \vartheta_0))| \geq a > 0,$$

resulting in inconsistent estimators. For instance,

$$\mathbb{E} \left( \sum_{i=1}^n (y_i - \beta_0 - \beta_1 \cdot X_i^*) \begin{pmatrix} 1 \\ X_i^* \end{pmatrix} \right) \neq \mathbb{E} \left( \sum_{i=1}^n (y_i - \beta_0 - \beta_1 \cdot X_i) \begin{pmatrix} 1 \\ X_i \end{pmatrix} \right) = 0$$

- Measurement error correction: Find an estimating function  $s^{\mathbf{X}^*}(\mathbf{Y}, \mathbf{X}^*, \vartheta)$  in the **error prone** data with

$$\mathbb{E}_{\vartheta_0} s^{\mathbf{X}^*}(\mathbf{Y}; \mathbf{X}^*; \vartheta_0) = \mathbf{0}.$$

## Corrected score functions

(Stefanski (1989, Comm. Stat. Theory Meth.), Nakamura (1990, Biometrika))

- Use the ideal score function as a building block,  $s^{X^*}(\mathbf{Y}, \mathbf{X}^*, \vartheta)$  such that

$$\mathbb{E}(s^{X^*}(\mathbf{Y}, \mathbf{X}^*, \vartheta) | \mathbf{X}, \tilde{\mathbf{Y}}) = s^{\mathbf{X}}(\mathbf{Y}, \mathbf{X}, \vartheta),$$

- Then (via iterated expectation!), indeed:

$$\begin{aligned}\mathbb{E}_{\vartheta_0} \left( s^{X^*}(\mathbf{Y}; \mathbf{X}^*; \vartheta_0) \right) &= \mathbb{E}_{\vartheta_0} \left( \mathbb{E}_y \left( s^{X^*}(\mathbf{Y}; \mathbf{X}^*; \vartheta_0) | \mathbf{X}, \mathbf{Y} \right) \right) \\ &= \mathbb{E}_{\vartheta_0} \left( s^{\mathbf{X}}(\mathbf{Y}; \mathbf{X}; \vartheta_0) \right) = \mathbf{0}.\end{aligned}$$

- Sometimes indirect proceeding: **corrected log-likelihood**  $l^{X^*}(\mathbf{Y}, \mathbf{X}, \vartheta)$  with

$$\mathbb{E}(l^{X^*}(\mathbf{Y}, \mathbf{X}^*, \vartheta) | \mathbf{X}, \mathbf{Y}) = l^{\mathbf{X}}(\mathbf{Y}, \mathbf{X}, \vartheta).$$

Under regularity conditions corrected score function by taking the derivative.

- + Functional method: no (unjustified !?) assumptions on the distribution of  $X$
- + Successful for generalized linear models, polynomial regression, etc. (Survey: Schneeweiß & Augustin, 2006, ASTA, Hübler & Frohn (eds.))
- + Extensions to misclassification (Akazawa, Kinukawa, Nakamura, 1998, J. Jap. Stat. Soc.; Zucker, Spiegelman, 2008, Stat. Med.)
- + Quite general error distribution can often be handled (only existing moment generating function needed)
  - Numerical difficulties for small samples
  - Handling of transformations (e.g.  $\ln X$ ) complicated or impossible
  - Non-existence of corrected score functions for some models





Tabelle 1:  $\beta = 1, \nu = 1.2, \lambda = 1, \text{size} = 1000$ ,  
 berechnet aus 100 Schätzern, Varianz des wahren  
 Prädiktors: 1.0, relative Median Fehler (in Pro-  
 zent)

Fehler- varianz	Zensierung (in % )	rel. Median Fehler für $\beta$	
		korrigiert	naiv
0.0	10	0.00	0.00
0.0	40	0.00	0.00
0.0	70	0.00	0.00
0.1	10	0.36	-13.20
0.1	40	1.59	-11.53
0.1	70	1.36	-10.38
0.3	10	0.43	-31.88
0.3	40	1.92	-29.61
0.3	70	1.46	-27.06
0.5	10	-0.76	-44.61
0.5	40	4.75	-41.41
0.5	70	0.96	-38.57
0.7	10	-0.50	-52.42
0.7	40	1.68	-49.89
0.7	70	2.69	-46.89

# Example: An Exact Correction for the Breslow Likelihood

(Augustin (2004, Scand. J.Stat.))

- Nakamura (1992, Biometrics): method not applicable to partial likelihood, seemingly approximate estimator
- The Breslow loglikelihood (Breslow (1972, JRSS SerB; 1974, Biometrics)), based on  $\lambda_0(t) := \lambda_j, \tau_{j-1} \leq t \leq \tau_j$ ,

$$\ln(\mathcal{L}(\lambda_1, \dots, \lambda_k, \beta)) = \sum_{j=1}^k \left( \ln \lambda_j + \beta' \mathbf{X}_{(j)} - \lambda_j (\tau_j - \tau_{j-1}) \sum_{i \in \mathcal{R}(\tau_j)} \exp(\beta' \mathbf{X}_i) \right)$$

does not have singularities and

$$\sum_{j=1}^k \left( \ln \lambda_j + \beta' \mathbf{X}_{(j)}^* - \lambda_j (\tau_j - \tau_{j-1}) \sum_{i \in \mathcal{R}(\tau_j)} \frac{\exp(\beta' \mathbf{X}_i^*)}{M_{U_i}(\beta)} \right)$$

is a corrected log-likelihood.

Proof / construction principle:

- Start with the naive log-likelihood
- It is sufficient to find constants  $c_1, c_2$  such that

$$\mathbb{E} \left( \sum_{j=1}^k \left( \ln \lambda_j + c_1 \cdot \beta' \mathbf{X}_{(j)}^* - \lambda_j (\tau_j - \tau_{j-1}) \sum_{i \in \mathcal{R}(\tau_j)} \frac{\exp(\beta' \mathbf{X}_i^*)}{c_2} \right) \mid \mathbf{X}, \mathbf{Y} = \right)$$

$$\sum_{j=1}^k \left( \ln \lambda_j + \beta' \mathbf{X}_{(j)} - \lambda_j (\tau_j - \tau_{j-1}) \sum_{i \in \mathcal{R}(\tau_j)} \exp(\beta' \mathbf{X}_i) \right)$$

- For that purpose note that:  $\mathbb{E}(\beta' \mathbf{X}_i^* \mid \mathbf{X}, \mathbf{Y}) = \beta' \mathbf{X}_i$ ;  
 $\mathbb{E}(\exp(\beta' \mathbf{X}_i^*) \mid \mathbf{X}, \mathbf{Y}) = \exp(\beta' \mathbf{X}_i) \cdot M_{U_i}(\beta)$

# 2.3. Extended Corrected Score Functions – A Unified View at Measurement Error and Censoring

# Extended Corrected Score Functions (and Parametric Survival Models)

- Parametric survival models: exponential, log-normal, Weibull, Gamma, log-logistic
- Superstructure: accelerated failure time models

$$\ln T_i = \beta_0 + \beta_{\mathbf{X}'} \cdot \mathbf{X}_i + \psi \cdot \epsilon, \quad \psi > 0,$$

- $\mathbb{E}(T_i) = \exp(\beta_0) \cdot \exp(\beta_{\mathbf{X}'} \cdot \mathbf{X}_i) \cdot \mathbb{E}(\exp(\psi \cdot \epsilon))$
- Use quasi-likelihood approach with quasi-score function

$$s^{\mathbf{X}}(\vec{T}; \mathbf{X}; \beta) = \sum_{i=1}^n \frac{\partial \mathbb{E}[T_i | \mathbf{X}_i; \beta]}{\partial \beta} \cdot \frac{\{T_i - \mathbb{E}[T_i | \mathbf{X}_i; \beta]\}}{\mathbb{V}[T_i | \mathbf{X}_i; \beta, \zeta]} = 0$$

- Use the accelerated failure time model as a **superstructure!**

↓ **General** form of the **'ideal'** quasi-score equation based on

$$\mathbb{E}[T_i^r | X_i; \beta, \psi, \gamma] = \underbrace{c_r(\psi, \gamma)}_{\text{distribution}} \cdot \underbrace{\exp(r \cdot (\beta_0 + \beta'_X \cdot X_i))}_{\text{covariates}}$$

↓ **General** form of the **corrected** quasi-score equation

↓ Measurement error correction in a concrete model through appropriate **choice of  $c_r(\cdot)$**

↓ Adopt for deficient dependent variables (**measurement error/censoring**)

# Handling Deficient Dependent Variables

- **Censoring:** Instead of  $T_i$  one observes

$$T_i^* = \min(T_i, C_i) \quad \text{and} \quad \Delta_i := I\{T_i \leq C_i\}$$

- **Measurement error:** (with  $V_i \perp \text{rest}$ ,  $\mathbb{E}(V_i) = m$ ,  $\mathbb{V}(V_i) = v$ )

- \* additive ( $T_i^* = T_i + V_i$ )

- \* multiplicative ( $T_i^* = T_i \cdot V_i$ )

response effects, memory effects (e.g., Skinner & Humphreys (1999, Lifetime Data Analysis))

Unifying approach:  $\Psi^{\mathbf{T}^*, \mathbf{X}^*}(\mathbf{T}^*; \mathbf{X}^*, \vartheta)$  extended corrected score function:

$$\mathbb{E}(\Psi^{\mathbf{T}^*, \mathbf{X}^*}(\mathbf{T}^*; \mathbf{X}^*, \vartheta) | \mathbf{T}, \mathbf{X}) = \Psi^{\mathbf{T}; \mathbf{X}}(\mathbf{T}, \mathbf{X}, \vartheta),$$

where  $\Psi^{\mathbf{T}, \mathbf{X}}(\mathbf{T}; \mathbf{X}; \vartheta)$  is an ideal unbiased estimating function.

It leads, again by the law of iterated expectation, to unbiased estimating function.



**Theorem:** Let  $y \mathbf{W} \in \{X_i^*, X_i\}$ , and  $\Psi^{\mathbf{T}}(\mathbf{T}_i, \mathbf{W}, \vartheta)$  be an unbiased estimating function of the very general form (including ML estimation in exponential families or Weibull type models, corrected QL estimation)

$$\psi^{\mathbf{T}}(\mathbf{T}_i, \mathbf{W}, \vartheta) = \sum_{i=1}^n \sum_{l=0}^q c_l(W_i, \vartheta) \cdot T_i^{\alpha_l}.$$

Then, under multiplicative measurement error,

$$\psi^{\mathbf{T}^*}(\mathbf{T}^*, \mathbf{W}, \vartheta) = \sum_{i=1}^n \sum_{l=0}^q c_l(W_i, \vartheta) \cdot \frac{\mathbf{T}_i^{*\alpha_l}}{\mathbb{E}(V_i^{\alpha_l})}$$

is an extended corrected score function. Similar results hold for the additive model, if  $\alpha_l \in \mathbb{N}$ .

**Corollary:**  $q = 1, \alpha = 1, \mathbb{E}(V_i) = 1 \implies \hat{\vartheta}_{ML,naive} \longrightarrow \vartheta$  for exponential families under unbiased measurement error in the *dependent* variable.

**Censoring** Instead of  $T_i$  one observes

$$T_i^* := \min(T_i, C_i) \quad \text{and} \quad \Delta_i := I(\{T_i \leq C_i\})$$

with  $C_i$  as the censoring variable.

Solve censoring by inverse probability weighting (e.g. Zhou (1992, Biometrika), Graf et al. (1999, Stat. Med.), Augustin (2002, Habil.), van der Laan & Robins, (2003, Springer), Hothorn et al (2006, Biostatistics), Gerds & Schuhmacher (2006, Biom. J.))

Consider independent random censoring,  $W \in \{X, X^*\}$ ,  $C_1, \dots, C_n$  i.i.d. with  $G(t) := P(C_i \geq t | W_i) = P(C_i \geq t) > 0$ ,  $\forall t \in \mathbb{R}_+$ . Then for "every"  $g(\cdot)$

$$\mathbb{E}(\Delta \cdot \frac{g(T_i^*, W, \vartheta)}{G(T_i^*)} | W_i, T_i) = g(T_i, W_i, \vartheta)$$

**Theorem:** Consider independent random censoring,  $W \in \{X, X^*\}$ ,  $C_1, \dots, C_n$  i.i.d. with  $G(t) := P(C_i \geq t | W_i) = P(C_i \geq t) > 0$ ,  $\forall t \in \mathbb{R}_+$ , and an unbiased estimation function for  $\vartheta$  of the form

$$\psi^*(\mathbf{T}, \mathbf{W}, \vartheta) = \sum_{i=1}^n a_0(W_i, \vartheta) + \sum_{l=1}^q a_l(W_i, \vartheta) \cdot T_i^{\alpha_l} \quad [C]$$

Then

$$\psi^{**}(\mathbf{T}^*, \mathbf{W}, \vartheta) = \sum_{i=1}^n \frac{\Delta_i}{G(T_i^*)} \cdot \psi_i^*(T_i^*, W_i, \vartheta)$$

and

$$\psi^{***}(\mathbf{T}^*, \mathbf{W}, \vartheta) = \sum_{i=1}^n \left( a_0(W_i, \vartheta) + \sum_{l=1}^q a_l(W_i, \vartheta) \cdot \frac{\Delta_i}{G(T_i^*)} \cdot T_i^{*\alpha_l} \right)$$

are extended corrected score functions.

**Proof:**  $\Delta_i = 0 \iff T_i^* \neq T_i$ . Therefore:

$$\frac{\Delta_i}{G(T_i^*)} \cdot g(T_i^*, W_i, \vartheta) = \frac{\Delta_i}{G(T_i)} \cdot g(T_i, W_i, \vartheta)$$

and

$$\begin{aligned} & \mathbb{E} \left( \frac{\Delta_i}{G(T_i^*)} \cdot g(T_i^*, W_i, \vartheta) \mid W_i, T_i \right) \\ &= \mathbb{E} \left( \frac{\Delta_i}{G(T_i)} \cdot g(T_i, W_i, \vartheta) \mid W_i, T_i \right) \\ &= 1 \cdot \frac{P(\Delta_i = 1 \mid W_i, T_i)}{G(T_i)} \cdot g(T_i, W_i, \vartheta) + 0 \\ &= g(T_i, W_i, \vartheta), \end{aligned}$$

due to  $P(\Delta_i = 1 \mid W_i, T_i) = P(C_i \geq T_i \mid W_i, T_i) = G(T_i)$ .

- 

$$P(\Delta_i = 1 | W_i, T_i) = P(C_i \geq T_i | W_i, T_i) = G(T_i)$$

Is it possible to allow dependence on  $W_i$  indeed?

Note,  $W$  is observable.

- Estimate  $G(\cdot)$  in a nonparametric way: Kaplan-Meier-estimator
- In the case of  $\psi^{**}(\cdot)$  condition [C] can be replaced by  $\psi^*(\mathbf{T}, \mathbf{W}, \vartheta) = \sum_{i=1}^n \psi_i^*(T_i, W_i, \vartheta)$ .
- Utilize the generality of the censoring lemma!

## Direct extensions:

### Utilize the generality of the censoring lemma

- Use soft weighting by truncating large pseudo-observation (Hothorn et al. (2006, Biostat.))
- Consistency under nonparametric estimation of  $G(\cdot)$  (with Kukush, Usoltseva ):
- Extend sophisticated methods to handle measurement error in the linear model (Cheng & van Ness (1999, Arnold); Wansbeek & Meijer (2000, Elsevier)) and the polynomial model (Cheng & Schneeweiss (1998, JRSS, Ser.B, 2002, TLS)), (Wansbeek & Meijer (2000, Elsevier, Chap. 11)) to the nonparametric AFT

$$\ln(T_i) = \beta' X_i + \epsilon_i \quad \text{with } \epsilon_i \text{ unspecified.}$$

- M-estimators under censorship
- handling of censored independent variables!?

- but be careful with standard application of weighting procedures: compare Basu's elephant in nonparametric statistics (Einbeck & Augustin, (2009, *Statistica Sinica*))

## 2.4. Corrected Score Functions for Berkson Models



- Up to now, the method of corrected score functions has not been applied to the Berkson model.

- Recall: **Corrected score function:**  $s^{X^*}(\mathbf{Y}, \mathbf{X}^*, \vartheta)$  such that

$$\mathbb{E}(s^{X^*}(\mathbf{Y}, \mathbf{X}^*, \vartheta) | \mathbf{X}, \tilde{\mathbf{Y}}) = \mathbf{s}^{\mathbf{X}}(\mathbf{Y}, \mathbf{X}, \vartheta)$$

- Only in the classical error model

$$[X_i^* | X_i; Y_i]$$

is directly available

- Idea: In the Berkson model

$$[X_i^* | X_i, Y_i]$$

can be obtained under additional assumptions, e.g., on the marginal distribution of  $X_i^*$ .

Let, for the moment, for  $i = 1, \dots, n$ ,

- $X_i^* \sim N(\mu^*, \sigma^{*2})$ ; w.l.o.g.  $\mu^* = 0$
- $\delta_i \sim N(0, \sigma_{\delta_i}^2)$

Berkson model  $X_i = X_i^* + \delta_i$ ;  $X_i^* \perp \delta_i$

Then

$$X_i^* | X_i, Y_i \sim N(\tilde{\mu}; \tilde{\sigma}^2)$$

with

$$\tilde{\mu} = X_i \left( 1 - \frac{\sigma_{\delta_i}^2}{\sigma_{X_i}^2} \right)$$

and

$$\tilde{\sigma}^2 = \frac{\sigma_{\delta_i}^2 \sigma_{X_i^*}^2}{\sigma_{X_i}^2}$$

# Theorem: Corrected loglikelihood under Berkson error:

- Poisson regression:

$$\ell(\mathbf{X}, Y, \beta) = \sum_{i=1}^n \beta X_i Y_i - \exp(\beta X_i)$$

$$\ell_{naive}(\mathbf{X}^*, Y, \beta) = \sum_{i=1}^n \beta X_i^* Y_i - \exp(\beta X_i^*)$$

$$\ell_{corr}(\mathbf{X}^*, Y^*, \beta) = \sum_{i=1}^n \frac{\sigma_X^2}{\sigma^{*2}} \cdot \beta X_i^* Y_i - \exp\left(\frac{\sigma_X^2}{\sigma^{*2}} \beta X_i^* - \frac{\sigma_{\delta_i}^2 \sigma_X^2}{2\sigma^{*2}} \beta^2\right)$$

- Cox model: Breslow log-likelihood, based on failure times  $\tau_j$ ,  $j = 1, \dots, k$ ,
  - \*  $\mathcal{D}(\tau_j)$  deaths at  $\tau_j$ ;  $d_j := |\mathcal{D}(\tau_j)|$ ,  $\mathcal{R}(\tau_j)$  risk set,
  - \*  $\lambda_j$  piecewise constant part of baseline hazard rate  $\lambda_0(t)$

$$\ell(\mathbf{X}, Y, \beta) = \sum_{i=1}^n d_j \ln \lambda_j + \sum_{i \in \mathcal{D}(\tau_j)} \beta \mathbf{X}_i - \lambda_j(\tau_{j+1} - \tau_j) \sum_{i \in \mathcal{R}(\tau_j)} \exp(\beta \mathbf{X}_i)$$

$$\ell_{naive}(\mathbf{X}^*, Y, \beta) = \sum_{i=1}^n d_j \ln \lambda_j + \sum_{i \in \mathcal{D}(\tau_j)} \beta \mathbf{X}_i^* - \lambda_j(\tau_{j+1} - \tau_j) \sum_{i \in \mathcal{R}(\tau_j)} \exp(\beta \mathbf{X}_i^*)$$

$$\ell_{corr}(\mathbf{X}^*, Y, \beta) = \sum_{j=1}^k \left( d_j \ln \lambda_j + \sum_{i \in \mathcal{D}(\tau_j)} \frac{\sigma_X^2}{\sigma^{*2}} \cdot \beta \mathbf{X}_i^* - \right. \\ \left. - \lambda_j(\tau_{j+1} - \tau_j) \sum_{i \in \mathcal{R}(\tau_j)} \exp \left( \frac{\sigma_X^2}{\sigma^{*2}} \beta \mathbf{X}_i^* - \frac{\sigma_{\delta_i}^2 \sigma_X^2}{2\sigma^{*2}} \beta^2 \right) \right),$$

# Construction

- Both models have a similar structure
  - \* linear term:  $\beta X_i$  and  $(\beta X_i \cdot Y_i)$ , resp.
  - \* exponential term:  $\exp(\beta X_i)$
- For finding a corrected log likelihood

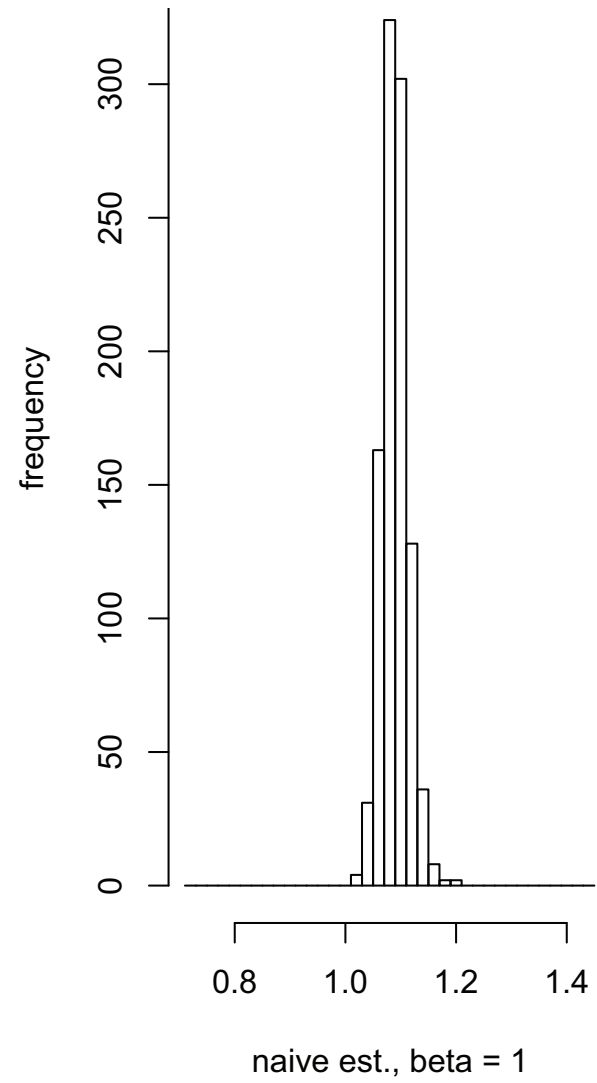
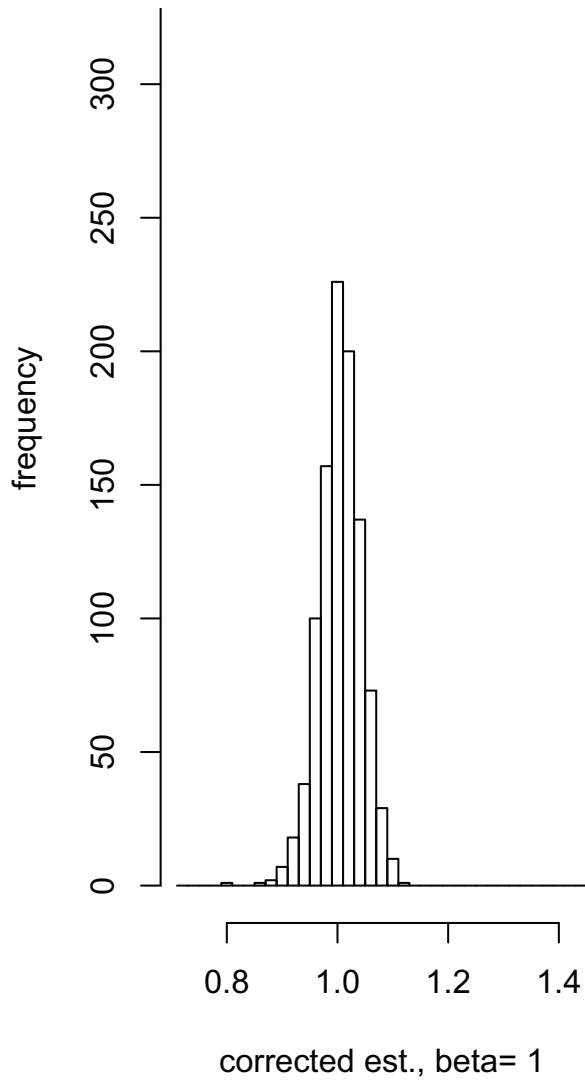
$$\mathbb{E}(l^{X^*}(\mathbf{Y}, \mathbf{X}^*, \vartheta) | \mathbf{X}, \mathbf{Y}) = l^X(\mathbf{Y}, \mathbf{X}, \vartheta).$$

it is sufficient that there exist  $c_1, c_2, c_3, c_4$  such that

$$\mathbb{E}\left( (c_1 \beta X_i^* + c_2) | X_i, Y_i \right) \stackrel{!}{=} \beta X_i$$

and

$$\mathbb{E}\left( (\exp(c_3 \beta X_i^* + c_4)) | X_i, Y_i \right) \stackrel{!}{=} \exp(\beta X_i).$$



- UNBIASEDNESS OF NAIVE SCORE TESTS ???