



# Advances in the analysis of aggregated data

Boris Beranger, Jaslene Lin Tom Whitaker, Scott A. Sisson,

**UNSW & ACEMS** 

Macquarie University, 22nd October, 2019

# Talk Outline



#### 1. What is Symbolic Data Analysis?

- Existing and new SDA models
- Statistical analyses using aggregates

Discussion

### Rise of non-standard data forms



Standard statistical methods analyse classical datasets

E.g.  $x_1, \ldots, x_n$  where  $x_i \in \mathcal{X} = \mathbb{R}^p$ 

However: Increasingly see non-standard data forms for analysis.

Simple non-standard forms:

- Can arise as result of measurement process
- Blood pressure naturally recorded as (low, high) interval
- Particulate matter directly recorded as counts within particle diameter ranges i.e. histogram



### Example: Discretised data = histogram



Scatterplot with loess line

▶ E.g. point (4.0, 0.0) actually lies within [3.95, 4.05) × [-0.05, 0.05)

Strong discretisation could have undesired inferential impact

# Symbolic Data Analysis



- Established by Diday & coauthors in 1990s.
- Basic unit of data is a distribution rather than usual datapoint.
  - interval (a, b)
  - p-dim hyper-rectangle
  - histogram
  - weighted list etc.
  - can be complicated by "rules"
- Classical data are special case of symbolic data:

E.g. symbolic interval s = (a, b)equivalent to classical data point x if x = a = b.

Or histogram  $\rightarrow \{x_i\}$  as # bins  $\rightarrow \infty$ .

So symbolic analyses must reduce to classical methods. 5/32

# How do symbolic data arise?



Big data  $\rightarrow$  small (symb) data Easier to analyse (hopefully!)

Possible use in data privacy? Individual can't be indentified.

#### Statistical question:

How to do statistical analysis for this form of data?

 Can arise naturally (measurement error):
 E.g. blood pressure, particulate histogram, truncation/rounding.

#### 'Big Data' context:

- Symbolic data points can summarise a complex & very large dataset in a compact manner.
- Retaining maximal relevant information in original dataset.
- Collapse over data not needed in detail for analysis.
- Summarised data have own internal structure, which must be taken into account in any analysis.

# How to analyse symbolic data?

#### A good idea in principle, however:

- Poorly developed in terms of inferential methods.
- Current approaches:
  - Descriptive statistics (means, covariances)
     ⇒ Methods based on 1<sup>st</sup>/2<sup>nd</sup> moments: clustering, PCA etc.
  - Ad-hoc approaches (e.g. regression)
     ⇒ Can be plain wrong for inference/prediction.
  - Single technique for constructing likelihood functions
     ⇒ Limited model-based inferences
- ► Over-prevalence of models for intervals (a, b) & assuming uniformity ⇒ Need to move beyond uniformity (Lynne Billard)

#### Current SDA research:

Developing practical model-based (e.g. likelihood-based) procedures for statistical inference using symbolic data for general symbols.

# Talk Outline



- What is Symbolic Data Analysis?
- Existing and new SDA models
- Statistical analyses using aggregates Discussion

### Existing models for symbols (Le Rademacher & Billard, 2011)

Symbol:  $S = (S^1, \dots, S^d)^\top$ 

E.g. For random intervals  $[a_i, b_i]$ , i = 1, ..., n:

• 
$$S_i = (a_i, b_i)^\top$$

• 
$$S_i = (m_i, \log r_i)^\top$$

Then specify a standard (classical data) model for  $S_1, \ldots, S_n$ . E.g.

 $(m_i, \log r_i)^\top \sim N(\mu, \Sigma)$ 

#### Problems:

- Model unstable/collapses as  $a_i \rightarrow b_i$  (classic data)
- How to fit equivalent models for classical data to symbols?
  - Fit to means? How to account for variation? etc.
- Symbols are summaries of classical data,  $S = \pi(X_1, \ldots, X_N)$ 
  - Model can only predict symbols
- Q: How to fit models and make predictions at the level of the classical data, based on observed symbols?

One possible approach (Beranger, Lin & Sisson, Submitted) Define  $S = \pi(X_{1:N}) : [\mathcal{X}]^N \to S$  such that  $x_{1:N} \mapsto \pi(x_{1:N})$  then,  $L(S|\theta) \propto \int_x g(S|x,\phi)L(x|\theta)dx$ 

where

- $L(x|\theta)$  standard, classical data likelihood
- $g(S|x, \phi)$  explains mapping to S given classical data x
- $L(S|\theta)$  new symbolic likelihood for parameters of classical model

Gist: Fitting the standard classical model, when the data are viewed only through symbols S as summaries

#### Example: No generative model $L(x|\theta)$

- $g(S|x,\phi) = g(S|\phi) \Rightarrow L(S|\theta) = g(S|\phi)$
- ► Directly modelling symbol = existing likelihood approach (Le Rademacher & Billard, 2011) √

### Modelling a random interval

 $\frac{\text{Aggregation: } S = \pi(X_{1:N}) : \mathbb{R}^N \to S = \{(a_1, a_2) \in \mathbb{R}^2 : a_1 \le a_2\} \times \mathbb{N}$ such that  $x_{1:N} \mapsto (x_{(l)}, x_{(u)}, N)$ .

Let  $s = (s_l, s_u, n)$  with  $s_l = X_{(\ell)}$ ,  $s_u = x_{(u)}$   $\ell < u$  and  $x_j \sim f(X|\theta)$ :

$$L(s|\theta) \propto \int_{x} g(s|x,\phi) L(x|\theta) dx$$
  
=  $\int I(X_{(l)} = s_{l} \& X_{(u)} = s_{u}) \prod_{j} f(X_{j}|\theta) dX_{1:n}$   
=  $\frac{n!}{(\ell-1)!(u-\ell-1)!(n-u)!} f(s_{l}|\theta) f(s_{u}|\theta) F(s_{l}|\theta)^{\ell-1} \times [F(s_{u}|\theta) - F(s_{l}|\theta)]^{u-\ell-1} [1 - F(s_{u}|\theta)]^{n-u}$ 

 $\Rightarrow$  the joint distribution of  $\ell$ -th and *u*-th order statistics from  $f(x|\theta)$ .  $\checkmark$ 

Symbolic  $\rightarrow$  Classical check: If  $s_I \rightarrow s_u = x$  and n = 1 then  $L(s|\theta) = f(x|\theta)$ .

- p: number of points involved in constructing the rectangle
- I(p) : locations of the points (taking values in T)

For 
$$s = (s_{\min}, s_{\max}, s_p, s_{I_p}, n)$$

$$L(s|\theta) = \frac{n!}{(n-s_p)!} \left[ \int_{s_{\min}}^{s_{\max}} f(z|\theta) dz \right]^{n-s_p} \times \ell_{s_p}.$$

• If  $s_p = 2$  then  $s_{l_p} = (s_{\min}, s_{\max})$  and  $\ell_2 = f(s_{\min}|\theta)f(s_{\max}|\theta)$ .

Aggregation: Marginal maxima and minima

 $\rho = 0.95$ 



Aggregation: Marginal maxima and minima

**ρ** = -0.6



Aggregation: Marginal maxima and minima

 $\rho = 0$ 



### Modelling a histogram with random counts <u>Aggregation</u>: $S = \pi(X_{1:N}) : \mathbb{R}^{d \times N} \to S = \{0, \dots, N\}^{B^1 \times \dots \times B^d}$ such that $x_{1:N} \mapsto (\sum_{i=1}^n \mathbb{I}\{x_i \in \mathcal{B}_1\}, \dots, \sum_{i=1}^n \mathbb{I}\{x_i \in \mathcal{B}_B\})$



- ► Assume some fixed bins  $\mathcal{B}_1, \dots, \mathcal{B}_B$  and let  $s = (s_1, \dots, s_B)^\top, \sum_b s_b = n$
- If the X<sub>i</sub> are *iid* then likelihood is multinomial:

$$L(s| heta) \propto rac{n!}{s_1! \dots s_B!} \prod_{b=1}^B p_b( heta)^{s_b}$$

where  $p_b(\theta) \propto \int_{\mathcal{B}_b} f(z|\theta) dz$  under the model.  $\checkmark$ 

 More complicated if data are not *iid* (Zhang, Beranger & Sisson, 2019)

### Modelling a histogram with random counts

• Can recover classical likelihood as  $B \to \infty$ 

$$\lim_{B\to\infty} L(S|\theta) \propto \lim_{B\to\infty} \frac{n!}{s_1!\dots s_B!} \prod_{b=1}^B \left[ \int_{D_b} f(z|\theta) dz \right]^{s_b} = L(X_1,\dots,X_n|\theta)$$

So recover classical analysis as we approach classical data.  $\checkmark$ 

- Consistency: Can show that with a sufficient number of histogram bins can perform analysis arbitrarily close to analysis with full dataset.
- Computationally scalable: Working with counts not computationally expensive latent data.
- Can consider histogram with random bins

# Talk Outline



- What is Symbolic Data Analysis?
- Some existing and new SDA models
- Statistical analyses using aggregatesDiscussion

# Fitting a GEV

n=1000. bins=11 Mean MSE  $\times 10^{-3}$  (1000 reps) 0.3 В  $\mu$  $\sigma$ 5 2.977 7.675 4.091 Density 0.2 10 1.385 1.030 0.916 20 1.278 0.762 0.682 <u>--</u> 1000 1.277 0.809 0.662 Standard 1.268 0.725 0.547 0.0 -2 ٥ 2 6 ν

- Use R's hist command to construct histograms, n = 1,000
- Use fgev command in evd package for standard approach
- Accuracy increases with more bins
- Accuracy close to using full dataset with only 20 bins (No real advantage to 1000 bins over 20)

È

# Fitting a GEV

#### Time in seconds

п	100	1K	10K	100K	1M	10M	100M
Standard	0.018	0.047	0.431	2.860	(*)	(*)	(*)
Symbolic (total)	0.060	0.062	0.062	0.107	0.247	2.217	42.994
Symbolic (hist)	0.055	0.057	0.059	0.104	0.243	2.209	42.943
Symbolic (mle)	0.005	0.005	0.004	0.003	0.004	0.007	0.051

- $\blacktriangleright$  Standard initially faster than symbolic for small datasets  $\sim 1 K$
- Symbolic scales much better > 1K
- \* = fgev crashed on my laptop!
- However, most time for symbolic on histogram construction
- Actual symbolic optimisation super fast (obviously)
- Possible laptop caching problems around 100M
- Faster ways to construct histogram counts than hist for really large datasets (e.g. map-reduce using DeltaRho)

### Spatial Extremes



#### Bureau of Meteorology, New South Wales 🤣 @BOM\_NSW

Fri marks peak day for some of #NSW most heavily populated areas. Temps in western #Sydney well into the 40's, regional western towns similar after many broke records this week, CBD likely to have 5th consecutive day above 30 for 1st time in 8 yrs ow.ly/E9QY50ke617 #heatwave





Bureau of Meteorology, Australia @BOM au

"Severe to extreme heatwave conditions across the southeast interior". Temperatures exceeding 45oC for many locations through western NSW and central Australia this afternoon. Latest at ow.ly/3W6s30n/rdY



- What is the maximum value that a process (Temperature) is expected to reach over some region of interest (NSW/Australia) within the next 20, 50 years?
- Whitaker, Beranger & Sisson (2019, Submitted)

### Spatial Extremes

- Max-stable processes are a useful tool to analyse Spatial Extremes
- ► For e.g. the d.f. of the Gaussian max-stable process model

$$P(Y_1(t) \leq y_1, \dots, Y_K(t) \leq y_K) = \exp\left\{-\sum_{j=1}^K \frac{1}{y_j} \Phi_{K-1}\left(c^{(j)}(\mathbf{y}); \Sigma^{(j)}\right)\right\}$$

- The d.f. of such models becomes rapidly intractable with the number of spatial locations
   Composite Likelihood methods (Padoan et al., 2010)
- Still unfeasible for a large number of locations and temporal observations!!

### Composite symbolic likelihood

- Consider we are only interested in a subset of size j of the K dimensions
- ► Let b<sup>i</sup> be the subset of b defining the coordinates of a *j*-dimensional histogram bin and let B<sup>i</sup> = (B<sup>i<sub>1</sub></sup>,..., B<sup>i<sub>j</sub></sup>) be the vector of the number of marginal bins.

The symbolic likelihood function associated with the vector of counts  $\mathbf{s}_{j}^{i} = (s_{1i}^{i}, \dots, s_{Bi}^{i})$  of length  $B^{i_{1}} \times \dots \times B^{i_{j}}$  is  $L(\mathbf{s}_{j}^{i}; \theta) = \frac{N!}{s_{1i}^{i_{1}} \cdots s_{Bi}^{i_{j}}} \prod_{b^{i}=1}^{B^{i}} P_{b^{i}}(\theta)^{s_{b^{i}}^{i}},$ where  $P_{b^{i}}(\theta) = \int_{\Upsilon_{b_{i_{1}}}^{i_{1}}} \dots \int_{\Upsilon_{b_{i_{j}}}^{i_{j}}} g_{X}(x; \theta) dx$  and  $g_{X}$  is a *j*-dim density.

### Composite symbolic likelihood

- s<sub>j</sub> = {s<sup>i</sup><sub>jt</sub>; t = 1,..., T, i = (i<sub>1</sub>,..., i<sub>j</sub>), i<sub>1</sub> < ... < i<sub>j</sub>} represents the set of j−dimensional observed histograms for the symbolic-valued random variable S<sub>j</sub>
- The symbolic *j*-wise composite likelihood function (SCL<sup>(j)</sup>) is given by

$$L_{SCL}^{(j)}(\mathbf{s}_{j};\theta) = \prod_{t=1}^{l} \prod_{i} L(\mathbf{s}_{jt}^{i};\theta)$$

Components of the Godambe matrix are given by

$$\hat{H}(\hat{\theta}_{SCL}^{(j)}) = -\frac{1}{N} \sum_{t=1}^{T} \sum_{i} \nabla^{2} I(\mathbf{s}_{jt}^{i}; \hat{\theta}_{SCL}^{(j)})$$
$$\hat{J}(\hat{\theta}_{SCL}^{(j)}) = \frac{1}{N} \sum_{t=1}^{T} \left( \sum_{i} \nabla I(\mathbf{s}_{jt}^{i}; \hat{\theta}_{SCL}^{(j)}) \right) \left( \sum_{i} \nabla I(\mathbf{s}_{jt}^{i}; \hat{\theta}_{SCL}^{(j)}) \right)^{\top}$$

### Spatial Extremes - Example

- ► Consider N = 1000 observations at K = 15 spatial locations and T = 1 random histogram
- Spatial dependence of Gaussian max-stable model is  $\sigma_{11} = 300$ ,  $\sigma_{12} = 0$  and  $\sigma_{22} = 300$

В	$\sigma_{11}$	$\sigma_{12}$	$\sigma_{22}$		
2	335.5 (585.5)	5.7 (232.2)	317.2 (125.1)		
3	301.0 (34.5)	-0.1 (16.9)	301.9 (33.5)		
5	299.1 (23.1)	-0.9 (13.2)	299.9 (24.1)		
10	299.8 (20.2)	-0.5 (11.1)	300.0 (20.9)		
15	299.8 (18.9)	-0.3 ( 10.4)	300.0 (19.5)		
25	299.7 (18.0)	-0.3 ( 10.0)	300.2 (18.9)		
Classic	300.76 (17.1)	-0.4 (9.7)	301.02 (18.1)		

Table: Mean (and standard errors) of the symbolic composite MLE  $\hat{\theta}_{SC_L}^{(2)}$  and composite MLE  $\hat{\theta}_{C_L}^{(2)}$  (Classic) from 1000 replications of the Gaussian max-stable process model, for  $B \times B$  histograms for varying values of B.

### Spatial Extremes - Example

► Consider B = 25 bins, K = 10,100 spatial locations and T = 1 random histogram. Repetitions = 10

N	K = 10					K	= 100	
14	tc	ts	t <sub>histDR</sub>	t <sub>histR</sub>	t <sub>c</sub>	ts	t <sub>hist</sub> DR	t <sub>histR</sub>
1 000	71.9	22.5	0.8	0.1	-	2238.0	78.8	12.0
5 000	291.8	19.0	0.8	0.3	-	2650.2	81.7	30.9
10 000	591.7	23.8	0.9	0.5	-	2356.6	85.8	54.1
50 000	2 6 2 6 . 8	24.2	1.7	2.1	-	2 300.6	131.6	237.0
100 000	5610.7	25.4	2.4	4.2	-	2766.9	188.2	461.8
500 000	31 083.1	23.2	7.5	20.6	-	3111.5	627.1	2243.5

Table: Mean computation times (seconds) for different components involved in computing  $\hat{\theta}_{CL}^{(2)}$  and  $\hat{\theta}_{SCL}^{(2)}$ .

#### Note: convergence of the variances

- ▶  $\hat{J}(\hat{\theta}_{SCL}^{(j)})$  requires  $T \to N$  and  $\mathbf{B} \to \infty$  for the convergence towards the classical Godambe matrices to occur.
- ▶ For *T* fixed, convergence still occurs as  $\mathbf{B} \to \infty$  towards a different expression

### Classification

- $Y \in \Omega = \{1, \dots, K\}$  (response),  $X \in \mathbb{R}^D$  (explanatory)
- ► Multinomial Logistic Regression: for realisations  $\mathbf{x} \in \mathbb{R}^{D \times N}$ ,  $y \in \Omega^N$ , parameters  $\boldsymbol{\beta} \in \mathbb{R}^{(D+1) \times K}$ , the likelihood is given by

$$L_{\mathrm{M}}(\mathbf{x}, y; \boldsymbol{\beta}) = \prod_{n=1}^{N} \prod_{k \in \Omega} P_{\mathrm{M}}(Y = k | X = x_n)^{\mathbb{1}\{y_n = k\}},$$

where

$$P_{\mathrm{M}}(Y=k|X) = rac{e^{eta_{k0}+eta_k^ op X}}{1+\sum_{j\in\Omegaackslash\{K\}}e^{eta_{j0}+eta_j^ op X}}.$$

- Other model: One-vs-rest
- Prediction:  $Y_n^{\text{Pred}} = \operatorname{argmax}_{k \in \Omega} P_{\text{Model}}(Y = k | X = X_n), \forall n$
- Prediction accuracy:  $PA^{\text{Model}} = \frac{1}{N} \sum_{n=1}^{N} \mathbb{1}\{Y_n^{\text{Pred}} = Y_n\}$

### Classification

- Let  $\mathbf{X}^{(k)} = (X_n | Y_n = k, n = 1, \dots, N) \in \mathbb{R}^{D \times N_k}$
- If  $N_k = \sum_{n=1}^N \mathbb{1}\{Y_n = k\}$  is huge then  $\mathbf{X}^{(k)}$  can be aggregated
- Histogram-valued symbol leads to likelihood

$$L_{ ext{SM}}(\mathbf{s};eta) \propto \prod_{k\in\Omega} \prod_{b_k=\mathbf{1}_k}^{\mathbf{B}_k} \left(\int_{oldsymbol{\Upsilon}_{\mathbf{b}_k}} P_{ ext{M}}(Y=k|X=x) dx
ight)^{s_{\mathbf{b}_k}}$$

- Statistical improvement: mixture symbolic and classical contributions
- Computational improvements: Composite Likelihood (again!) but requires some adjustment.

### Classification - Example

► Use a Supersymmetric (SUSY) benchmark dataset which consists of:

- Binary response (*K* = 2): signal process (which produces supersymmetric particles) vs background process
- N = 5 million observations
- D = 18 features (8 kinematic properties, 10 functions)
- Comparison with optimal sub-sampling method (Wang et al., 2018)
- Training data: 4 500 000 obs.
- Test data: 500 000 obs.
- We consider the following:
  - One-vs-Rest model
  - Marginal composite likelihood
  - Histogram with random bins  $L_{OO}^{(1)}$
  - Histogram with random counts  $L_{\rm SO}^{(1)}$

### Classification - Example

				Bins			
Likelihood	6	8	10	12	15	20	25
$L_{00}^{(1)}$	74.9	75.9	76.6	77.7	78.1	77.9	78.1
	(11.7)	(14.5)	(12.2)	(15.0)	(18.9)	(21.3)	(27.6)
$L_{SO}^{(1)}$	74.4	73.5	75.8	77.8	77.4	78.0	78.0
~~~	(13.3)	(12.6)	(11.5)	(13.9)	(16.8)	(18.0)	(21.4)

Table: Prediction accuracies percentage (computing time in seconds) on the Supersymmetric dataset using histograms with B bins per margins.

- ▶ Wang et al. (2018) obtain a prediction accuracy of 78.2 with a computation time of 86.1 seconds.
- Simulation study: as good or better prediction accuracy, shorter computation time
- Sub-sampling will produce better MSE of the regression coefficients.

# Talk Outline



What is Symbolic Data Analysis?
Some existing and new SDA models
Statistical analyses using aggregates
Discussion

# Summary

Completely new approach to SDA:

- Based on fitting underlying (classical) model
  - Radically different approach to existing SDA methods
  - Ours is much better!
- Views latent (classical) data through symbols
- Recovers known existing models for symbols but is more general
- Works for more general symbols than currently in use

#### Still to do/Working on:

- Implement more sophisticated statistical techniques using Symbols (Tom's PhD)
- Characterise impact of using symbols on accuracy
  - Trade-off of accuracy vs computation
- Design of symbols for best performance
  - Histogram setting: How many bins? Bin locations?

# How to design symbolic data?



How to design symbols to most efficiently represent dataset without (much) loss of critical information?

E. g. Linear regression with 10 million datapoints.



# **THANK YOU**

#### Manuscripts:

- New models for symbolic data. Beranger, Lin & Sisson. https://arxiv.org/pdf/1805.03316.pdf.
- Composite likelihood methods for histogram-valued random variables. Whitaker, Beranger & Sisson. https://arxiv.org/pdf/1908.11548.pdf.
- Logistic regression models using aggregated data. Whitaker, Beranger & Sisson. In prep.

#### Contact:

B.Beranger@unsw.edu.au