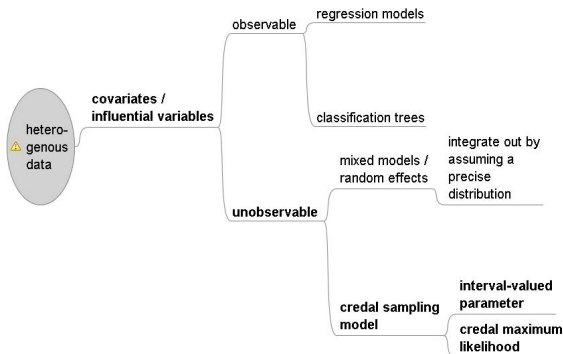


Credal maximum likelihood – an imprecise probability alternative to mixed models?

Thomas Augustin
LMU Munich

What can IP contribute to the reliable handling of unobserved data heterogeneity?



1) Traditional Maximum Likelihood

- ▶ THE frequentist estimation method
 - consistency
 - asymptotic normality
 - asymptotic efficiency
 - universally applicable
 - gives immediately confidence regions and tests
- ▶ Observation i , $i = 1, \dots, n$
- ▶ $Y_1, \dots, Y_n := \mathbf{Y}$ outcome
- ▶ $X_1, \dots, X_n := \mathbf{X}$ covariates
- ▶ $Y_i|X_i \sim p_{\vartheta, X_i}$ with density f_{ϑ, X_i}
- ▶ Estimate ϑ from observations of Y_1, \dots, Y_n
- ▶ After having observed y_1, \dots, y_n , the higher

$$\prod_{i=1}^n P_{\vartheta}(Y_i = y_i|X_i) \quad \text{or} \quad \prod_{i=1}^n f_{\vartheta, X_i}, \quad (*)$$

the more plausible the conclusion that ϑ is the true parameter.

- ▶ So estimate ϑ by maximizing (*) with respect to ϑ

→ *maximum likelihood estimator*

Examples

1. Y_1, \dots, Y_n normally distributed with unknown mean μ and given variance σ^2 :

$$\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma}} \cdot \exp\left(-\frac{1}{2\sigma^2}(y_i - \mu)^2\right) \rightarrow \max_{\mu}$$

$$\iff \sum_{i=1}^n (y_i - \mu)^2 \rightarrow \min_{\mu}$$

Least square problem! (Solution $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n y_i$)

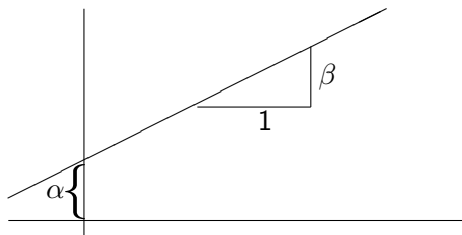
2. $Y_1, \dots, Y_n \sim \text{Poisson}(\lambda)$

$$\prod_{i=1}^n \frac{\lambda^{y_i}}{y_i!} \exp(-\lambda) \rightarrow \max_{\lambda}$$

$$\iff \sum_{i=1}^n (y_i \ln \lambda - \lambda) \rightarrow \max_{\lambda} \Rightarrow \hat{\lambda} = \frac{1}{n} \sum_{i=1}^n y_i$$

3. Linear regression

$$Y_i = \alpha + \beta \cdot x_i + \varepsilon_i$$



$$y_i | x_i \sim N(x_i' \beta, \sigma^2)$$

Again maximum likelihood principle and least squares principle coincide

$$\hat{\alpha}, \hat{\beta} \quad \text{by} \quad \sum_{i=1}^n (y_i - \alpha - x_i' \beta)^2 \rightarrow \min_{\alpha, \beta}$$

2) Credal (Parametric) Sample Models

- ▶ Let $\Theta \subseteq \mathbb{R}$, parametric family of classical distributions $(p_{\vartheta, X_i})_{\vartheta \in \Theta}$.
- ▶ Credal parametric sampling model (imprecise model, not just imprecise data!).

Parameter interval-valued

$$[\underline{\vartheta}, \overline{\vartheta}]$$

Credal set

$$\mathcal{M}_{X_i} = \{P_{\vartheta, X_i} \mid \vartheta \in [\underline{\vartheta}, \overline{\vartheta}]\}$$

Strongly independent observations

$$\prod_{i=1}^n \mathcal{M}_{X_i} = \left\{ \prod_{i=1}^n P_{\vartheta_i, X_i} \mid \vartheta_i \in [\underline{\vartheta}, \overline{\vartheta}] \right\}$$

What is it good for?

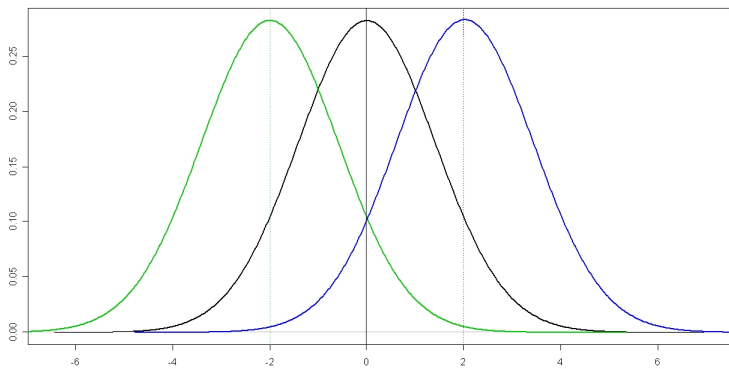
Heterogeneity interpretation:

Overall parameter + individual parameter:

$$\vartheta_i = \vartheta_{overall} + \nu_i$$

			unobserved
biometrics	overall effect	treatment	hospital-, patient-specific
insurance	overall risk		individual risk attitude
dynamical econometrical model	overall chance		individual characteristics

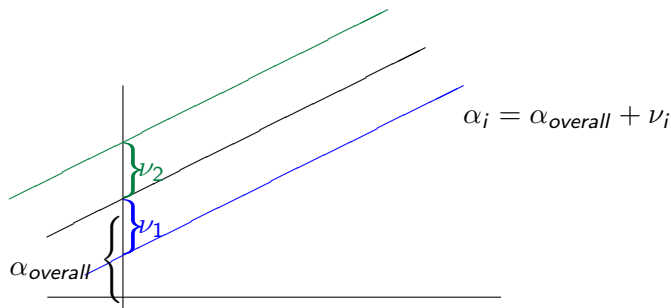
Credal (Parametric) Sample Models



In linear regression analysis “set of “true” regression lines”

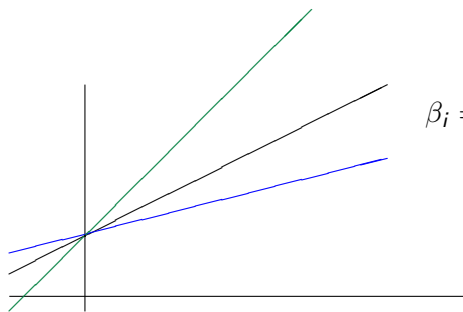
simple linear regression: x_i one-dimensional

$$i) \quad y_i = \alpha_i + \beta x_i + \varepsilon$$



x_i dummy variables: Analysis of variance with random effects

$$\text{ii) } y_i = \alpha + \beta_i x_i + \varepsilon$$



$$\beta_i = \beta_{\text{overall}} + \nu_i$$

3) Traditional solution: random effects model

Assume certain distribution, described by $\tilde{f}(\cdot)$, for ν_i , typically $\nu_i \sim N(0, \sigma_\nu^2)$

Consider likelihood

$$\prod_{i=1}^n f_{\vartheta_{\text{overall}}}(x_i) = \prod_{i=1}^n \int f_{\vartheta_{\text{overall}}}(x_i | \nu_i) \cdot \tilde{f}(\nu_i) d\nu_i$$

to estimate $\vartheta_{\text{overall}}$

- point estimator, irrespectively of amount of heterogeneity
- depends, of course, strongly on $\tilde{f}(\cdot)$

4) Level δ – Credal Maximum Likelihood Estimation

Definition:

Let $\delta \geq 0$ be fixed and let, for given data y_1, y_2, \dots, y_n ,

$$\hat{\vartheta}_1, \dots, \hat{\vartheta}_n, \quad \widehat{L\vartheta}, \widehat{U\vartheta},$$

be an optimal solution of

$$\prod_{i=1}^n f_{\vartheta_i}(y_i) \rightarrow \max_{\vartheta_1, \dots, \vartheta_n, L\vartheta, U\vartheta}$$

subject to

$$\begin{aligned} L\vartheta \leq \vartheta_i &\leq U\vartheta, \quad i = 1, \dots, n \\ U\vartheta - L\vartheta &\leq \delta, \end{aligned}$$

then

$$\left[\widehat{L\vartheta}, \widehat{U\vartheta} \right]$$

is called level- δ credal maximum likelihood estimator.

Remarks:

i) Obviously

$$\delta = 0 \Rightarrow \widehat{L}\vartheta = \widehat{U}\vartheta = \widehat{\vartheta}_{ML}$$

(the traditional ML estimator)

ii) Of course, it is much more convenient to replace the objective function by the equivalent objective function

$$\sum_{i=1}^n \ln f_{\vartheta_i}(y_i) \rightarrow \max$$

5) Examples: Least Squares Problems

Example I: normal model: Normal distribution (ML and Least Squares coincide), parameter μ_j .

We have to consider the quadratic optimization problem

$$\sum_{i=1}^n (y_i - \mu_i)^2 \rightarrow \min$$

subject to

$$L\mu \leq \mu_i \leq U\mu \quad \text{and} \quad U\mu - L\mu \leq \delta,$$

which can be solved by standard software.

- i) At least in the case of normal distribution with unknown location parameter

$$\delta \rightarrow \infty : \widehat{L\vartheta} = \min_{i=1, \dots, n} y_i$$

- ii) The problem can be viewed as a function of the lower interval limit T of the estimator

$$\mathcal{E}(y_i, T) = (y_i - T)^2 \cdot I\{y_i \leq T\} + (y_i - (T + \delta))^2 \cdot I\{y_i \geq T + \delta\}$$

Some numerical toy examples ($n = 4$, MAPLE)i) $y_1 = 1; y_2 = 2; y_3 = 3; y_4 = 4$

δ	$[\widehat{L\mu}, \widehat{U\mu}]$
0:	2.5 ...
0.1:	[2.45; 2.55]
0.5:	[2.25; 2.75]
1:	[2; 3]

ii) Note $[\widehat{L\mu}, \widehat{U\mu}]$ is not just $\hat{\mu} \pm$ something $y_1 = 1; y_2 = 2; y_3 = 3; y_4 = 14$

δ	$[\widehat{L\mu}, \widehat{U\mu}]$
0:	5
0.1:	[4.975; 5.075]
0.5:	[4.875; 5.375]
1:	[4.75; 5.75]

Example II: simple linear regression In the regression context we have to consider

$$\sum_{i=1}^n (y_i - \alpha_i - \beta x_i) \rightarrow \min$$

or

$$\sum_{i=1}^n (y_i - \alpha - \beta_i x_i) \rightarrow \min$$

subject to the restrictions

$$\alpha_i \in [\widehat{L}_\alpha, \widehat{U}_\alpha], \quad \widehat{U}_\alpha - \widehat{L}_\alpha \leq \delta$$

and

$$\beta_i \in [\widehat{L}_\beta, \widehat{U}_\beta], \quad \widehat{U}_\beta - \widehat{L}_\beta \leq \delta$$

respectively

6) Outlook

Conjecture: Objective function convex then

$$\delta_1 \leq \delta_2 \Rightarrow [\widehat{L}_{\delta_1} \vartheta, \widehat{U}_{\delta_1} \vartheta] \subseteq [\widehat{L}_{\delta_2} \vartheta, \widehat{U}_{\delta_2} \vartheta]$$

→ Note special case $\delta = 0$ (tradit. ML)

Then:

Under i.i.d ($\vartheta_i \equiv \vartheta$)

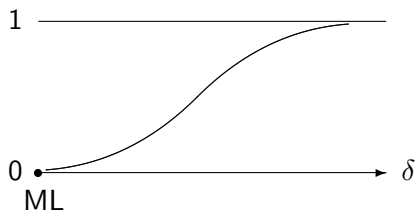
$$\lim_{n \rightarrow \infty} P_{\vartheta} \left([\widehat{L}_{\delta} \vartheta^{(n)}, \widehat{U}_{\delta} \vartheta^{(n)}] \ni \vartheta \right) = 1$$

“i.i.d consistency of level δ -ML estimation”

(Proof: traditional consistency of ML; conjecture above)

On the Choice of δ

- a) Look at the objective function as a function of δ :
Use that δ where a further increase does not improve the objective function substantially (cp. elbow criterion in principal component analysis)
- b) fuzzy set interpretation; estimator as a fuzzy set, membership function increasing in δ



- c) Penalization (like in nonparametric statistic)
look at the objective function

$$\sum_{i=1}^n \ln f_{\vartheta}(y_i) - \lambda \cdot \delta \rightarrow \max$$

λ for instance by cross-validation

Additional aspects

- ▶ What can we learn when $(P_{\vartheta})_{\vartheta \in \Theta}$ is stochastically ordered?
- ▶ Comparison to traditional random effect models
- ▶ Method can be extended to robust objective functions \rightarrow credal M-estimators