



# Estimating Equations for data summaries

Tom Whitaker, Boris Beranger, Scott A. Sisson,

**UNSW & ACEMS** 

New Perspectives in Data Science IHP, 6th March, 2020



SDA Seminal contributor Citations: 9282 H-index: 41

1987: Edwin's first SDA paper and ... I was born!

- 2016: I started research in SDA
- 2017: Edwin and I met at the SDA workshop in Ljubljana



SDA Seminal contributor Citations: 9282 H-index: 41

1987: Edwin's first SDA paper and ... I was born!

2016: I started research in SDA

2017: Edwin and I met at the SDA workshop in Ljubljana



SDA Seminal contributor Citations: 9282 H-index: 41

1987: Edwin's first SDA paper and ... I was born!

2016: I started research in SDA

2017: Edwin and I met at the SDA workshop in Ljubljana



SDA Seminal contributor Citations: 9282 H-index: 41

- 1987: Edwin's first SDA paper and ... I was born!
- 2016: I started research in SDA
- 2017: Edwin and I met at the SDA workshop in Ljubljana



SDA Seminal contributor Citations: 9282 H-index: 41

- 1987: Edwin's first SDA paper and ... I was born!
- 2016: I started research in SDA
- 2017: Edwin and I met at the SDA workshop in Ljubljana



SDA Seminal contributor Citations: 9282 H-index: 41

- 1987: Edwin's first SDA paper and ... I was born!
- 2016: I started research in SDA
- 2017: Edwin and I met at the SDA workshop in Ljubljana
- 2020: HAPPY BIRTHDAY EDWIN!

# Talk Outline



- Background Information
- Estimating Equations for symbolic data
- . Examples
  - Discussion

# Background Information (1)

▶ Let  $X = (X_{[1]}, ..., X_{[D]}) \in \mathcal{D}_X \subset \mathbb{R}^D$  with d.f.  $F_X$  (unknown).

Let X = (X<sub>1</sub>,...,X<sub>N</sub>) be the collection of N i.i.d. replicates of X with realisation given by x = (x<sub>1</sub>,...,x<sub>N</sub>)

Without any parametric assumption about  $F_X$ :

Make statistical inference on  $\theta \in \mathcal{D}_{\theta} \subset \mathbb{R}^{M}$  using  $R \geq M$  functionally independent estimating equations (EE):

 $g(X,\theta) = (g_1(X,\theta),\ldots,g_R(X,\theta))^{\top},$ 

with the condition

 $\mathbb{E}_{F_X}[g_r(X,\hat{\theta})] = 0$ , for all  $r = 1, \ldots, R$ .

# Background Information (2)

• No assumption on  $F_X \Rightarrow$  Empirical alternative

The **empirical likelihood (EL)** associated with the observed sample **x** is

$$L(F_X; \mathbf{x}) = \prod_{n=1}^N \mathrm{d} F_X(x_n) = \prod_{n=1}^N P(X = x_n).$$

 $F_X$  can be seen as a discrete distribution on  $\{x_1, \ldots, x_N\}$  with probability vector  $p = (p_1, \ldots, p_N)$  defined such that:

(C1): 
$$p_n = P(X = x_n) > 0, n = 1, ..., N.$$
  
(C2):  $\sum_{n=1}^{N} p_n = 1$ 

If there are no other conditions on x other than (C1) and (C2) then

$$\hat{p}_n = \operatorname*{arg\,max}_{p_n \mid C1, C2} \prod_{n=1}^N p_n = \frac{1}{N},$$

and the estimating equation becomes:

$$\frac{1}{N}\sum_{n=1}^{N}g_r(x_n,\hat{\theta})=0, \text{ for all } r=1,\ldots,R.$$

# Background Information (3)

#### Example:

Assume we are interested in estimating the mean and variance, i.e. for M = 2:

$$\theta = (\mu, \sigma^2) = (\mathbb{E}(X), \mathbb{V}(X)).$$

A natural choice of estimating function with R = 2 is then

$$g(X,\theta) = \left(X - \mu, (X - \mu)^2 - \sigma^2\right),$$

for which the estimating equations give

$$\begin{cases} \frac{1}{N} \sum_{n=1}^{N} (x_n - \hat{\mu}) = 0 \\ \frac{1}{N} \sum_{n=1}^{N} ((x_n - \hat{\mu})^2 - \hat{\sigma}^2) = 0 \end{cases}$$

$$\Leftrightarrow \left\{ \begin{array}{l} \hat{\mu} = \frac{1}{N} \sum_{n=1}^{N} x_n \\ \\ \hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^{N} (x_n - \hat{\mu})^2 \end{array} \right.$$

#### Symbolic Data

For a <u>class</u> c = 1, ..., C, let  $\mathbf{X}^{(c)} = \{X_1^{(c)}, ..., X_{n_c}^{(c)}\}$  denote the *c*-th subset of  $\mathbf{X}$  of size  $n_c$  s.t.

$$\begin{cases} \bigcup_{c=1}^{C} \boldsymbol{X}^{(c)} = \boldsymbol{X} \\ \sum_{c=1}^{C} n_{c} = N \end{cases}$$

The *Symbolic object*  $S_c \in \mathcal{D}_{S_c}$  is a summary of the information contained in  $\mathbf{X}^{(c)}$ , obtained through an aggregation function  $\pi(\cdot)$  which may contain some deterministic elements  $\vartheta$ 

For observations  $\mathbf{x}^{(c)}$ , the information is summarised in  $\mathbf{s}_c = \{\mathbf{n}_c, \Upsilon_c, \alpha_c\}$ , where  $\Upsilon_c = \mathcal{D}(\mathbf{X}^{(c)})$  and  $\alpha_c$  contains the summary statistics specific to  $\pi(\cdot)$ .

# Estimating Equations for Symbolic Data (1)

- Let  $\phi_c := f_{X|S_c=s_c}$  denote the density of X given a summary  $s_c$ .
- Let  $\phi := f_{X|S=s}$  denote the density of X given the set of all summaries s.
- These densities are linked through

$$\phi_c(x) = \frac{\mathbb{1}(x \in s_c)\phi(x)}{\mathbb{P}(s_c)},$$

where the indicator restricts to the *c*-th symbol and the denominator corresponds its probability of occurrence  $\mathbb{P}(s_c) = \int \phi(y) \mathbb{1}\{y \in s_c\} dy$ .

• **Example:** If classes/symbols are independent then we can take  $\overline{\mathbb{P}(s_c)} = \frac{n_c}{N}$  such that  $\sum_{c=1}^{C} \mathbb{P}(s_c) = 1$  which yields

$$\phi(x) = \sum_{c=1}^{C} \frac{n_c}{N} \phi_c(x).$$

# Estimating Equations for Symbolic Data (2)

We define the estimating equations for symbolic inputs as

$$\mathbb{E}_{F_S}[g'_r(S,\theta,\vartheta)] = 0, \text{ for all } r = 1,\ldots,R$$

where  $g'_r$  is a transformation of  $g_r$  due to the aggregation and  $F_S$  is the symbolic distribution function.

1. Need to define an empirical alternative for  $F_S$ 

2. Need to derive  $g'_r$  (from  $g_r$ )

## Estimating Equations for Symbolic Data (3)

The symbolic empirical likelihood associated with the observed symbol s is

$$L(F_S; \boldsymbol{s}) = \prod_{c=1}^{C} \mathbb{P}(S = \boldsymbol{s}_c) = \prod_{c=1}^{C} \prod_{n=1}^{n_c} \mathbb{P}(X = \boldsymbol{x}_n^{(c)}),$$

 $F_S$  can be seen as a discrete distribution on  $s_1, \ldots, s_C$  with probabilities  $p_1, \ldots, p_C$ , such that:

$$\begin{cases} (C3): p_c = \mathbb{P}(S_c = s_c) > 0, c = 1, \dots, C \\ (C4): \sum_{c=1}^{C} p_c = 1 \end{cases}$$

Each observation  $x_n^{(c)}$ ,  $n = 1, ..., n_c$ , c = 1, ..., C has probability  $q_n^{(c)} = \frac{p_c}{n_c}$ , such that  $\sum_{c=1}^{C} \sum_{n=1}^{n_c} q_n^{(c)} = 1$ . If there are no other conditions on s other than **(C3)** and **(C4)** then

$$\hat{q}_{n}^{(c)} = \operatorname*{arg\,max}_{q_{n}^{(c)}|C3,C4} \prod_{c=1}^{C} \prod_{n=1}^{n_{c}} q_{n}^{(c)} = \frac{1}{N}, \implies \hat{p}_{c} = \frac{n_{c}}{N}$$

and the estimating equation becomes:

$$\sum_{c=1}^{C} \frac{n_c}{N} g'_r(s_c, \hat{\theta}, \vartheta) = 0, \text{ for all } r = 1, \dots, R.$$

## Estimating Equations for Symbolic Data (4)

The symbolic estimating equations are obtained by integrating the classical estimating equation over all the data points x from which the symbols s are produced with their corresponding weights, i.e.

 $\mathbb{E}_{F_{S}}[g_{r}'(S,\theta,\vartheta)] = \int_{\mathcal{D}_{Y}^{N}} \mathbb{E}_{F_{X}}[g_{r}(X,\theta)]f_{S|X}(s|x,\vartheta)\phi(x)dx.$ 

$$\begin{split} \mathbb{E}_{F_{S}}[g_{r}'(S,\theta,\vartheta)] &= \int_{\mathcal{D}_{X}^{N}} \mathbb{1}\left\{\pi\left(\mathbf{x}^{(c)}\right) = \mathbf{s}_{c}; c = 1, \dots, C\right\} \\ &\times \left\{\frac{1}{N}\sum_{n=1}^{N}\sum_{c=1}^{C}\mathbb{1}\left\{x_{n} \in \mathbf{x}^{(c)}\right\}g_{r}(x_{n};\theta)\right\}\phi(\mathbf{x})\mathrm{d}\mathbf{x} \\ &= \frac{1}{N}\sum_{n=1}^{N}\sum_{c=1}^{C}\int_{\mathcal{D}_{X}^{N}}\mathbb{1}\left\{x_{n} \in \mathbf{x}^{(c)}\right\}\mathbb{1}\left\{\pi\left(\mathbf{x}^{(c)}\right) = \mathbf{s}_{c}\right\}g_{r}(x_{n};\theta)\phi(\mathbf{x})\mathrm{d}\mathbf{x} \\ &= \sum_{c=1}^{C}\frac{n_{c}}{N}\int_{\Upsilon_{c}}\phi_{c}(x)g_{r}(x;\theta)\mathrm{d}x, \end{split}$$

And we conclude

$$g_r'(s_c, \theta, \vartheta) = \int_{\Upsilon_c} \phi_c(x) g_r(x, \theta) \mathrm{d}x.$$

#### Estimating the within-symbol density $\phi_c$ (1)

Back to the Example:  $\theta = (\mu, \sigma^2)$ , we can rewrite

$$\mathbb{E}_{F_X}\left[g\left(\boldsymbol{X},\hat{\theta}\right)\right] = \left(\frac{1}{N}\sum_{n=1}^{N}(x_n-\hat{\mu}), \frac{1}{N}\sum_{n=1}^{N}\left((x_n-\hat{\mu})^2-\hat{\sigma}^2\right)\right)$$
$$= \left(\sum_{c=1}^{C}\frac{n_c}{N}(\mu_c-\hat{\mu}), \sum_{c=1}^{C}\frac{n_c}{N}\left((\mu_c-\hat{\mu})^2+\sigma_c^2-\hat{\sigma}^2\right)\right) = 0$$

 $\Rightarrow$  Need  $\hat{\mu}_c = \mathbb{E}_{\phi_c}(X)$  and  $\hat{\sigma}_b^2 = \mathbb{V}_{\phi_c}(X)$ .

$$\begin{array}{ll} \text{For skewness:} & \gamma_{c} = \mathbb{E}_{F_{X^{c}}}[((X^{(c)} - \mu_{b})/\sigma_{b})^{3}] \\ \text{For correlation:} & \rho_{cde} = \mathbb{E}_{F_{X^{c}}}[(X^{(c)}_{[d]} - \mu_{c[d]})(X^{(c)}_{[e]} - \mu_{c[e]})]\sigma^{-2}_{c[d]}\sigma^{-2}_{c[e]} \end{array}$$

#### We first need to estimate $\phi$

Common approach: (Bertrand and Goupil, 2000; Billard and Diday, 2003)

- Assume uniformity within classes:  $\phi_c$  is uniform density;
- **Example:**  $\alpha_c = (\alpha_{c,l}, \alpha_{c,u})$  the upper and lower bounds of an interval

• 
$$\hat{\mu}_c = \frac{\alpha_{c,l} + \alpha_{c,u}}{2}$$
,  $\hat{\sigma}_c^2 = \frac{(\alpha_{c,u} - \alpha_{c,l})^2}{12}$ 

## Estimating the within-symbol density $\phi_c$ (2)

Our approach: Borrowing information from adjacent symbols

- Attribution of a realisation to a class is random:
  Λ<sub>x</sub> ∈ {1,..., C} with d.f. H<sub>x</sub>, indicates the list of symbols in which an observation x could have been aggregated in.
- Redefine \(\phi\) (density of X given \(\boldsymbol{S}\)) as

$$\phi(x) = \int_{\mathcal{D}(\Lambda_x)} f_{X|\mathbf{S}=\mathbf{s}'}(x) \mathrm{d}H_x(\lambda) \mathrm{d}\lambda$$

where s' denotes the set of symbols given that the observation x is assumed to be grouped in the  $\lambda$ -th class.

If an observation x can only be associated to a unique symbol  $s_c$  then  $H_x$  is a Dirac delta function and  $\phi(x) = f_{X|S=s'}(x)$ .

• The density of an observation given a symbol  $s_c$  is given by

$$\phi_c(x) = \frac{\mathbb{1}\{x \in s_c\}\phi(x)}{P_H^{s_c}}, \quad P_H^{s_c} = \int_{\Upsilon_c} \phi(x) \mathrm{d}H_x(c) \mathrm{d}x.$$

#### $\phi_c$ for interval-valued data

- ▶ s<sub>c</sub>: D-dimensional intervals with  $\alpha_c = (\alpha_{c,l}, \alpha_{c,u})$  and  $\Upsilon_c = [\alpha_{c,l}, \alpha_{c,u}]$
- If x is in a region where some intervals from the set s overlap, then Λ<sub>x</sub> has a discrete outcome λ, a subset of {1,..., C}

$$H_x(\lambda, c) = rac{n_c}{\sum_{a \in \lambda} n_a}, ext{ for all } c \in \lambda.$$

We can derive that:

$$\phi(x) = \sum_{c=1}^{C} \frac{n_c}{N|\Upsilon_c|} \mathbb{1}\{x \in \Upsilon_c\},$$

Split the range of x into a grid of subintervals denoted by  $v_b$  with  $b = (b_1, \ldots, b_D)$  and  $b_d = 1, \ldots, (2C + 1); d = 1, \ldots D$ 

The normalising term  $P_H^{s_c} = m_c(1)$  where

$$\begin{split} m_{c}(f(x)) &= \int_{\Upsilon_{c}} f(x)\phi(x)\mathrm{d}H_{x}(c)\mathrm{d}x\\ &= \frac{1}{N}\sum_{c'=1}^{C}\frac{n_{c'}}{|\Upsilon_{c'}|}\left(\sum_{\boldsymbol{b}}\mathbbm{1}\{\upsilon_{\boldsymbol{b}}\subset\Upsilon_{c'}\}H_{\boldsymbol{b}}(\lambda,c)\int_{\upsilon_{\boldsymbol{b}}}f(y)\mathrm{d}y\right), \end{split}$$

## $\phi_{\textit{c}}$ for interval-valued data

The estimates of the mean, variance, skewness and correlation within an interval c are obtained using the density  $\phi_c$  and are given as follows

$$\hat{\mu}_{cd} = \frac{m_c \left( \mathbf{x}_{[d]} \right)}{m_c(1)}, \qquad \qquad \hat{\sigma}_{cd}^2 = \frac{m_c \left( \mathbf{x}_{[d]} \right)}{m_c(1)} - \hat{\mu}_{cd}^2,$$

$$\hat{\gamma}_{cd} = \frac{m_c \left( \mathbf{x}_{[d]} \right)}{m_c(1)\hat{\sigma}_{cd}^3} - \frac{\hat{\mu}_{cd}^3}{\hat{\sigma}_{cd}^3} - 3\frac{\hat{\mu}_{cd}}{\hat{\sigma}_{cd}}, \qquad \hat{\rho}_{cde} = \frac{m_c \left( \mathbf{x}_{[d]} \right)}{m_c(1)\hat{\sigma}_{cd}\hat{\sigma}_{ce}} - \frac{\hat{\mu}_{cd}\hat{\mu}_{ce}}{\hat{\sigma}_{cd}\hat{\sigma}_{ce}},$$
for  $d = 1, \dots, D.$ 

## $\phi_{c}$ for histogram-valued data

- **s**<sub>c</sub>: D-dimensional histogram bins with  $n_c = \alpha_c = \sum_{n=1}^N \mathbb{1}\{x_n \in \Upsilon_c\}$  and  $\Upsilon_c = ((y_{c_1-1}^1, y_{c_1}^1) \times \cdots \times (y_{c_D-1}^D, y_{c_D}^D)).$
- Assume the marginal bin width to be equal and given by  $\delta_d = y_{c_d}^d y_{c_d-1}^d \ \forall c_d, d \text{ s.t.}$  the area of each bin is  $|\Upsilon_c| = \prod_{d=1}^D \delta_d$ .
- Choice of the bin locations  $\Upsilon_c$  is arbitrary and other histograms could have arisen by shifting the bin locations by up to half of their width to the left or to the right.
- New bin locations  $\Upsilon' = \Upsilon + \boldsymbol{u}$  where  $\boldsymbol{u} = (u_1, \dots, u_D)$ ,  $u_d \sim \mathcal{U}\left(-\frac{\delta_d}{2}, \frac{\delta_d}{2}\right)$ .
- Λ<sub>x</sub>: the bin Υ in which x could have been aggregated to, is now a continuous as the set of outcomes contains all the shifted bins Υ' that will include x.

$$H_{x}(\Upsilon_{c}') = \frac{1}{\prod_{d=1}^{D} \delta_{d}} \mathbb{1}\{x \in \Upsilon_{c}'\}$$

 Assuming that histogram bins can be shifted relates to the ideas of Heitjan and Rubin (1991)

#### **Examples**

N = 2000 draws from the <u>Skew-Normal</u> distribution.C min/max intervals

- Classical estimates and 95% CI (Solid black)
- Borrowing/Not borrowing information (Solid grey/Dotted black)
- Billard and Didday (2003) histogram estimator of intervals by partitioning, borrowing (Dotted grey) and not borrowing information (Dashed grey)



#### **Examples**

N = 2000 draws from the <u>Skew-Normal</u> distribution.• 1 Histogram with C equal bins

- Classical estimates and 95% CI (Solid black)
- Borrowing/Not borrowing information (Solid grey/Dotted black)



## Discussion

- A new non-parametric approach to estimate statistics presented in the general framework of Estimating Equations
- Proposal to borrow information from neighbours to get a sense of the dependence:
  - Borrow lots of information = low dependence
  - Borrow no information = dependence
- Not always necessary: e.g. mushroom dataset analysed in de A Lima Neto et al. (2011)



# THANK YOU

#### Manuscripts:

- Whitaker, Beranger & Sisson (2020). Estimating Equations for data summaries. Available very soon!
- Beranger, Lin & Sisson (2018). New models for symbolic data. https://arxiv.org/pdf/1805.03316.pdf
- Whitaker, Beranger & Sisson (2019). Composite likelihood methods for histogram-valued random variables.

https://arxiv.org/pdf/1908.11548.pdf

#### Contact:

B.Beranger@unsw.edu.au