

Globular cluster detection in the Gaia survey

M. Mohammadi¹, R. F. Peletier², F.M. Schleif³, N. Petkov¹, and K. Bunte¹ *

1- Johann Bernoulli Institute, University of Groningen

2- Kapteyn Institute, University of Groningen

3- Univ. of Appl. Sc. Würzburg-Schweinfurt, Dept. of CS, Würzburg, Germany

Abstract. Existing algorithms for the detection of stellar structures in the Milky Way are most efficient when full phase-space and color information is available. This is rarely the case. Since recently, the Gaia satellite surveys the whole sky and is providing highly accurate positions for more than one billion sources. In this contribution we propose two independent strategies to find globular clusters in this database, based on magnitude distributions only. One approach is a nearest neighbor retrieval and the other an anomaly detection. Both techniques are able to find known globular clusters within our observation frame consistently, as well as additional candidates for further investigation.

1 Introduction

A major interest in astronomy is finding interesting stellar structures in the sky, that are present in large modern astronomical databases, but are hidden by structures in front, like for example dust. A globular cluster (GC) is a spherically shaped collection of stars bound by gravity with high stellar densities toward its center (see e.g. Fig. 1 upper left corner). It is assumed that its stars are formed almost at the same time, and therefore being of similar age and chemical compositions, making them effective tracers of Galaxy formation and evolution. The Λ *cold dark matter model* suggests that big systems like galaxies are built by merging smaller systems. GC and dwarf galaxies are considered the building blocks for such merging events. Moreover, the dynamical interactions of GCs with the galactic potential inferred from effects on stellar density and velocity profiles, can be used to study the mass distribution of the Galaxies[1].

A few algorithms like ROCKSTAR[2] and OPTICS[3] identify substructures, such as GCs, from full phase-space information. However, if only the position and magnitudes of the stars on the sky and color information is available often matched filter methods are applied[4]. In this paper we focus on the detection of GCs using the DR1 catalog collected by the Gaia satellite. Gaia was launched at the end of 2013 to survey the whole sky and determine highly accurate positions, parallaxes, and proper motions for more than one billion sources brighter than magnitude 20.7 in the white-light photometric band G[5]. However, neither the full phase-space information nor the color information of the stars is available in the Gaia DR1 catalog, and therefore above mentioned methods cannot be used. In this contribution we investigate the detection of GCs in the Gaia survey in the region $120^\circ < \alpha < 246^\circ$ right ascension (RA) and $-2^\circ < \delta < 60^\circ$ declination (Dec), respectively. We follow two strategies: a retrieval technique using known galaxy clusters to detect similar structures and a novelty detection technique.

*This work was supported by the European H2020-MSCA-ITN *Survey Network for Deep Imaging Analysis and Learning* (SUNDIAL), project ID 721463.

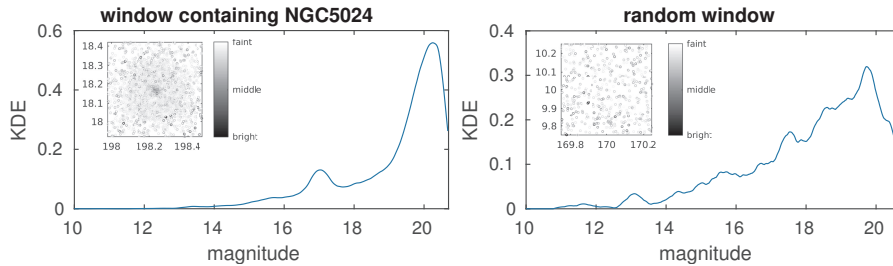


Fig. 1: Left: Magnitude distribution of a window of 0.5 degree centering at globular cluster NGC5024. Right: Magnitude distribution of a random window.

2 Data and Preprocessing

To find globular clusters we investigate the magnitude distribution of the stars contained in rectangular windows of a certain size. As proposed in [4] we consider windows of 0.5 degree length in both right ascension (RA) and declination (Dec). The left panel in Fig. 1 shows the magnitude distribution as estimated via Kernel Density Estimation (KDE) based on the magnitudes extracted from a window centering at position (198.23 RA, 18.17 Dec) of globular cluster NGC5024 in Gaia. In contrast to random windows (Fig. 1 right panel) the distribution of a window containing a globular cluster will exhibit a bump in the luminosity function. For NGC5024 we see a clear bump at a characteristic magnitude 17, which comes from the globular cluster and a sharp increase as soon as fainter stars are included. The latter indicates magnitude limit of 20.7 on the magnitude of the visible stars in the Gaia survey. This bump consists of stars on the horizontal branch, which happen to all have roughly the same absolute magnitude of about $V = 0.5$. Depending on the distance of the cluster, the bump will shift in magnitude, which can be described by the distance modulus formula:

$$m - M = 5 \log_{10} d - 5, \quad (1)$$

where m , M and d are apparent magnitude, absolute magnitude and the distance (in pc), respectively. Here the apparent magnitude indicate how bright a star appears to an observer, and the absolute magnitude is its intrinsic brightness.

To approximate the magnitude distribution of each window, we use KDE for magnitude values (x_1, \dots, x_n) :

$$f_h(x) = \frac{1}{n} \sum_{i=1}^n K_h(x - x_i) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right) \quad (2)$$

with bandwidth h and the Epanechnikov function as kernel K .

3 Methods

To compare magnitude densities extracted from a window we regard each magnitude distribution as temporal sequence and deploy Dynamic Time Warping (DTW) [6] to measure their similarity.

DTW aligns two sequences non-linearly along the temporal axis (see Fig. 2) leading to a distance like measure which, however, does not fulfill all metric properties. This measure has the potential to detect similar structures in the densities, such as for bumps, despite their exact location on the magnitude axis.

3.1 Nearest Neighbor Retrieval

One method to find similar objects of interest in a data base is Nearest Neighbor Search (NNS). Given a query a nearest neighbor retrieval system returns a desired number of most similar objects from the data base according to the distance measure used. Following this strategy we aim to find interesting GCs in the Gaia survey by retrieving similar densities to those extracted from the known globular clusters using DTW as similarity measure. We compare the extracted distribution with a k nearest neighbors of known structures, by counting the number of common GCs for each window. The higher the obtained count the more likely it is that there is something interesting in it. This technique can be easily extended to deal with even larger data sets by using approximate nearest neighbor search approaches, as for example locality-sensitive hashing[8], best bin first[9] and balanced box-decomposition tree based search[10].

3.2 Novelty Detection

Since the sky is mostly empty we can regard the search for interesting structures as novelty or outlier detection. We use an approach which was inspired by the Support Vector Classifier called Support Vector Data Description (SVDD) [11]. Assume given data $\mathbf{x}_i \in \mathbb{R}^n$ and a fixed chosen kernel $K : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ which is associated to the feature map $\Phi : K(\mathbf{x}, \mathbf{y}) = \Phi(\mathbf{x})^t \Phi(\mathbf{y})$. The goal of SVDD is to find a generalized linear mapping $\mathbf{x} \mapsto \text{sgn}(f(\mathbf{x})) = \text{sgn}(\mathbf{w}^t \Phi(\mathbf{x}) - \rho)$, which defines a separation of the given data to outliers by means of its sign. Where the separation boundary corresponds to a linear separation in the feature space induced by Φ . The problem to find suitable parameters \mathbf{w} and ρ for a given data set of typical points \mathbf{x}_i can be formalized as an optimization problem which aims at a separation of the given data from the origin with maximum margin. This leads to the following primal optimization problem SVDD (primal): $\min_{\mathbf{w}, \rho, \xi_i} \frac{1}{2} \cdot \|\mathbf{w}\|^2 - \rho + \frac{C}{2} \cdot \sum_i \xi_i^2$, such that $\mathbf{w}^t \Phi(\mathbf{x}_i) \geq \rho - \xi_i \quad \forall i$, where $C > 0$ is a fixed constant, and the parameters ξ_i refer to the slack variables to allow for some error tolerance. Using the Karush-Kuhn-Tucker conditions, the Lagrange dual problem becomes

$$\begin{aligned} \text{SVDD(dual)} \\ \max_{\alpha_i} & -\frac{1}{2} \cdot \sum_{ij} K(\mathbf{x}_i, \mathbf{x}_j) \alpha_i \alpha_j + \frac{1}{C} \cdot \sum_i \alpha_i^2 \\ \text{such that} & \alpha_i \geq 0 \quad \forall i \\ & \sum_i \alpha_i = 1 \end{aligned}$$

This dual problem can be directly optimized relying on linearly constraint convex quadratic optimization. The solution \mathbf{w} and ρ of the primal problem can then

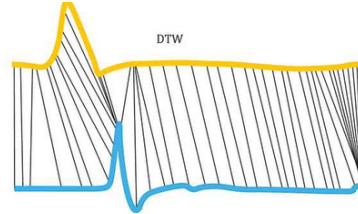


Fig. 2: DTW dissimilarity [7].

be recovered from the dual variables α_i . These are non-vanishing for support vectors only and hence result in a sparse description.

The tightness of the data description can be controlled by setting an error tolerance or using outlier examples one desires to detect. However, for this approach we require a positive definite kernel also known as Mercer kernel [12], which we cannot guarantee if the DTW dissimilarity is used due to missing metric properties. Therefore, we approximate a positive definite kernel (psd) based on the DTW similarity following the strategy discussed in [6] using the *square approach*. *Square* changes the eigenspectrum of a given similarity matrix such that all eigenvalues are non-negative and a psd kernel is obtained. Therefore we follow two steps to convert the DTW pairwise dissimilarity matrix D between magnitude distributions into an appropriate kernel as input to the SVDD: 1) compute similarity matrix S from D via double centering and 2) convert S to a psd kernel by $K = S \cdot S^\top$.

4 Experiments

To assess the different strategies for the detection of globular clusters we investigate the part of the Gaia survey ranging from $120^\circ < RA < 246^\circ$ and $-2^\circ < Dec < 60^\circ$ respectively. We extract windows of 0.5 degree bin width from this part of the sky and approximate the magnitude distributions using KDE with bandwidth $h = 0.15$. There are seven known globular clusters at position (α, δ) for our investigation [13]: 1) NGC4147 (182.53, 18.54), 2) NGC 5024 (198.23, 18.17), 3) NGC5053 (199.11, 17.7), 4) NGC 5272 (205.55, 28.38), 5) NGC5466 (211.36, 28.53), 6) NGC5904 (229.02, -0.11) and 7) Palomar5 (229.64, 2.08). Since we have a limited magnitude of 20.7 of stars in the Gaia DR1 survey, we do not include the more distant known GCs.

Nearest Neighbor Retrieval: For the nearest neighbor retrieval we do not only use the windows centered at the 7 known GCs as example structures as mentioned before, but we also include sliding windows around it which are shifted horizontally and/or vertically by 0.25 degree. The globular cluster will still at least partly be contained in those. All these examples are then compared with DTW to windows of equal size extracted by sliding with 0.25 degree through the part of the sky under observation. Using the counting strategy discussed above the promising candidates are identified.

Novelty Detection: For the outlier detection we build 9 psd kernels for cross validation based on the DTW dissimilarities as mentioned before. As a proof of concept we build kernels of 10000 samples each containing: a) a training subset of 8000 random windows taken from the investigated part of Gaia and b) a testing subset containing the windows extracted around the known GCs, potential candidates from the retrieval and ca. 2000 random windows. Therefore, we train 9 SVDD models on different random subsets. To evaluate our method, we use the testing set (b) and report the detection rate of known GCs as well as overlap with the retrieval method. In our experiments we investigate the threshold ρ for the decision boundary and therefore the tightness of the bound by an error tolerance E .

5 Results

The summary of the results is given in Table 1. Both methods retrieve an increasing number of potential candidates with increasing their corresponding hyper-parameter. The retrieval technique is more specific with fewer number of potential candidates in comparison to the outlier detection. This was expected since it uses examples of the structure we are looking for. However, Palomar5 or NGC5904 are more difficult. They can be detected, if we increase the hyper-parameters. The average vote for $k = 20$ and corresponding candidates are shown in Fig. 3. Windows with more than one fourth of the maximum vote are clustered based on Ra and Dec and hits closer than 0.75° denote the same candidate. In the outlier detection we transform the de-

Table 1: Mean hit number (std) for different hyper-parameter in retrieval and outlier detection.

HP	Candidates	GCs found
k 20	23.43 (12.14)	6.71 (0.49)
k 50	35.57 (11.86)	6.71 (0.49)
k 100	77.57 (5.56)	7.00 (0.00)
E 1%	22.67 (7.37)	5.00 (0.00)
E 5%	107.89 (31.21)	5.00 (0.00)
E 10%	139.22 (47.96)	5.11 (0.33)
E 20%	222.33 (68.01)	5.33 (0.50)
E 30%	256.44 (41.22)	5.11 (0.33)
E 40%	377.33 (118.97)	5.44 (0.53)
E 50%	374.44 (107.56)	6.00 (0.00)

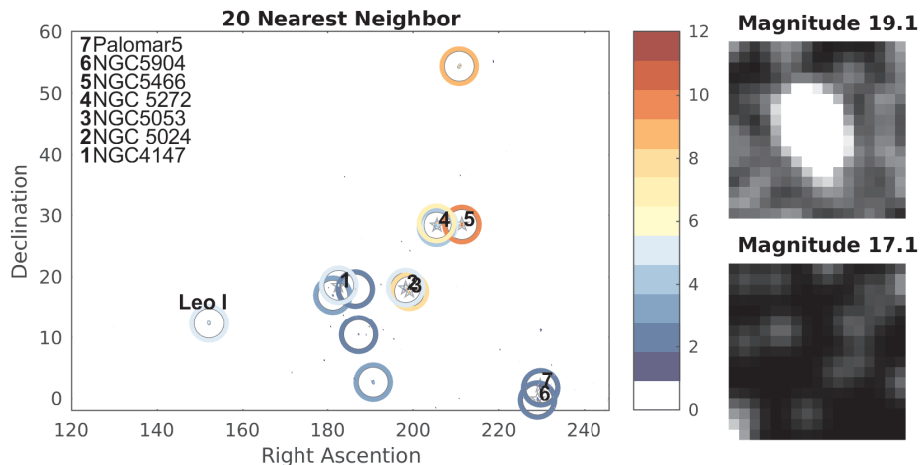


Fig. 3: Average vote for leave-one-out cross-validation using 20 NN retrieval. The circles mark interesting regions with the color indicating the similarity to known GCs. Numbers and text flag positions of known GCs. On the right we also include the response of spatial Gaussian filters at magnitude 19.1 and 17.1 respectively. The former shows clearly the globular stellar structure.

viation from the error tolerance threshold for every hit and model to achieve a score. Every score which is smaller than 80% of the minimum score is clustered as in the retrieval and considered as candidate. The outlier models were trained solely on random windows and not provided any example of the known structures but still detect between 5 and 6 of the 7 known GCs used for testing. Furthermore the outliers overlap with candidates detected using the retrieval method. To determine the exact position of a GC and to reduce false positives

one can use circular Gaussian filters [14] to judge the candidates. The right side of Fig. 3 shows Gaussian filter responses to the window containing Leo I at different magnitude. There is a high response in the center (i.e. a circular shaped group of stars) at magnitude around 19.1, but a weak response (i.e. uniformly distributed stars) at magnitude around 17.1.

6 Conclusion and Outlook

In this contribution we propose two independent strategies to find globular clusters in the Gaia survey. Potentially interesting structures are found by a nearest neighbor retrieval method comparing examples extracted from 7 known GCs and the support vector data description, a kernel-based anomaly detection technique. Both techniques find most of the 7 known GCs used for testing and also depict overlapping agreement on potential other candidates. The exact position of GCs as well as false positives can be determined by Gaussian Filters on the candidate windows. Future work include the detailed investigation of the promising candidates and efficient out-of-sample extension.

Acknowledgement: We thank the Center for Information Technology of the University of Groningen for providing access to the Peregrine high performance computing cluster.

References

- [1] CM. Rockosi, M. Odenkirchen, EK. Grebel, W. Dehnen, KM. Cudworth, JE. Gunn, DG. York, J. Brinkmann, GS. Hennessy, and Z. Ivezić. A matched-filter analysis of the tidal tails of the globular cluster palomar 5. *The Astronomical Journ.*, 124(1):349, 2002.
- [2] PS. Behroozi, RH. Wechsler, and H-Y. Wu. The rockstar phase-space temporal halo finder and the velocity offsets of cluster cores. *The Astrophysical Journ.*, 762(2):109, 2012.
- [3] SA. Sans Fuentes, J. De Ridder, and J. Debosscher. Stellar halo hierarchical density structure identification using (f) optics. *Astronomy & Astrophysics*, 599:A143, 2017.
- [4] H. Tian. *Hide and Seek in the Halo of the MW*. PhD thesis, Univ. of Groningen, 2017.
- [5] AGA. Brown, A. Vallenari, T. Prusti, JHJ. De Bruijne, F. Mignard, R. Drimmel, C. Babusaix, CAL. Bailer-Jones, U. Bastian, et al. Gaia data release 1-summary of the astrometric, photometric, and survey properties. *Astronomy & Astrophysics*, 595:A2, 2016.
- [6] F-M. Schleif and P. Tiño. Indefinite proximity learning: A review. *Neural Computation*, 27(10):2039–2096, 2015.
- [7] P. Ranacher and K. Tzavella. How to compare movement? a review of physical movement similarity measures in geographic information science and beyond. *Cartography and geographic information science*, 41(3):286–307, 2014.
- [8] J. Leskovec, A. Rajaraman, and JD. Ullman. *Mining of massive datasets, Ch. 3*. Cambridge university press, 2014.
- [9] JS. Beis and DG. Lowe. Shape indexing using approximate nearest-neighbour search in high-dimensional spaces. In *Computer Vision and Pattern Recognition, 1997. Proceedings., 1997 IEEE Computer Society Conference on*, pages 1000–1006. IEEE, 1997.
- [10] S. Arya, DM. Mount, NS. Netanyahu, R. Silverman, and AY. Wu. An optimal algorithm for approximate nearest neighbor searching fixed dimensions. *Journal of the ACM (JACM)*, 45(6):891–923, 1998.
- [11] DMJ. Tax and RPW. Duin. Support vector domain description. *Pattern Recognition Letters*, 20(11):1191 – 1199, 1999.
- [12] E. Pekalska and RPW. Duin. *The dissimilarity representation for pattern recognition: foundations and applications, Ch. 2*. World Scientific, 2005.
- [13] WE. Harris. A catalog of parameters for globular clusters in the milky way. *The Astronomical Journal*, 112:1487, 1996.
- [14] B. Gecer, G. Azzopardi, and N. Petkov. Color-blob-based cosfire filters for object recognition. *Image and Vision Computing*, 57:165–174, 2017.