

Harnessing Tailored Statistical Techniques to Discover Star Clusters

KICC Focus Meeting on Astrostatistics & Astro-ML

Rafael S. de Souza

Chair: The Cosmostatistics Initiative

University of Hertfordshire



Astrophysical Outline

- Why Star Clusters?
- The Old & The Young
- Optimizing discoveries
- How star clusters can map the galaxy and beyond

Methodological Outline

Some tailored solutions for everyday statistical problems

- Missing data imputation
- Low-rank Heteroskedastic data-denoising
- GPU Scalability
- And old-school Hierarchical Bayesian Models

Types of Star Clusters

- Young stellar objects clusters
 - Offers a glimpse into early star and planet formation processes.
 - They are independent tracers of the galactic spiral arms structure.
- Open Clusters (OCs)
 - Comprised of stars of mixed ages and higher metallicity, OCs map galactic chemical enrichment.
 - Their location helps tracing the galaxy's spiral structure and star formation history.
- Globular Clusters (GCs)
 - Old, metal-poor stars, they are relics of the early Universe, shedding light on the formation and evolution of the Milky Way.
 - Their dynamics provide constraints on dark matter.

Mapping Young Stellar Objects in the Milky Way

- YSOs live in regions of intense star formation.
- They enable to map of the galactic structure. Because they are close to the place they are born.
- Challenge is to identify them among 10⁸⁻⁹ objects observed by the Gaia space mission. With upcoming surveys, those numbers will be at least ten times larger.



Images of range of and charms, sharp with NGC 1312, an object provinsingl chandled in a charm' but now known to be an intranse. The field of view in a fitter sharp of a particular sharp of the sharp o

Mapping Young Stellar Objects in the Milky Way

- YSOs live in regions of intense star formation.
- They enable to map of the galactic structure. Because they are close to the place they are born.
- Challenge is to identify them among 10⁸⁻⁹ objects observed by the Gaia space mission. With upcoming surveys, those numbers will be at least ten times larger.



Images of range of and charters, holge with NGC 1217, and other processing distuffed in a charter but now known to the an interime. The field of view in a Hinnes is Jac v 5 Jacon Morthin jace, patier havin as relateded by a hole. The Art a staget of time Roberts or et al. 2013, particle from NASA and ISA, part of them NASA and ISA, Davids Dr. Martin (ESA/Haldels and Zakar W Chinessis) (University of Annues, USA), and el Born NASA, RSA, RSA (STSE), and el holdes: Einstrugt of TSE/A/RAAF SA/Haldel SA/Haldels and NASA. Chain: Capacitance particle phone NASA, RSA, and the Haldes: Einstrugt of TSE/A/RAAF SA/Haldel SA/Haldels and NASA. Chain: Capacitance particle phone NASA, RSA, and the Haldes: Einstrugt of TSE/A/RAAF SA/Haldel and the safety of the saf

Mapping Young Stellar Objects in the Milky Way

- YSOs live in regions of intense star formation.
- They enable to map of the galactic structure. Because they are close to the place they are born.
- Challenge is to identify them among 10⁸⁻⁹ objects observed by the Gaia space mission. With upcoming surveys, those numbers will be at least ten times larger.



Images of a range of and charters, Kalang with NGC 1217, and shops provinsing dranked in a charter but now known to its an internation field of vision in Binnois 19, as v 19, and Northin jung anglen thins are indicated by an italian. The Narth Agend Binnois 19, as v 19, and 19, and 19, but with 19, but with

YSO data: Spectral Energy Distribution





First issue: Missing data

- Most off-the-shelf approaches assume missingness at random:
- An alternative is to learn the joint distribution from the complete data, which often requires assumptions about the joint density



First issue: Missing data

 Astronomical data shows non-trivial missing patterns



First issue: Missing data

• How can we take advantage of the data's correlated structure for arbitrary marginal distributions?



Sklar's Theorem: Let *F* be a *p*-dimensional joint distribution function with marginals F_1, \ldots, F_p . Then there exists a copula *C* with uniform marginals such that

$$F(x_1,\ldots,x_p)=C(F_1(x_1),\ldots,F_p(x_p))$$



Multiple Imputation via Generative Adversarial Networks



MIGAN employs a self-attention mechanism, which learns a sparse representation of the relevant features for a given task (de Souza et al, in prep). Initially used for images, can be adapted to Astronomical catalogues.

Multiple Imputation via Generative Adversarial Networks



MIGAN employs a self-attention mechanism, which learns a non-local sparse representation of the data.

The MICE Algorithm

Missing data is in red. There is a strong correlation between A and B, so let's try to impute A using B and C.

в C ٨ 0.93 1.40 1.53 0.24 0.46 0.76 0.80 0.95 1.24 1.46 0.23 0.57 0.90 1.28 0.15 0.42 0.47 0.54 0.63 1.14 0.89 1.23 1.45



Missing data is filled in randomly. This dillutes the correlations, but allows us to impute using all available data.

Α в с 0.93 1.40 1.53 0.24 0.46 0.76 0.90 0.80 0.95 1.24 1.46 0.23 0.57 0.90 0.46 1.28 0.15 0.42 0.47 0.54 0.63 1.14 0.89 1.23 1.45

 $R^2 = 0.4106$

2.0

A random forest is used to predict A with B and C. Notice the correlation between A and B improved.

٨

0.93

0.24

0.24

0.95

0.23

0.90

0.15

0.47

0.89

0.89

B

1.40

0.46

0.80

1.24

0.57

0.42

0.54

1.14

1.23

 $R^2 = 0.5311$

c

1.53

0.76

1.46

1.28

0.63

1.45

After Imputing B using A and C, we have achieved a correlation between A and B much closer to the original data.

	Α	В	C
	0.93	1.40	1
	0.24	0.46	0.
	0.24	0.80	1
	0.95	1.24	1
5	0.23	0.57	1
	0.90	1.24	1
	0.15	0.42	1
	0.47	0.54	0.
	0.89	1.14	1
	0.89	1.23	1





MIGAN as Emulator



MIGAN also enables to user to mimic a particular model of choice as e.g. Multiple Imputation via Chained Equations.

YSO search pipeline





The SPitzer/IRAC Candidate YSO Catalog



The largest catalogue of YSOs (\sim 200,000) in the Milky Way midplane



For each YSO association

For star *i* of a cluster, the probability distribution is,

$$\begin{aligned} p_{\text{clust}}(\varpi_i, \mu_{\ell^\star, i}, \mu_{b, i} | \varpi_0, \mu_{\ell^\star, 0}, \mu_{b, 0}) &= \\ \phi(\varpi_i | \varpi_0, \sigma_{\varpi_i}^2) \cdot f(\mu_{\ell^\star, i} | \mu_{\ell^\star, 0}, \sigma_{\mu_{\ell^\star, 0}}^2, \nu_{\mu}) \cdot \\ f(\mu_{b, i} | \mu_{b, 0}, \sigma_{\mu_{b, 0}}^2, \nu_{\mu}), \end{aligned}$$

where $\theta = (\varpi_0, \mu_{\ell^*,0}, \mu_{b,0})$ are the mean astrometric values for the cluster, $x_i = (\varpi_i, \mu_{\ell^*,i}, \mu_{b,i})$ are the measured values for the *i*th star, σ_i are corresponding uncertainties, ϕ denotes a Gaussian distribution, and *f* denotes a *t*-distribution.



• YSOs are independent tracers of Spiral Arm Structure



• We have identified a new structure near the Sagittarius arm



• We then compared it with other independent tracers such as dust maps and masers to confirm the structure was not an artifact



• Our analysis provided the first evidence of a high-pitch angle structure in the galactic spiral arms



SPICY byproducts



- Hundreds of thousands Light-curves (Time-Series)
- The light curve of Gaia23bab (=SPICY 97589) suggests the presence of an accretion outburst.
- These still scarce class of objects play a significant role in our understanding of star and planetary system formation.



SPICY byproducts



- 117,224 stamps of star forming regions
 - → Computer vision
 - → Fourier and Wavelets Analysis
 - → Marked Point Process



Searching for Extragalactic Globular Clusters

- Approximate figures
- Dwarf galaxies: 0 10 GCs
- Disk Galaxies 10s 100s GCs
- Elliptical Galaxies 100s 10k GCs
- Unsurprisingly GCs are usually targeted around E/SO galaxies, because of large numbers and easier detection



Searching for Extragalactic Globular Clusters

- To help mitigate this bias, we start a campaign to search for GCs around Spirals
- Only 105 confirmed GCs around the region (spectroscopic + HST data)



Searching for Extragalactic Globular Clusters

- Data from S-PLUS a ongoing survey mapping about 9300 square degrees of the southern sky with an optical 12-bands.
- The figure shows a typical GC SED and Spectra



Photometric Selection - 7.2K point sources

 A traditional GC selection would apply color-magnitude cuts around regions of known GCs



Photometric Selection - 7.2K point sources

 Going a bit further, we can just apply a Principal Components Analysis



Photometric Selection - 7.2K point sources

- But what about handling heteroskedastic errors with known variance?
- Off-the-shelf packages often don't account for errors in measurements



Yonder: Low-rank data denoising

RNAAS RESEARCH NOTES OF THE AAS

OPEN ACCESS

Yonder: A Python Package for Data Denoising and Reconstruction

Peng Chen (形況)¹ and Rafael S. de Souza² Published March 2022 - 0 2022. The Author(s). Published by the American Astronomical Society. Research Notes of the AAS, Volume B, Number 3 Citation Peng Chen and Rafael S. de Souza 2022 Res. Notes AAS 6 51

Figures . References .

- Uncertainty aware PCA
- Data-denoising



PCA Scalability

	comment and succession of the succession	
-	Comente lists available at SummOnut	The second
CALL OF THE OWNER	Astronomy and Computing	10-00
ET SPOITS	journal formegrager and it	

Full length article

qrpca: A package for fast principal component analysis with GPU acceleration

.

R. S. de Souza^{1,1}, X. Quanferng¹, S. Shen^{1,0}, C. Peng^{1,1}, Z. Mu¹ ¹EV (absorb) for Remote Ta Galaxies and Causting, Shinglar Advanced Discreases, Oliver Nature of Sources III Insure ¹EV (absorb Nature 2), Olive ¹EV (2004), Cline ¹EV (absorb Nature 2), Olive ¹EV (2004), Cline ¹EV (absorb Nature 2), Oliver ¹EV (2004), Cline ¹EV (absorb Nature 2), Oliver ¹EV (2004), Olive¹ ¹EV (absorb Nature 2), Olive¹ ¹EV (ab

- I was somewhat dissatisfied with the standard Python and R implementations of PCA, particularly when applied to IFUs (data cubes).
- I developed a QR-based PCA package.



PCA Scalability



Full length article

qrpca: A package for fast principal component analysis with GPU acceleration

R. S. de Souza⁺⁺⁺, X. Quanfeng⁺, S. Shen⁺⁺⁺, C. Peng⁺⁺⁺, Z. Mu⁺⁺ ⁺ Exy advances by the feature in a Gamma di Caming, Shingher Advances di Denmary, Dinner Assistrary of Sources, in the ⁺ Part and in American Shingha 2004. Cone

Shorphan Institute of Technology, 100 Manyuov Rd., Sharphan 201418, Otiou-

- It utilizes Torch and Pytorch for GPU acceleration.
- QRPCA behaves similarly to standard implementations in R and python
- It is 10-20 \times faster then sklearn and prcomp



Back to Extragalactic GCs

 After employing our customized pre-processing, including imputation, denoising, proper motion cuts, and a propensity score matching, we compiled an initial list of 640 GC candidates out of 7k sources.



Back to Extragalactic GCs

- The first compilation of extragalactic GCs around the triplet.
- In the figure orange stands for GCs with lower proper motions, while cyans are higher in comparision to the known GC.
- We are systematically performing spectroscopic follow-up, which has borne fruit so far



Back to Extragalactic GCs

 An analysis of their spatial distribution suggests possible evidence for a bridge between M81 and M82, which is currently under investigation

