Two problems in Astrostatistics:

Identification of Shell and Stream Substructures in Galaxy Debris¹ &

Nonparametric Estimation in Contamination Models²

Bodhisattva Sen³ Department of Statistics Columbia University, New York

Statistical Challenges in Modern Astronomy – VI Carnegie Mellon University, Pittsburgh

9 June, 2016

 1 Joint work with R. Patra, R. Biswas, D. Hendel, K. Johnston 2 Joint work with R. Patra

³Supported by NSF Grant DMS-1150435 (CAREER)

I. Identification of *shell* and *stream* substructures

in galaxy debris

<ロ > < 団 > < 豆 > < 豆 > < 豆 > < 豆 > < 豆 > < 豆 > < 2/32

Tidal debris morphology: The *debris* from minor *mergers* can be broadly divided into *two* morphological categories – *streams* and *shells*



Figure : Galaxy NGC 5907 (left) and Galaxy NGC 474 (right)

- The goal of this collaborative project is to *identify* and *distinguish* different morphologies (i.e., shells and streams)
- We hope that this analysis will provide new insights into our understanding of *structure formation*

- *Stream*-like substructures stretch approximately along the progenitor's orbit, sometimes *wrapping* around the host multiple times
- *Shell*-like structures may extend both along and perpendicular to the path of their disrupted parent, forming an *umbrella*-shaped distribution of stars and/or *sharp* edges in the light distribution
- We study the outputs of *N-body* simulations (across a wide range of orbital and galactic parameters); Hendel and Johnston (2015)



We will work with these two simulations



Shells versus Streams

• Shells have a sharp edge in surface brightness on one side and a smooth decay on the other

(日) (同) (日) (日)

• Streams appear more symmetric around the 'ridge'

400 400 200 200 0 0 -200 200 450.5::0.45::1 63.1::0.5::2 400 296.6::0.7::2 400 320.8::0.95::2 338.2::1.2::2 420.1::0.25::1 525.4::0.5::1 -200 -400 200 400 600 -600 -400 -200 200 400 0 0



Detected shells for simul 2

600

Shells versus Streams

- *Shells* have a *sharp edge* in surface brightness on one side and *smooth decay* on the *other*
- Streams appear more symmetric around the 'ridge'

A ridge a *local maxima* in *one* direction (like a *mode*)



(日)

- **Data**: 2-D position of particles $(\{\vec{x_i} := (x_{i1}, x_{i2})\}_{i=1}^n)$
- Construct kernel density estimator $p(\vec{x})$ to obtain the surface density, i.e.,

$$p(\vec{x}) = rac{1}{nh^2} \sum_{i=1}^n K\left(rac{x_{i1} - x_1}{h}, rac{x_{i2} - x_2}{h}
ight)$$



▲ロト ▲舂 ト ▲ 臣 ト ▲ 臣 ト ○臣 … 釣ん(で

- Hessian matrix: $H(\vec{x}) = \frac{\partial^2}{\partial x_1 \partial x_2} p(\vec{x})$ at each point \vec{x}
- *Eigenvector* of the *smallest eigenvalue* of the $H(\vec{x})$: $\vec{v}(\vec{x})$
- The *ridge* is the set of points $\vec{x_r}$ such that

 $ec{v}(ec{x_{\mathrm{r}}})^{ op}
abla p(ec{x_{\mathrm{r}}}) = 0 \quad \text{and} \quad ec{v}(ec{x_{\mathrm{r}}})^{ op} H(ec{x_{\mathrm{r}}}) ec{v}(ec{x_{\mathrm{r}}}) < 0,$

local maxima in p along the eigendirection of the smallest eigenvalue

• Subspace Constrained Mean Shift [Ozertem, U. & Erdogmus, D. (2011)]





- Fix a *window* perpendicular to the ridge at a point (in the ridge)
- Streams are classified by looking for a symmetric peak in the window





< • • • **•**

≣ ▶ ≣ *∙*0.<?



If a *peak* is not detected, then we try to look for a *shell* nearby

Alpha= 0.045



(目) ▲目) 目 のへで



(ロトメ都トメミトメミト、ミーク

Simul 1

Detected shells for simul 1



х

х

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ 三臣 - のへ(?)

Simul 1

Detected streams for Simul 2







イロト イヨト イヨト イヨト





Multi-scale idea

・ロト・西・・川・・ 小田・ ・日・





















Detected streams for Simul 2

х



Detected shell and streams for simul 1 at h= 3



Detected shell and streams for simul 1 at h= 5



Detected shell and streams for simul 1 at h= 7



Detected shell and streams for simul 1 at h= 10



Detected shell and streams for simul 1 at h= 12



Detected shell and streams for simul 1 at h= 15



Detected shell and streams for simul 1 at h= 17



Detected shell and streams for simul 1 at h= 20



Detected shell and streams for simul 1 at h= 22



Detected shells for simul 1

Detected streams for Simul 2

23/32

Issues

- How to choose the *bandwidth* (crucial)? Use a *multi-scale* idea
- Peaks are not always well-defined
- How to detect the *boundary* of the shell? Look at the *drop* in density
- How should we choose the width/length of the window?
- *Dip test* (for unimodality) used to choose the *length* of the window

II. Contamination Models

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ 三重 めんの

- Data: Radial velocity (RV) of stars (n = 1215) from Carina (dSph), contaminated with Milky Way stars in the field of view.
- **Question**: Find the distribution of RV of stars in *Carina* F_s and the *proportion* of stars from Carina α
- *F_b*: The distribution of RV of the contaminating stars; known from the Besancon Milky Way model (Robin et. al, 2003)



Plot of histogram of F_b (blue) overlaid with the (scaled) KDE of data

・ロト・四ト・ヨト・ヨー シタウ

Mixture model with two components

$$F(x) = \alpha F_s(x) + (1 - \alpha)F_b(x)$$

- F_b is a *known* distribution function (DF)
- Unknowns: Mixing proportion $\alpha \in (0,1)$ and DF F_s $(\neq F_b)$
- Problem: Given a random sample from $X_1, X_2, \ldots, X_n \stackrel{i.i.d.}{\sim} F$ (and F_b), we wish to (nonparametrically) estimate F_s and the parameter α

Applications

- In contamination problems application in astronomy
- In *multiple* testing problems the *p*-values are *uniformly* distributed on [0,1], under *H*₀, while their distribution associated with *H*₁ is *unknown*; see e.g., Efron (2010)

Identifiability

- When α is *unknown*, the problem is *ill-posed*
- If $F = \alpha F_s + (1 \alpha)F_b$ for some F_b (known) and α (unknown), then the mixture model can be re-written as

$$F = (\alpha + \eta) \left(\frac{\alpha}{\alpha + \eta} F_s + \frac{\eta}{\alpha + \eta} F_b \right) + (1 - \alpha - \eta) F_b,$$

for $0 \le \eta \le 1 - \alpha$, and the term $(\alpha F_s + \eta F_b)/(\alpha + \eta)$ can be thought of as the *unknown* DF. Thus, if the model holds for α then it holds for $\alpha + \eta$, for every $\eta \le 1 - \alpha$.

Identifiable parameter

We redefine the mixing proportion as

$$\alpha_{0} := \inf \left\{ \gamma \in (0,1] : \frac{F - (1 - \gamma)F_{b}}{\gamma} \text{ is a valid } DF \right\}$$

Intuitively, this definition makes sure that the "signal" distribution F_s does not include any contribution from the known "background" F_b

・ロット 全部 マート・ キョン

= na<</p>

Estimation of α_0 and F_s

- Can develop a *consistent* estimator $\hat{\alpha}_n$ of α_0
- Can construct *tuning parameter free* nonparametric estimators of *F_s* and *f_s* the density of *F_s* (under some conditions)

Result

• If the model is *identifiable*, then $\hat{F}_{s,n}$ and $\hat{f}_{s,n}$ are consistent



Plot of $\hat{F}_{s,n}$ (in dotted red) and F_s (in dashed black) when n = 300

Lower confidence bound for α_0

• We can construct a *finite sample* (honest) *lower confidence bound* $\hat{\alpha}_L$ with the property

$$\mathbb{P}(\alpha_0 \geq \hat{\alpha}_L) \geq 1 - \beta, \quad \text{for all } n,$$

for a specified confidence level $(1 - \beta)$, $0 < \beta < 1$

- Would allow one to assert, with a specified level of confidence, that the proportion of "signal" is *at least* $\hat{\alpha}_L$
- $\hat{\alpha}_L$ is defined in terms of the (1β) -quantile of a distribution H_n
- H_n is *distribution-free*, when F is continuous, and can be simulated
- Requires no tuning parameters; lower bound holds for all n

Astronomy application

 Data: Radial velocity (RV) of stars (n = 1215) from Carina (dSph), contaminated with Milky Way stars in the field of view.



- Astronomers usually assume the distribution of the radial velocities for these dSph galaxies to be *Gaussian* in nature.
- The right panel shows $\hat{F}_{s,n}$ (in dashed red) overlaid with the *closest* Gaussian distribution (in blue).
- Our estimate $\hat{\alpha}_n$ of α_0 turns out to be 0.356, while the lower bound $\hat{\alpha}_L$ (at level 0.05) is found to be 0.322.

References

- Biswas, R., Hendel, D., Johnston, K. V., Patra, R. K., and Sen, B. (2016). Statistical tools to categorize debris morphologies in a galaxy. in preparation.
- Hendel, David and Johnston, Kathryn V. (2015). Tidal debris morphology and the orbits of satellite galaxies, Monthly Notices of the Royal Astronomical Society, 454,2472–2485.
- [3] Ozertem, U. and Erdogmus, D. (2011). Locally defined principal curves and surfaces. The Journal of Machine Learning Research, 12, 1249–1286.
- [4] Patra, R. K. and Sen, B. (2015). Estimation of a two-component mixture model with ap- plications to multiple testing. J. Roy. Statist. Soc. Ser. B (to appear).

Thank You! Questions?

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで