





Composite likelihood and logistic regression models for aggregated data

Tom Whitaker, Boris Beranger, Scott A. Sisson,

UNSW & ACEMS

Mathematics Colloquium, UTS, August 14th

Talk Outline



- Symbolic likelihood
- ... and its limitations
- 2. Symbolic composite likelihoods
 - Methods
 - Applications to spatial extremes
- 3. SCL for logistic regression
 - Methods
 - Application to satellite crop prediction
- 4. Discussion



Main idea

Standard statistical methods analyse classical datasets

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

However there is a rise of *non-standard* data forms:

- As a result of measurement process;
- Blood pressure recorded as (low, high) interval;
- Particulate matter recorded as counts within particle diameter ranges i.e. histogram;



Main idea

'Big Data' context:

Symbolic data points to summarise a complex & very large dataset in a compact manner.

 $S = \pi(X_{1:N}) : [\mathcal{X}]^N \to S$ such that $x_{1:N} \mapsto \pi(x_{1:N})$

- 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.
- \bigoplus Big data ightarrow small (symb) data
- \oplus Possible use in data privacy? Individual can't be identified

Statistical question: How to do statistical analysis for this form of data?

Likelihood-based SDA (Beranger, Lin & Sisson, 2018)

The general approach:

$$L(S|\theta,\phi) \propto \int_{x} g(S|x,\phi) L(x|\theta) dx$$

where

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

Gist: Fitting the standard classical model $L(x|\theta)$, when the data are viewed only through symbols S as summaries.

As
$$S_i o x_i$$
, then $g(S_i|x,\phi) = g(x_i|x) = \delta_{x_i}(x)$ and so

 $L(S_i|\theta,\phi) \propto \int_x \delta_{x_i}(x)L(x|\theta)dx = L(x_i|\theta)$ (classical likelihood)

Different symbols give different forms for $g(S|x, \phi)$ (and $\therefore L(S|\theta, \phi)$).

Specific case: Random histograms

Underlying data $X_1, \ldots, X_N \in \mathbb{R}^d \sim g(x|\theta)$ collected into random counts histogram, with fixed bins $\mathcal{B}_1, \ldots, \mathcal{B}_B$.

Aggregation:

 $S = \pi(X_{1:N}) : \mathbb{R}^{d \times N} \to S = \{0, \dots, N\}^{B^1 \times \dots \times B^d} \text{ such that} \\ x_{1:N} \mapsto \left(\sum_{i=1}^n \mathbb{I}\{x_i \in \mathcal{B}_1\}, \dots, \sum_{i=1}^n \mathbb{I}\{x_i \in \mathcal{B}_B\}\right).$

 $g(S|x,\phi) = \begin{cases} 1 & \text{if } s_b \text{ observations in bin } b; \text{ for each } b = 1, \dots, B \\ 0 & \text{else} \end{cases}$

The symbolic likelihood is then (multinomial):

$$L(S|\theta) \propto \int_{x} g(S|x) \prod_{k=1}^{n} g(x_{k}|\theta) dx \propto \prod_{b=1}^{B} \left(\int_{B_{b}} g(z|\theta) dz \right)^{s_{b}}$$

 \Rightarrow generalises univariate result of McLachlan & Jones (1988). \checkmark

Specific case: Random histograms

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

$$\lim_{B\to\infty} L(S|\theta) \propto \lim_{B\to\infty} \prod_{b=1}^{B} \left[\int_{B_b} g(z|\theta) dz \right]^{s_b}$$
$$= L(X_1, \dots, X_n|\theta)$$

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.
- Some approximation of L(S|θ) to L(x|θ) depending on level of discretisation. Work needed to quantify this.
- More complicated if data are not *iid* but exchangeable (Zhang, Beranger & Sisson (2020), SJS)

Fitting a GEV

Suppose $X_1, \ldots, X_n \sim \text{GEV}(\mu, \sigma, \xi)$ for large *n*. Create histogram of counts $s = (s_1, \ldots, s_B)$. Symbolic log-likelihood function is then

$$\ell(s|\mu,\sigma,\xi) \propto \sum_{b=1}^{B} s_b \log \left[G(a_{b+1}|\mu,\sigma,\xi) - G(a_b|\mu,\sigma,\xi)\right]$$



Note that $\ell(s|\mu, \sigma, \xi)$ tends to standard likelihood as # bins gets large (so 1 or 0 observations per bin)

Computation:

- ▶ Optimisation of ℓ (v. quick)
- Creation of histogram s (slower)
- \leftarrow good fits with moderate bin numbers

Limitations: Calculating $\int_{B_b} g(z|\theta) dz$



1D: Probability of X_i falling in bin $(L_{s_1}, L_{s_1-1}]$ is $F(L_{s_1}|\theta) - F(L_{s_1-1}|\theta)$.

2D: Probability of X_i falling in bin $(L_{s_1}^{(1)}, L_{s_1-1}^{(1)}] \times (L_{s_1}^{(2)}, L_{s_1-1}^{(2)}]$ is

 $F\left(L_{s_{1}-1}^{(1)},L_{s_{2}-1}^{(2)}|\theta\right)-F\left(L_{s_{1}-1}^{(1)},L_{s_{2}}^{(2)}|\theta\right)-F\left(L_{s_{1}}^{(1)},L_{s_{2}-1}^{(2)}|\theta\right)+F\left(L_{s_{1}}^{(1)},L_{s_{2}}^{(2)}|\theta\right).$

d-D: Has 2^d components – viable for low d.

Limitations for histograms

For the symbolic log likelihood (using multivariate histograms)

$$\log L(s| heta) \propto \sum_{b=1}^{B} s_b \log \left[\int_{B_b} g(z| heta) dz
ight]$$

there are some limitations:

- Multivariate histograms become inefficient as d gets large

 number of bins to cover d dimensions accurately gets large fast.
- $\int_{B_b} g(z|\theta) dz$ has 2^d components for each bin only viable for low d.
- ▶ This means that we are limited to low-dimensional symbols.
 - SDA traditionally uses 1- or 2-dimensional symbols, so maybe this is ok;
 - But in principle there will be an analysis that requires higher-D information
 So should resolve this problem if possible.

One option: Composite likelihoods.

Talk Outline



- Symbolic likelihood
- ...and its limitations
- 2. Symbolic composite likelihoods
 - Methods
 - Applications to spatial extremes
- 3. SCL for logistic regression
 - Methods
 - Application to satellite crop prediction
- 4. Discussion



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 (2020, Stat. Comput.)

Modelling Australian temperature spatial extremes



- ▶ 105 spatial locations with temperature observation, over time
- Want to fit spatial model to temperature extremes. Spatial (multivariate) information is important!
- Lots of data can form 105-dimensional histogram(!)
- Can't fit this using $L(S|\theta)$. What can we do?

Modelling Australian temperature spatial extremes



- ▶ 105-dimensional histograms are completely infeasible.
- ► But lower-dimensional histograms could still be very informative.
- E.g. 2-dimensional.
- ▶ But how to do inference with 105×104/2 bivariate histograms?
- One answer: Composite likelihoods

Composite likelihoods

Standard likelihood $L(\mathbf{x}|\theta)$ where $\mathbf{x} = (x_1, \dots, x_d)$ where e.g. $x_i = \text{data at } i\text{-th spatial location.}$

Suppose $L(\mathbf{x}|\theta)$ is computationally intractable except for e.g. d = 2 (as for the spatial extremes model we are using).

Then can construct (say) **pairwise composite likelihood** $L_{CL}^{(2)}(\mathbf{x}|\theta) \propto \prod_{i} \prod_{j>i} L(x_i, x_j|\theta)$ from all bivariate marginal events.

Works if $\ell(x_i, x_j | \theta)$ is an unbiased estimating equation for θ (as then log likelihood is a sum of these and so is also an unbiased estimating equation for θ). So works well for spatial models.

Similarly *j*-wise composite likelihoods.

Composite likelihoods

Behaviour of composite MLE

 $\hat{ heta}_{CL}^{(j)}$ is asymptotically ($N
ightarrow \infty$) consistent and distributed as

$$\sqrt{N}\left(\hat{\theta}_{CL}^{(j)}-\theta\right) \to N\left(0,\ G^{(j)}(\theta)^{-1}\right)$$

where

•
$$G^{(j)}(\theta)^{-1} = H^{(j)}(\theta)J^{(j)}(\theta)^{-1}H^{(j)}(\theta)$$
 is Godambe information matrix
• $H^{(j)}(\theta) = -\mathbb{E}(\nabla^2 \ell_{CL}^{(j)}(\theta; \mathbf{x}))$ is the sensitivity matrix
• $J^{(j)}(\theta) = \mathbb{V}(\nabla \ell_{CL}^{(j)}(\theta; \mathbf{x}))$ is the variability matrix.

- ► For standard likelihoods j = d and $H(\theta) = J(\theta)$ and so $G(\theta) = H(\theta) = I(\theta)$ is the Fisher information matrix.
- How can this help us with L(S|θ) when S is a 105-dimensional histogram?

Composite symbolic likelihoods

As with $L_{CL}^{(2)}(\mathbf{x}|\theta) \propto \prod_{i} \prod_{j>i} L(x_i, x_j|\theta)$ we may similarly have $L_{CL}^{(2)}(\mathbf{S}|\theta) \propto \prod_{i} \prod_{j>i} L(S_{ij}|\theta)$

where S_{ij} is the bivariate marginal histogram for dimensions (i, j).

In this setting

$$L(S_{ij}|\theta) \propto \prod_{b} \left(\int_{B_b} g(z_1, z_2|\theta) dz_1 dz_2 \right)^{s_b}$$

as before.

Composite symbolic likelihoods



But when the bins (number and volume) are fixed then

$$\sqrt{N}\left(\hat{\theta}-\theta\right)
ightarrow N\left(??(\theta,\textit{bins}),??(\theta,\textit{Bins})^{-1}
ight).$$

Currently working on non-asymptotic (in bins) distribution of MLE

Composite symbolic likelihoods



But when the bins (number and volume) are fixed then, as before

 $\sqrt{N}\left(\hat{\theta}_{SCL}^{(j)}-\theta\right)
ightarrow N\left(??(\theta,\textit{bins}),??(\theta,\textit{Bins})^{-1}
ight).$

Similarly work in progress.

Simulated spatial extremes

(Mean) Pairwise symbolic composite likelihood estimates ($\hat{\theta}_{SCL}^{(2)}$):

- Consider $N = 1\,000$ observations at K = 15 spatial locations and T = 1 random histogram
- ▶ Spatial dependence of Gaussian max-stable model is $\sigma_{11} = 300$, $\sigma_{12} = 150$ and $\sigma_{22} = 200$

| В | σ_{11} | σ_{12} | σ_{22} |
|---------|---------------|---------------|-----------------|
| 2 | 321.6 (360.0) | 162.3 (210.6) | 210.8 (131.2)) |
| 3 | 296.1 (30.6) | 147.4 (20.1) | 197.9 (19.9) |
| 5 | 298.8 (23.3) | 149.4 (15.3) | 199.6 (15.4) |
| 10 | 299.0 (19.3) | 149.6 (12.3) | 199.7 (12.9) |
| 15 | 299.5 (18.7) | 149.8 (11.6) | 199.8 (12.1) |
| 25 | 299.7 (17.8) | 150.0 (11.2) | 200.0 (11.8) |
| Classic | 300.7 (16.4) | 150.6 (10.2) | 200.6 (10.9) |

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.

- As "bins → ∞" performance approaches classical composite likelihood (also estimated the marginal parameters).
- "Acceptable" results for B = 10

Simulated spatial extremes

(Mean) Time comparisons for increasing N

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

| N | K = 10 | | | K = 100 | | | | |
|---------|----------------|------|---------------------|--------------------|----------------|---------|---------------------|--------------------|
| /\ | t _c | ts | t _{histDR} | t _{histR} | t _c | ts | t _{histDR} | t _{histR} |
| 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 | 2 243.5 |

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

- Classical composite likelihood rapidly not feasible as spatial dimensions increases (K = 20)
- Symbolic approach much more efficient

Simulated spatial extremes

Plenty more simulations and interesting results in the paper¹

- Including effect of number of histograms
- ...and allocation of micro-data between them;
- Comparing bivariate SCL and trivariate SCL.

Now consider a different problem with "high" dimensional histograms: \longrightarrow logistic regression \longleftarrow

¹Whitaker T., B. Beranger and S. A. Sisson (2020). Composite likelihood methods for histogram-valued random variables. Stat. Comput., In Press.

Talk Outline



- Symbolic likelihood
- ... and its limitations
- 2. Symbolic composite likelihoods
 - Methods
 - Applications to spatial extremes
- 3. SCL for logistic regression
 - Methods
 - Application to satellite crop prediction
- 4. Discussion







- \sim 250K pixels with 7-dimensional predictor variable x
- 7 response categories with known ground truth
- ▶ Multinomial logistic regression, but computational to fit (8+ hours)
- Can we use SDA to speed things up while maintaining prediction quality?

 (Y_i, X_i) pairs, $Y_i \in \Omega = \{1, \dots, K\}$ and $X \in \mathbb{R}^d$, $i = 1, \dots, N$

Want to predict category Y_i given vector X_i .



Odds ratio linear model

$$\log\left(\frac{P(Y=k|X)}{P(Y\neq k|X)}\right) = \beta_{0k} + \beta_k^\top X$$

so that

$$P(Y = k | X) = rac{\exp\{eta_{0k} + eta_k^\top X\}}{1 + \exp\{eta_{0k} + eta_k^\top X\}}$$

and so the standard classical likelihood is

$$L(\mathbf{Y}, \mathbf{X}|\boldsymbol{\beta}) \propto \prod_{i=1}^{N} \left(P(\mathbf{Y} = Y_i | X = X_i) \prod_{j \in \Omega \setminus \{Y_i\}} P(\mathbf{Y} \neq j | X = X_i) \right).$$

 (Y_i, X_i) pairs, $Y_i \in \Omega = \{1, \dots, K\}$ and $X \in \mathbb{R}^d$, $i = 1, \dots, N$

Want to predict category Y_i given vector X_i . Construct 7-dimensional predictor histogram S^j from $\{X_i : Y_i = j\}$, for each crop type j = 1, ..., 7.

The histogram-based likelihood for $\boldsymbol{S} = (\boldsymbol{S}^1, \dots, \boldsymbol{S}^7)$ is then $L_{\boldsymbol{S}}(\boldsymbol{S}|\boldsymbol{\beta}) \propto \prod_{k \in \Omega} \prod_{b_k} \left(\int_{B_{b_k}} P(Y = k | X = x) dx \prod_{j \in \Omega \setminus \{k\}} \int_{B_{b_k}} P(Y \neq k | X = x) dx \right)^{s_{b_k}}.$

- Standard application of symbolic likelihood
- Can only do this integral if each category has only one predictor X. (Will return to this shortly...)

Interesting result on existence of MLEs

Standard logistic regression:

 Â = arg max L(Y, X|β) exists and is unique if their is neither complete nor quasi-complete separation of the data (Albert and Anderson, 1984).

 Histogram-based logistic regression:

2) $\hat{\boldsymbol{\beta}}_{S} = \arg \max L_{S}(\boldsymbol{S}|\boldsymbol{\beta})$ exists and is unique if the set of histograms $(\boldsymbol{S}^{1}, \dots, \boldsymbol{S}^{K})$ does not exhibit complete nor quasi-complete separation² of the data (Whitaker, Beranger & Sisson, 2019; arxiv).

- 2) is stronger than 1), so 2) ightarrow 1)
- So if $\nexists \hat{\beta}_S \Rightarrow \nexists \hat{\beta}$ (i.e. if $\exists \hat{\beta} \Rightarrow \exists \hat{\beta}_S$)
- However $\hat{\boldsymbol{\beta}}_{\mathcal{S}}$ can exist where $\hat{\boldsymbol{\beta}}$ does not
- So can do something in SDA that you can't with classical data³
- (Gets a bit crazy when $\boldsymbol{S} \to \boldsymbol{X}$... as then $\exists \hat{\boldsymbol{\beta}}_S \to \nexists \hat{\boldsymbol{\beta}}!!$)

²For modified definitions of separation compared to Albert and Anderson (1984)

³Not sure this is useful though!

Recall $L_{S}(\boldsymbol{S}|\boldsymbol{\beta})$ is proportional to

$$\prod_{k\in\Omega}\prod_{b_k}\left(\int_{B_{b_k}}P(Y=k|X=x)dx\prod_{j\in\Omega\setminus\{k\}}\int_{B_{b_k}}P(Y\neq k|X=x)dx\right)^{s_{b_k}}$$

- This works as standard application of symbolic likelihood
- ► However, can only do this integral if each category has only one predictor X.
- ► So either need to do *d*-dimensional computational integration or
- ...abuse ideas from composite symbolic likelihoods

... Abuse ideas from composite symbolic likelihoods

- As we can integrate $L_S(\boldsymbol{S}|\boldsymbol{\beta})$ for univariate predictor
- Construct composite likelihood over all univariate predictor likelihoods
- Or over all 2-dimensional predictor likelihoods

• Or . . . etc.

The Good:

Gets around high-dimensional integration (1-d is particularly good)

The Bad:

- Each marginal event is not an unbiased estimating equation
- So this is not a "true" composite likelihood
- The estimates of β will be biased
- All parameters depressed $\beta \downarrow 0$ (known result)

The (partial) Fix:

- Can reduce the bias using some modifications to this composite likelihood following ideas in a related context by Cramer (2007)
- Does not eliminate it
- ► However prediction can still be good if reduction in β ↓ 0 is similar for all parameters
- This is what we found to happen in practice



- \sim 250K pixels with 7-dimensional predictor variable x
- 7 response categories with known ground truth
- ► Use (modified) symbolic composite likelihood over all 1-D predictors
- (Lasso regularisation included)





Histograms of band3

Histograms of band4



Histograms of band5



band5



Some predictors clearly identify crops even in 1-d (e.g. Bare soil)

| | | | | | Bins | | | |
|-----------------|----------------|-------|-------|-------|-------|-------|-------|--|
| Crop type | N _k | 6 | 8 | 10 | 12 | 15 | 20 | $L_M(\boldsymbol{X}, \boldsymbol{Y} \boldsymbol{\beta})$ |
| Cotton | 72 450 | 90.5 | 90.6 | 92.8 | 93.6 | 94.0 | 94.1 | 92.2 |
| Sorghum | 66 751 | 74.6 | 74.8 | 75.7 | 76.4 | 76.2 | 76.3 | 80.3 |
| Pasture Natural | 27 479 | 75.7 | 75.4 | 76.0 | 76.8 | 77.0 | 77.1 | 77.6 |
| Bare Soil | 26 173 | 88.0 | 89.6 | 89.2 | 90.0 | 89.5 | 90.1 | 91.0 |
| Peanut | 17 868 | 81.2 | 81.3 | 81.5 | 81.5 | 81.9 | 81.6 | 82.9 |
| Maize | 12 986 | 9.7 | 9.9 | 10.2 | 10.4 | 10.3 | 10.4 | 14.2 |
| Wheat | 10 778 | 3.4 | 4.0 | 4.8 | 5.0 | 5.2 | 5.7 | 10.3 |
| Overall | 234 485 | 74.6 | 75.5 | 76.4 | 77.1 | 77.2 | 77.2 | 78.1 |
| Time (secs) | | (164) | (162) | (221) | (229) | (276) | (508) | (6071) |

Table: Crop specific and overall prediction accuracies (%) using univariate marginal histograms with *B* bins. The likelihood optimisation times (in seconds) are reported in the last row. The full model is the standard multinomial likelihood $L_M(X, Y|\beta)$ with LASSO regularisation.

Pretty good results with 10 bins

- Overall accuracy 76.4% (histogram) versus 78.1% (classical)
- Poorer performance for less numerous crops (wheat, maize)

\blacktriangleright ... and 27× faster

More simulations and details in the paper.

Talk Outline



- Symbolic likelihood
- ...and its limitations
- 2. Symbolic composite likelihoods
 - Methods
 - Applications to spatial extremes
- 3. SCL for logistic regression
 - Methods
 - Application to satellite crop prediction
- 4. Discussion



Summary

Summary

- Symbolic composite likelihoods a natural extension of symbolic likelihoods
- Can avoid likelihood integration issues for some models
- Was useful for prediction even in models for which composite likelihoods are not suited (logistic regression)

Questions

- Other ways to avoid problematic integration in high dimensional histograms?
- Other ways to do high-dimensional regression (general) with symbolic likelihood?
- Etc.







THANK YOU

Relevent Manuscripts:

Beranger, Lin & Sisson (2018). New models for symbolic data analysis.

https://arxiv.org/abs/1809.03659

Whitaker T., B. Beranger and S. A. Sisson (2020). Composite likelihood functions for histogram-valued random variables. Stat. Comput., In press.

Whitaker T., B. Beranger and S. A. Sisson (2019). Logistic regression models for aggregated data.

https://arxiv.org/abs/1912.03805

Contact:

B.Beranger@unsw.edu.au www.borisberanger.com