Significance Mode Analysis for hierarchical structures **Extracting stellar populations from large-scale surveys**

Sebastian Ratzenböck @ Uni Vienna

ML-IAP/CCA-2023



Gaia data





Stellar populations Born from same molecular cloud

- Thought to be birthplace of most stars (Lada & Lada 2003; Parker & Goodwin 2007)
- Structure formation and evolution
- Chemical composition of Milky Way
- Exoplanet formation and evolution
- Stellar initial mass function

Probe for

Pleiades (c) ESO/S. Brunier





ESA Gaia, DPAC - Moitinho+2017

Pleiades (c) ESO/S. Brunier



 Low dimensional feature space 3 positional axes + 2 tangential velocities Stars that move together were born together (Kamdar+2019)



- Low dimensional feature space
- Projection effects in velocities















- Low dimensional feature space
- Projection effects in velocities
- Millions to billions of data points

S

- Low dimensional feature space
- Projection effects in velocities
- Millions to billions of data points
- 95 99% noise

ts

- Low dimensional feature space
- Projection effects in velocities
- Millions to billions of data points
- 95 99% noise
- Wide variety of (non-convex) cluster morphologies

Tidal tails (Meingast+2019a), Streams (Meingast+2019b), Strings (Kounkel+2019), Rings (Cantat-Gaudin+2019), Snakes (Tian+2020), Pearls (Coronado+2021), ...







- Low dimensional feature space
- Projection effects in velocities
- Millions to billions of data points
- 95 99% noise
- Wide variety of (non-convex) cluster morphologies No accurate simulations / forward models



- Low dimensional feature space
- Projection effects in velocities
- Millions to billions of data points
- 95 99% noise
- Wide variety of (non-convex) cluster morphologies No accurate simulations / forward models



Nonparametric, density based clustering

Recap: Density based clustering

• Data set $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, x_i \in \mathbb{R}^p$



- Data set $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, x_i \in \mathbb{R}^p$
- Data generated from density: $X \sim f$





- Data set $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, x_i \in \mathbb{R}^p$
- Data generated from density: $X \sim f$
- Wishart (1969) cluster definition
 - \mathbf{X}_i associated with modes of f
 - Propagate \mathbf{x}_i along ∇f

50

40 30 20 10 10 20 30 40 0 Х



- Data set $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, x_i \in \mathbb{R}^p$
- Data generated from density: $X \sim f$
- Wishart (1969) cluster definition
 - \mathbf{X}_i associated with modes of f
 - Propagate \mathbf{x}_i along ∇f

50

40 30 20 10 10 40 20 30 0 X



- Data set $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, x_i \in \mathbb{R}^p$
- Data generated from density: $X \sim f$
- Wishart (1969) cluster definition
 - \mathbf{X}_i associated with modes of f
 - Propagate \mathbf{x}_i along ∇f



- Data set $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, x_i \in \mathbb{R}^p$
- Data generated from density: $X \sim f$
- Wishart (1969) cluster definition
 - \mathbf{X}_i associated with modes of f
 - Propagate \mathbf{x}_i along ∇f



- Data set $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, x_i \in \mathbb{R}^p$
- Data generated from density: $X \sim f$
- Wishart (1969) cluster definition
 - \mathbf{X}_i associated with modes of f
 - Propagate \mathbf{x}_i along ∇f



- Data set $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, x_i \in \mathbb{R}^p$
- Data generated from density: $X \sim f$
- Wishart (1969) cluster definition
 - \mathbf{X}_i associated with modes of f
 - Propagate \mathbf{x}_i along ∇f





- Data set $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, x_i \in \mathbb{R}^p$
- Data generated from density: $X \sim f$
- Wishart (1969) cluster definition
 - \mathbf{X}_i associated with modes of f
 - Propagate \mathbf{x}_i along ∇f





20

Х

30

- Data set $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, x_i \in \mathbb{R}^p$
- Data generated from density: $X \sim f$
- Wishart (1969) cluster definition
 - \mathbf{X}_i associated with modes of f
 - Propagate \mathbf{x}_i along ∇f

50

40 30 20 10 40 10 20 30 0 X



- Level set: $L(\lambda) = \{f(\mathbf{x}) \ge \lambda\}$
- Hartigan (1975) cluster definition
 - Connected components of $L(\lambda)$
 - Cluster tree: vary $\lambda: \infty \to -\infty$





- Estimate density \hat{f} from data X
- produces spurious clusters



- Estimate density \hat{f} from data X
- produces spurious clusters



Estimate density *f̂* from data *X* ➡ produces spurious clusters





• Estimate density \hat{f} from data X \Rightarrow produces spurious clusters



Pruning cluster tree Current strategies

- Density difference $\Delta \hat{f}$ (Chazal+2013)
- Normalised $\Delta \hat{f}$ (Ding+2016)
- Distance based (Stuetzle+2010; Kpotufe+2011; Chaudhuri+2014)
- Relative excess of mass (HDBSCAN; Campello+2013)

Pruning cluster tree Current strategies

Density difference $\Delta \hat{f}$ (Chazal+2013)

Normalised $\Delta \hat{f}$ (Ding+2016)

Hard to determine threshold for $N \gg 1$

Distance based (Stuetzle+2010; Kpotufe+2011; Chaudhuri+2014)

Relative excess of mass (HDBSCAN; Campello+2013)

Typically over-merges

Going back to Wishart (1969) Clusters are modes of f

What constitutes a cluster?

Clusters are modal regions of fTest for multimodality

What constitutes a cluster?

Clusters are modal regions of *f* → Test for multimodality

 H_0 : Points belong to single mode 20- H_1 : Points belong to multiple modes 10-

Modality along paths





Nultimodality test statistic

$H_0: \nexists t | f(r_t) < \min\{f(r_0), f(r_1)\}$

$T(t) := \min\{\log f(r_0), \log f(r_1)\} - \log f(r_t)\}$

Multimodality test statistic

$H_0: \nexists t \mid f(r_t) < \min\{f(r_0), f(r_1)\}$

$T(t) := \min\{\log f(r_0), \log f(r_1)\} - \log f(r_t)$

 $H_0: T(t) \le 0 \quad \forall t \in (0,1)$



$H_0: T(t) \le 0 \quad \forall t \in (0,1)$









On estimated density?

Let's apply: \hat{f}



$H_0: T(t) \le 0 \quad \forall t \in (0,1)$









 $f \to \hat{f} : T(t) \to \hat{T}(t)$



Multimodality test statistic: $\hat{T}(t)$

$T(t) := \min\{\log f(r_0), \log f(r_1)\} - \log f(r_t)\}$





Multimodality test statistic: $\hat{T}(t)$ $T(t) := \min\{\log f(r_0), \log f(r_1)\} - \log f(r_t)\}$ $\hat{T}(t) := -p \max\{\log d_k(r_0), \log d_k(r_1)\} + p \log d_k(r_t)$

Multimodality test statistic: $\hat{T}(t)$ $T(t) := \min\{\log f(r_0), \log f(r_1)\} - \log f(r_t)$ $\hat{f}(x) \propto \frac{1}{d_k^p(x)}$ k-NN density estimator $\hat{T}(t) := -p \max\{\log d_k(r_0), \log d_k(r_1)\} + p \log d_k(r_t)$

Burman & Polonik (2009) show $H_0: \hat{T}(t) \sim \mathcal{N}(0,1) \times c$



Let's test it!



 $H_1: \hat{T}_n(t) > \Phi^{-1}(1-\alpha)$?







 $H_1: \hat{T}_n(t) > \Phi^{-1}(1-\alpha)$







 $H_0: \hat{T}_n(t) \le \Phi^{-1}(1-\alpha)$?







 $H_0: \hat{T}_n(t) \le \Phi^{-1}(1-\alpha)$





Putting it all together

1. Gradient ascent step — cluster tree





- 1. Gradient ascent step
- 2. Scan saddle points: $\max \hat{f} \rightarrow \min \hat{f}$





- 1. Gradient ascent step
- 2. Scan saddle points: $\max \hat{f} \rightarrow \min \hat{f}$
 - A. Test modality between modes



Min. energy path Arbitrary path Saddle point

- 1. Gradient ascent step
- 2. Scan saddle points: $\max \hat{f} \rightarrow \min \hat{f}$
 - A. Test modality between modes
 - B. If H_0 cannot be rejected merge



Min. energy path --- Arbitrary path Saddle point

- Gradient ascent step
- 2. Scan saddle points: $\max \hat{f} \rightarrow \min \hat{f}$
 - A. Test modality between modes
 - B. If H_0 cannot be rejected merge





Next saddle point



- 1. Gradient ascent step



How to set parameters?

 $SigMA(k, \alpha)$

Choosing k $\hat{T}_n(t) \sim \mathcal{N}(0,1) \iff \log N < k < N^{4/(4+p)}$





Choosing a

- Many hypotheses tests increases chance of false positives
- Limit proportion of false positives among all positives
 - Apply Benjamini & Hochberg procedure
 - \blacktriangleright Data driven way of choosing significance α

Results on Sco-Cen







30°

20°

10°

10°

-20

Consists of 37 groups • Unseen substructure Validated via • narrow HRD

• B stars in center





40°

30°

20°

10°

0°

-10°

-20

- High spatial resolution age map
 - Investigate SF history







Backup

I me complexity

mode & saddle Density computation search (union find) (k-d tree) $\mathcal{O}(pN\log N) + \mathcal{O}(pN\log N) + \mathcal{O}(Nk) + \mathcal{O}(|\mathcal{S}|)$ Graph construction Cluster tree prunina



Robustness of $\hat{T}_n(t)$

Graph

N

k

 β -Skeleton

Feature scaling



-2 -1 0 1 2 -2 -1 0 1 2 -2 -1 0 1 2

- Consists of 37 groups
- Unseen substructure
- Validated via
- narrow HRD
- B stars in center
- Age gradients



Background reduction



 $\rho \rightarrow$

Background reduction



 $\rho \rightarrow$



Likelihood of samples under model k-NN density estimation

