

## H $\alpha$ emitters from the Southern Photometric Local Universe Survey (S-PLUS)

L. A. Gutiérrez-Soto<sup>1\*</sup>, R. Lopes de Oliveira<sup>1,2,3</sup>, S. Akras<sup>4</sup>, D. R. Gonçalves<sup>5</sup>, C. Mendes de Oliveira<sup>1</sup>, F. Almeida-Fernandes<sup>1,6</sup>, F. R. Herpich<sup>1</sup>, C. Cheng<sup>7</sup>, T. S. Gonçalves<sup>5</sup>, L. Nakazono<sup>1</sup>, E. Telles<sup>3</sup>, A. Alvarez-Candal<sup>3,8,9</sup>, A. Kanaan<sup>10</sup>, T. Ribeiro<sup>11</sup>, W. Schoenell<sup>12</sup>

<sup>1</sup>Departamento de Astronomia, IAG, Universidade de São Paulo, Rua do Matão, 1226, 05509-900, São Paulo, Brazil

<sup>2</sup>Departamento de Física, Universidade Federal de Sergipe, Av. Marechal Rondon, S/N, 49100-000, São Cristóvão, SE, Brazil

<sup>3</sup>Observatório Nacional, Rue Gal. José Cristino 77, 20921-400, Rio de Janeiro, RJ, Brazil

<sup>4</sup>Institute for Astronomy, Astrophysics, Space Application & Remote Sensing, National Observatory Athens, GR-15236, Athens, Greece

<sup>5</sup>Observatório do Valongo, Universidade Federal do Rio de Janeiro, Ladeira Pedro Antonio 43, 20080-090, Rio de Janeiro, Brazil

<sup>6</sup>Community Science and Data Center/NSF NOIRLab, 950 N. Cherry Ave., Tucson, AZ 85719, USA

<sup>7</sup>Chinese Academy of Sciences South America Center for Astronomy, National Astronomical Observatories, CAS, Beijing 100101, China

<sup>8</sup>Instituto de Astrofísica de Andalucía, CSIC, Ap. 3004, E18080 Granada, Spain

<sup>9</sup>Instituto de Física Aplicada a las Ciencias y las Tecnologías, Universidad de Alicante, San Vicent del Raspeig, E03080, Alicante, Spain

<sup>10</sup>Departamento de Física, Universidade Federal de Santa Catarina, Florianópolis, SC, 88040-900, Brazil

<sup>11</sup>NOAO, P.O. Box 26732, Tucson, AZ 85726

<sup>12</sup>GMTO Corporation 465 N. Halsted Street, Suite 250 Pasadena, CA 91107

Accepted XXX. Received YYY; in original form ZZZ

### ABSTRACT

The ongoing multi-band survey performed by the S-PLUS project will have covered 9300 deg<sup>2</sup> of the Southern skies by the time it is completed. S-PLUS has a crucial feature: images over the whole field taken in the H $\alpha$  narrow-band. The H $\alpha$  transition provides a superb tool for the study of a number of important astrophysical processes and, in particular, it allows the classification of different types of astrophysical sources. Here we explore the S-PLUS data release 3, which covers 2000 deg<sup>2</sup>, including the Stripe-82 area, to highlight the potential of the survey for finding H $\alpha$  emitters using the (r - J0660) versus (r - i) colour-colour diagram and in distinguishing the red from the blue sources based on the (r - z) versus (g - r) diagram. Our H $\alpha$ -emitter catalog contains 3,187 objects that exhibit excess in the narrow J0660 band. For 2,869 of them (or 90%), the excess is thought to be due to the H $\alpha$  emission line, while for the remaining, the excess may be due to redshifted lines. Unsupervised clustering machine learning approach reveals two distinct populations: one with an intense blue continuum and another one with an intense red one. The hierarchical agglomerative clustering algorithm (HAC) was compared with the hierarchical density-based cluster selection (HDBSCAN) in order to reinforce the robustness of the red and blue populations' classification. By adopting a so-called "soft" clustering approach, we assigned the probability of each emitter belonging to a given population, blue or red. Around 87% of emitters were successfully classified as blue or red sources. We use synthetic and observed spectra to emphasize the potential of colour-colour diagrams in distinguishing several classes of H $\alpha$  emission-line emitters that include planetary nebulae, H II regions, young stellar objects, symbiotic stellar systems, cataclysmic variables, blue compact galaxies, star-forming galaxies, and quasars. In summary, the method described in detail in this paper is shown to be an efficient tool to find new emitters and to classify them, using multi-colour data.

**Key words:** surveys – techniques: photometric – stars: novae, cataclysmic variables – galaxies: dwarf – quasars: emission lines

### 1 INTRODUCTION

Atomic excitation followed by recombination in Balmer hydrogen emission lines may be ignited in different ways, thermal and non-thermal collisional excitation in shock-heated gas and energetic pho-

tons acting over a diffuse gas. As a practical result, and the Universe being hydrogen abundant, the observation of those electronic transitions offer an important window into the study of astrophysical objects. Among all the possible electronic transitions, the Balmer series represent extremely useful tools in Astronomy. Particularly, the H $\alpha$  emission line – rest-frame wavelength of 6564.614 Å at vacuum – that corresponds to the electron transition from the  $n = 3$  to the  $n = 2$

\* E-mail: gsoto.angel@gmail.com

## Summary of comments: S\_PLUS\_Emission\_line\_objects\_Eduv2.pdf

No Comments.

energy level, is the strongest one, in both emission or absorption, and the most widely used to identify various types of objects (e.g. star-forming regions, H II regions, planetary nebulae (PNe), supernovae, novae, young stellar objects (YSO), Herbig-Haro objects, circumstellar disks, post-asymptotic and asymptotic giant stars (AGB), red giant stars (RGB), active late-type dwarfs). Amongst massive stars, emission lines are observed in Be stars with decretion disks, Wolf-Rayet (WR) stars, interacting binary systems that experiencing mass exchange like symbiotic stars (SySt), cataclysmic Variables (CVs), among others.

At much larger scales, H $\alpha$  emission is also detected in extended no-point sources like PNe, H II regions, supernova remnants, as well as star-forming regions in galaxies, among others. In the case of high-redshifted sources like starburst galaxies and quasi-stellar object (QSOs), the detection of an emission at 6563 Å is not associated with the recombination of H $\alpha$  but with other UV emission lines.

Most of the aforementioned classes of objects are not homogeneous and far from complete even in the local Universe, with some being highly populated while others being highly underrepresented. For example, there are  $\sim 320$  known SySt, with only  $\sim 65$  of those located in galaxies other than the Milky Way (Akras et al. 2019a; Merc et al. 2019). The number of known PNe in our Galaxy is of the order of  $\sim 3500$  (Parker et al. 2016), which may represent only 15–30% of the total population (Frew 2008; Jacoby 2010).

H $\alpha$  surveys in a variety of angular resolutions, sky coverage, and sensitivity were carried out in the past. Some of them, with modest spatial resolutions, revealed spatially resolved, extended nebular emission to study supernova remnants, galaxy groups, and star-forming regions (e.g. Davies et al. 1976). Others with higher spatial resolution disclosed compact emission-line sources in the Milky Way and nearby galaxies. Examples are the INT Photometric H $\alpha$  survey (IPHAS; Drew et al. 2005; Barentsen et al. 2014), the SuperCOSMOS H $\alpha$  Survey with the UK Schmidt Telescope (UKST) of the Anglo-Australian Observatory (Parker et al. 2005), and the VST Photometric H $\alpha$  Survey (VPHAS+; Drew et al. 2014).

Colour-colour diagrams from photometric surveys are also used to identify possible H $\alpha$  emitters. For example, the  $(r - H\alpha)$  versus  $(r - i)$  colour-colour diagram, and similar diagrams, has been used to find CVs (Witham et al. 2006, 2007), YSOs (Vink et al. 2008), SySt (Corradi et al. 2008; Corradi & Giannamanco 2010; Corradi et al. 2011; Akras et al. 2019b), early-type emission-line stars (Drew et al. 2008), and PNe (Viironen et al. 2009; Sabin et al. 2010; Akras et al. 2019c).

There are two ongoing multi-band surveys observing the sky in a systematic, complementary way, with 5 broad and 7 narrow-band filters, including H $\alpha$ : the Javalambre Photometric Local Universe Survey (J-PLUS)<sup>1</sup>; Cenarro et al. 2019), covering the Northern celestial hemisphere, and the Southern-Photometric Local Universe Survey (S-PLUS<sup>2</sup>; Mendes de Oliveira et al. 2019), covering the southern sky with a twin 83 cm telescope and filter system. These are paving the way for an even more ambitious survey, the Javalambre Physics of the Accelerating Universe Astrophysical Survey (J-PAS; Benítez et al. 2014 and mini-JPAS; Bonoli et al. 2021), which will observe the Northern sky with 56 narrow-band filters. As source hunters, the spectral energy distribution<sup>3</sup> provided by these surveys enable an unprecedented source classification using photometry only. However, in the Big Data era, efficient investigation tools are required to deal with their massive imaging and catalogues production and machine

learning techniques have been increasingly used to explore these data sets.

Here we present a census of H $\alpha$  emitters from the S-PLUS DR3 by employing the  $(r - J0660)$  versus  $(r - i)$  colour-colour diagram and unsupervised machine learning techniques to classify them as blue or red sources. Section 2 describes the observations related to the S-PLUS project, as well as important information on the third data release. It also presents the technique implemented to select the H $\alpha$  emitters and machine learning approaches used to divide the sample into two populations based on their colours. In section 3 our findings are described and finally section 4 discusses our main results and conclusions.

## 2 METHODOLOGY

### 2.1 Observations: the S-PLUS project

This paper uses data from the S-PLUS DR3 (Buzzo et al., in prep), which covers 2,000 square degrees. The S-PLUS DR3 can be accessed in the database of the project, S-PLUS Cloud<sup>3</sup>. S-PLUS is being carried out by a dedicated 0.83m robotic telescope located at Cerro Tololo, Chile (Mendes de Oliveira et al. 2019). The project is surveying the southern sky using the 12 filters from the so-called Javalambre filter system (Marín-Franch et al. 2012), that spans the wavelength range from 3000 Å to 10000 Å. The system includes seven narrow-band filters ( $J0378$ ,  $J0395$ ,  $J0410$ ,  $J0430$ ,  $J0515$ ,  $J0660$ , and  $J0861$ ) and five broad-band Sloan-like (Fukugita et al. 1996) filters (see Fig. 1). The narrow-band  $J0660$  filter used in S-PLUS is centered at lambda 6614 Å and has a width of about 147 Å (Table 2 of Mendes de Oliveira et al. 2019), and therefore it covers both the H $\alpha$  and the doublet [N II]  $\lambda\lambda$ 6548, 6584 spectral lines for sources up to a redshift of approximately 0.02.

The data set used for this study, DR3, includes about 60 million objects distributed over  $\sim 2,000 \text{ deg}^2$  (of the total of  $\sim 8,000 \text{ deg}^2$  of high Galactic latitudes fields with  $b > 30^\circ$  planned to be covered when the survey is complete). The galactic disk and bulge are not included in DR3 despite S-PLUS plans to cover an area  $\sim 1,300 \text{ deg}^2$  of them and will be available in DR4. Amongst the different aperture photometer available in the catalog, the  $P_{\text{total}}$  photometry is used, which is a 3-arcsec aperture corrected magnitudes (Almeida-Fernandes et al. 2022). In order to ensure that high-quality data are used in the present analysis, only objects detected in at least the  $r$ ,  $i$  and  $J0660$  bands, simultaneously, with errors less than 0.2 mag, are considered. We also selected objects with the probability of being point sources (with  $\text{CLASS\_STAR} > 0.5$ ). Following Almeida-Fernandes et al. (2022), we implemented  $\text{PhotoFlag} = 0$  in the filters  $r$ ,  $J0660$  and  $i$  for the selection of targets with good photometry in these three filters.

The first goal of this paper is the identification of H $\alpha$  emitters in the S-PLUS DR3. For this, we applied an iterative and automatic technique to select objects with an excess in the  $J0660$  band, which is consistent with the detection of the H $\alpha$  line in emission. Next, the sample of H $\alpha$  sources is divided into two subgroups: the blue and red one. This classification was made by employing optical colours in combination with unsupervised machine learning/statistical tools. These procedures are described in the following subsections.

## Page:2

Author: josee Subject: Note Date: 2022-08-08 10:30:41  
not sure what you mean. Can you elaborate?

Author: josee Subject: Squiggly Underline Date: 2022-08-08 10:30:04

Author: josee Subject: Note Date: 2022-08-08 10:32:11  
in our Galaxy.

Author: josee Subject: Squiggly Underline Date: 2022-08-08 10:31:58

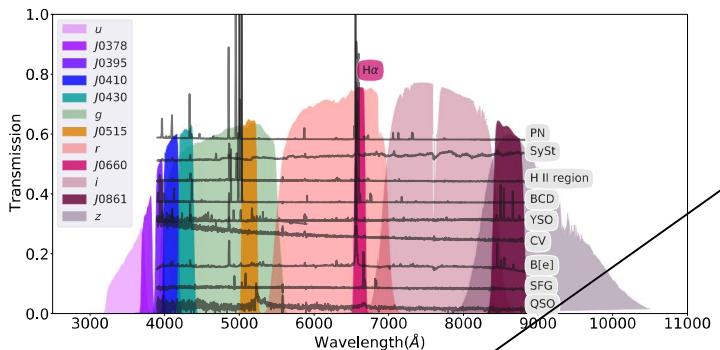
Author: josee Subject: Note Date: 2022-08-08 10:37:16  
or the spectral energy distribution of the detected sources provided by the filter system of these surveys

Author: josee Subject: Replace Text Date: 2022-08-08 10:35:08  
you mean  
wavelength coverage?

<sup>1</sup> <https://www.j-plus.es>

<sup>2</sup> <http://www.splus.iag.usp.br>

<sup>3</sup> <https://splus.cloud/>



**Figure 1.** Transmission curves of the S-PLUS filter set. The narrow-band filter  $J_{0660}$  includes the H $\alpha$  emission line. Over-plotted are spectra of different classes of emission line objects. From top to bottom: a PN, a symbiotic star, an extragalactic H II region, a blue compact/H II galaxy, a YSO, a CV star, a B[e] star, a star forming galaxy and a QSO at a redshift of  $\sim 3.31$ .

## 2.2 Selection of H $\alpha$ emitters

Before search for the H $\alpha$  emitters, potential hidden in the S-PLUS DR3 footprint, we first divided our sample into four sub-samples based on their magnitudes in the  $r$ -bands: (i)  $r$ -band  $< 16$ , (ii)  $16 \leq r < 18$ , (iii)  $18 \leq r < 20$ , and (iv)  $20 \leq r < 21$ . In this way, we avoided mixing up bright and faint sources with low and high uncertainties, respectively. Otherwise, the selection criteria could be affected by the intrinsic scatter in the measurement of faint objects.

The identification of H $\alpha$  emitters is based on the method successfully applied by [Witham et al. \(2008\)](#) to the IPHAS project, given that similar filters are also available in S-PLUS:  $r$ ,  $J_{0660}$ , and  $i$  filters. The same technique was also used by [Scaringi et al. \(2013\)](#) and [Wevers et al. \(2017\)](#) to reveal H $\alpha$  emitters.

We first generated the  $(r - J_{0660})$  versus  $(r - i)$  diagram for each sub-sample and attempted to fit the loci mainly occupied by main-sequence and giant stars with a linear regression. We then implemented an iterative  $\sigma$ -clipping technique so that, by construction, H $\alpha$  emitters should satisfy the condition:

$$(r - J_{0660})_{\text{obs}} - (r - J_{0660})_{\text{fit}} \geq C \times \sqrt{\sigma_s^2 + \sigma_{\text{phot}}^2} \quad (1)$$

where  $\sigma_s$  is the root mean squared value of the residuals around the fit and  $\sigma_{\text{phot}}$  is the error on the observed  $(r - J_{0660})$  colour index.  $C$  is a constant parameter with a value of 4, following [Wevers et al. \(2017\)](#). The fits were made by employing `astropy.modeling`<sup>4</sup>.

Fig. 2 illustrates the procedure applied. The solid black lines indicate the initial fit and the dashed lines show the 4- $\sigma$  clipping fit. The dotted lines correspond to the selection criteria for the H $\alpha$  emitters, 4- $\sigma$  above of the final fit. It should be noted that these cut-off lines are just approximations, as they only represent the residual around the fit. The photometric uncertainty of the  $(r - J_{0660})$  colour index for each individual point is also taken into account (see Equation 1).

Once the list of H $\alpha$  emitters was obtained, we proceeded with a visual inspection of their false-colour images and their spectral energy distributions, constructed with 12 points corresponding to

<sup>4</sup> <https://docs.astropy.org/en/stable/modeling/index.html>

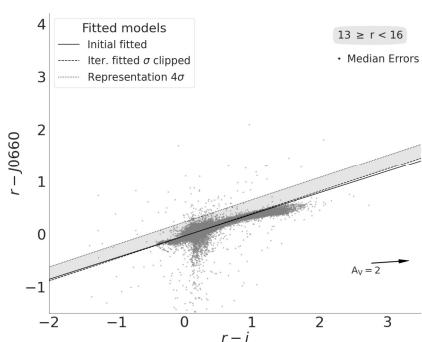
the 12 S-PLUS filter mean magnitudes for each source, hereafter called the S-spectra, to remove artefacts from the list. We repeat the methodology explained above without considering `PhotoFlag = 0`, to recover objects with good photometry but that were labeled with a flag different from 0. The upper panel of Fig. 3 shows an example of what the S-spectrum of an H $\alpha$  emitter looks like, while the bottom panel presents the SDSS spectrum of the same source. It is evident from the comparison of the two spectra that the excess in the  $J_{0660}$  band is linked with the H $\alpha$  emission line. The spectroscopic redshift (SDSS) of this object is 0.009 confirming that the excess in the  $J_{0660}$ -band is due to the H $\alpha$  emission line.

The distribution of the H $\alpha$  emitters in the  $(r - J_{0660})$  versus  $(r - i)$  colour-colour plane is shown in Fig. 4. The loci of the main-sequence and giant stars derived from synthetic spectra ([Pickles 1998](#)) convolved with the transmission of the S-PLUS filters in the AB magnitude system ([Oke & Gunn 1983](#)) are also plotted. All sources located above the locus of the main and giant stars exhibit an excess in  $J_{0660}$  filter and it is attributed to H $\alpha$  line. The wide distribution of sources across the  $(r - J_{0660})$  and  $(r - i)$  colour-colour diagram indicates that several types of H $\alpha$  emitters are selected. Sources with high  $(r - J_{0660})$  colour index are likely associated with PNe, H II regions, SySt or blue compact galaxies. On the other hand, the  $(r - i)$  colour index indicates redder sources such as SySt and YSO, while sources with strong blue continuum such as CVs and QSOs exhibit lower  $(r - i)$  values (see Fig. 2 of [Gutiérrez-Soto et al. 2020](#)).

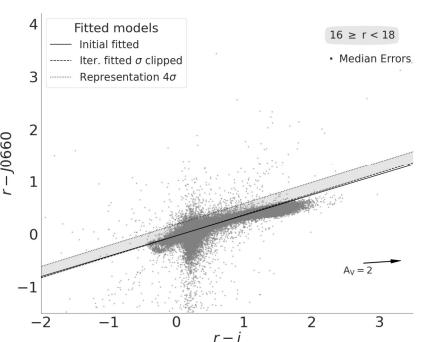
Fig. 5 displays the distribution of all H $\alpha$  emitters in Galactic latitude and longitude. The density map regions represent the spatial positions of the objects on the sky. The surface density of  $J_{0660}$ -excess objects is highest near the Galactic plane. In fact, the distribution in function of the latitude present two peaks, one at  $\sim 16^\circ$ , which corresponds to the Stripe-82 region and another at  $\sim 41^\circ$  (see inset figure).

Our list of H $\alpha$  emitters includes 3,187 sources. We now proceed to their classification into blue and red populations.

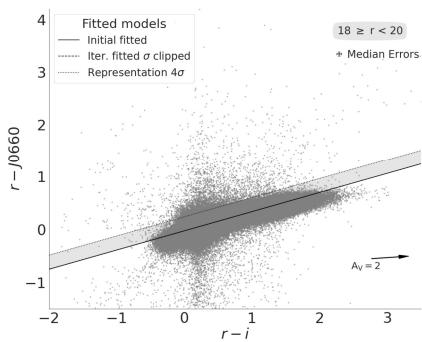
(a)



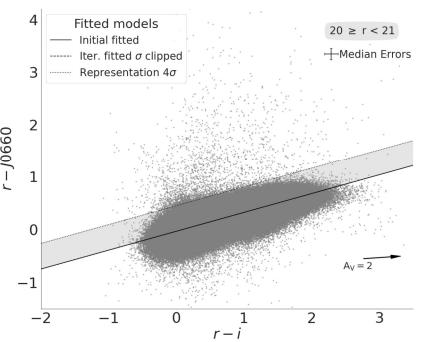
(b)



(c)



(d)



**Figure 2.** An illustration of the selection criteria used to identify strong emission-line objects via colour-colour plots. The data shown here are all from the S-PLUS DR3. The data are split up into four magnitude bins, as shown in the four panels. Objects with H $\alpha$  excess should be located near the top of the colour-colour diagrams. The thin continuous lines illustrate the original linear fit to all the data (grey points). The dashed lines represent the final fits to the stellar locus of points which were obtained by applying an iterative  $\sigma$ -clipping technique to the initial fit. The actual cuts used to select H $\alpha$  emitters are shown by the dotted lines. These correspond to  $4\sigma$  above of the final fit. Objects selected as H $\alpha$  emitters must be located above the dotted line. Note that the position of these lines (selection criteria) shown in the figure are approximated, given that the actual selection criterion also considers the errors on each source.

### 2.3 Unsupervised machine learning clustering approach

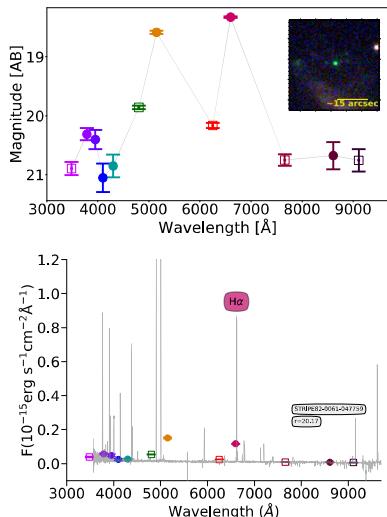
#### 2.3.1 Hierarchical agglomerative clustering

Hierarchical clustering (HC) belongs to the family of clustering algorithms of which clusters are constructed by recursively grouping and splitting the sources. Being an unsupervised algorithm, HC does not require a training sample or pre-conceived hypotheses. Data elements are grouped based on patterns in a given space of parameters and on the levels of similarity at which the groupings change (Jain et al. 1999). In the end, HC returns a diagrammatic representation of the groups as a tree – a dendrogram that follows an hierarchical structure.

There are two types of hierarchical clustering: the *hierarchical*

For the split of the sample of H $\alpha$  emitters into two classes, the blue and the red populations, we follow an unsupervised machine learning approach implementing two clustering techniques: hierarchical agglomerative clustering and hierarchical density-based cluster selection, both based on the ( $r - z$ ) and ( $g - r$ ) colours, whose results are mutually compared.

No Comments.



**Figure 3.** Top panel: S-spectrum of a random emitting object found by the method explained in Section 2.2. Open squares represent the SDSS-like broad-band filters. From left to right:  $u$ ,  $g$ ,  $r$ ,  $i$  and  $z$  magnitudes. Circles represent the narrow-band filters, which from left to right correspond to J0378, J0395, J0410, J0515, J0660 and J0861. The inset figure is the coloured image of the object which was produced by combining all twelve bands. Bottom panel: SDSS spectrum and S-PLUS photometry in flux unity of the object - the H $\alpha$  line is marked. The spectroscopic redshift of this object is  $z = 0.009$ .

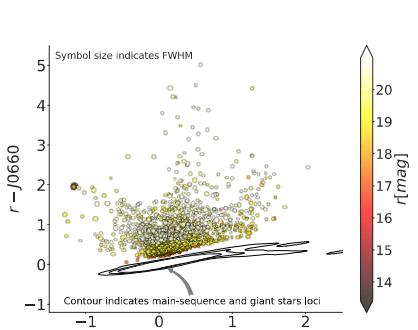
agglomerative clustering (HAC; the one used in this work), which follows “bottom-up” approach, and the *hierarchical divisive clustering* that follows “top-down” approach. The HAC consists of building a binary merge tree, starting from each data element stored at the leaves (interpreted as individual clusters) and proceeds by merging two by two the “closest” sub-sets (stored as nodes) until the root – unique cluster – of the tree that contains all the elements of the data set is reached. The term “agglomerative” is used to point out that data elements are successively agglomerated into higher-levels. In each iteration, two “nearby” clusters are collapsed into a new, more populated group (Mann & Kaur 2013; Aggarwal 2015). Hence, each step reduces the number of clusters. The procedure may be summarized in three steps:

(i) Initially, each data element represents one cluster, i.e. “leaves of the tree”. This means that at the beginning, the total number of clusters/leaves is equal to the number of the elements in the data set.

(ii) Through a looping process, the clusters are merged into new ones that are described by the maximum similarity among them.

(iii) Finally, all the clusters belong to an unique cluster, “the root of the tree” structure.

On the other hand, the *hierarchical divisive clustering* algorithm follows a “top-down” approach. This means that the clustering starts from data element from only one cluster and then moves down recursively in the hierarchy to smaller groups. In simple words, hi-



**Figure 4.** Colour-colour diagram with all the emission-line objects selected from the S-PLUS DR3. The size of the symbols represent the measured FWHM assuming a Gaussian core (for more detail see Almeida-Fernandes et al. 2022). Coloured bar indicates the magnitude values of the  $r$ -band. The contours represent the S-PLUS synthetic photometry of main-sequence and giant stars loci from the library of stellar spectral energy distributions of Pickles (1998).

archical (agglomerative and divisive) clustering algorithms intend to gather similar objects into groups called clusters in the space of parameters which is investigated.

### 2.3.2 Hierarchical density-based cluster selection

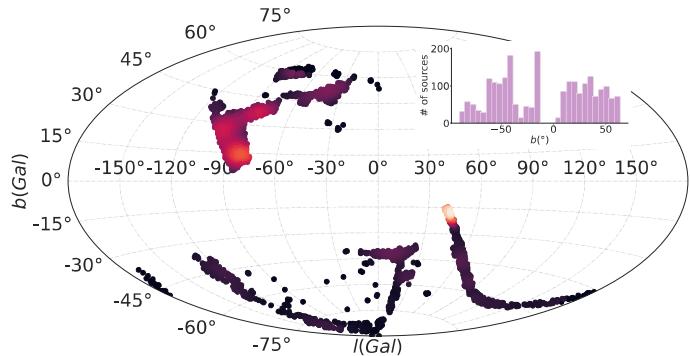
Hierarchical density-based cluster selection (hereafter HDBSCAN; Campello et al. 2013) is another unsupervised machine learning algorithm that relies on clustering. It is based on a slightly modified version of density-based spatial clustering of applications (DBSCAN; Ester et al. 1996) which declares data points as noise. It assumes that clusters are characterized by “islands” of high density in the sea of the parameter space. HDBSCAN takes the DBSCAN concept forward by introducing a hierarchy to the clustering, with “persistent” clusters finally extracted from the hierarchical tree. The main advantage of HDBSCAN in comparison with its predecessor consists in the possibility of finding clusters of variables densities and different shapes. Following Malzer & Baum (2021) and Ntawtsile & Geach (2021) it works as follows:

(i) HDBSCAN defines the “core” distance for a data point  $x$ ,  $\text{core}_k(x)$ , as the distance of an object to its  $k$ th nearest neighbour. This mean that lower values of  $\text{core}_k(x)$  represent higher densities and vice-versa.

(ii) The “mutual readability distance” between two points  $a$  and  $b$  is defined as  $d_m(a, b) = \min\{\text{core}_k(a), \text{core}_k(b), d(a, b)\}$ , where  $d(a, b)$  is the distance between  $a$  and  $b$  according, for instance, to Euclidean metric. The mutual readability distance allows data points in dense regions to stay close together and those that are in less dense regions to move away.

(iii) The mutual readability plot is used to construct the minimum spanning tree, and sorting its edges by the mutual readability distance resulting in a hierarchical tree structure. The hierarchy of connected components is defined by sorting the edges of the tree by distance

No Comments.



**Figure 5.** Distribution of the emission-line objects in galactic longitude and latitude coordinates. The inset figure represents the distribution of the objects in galactic latitude.

in reverse order, describing a dendrogram (the diagram explained in 2.3.1). This is the structure from which the cluster will be identified.

(iv) HDBSCAN allows extracting clusters of variable density by cutting the dendrogram at different levels of grouping.

(v) The cluster tree is condensed into a simpler structure (see, for instance, Figure A1 of Appendix A). Considering the single main trunk which contains all the data points, the tree splits into branches. A condensed cluster hierarchy can be described by considering the number of points that are kept in each branch as it splits. It is important to mention that there is a key parameter called minimum cluster size. If a given branch splits into two, with one branch containing fewer points than the minimum cluster size, the larger branch “persists” and the smaller split branch “falls out” of the cluster. If a branch splits into two with both branches exceeding the minimum cluster size, both new branches are preserved.

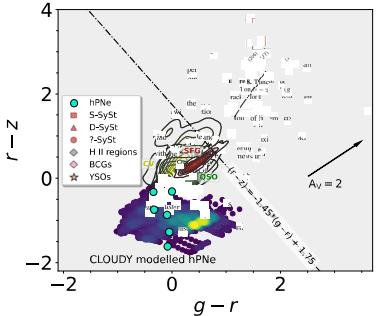
(vi) The clusters are extracted on the notion of persistence in the hierarchy. The parameter  $\lambda = d_m^{-1}$  is defined, and each cluster has a  $\lambda_{\text{birth}}$  (the point at which the cluster split off) and  $\lambda_{\text{death}}$  (the point when the cluster split into other clusters). In each cluster, we have  $\lambda_p$  describing when each point fell out of the cluster (or was split off into a new cluster), so that  $\lambda_{\text{birth}} \leq \lambda_p \leq \lambda_{\text{death}}$ . Cluster stability  $S$  is defined as the sum of  $\lambda_p - \lambda_{\text{birth}}$  for all points in the cluster. To extract clusters, the following procedure is implemented. First, each leave constitutes a cluster. Then, moving through the hierarchy, it is considered the stability of a parent cluster  $S_p$  and its  $n$  descendants  $S_d^{0,1,2,\dots,n}$ : if  $S_p > \sum_{i=0}^n S_d^i$ , we unselect all the descendants; otherwise, the cluster stability is set as  $S_p = \sum_{i=0}^n S_d^i$ . At the root node, we have our set of the selected clusters. Any data point in the sample that does not fall into one of the selected clusters is defined as noise.

(vii) By adopting the **soft clustering** or the **fuzzy clustering** technique it is possible to mitigate the need to establish or define the cluster membership limit. In fact, each source has a finite probability of belonging to every selected cluster. In this approach, all points (including noises) are not assigned to a cluster label, but are instead assigned to a vector of probabilities whose length is equal to the

number of clusters found. Such an approach aims solve the problem of noise classification.

The HDBSCAN algorithm starts off much the same as DBSCAN: transforming the space according to density, exactly as DBSCAN does, and perform single linkage clustering on the transformed space. Instead of taking an epsilon<sup>5</sup> value as a cut level for the dendrogram, a different approach is followed: the dendrogram is condensed by viewing splits that result in a small number of points splitting off as points “falling out of a cluster”. This results in a smaller tree with fewer clusters that “lose points”. That tree can then be used to select the most stable or persistent clusters. This process allows the tree to be cut at varying height, picking our varying density clusters based on cluster stability. The immediate advantage of this is that we can have varying density clusters; the second benefit is that we have eliminated the epsilon parameter as we no longer need it to choose a cut of the dendrogram. Instead we have a new parameter `min_cluster_size` which is used to determine whether points are “falling out of a cluster” or splitting to form two new clusters.

Over the last few years, HDBSCAN has been used for different tasks in Astronomy. HDBSCAN was used to identify IR bubbles from Spitzer images (Jayasinghe et al. 2019). Webb et al. (2020) implemented HDBSCAN for discovering transients. Logan & Fotopoulou (2020) presented HDBSCAN as a viable tool to separate stars, galaxies and QSOs using photometric data. Recently, Ntwaetsile & Geach (2021) employed HDBSCAN to group radio sources into a sequence of morphological classes, illustrating a simple methodology to classify and label new, unseen galaxies in large samples. This approach was also implemented to identify stellar groups in Canis Major OB1 (Santos-Silva et al. 2021).



**Figure 6.** The  $(r-z)$  versus  $(g-r)$  synthetic colour-colour diagram of several classes of emission line objects. Included in the diagrams, there are families of CLOUDY modelled halo PNe spanning a range of properties (density map region). Cyan circles represent S-PLUS photometry from observed spectra of PNe. Grey diamonds represent H II regions in NGC 55. Red boxes and triangles display S- and D-type symbiotic stars, respectively. Red circles are SySt with associated type. This group includes Galactic and external SySt from NGC 205 IC 10 and NGC 185. Yellow contours correspond to CVs from SDSS. Pink diamonds indicate blue compact galaxies (BCGs) from SDSS. Brown contours refer to SDSS star-forming galaxies (SDSS SFGs). SDSS QSOs at different redshift ranges are shown as green contours, and YSOs from Lupus and Sigma Orionis are represented by salmon stars. The diagonal dashed line represents a subjective criterion to separate the objects into two colour types. The arrow indicates the reddening vector with  $A_V \approx 2$  mag.

#### 2.4 Splitting the H $\alpha$ emitters into blue and red populations

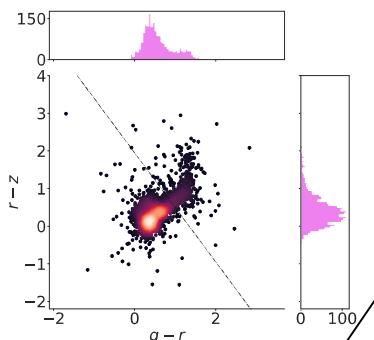
To select the blue and red populations in the sample of H $\alpha$  emitters, we first looked for the best colour-colour diagram by using the S-PLUS synthetic photometry of several classes of emission line objects. The  $(r-z)$  versus  $(g-r)$  diagram is displayed in Fig. 6. SySt and YSOs span a range on  $(g-r)$ , from 0.4 to 3.6 at the upper right. On the other hand, PNe, H II regions, CVs, QSOs, and emission line galaxies are located on the lower-left region of the diagram. The dashed line in Fig. 6 separates the blue and red zones.

Fig. 7 displays the  $(r-z)$  versus  $(g-r)$  diagram from the list of H $\alpha$  emitters in S-PLUS. Obviously only such emitters with detection in the  $g$  and  $z$  filters are considered for this colour classification by making a cut in the magnitude errors at 0.2, totaling 2,892 objects. A bi-modal distribution is found for both colour indices (see side and upper plots of the Fig. 7). The two peaks on the  $(g-r)$  and  $(r-z)$  distributions have immediate correspondence with the blue and red zones pointed out from the synthetic diagram (Fig. 6). One can also see that the fraction of blue objects is considerable larger than that of the red ones.

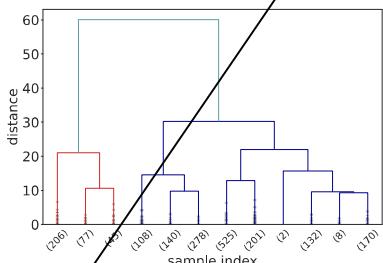
##### 2.4.1 Application of hierarchical agglomerative cluster

The ideal way to choose the number of clusters is by displaying the **dendrogram diagram**. Firstly, the hierarchical cluster output dendrogram can be used to obtain the desired clustering. Secondly, the dendrogram allows a convenient way to establish the entity-relationship at all levels of granularity.

Fig. 8 illustrates the dendrogram based on the  $(g-r)$  and  $(r-z)$



**Figure 7.** The  $(r-z)$  versus  $(g-r)$  colour-colour diagram with all the emission line objects selected in S-PLUS. The side and upper figures represent the  $(r-z)$  and  $(g-r)$  colour distributions, respectively.



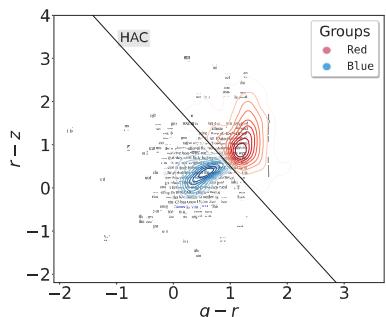
**Figure 8.** Truncated dendrogram of complete-linkage hierarchical clustered based on  $(r-z)$  and  $(g-r)$  colours. The cluster sizes are exposed in the brackets for the 12 truncated clusters.

colours of H $\alpha$  emitters, and it highlights the order and distances of the groups in the hierarchical clustering, stopping at 12 nodes:

- The x-axis specifies the population in the nodes in a given level of grouping – that summed up correspond to the total number of elements under investigation.
- The y-axis represents the “distance”, which is a measurement of the closeness of the clusters or data points in different levels of clustering.

Reading the diagram from the top to the bottom, we see that all systems are divided after the very first level from the top already into (only) two groups: as expected, they correspond to the red and blue populations of H $\alpha$  emitters presented in Fig. 7 as it is showed in Fig. 9. From that point on, the groups were subdivided without evident distinction, and truncation was thus assumed when the 12-node level was reached. The truncation is an usual procedure when dealing with big data.

In this work, HAC was employed by using the library



**Figure 9.** The  $(r - z)$  versus  $(g - r)$  colour-colour diagram with the two population found by implementing HAC algorithm. The blue and red symbols represent the sources with intense blue continuum and those with intense red continuum, respectively. The straight line is the same line in Fig. 6.

Scikit-learn<sup>6</sup> (Pedregosa et al. 2011). Then, the `DENDROGRAM()` function, which is included in the package `scipy`<sup>7</sup>, was used in conjunction with the task `Dendrogram Truncation` to generate the (truncated) dendrogram. The following parameters must be taken into account when the algorithm is applied to data: `n_clusters`, `Affinity`, and `Linkage`. The parameter `n_clusters` defines the number of clusters expected by the user. Given that our goal is to divide our sample into two groups, `n_clusters` is set to "2". `Affinity` determines the "metric" that compute the linkage. We have found that a simplistic metric, the "Euclidean", is effective for our purpose. `Linkage` determines which distance to use between sets of observation. `Linkage` defines how the similarity between two clusters is calculated, by determining the distance between sets of observations as a function of the pairwise distances between elements. The algorithm merges the pairs of cluster that minimize this criterion. Ward's method minimizes the variance of the clusters that are being merged (Ward 1963). To implement this method, find the pair of clusters that leads to a minimum increase in total within-cluster variance after merging. Ward procedure uses the error sum of squares to measure this variance. The two clusters with the smallest error sum of squares will eventually form a new cluster.

At this point, our list of  $\text{H}\alpha$  emitters is divided into two populations based on their `continuum`, with the blue population (1,564 objects) being larger than the red one (328).

#### 2.4.2 Application of HDBSCAN

For the sake of comparison with the results from HAC, we also used HDBSCAN to distinguish the blue and the red sources. The main difference between these two algorithms is that HDBSCAN is more conservative in the sense that several data points are classified as noise. For this task, the Python implementation of HDBSCAN<sup>8</sup> (McInnes et al. 2017) was adopted.

Similarly to HAC, there are some key parameters that should be

considered when the algorithm is applied. Regarding the metric, the "Euclidian" one is assumed. The two most critical parameters are the "minimum cluster size" and "minimum number of samples". The former refers to the smallest size of a group that it is considered as a cluster. The value of "60" has been adopted for the "minimum cluster size". The "minimum number of samples" provides a measure of how conservative our clustering method will be, expressed as the fraction of data classified as noise, and the value of "15" was adopted. With this model configuration, two clusters were identified.

Left panel of Fig. 10 shows the two clusters found with HDBSCAN. One cluster contains 1,413 blue sources and the other one 131 red sources. The number of objects classified as noise is 348. This result is overall consistent with those obtained with HAC, although being more restrictive when classifying members of the red group. The two main clusters obtained with HDBSCAN are located in the same region in the  $(r - z)$  versus  $(g - r)$  diagram as those groups found based on the HAC. About 98% of the blue sources selected by HDBSCAN are in the list of blue objects identified by HAC. All the red sources selected by HDBSCAN were also classified by HAC as red objects. In fact, by applying the condensed tree to the data colours two clusters are selected. The `condensed_tree_` attribute is the equivalent dendrogram plot for HDBSCAN which displays the cluster tree mentioned in the section 2.3.2. (see Appendix A for more details about `condensed_tree_` attribute).

#### 2.4.3 Soft clustering for HDBSCAN

The main disadvantage of HDBSCAN is that several sources are labelled as "noise", so that they are not assigned to any cluster. As mentioned earlier, this comes from the conservative nature of HDBSCAN and the fact that these data points (data noise) are located far away of the clusters' cores. An alternative way to avoid outliers (data noises) classifications is the implementation of the "soft clustering" (see section 2.3.2). Soft clustering from HDBSCAN<sup>9</sup> was used to assign every object to a cluster that they most likely belong to. According to this approach, data points are not assigned in a deterministic way to a cluster but to a vector of probabilities as a measure of belonging to different clusters; the probability value at the  $i$ th entry of the vector is the probability that a data point is a member of the  $i$ th cluster. We can, then, simply assign cluster labels for every data point by taking the most likely cluster it belongs to, using probability thresholds. Therefore, soft clustering for HDBSCAN is achieved through an outlier score modification to consider how distant an outlier is from each cluster, which is based on the Global-Local Outlier Score from Hierarchies (GLOSH) algorithm (Campello et al. 2015). This is combined with a measure of distance from a given cluster to estimate the probability that a given data point belongs to any of the fixed groups drawn from the condensed tree.

The right panel of Fig. 10 shows which cluster the data points classified as the noise by HDBSCAN belong to. Blue and red points indicate those with the highest probability of being in the blue and red groups, respectively. This procedure fills out the clusters nicely. There were many noise points that most likely belong to the expected clusters in very good agreement with the results obtained from HAC. Indeed, our separation of the  $\text{H}\alpha$  emitters into blue and red sources has been improved. Instead of forcing the algorithm to make a decision to which group a data point belongs to, as HAC does, we have quantified the likelihood of a given observation to belong to any of

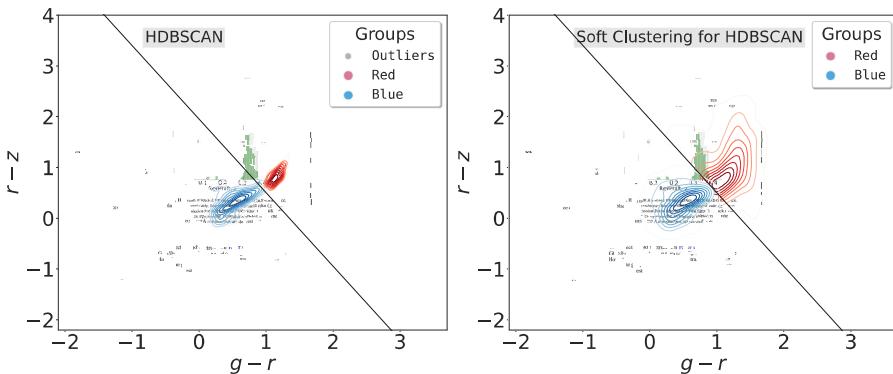
<sup>6</sup> <https://scikit-learn.org/stable/>

<sup>7</sup> <https://www.scipy.org/>

<sup>8</sup> <https://hdbscan.readthedocs.io/en/latest/>

obviously the continuum dominates the broad band, but you measured the integrated flux

so I would say: integrated flux



**Figure 10.** Left panel: as Figure 9 but with the results after to apply HDBSCAN to the sample of emission-line sources. Blue and red symbols correspond to blue and red sources, respectively, with gray symbols representing sources classified as noise by HDBSCAN. The straight line is the same line as Figures 6 and 9. Right panel: results after apply a soft clustering to the HDBSCAN results.

the two clusters found in our data set (see, for instance, the two last columns of Table B1).

### 3 RESULTS AND DISCUSSION

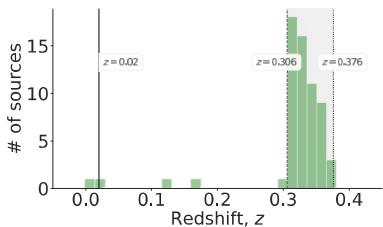
Our strategy that is focused on the identification of H $\alpha$  emitters in the S-PLUS footprint, exploiting the unique filter system of the survey returned 3,187 objects with excess in the J0660 band. The fractional contribution of different classes of H $\alpha$  emitters to the overall sample was evaluated by cross-matching the objects' list with the SIMBAD database<sup>10</sup>. Optical spectra available in the SDSS DR16 (Ahumada et al. 2020) and in the Large Sky Area Multi-Object Fiber Spectroscopic Telescope (LAMOST DR7; Wu et al. 2011) were also explored. In all cases, we assumed the angular distance on the sky-plane between sources considering positive matches those of mutually closest sources to each other within a given limit ( $d_{max,proj}$ ).

#### 3.1 Matches with SIMBAD sources

We found 393 positive matches between our catalog of H $\alpha$  emitters and SIMBAD database considering a radius of  $d_{max,proj} = 2$  arcsec. The results are described below and are listed in Table 1.

##### 3.1.1 Ionized nebulae

As it was mentioned, several classes of objects with diffuse appearance and/or nebular lines in our Galaxy as well as in nearby galaxies are listed in our sample, most being PNe. Planetary nebulae represent the final stages of low- and intermediate-mass stars from which the material has been previously ejected in the phases of AGB and post-AGB and is ionized by the high energetic radiation from a hot stellar remnant core. In our list of H $\alpha$  emitters one compact PN is catalogued in SIMBAD. This object was classified as a blue source.



**Figure 11.** Histogram of redshift for the galaxies with SIMBAD coincidence. The black vertical continuous line indicates the threshold value ( $z \sim 0.02$ ) on which the H $\alpha$  emission-line is detected in J0660 filter. The black vertical dashed and dotted lines represent the accumulative redshift range where the H $\beta$  and [O III] 4959, 5007 Å emission lines are detected on the J0660-band.

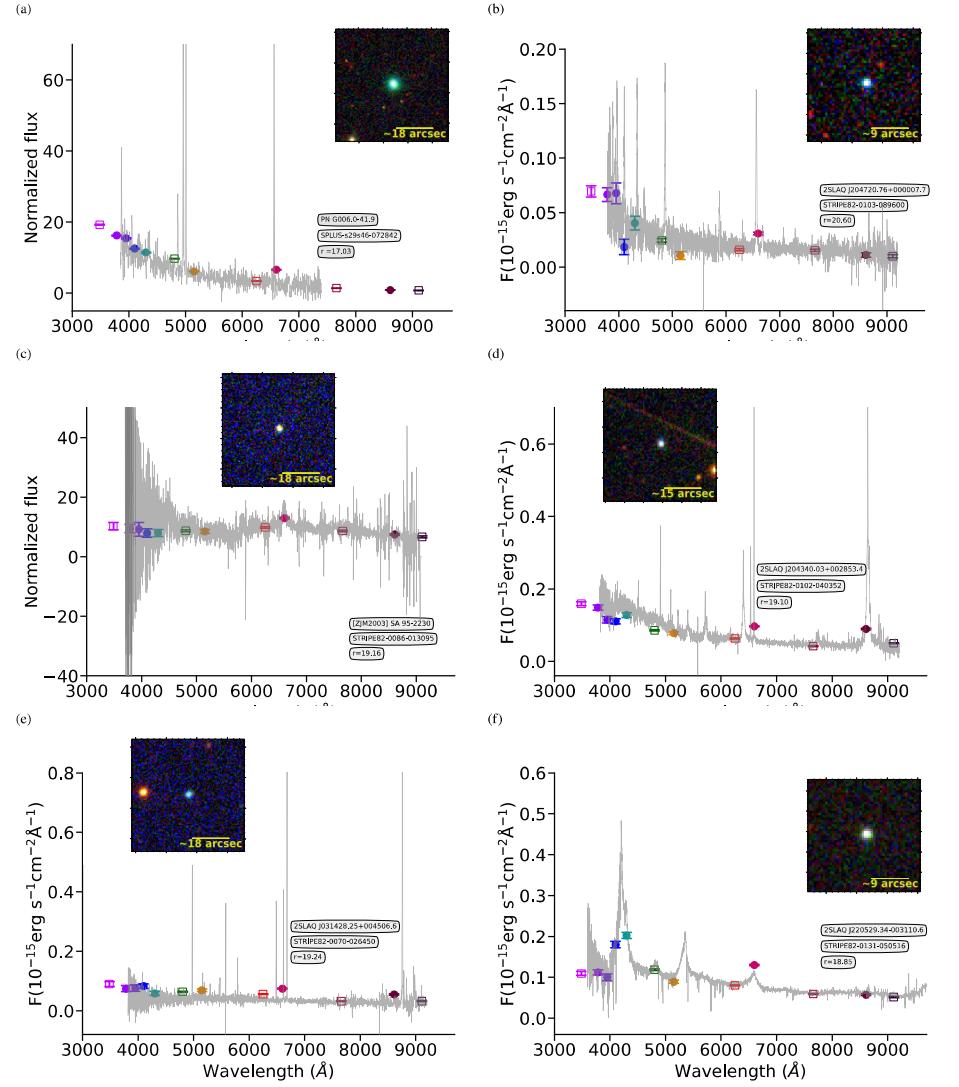
The S-spectrum and spectrum from Parker et al. (2016) of the PN G006.0-41.9 are displayed on panel (a) of Fig 12. Emission lines like H $\alpha$  and [N II] are clearly visible in its spectra. PN G006.0-41.9 belongs to the group of Galactic halo PNe. This rare group of PNe are of particular interest because they are characterized by low metallicity and present large velocities.

##### 3.1.2 Supernovae

Supernovae also display emission line spectra from evolved stars through multiple channels. However, the energy-input mechanism is quite different in each case. One confirmed SN is on the list. This SN is classified as a blue source by HAC but the probability of being blue and red of 0.16 and 0.29, respectively, according to HDBSCAN.

<sup>10</sup> <http://simbad.u-strasbg.fr/simbad/>

No Comments.



**Figure 12.** Summary of our selection results showing the spectrum (gray line in each panel) of different classes of emission line sources identified in our target list. A spectrum of a PN (*a*) from Parker et al. (2016). The SDSS spectra of a Seyfert (*b*) and a cataclysmic variable star (*c*). SDSS spectra of a Seyfert 1 (*d*) and an emission line galaxy (*e*) with  $z = 0.317$  and  $z = 0.334$ , respectively. For two objects, the excess on the  $J0660$  is due to the [O III] 4959, 5007 Å emission lines, and a QSO (*f*) with redshift of  $\sim 2.45$ , meaning that excess on the narrow-band filter is produced by the C III] emission line. As in Figure 3, coloured square and circles symbols represent the S-PLUS photometry. All these objects show a significance excesses on the  $J0660$  filter in comparison with the broad-bands.

No Comments.

**Table 1.** A summary of the results obtained of the positional cross-match between the S-PLUS list of emission line objects and the SIMBAD database. We used a search radius of 2 arcsec. SIMBAD categories of objects are listed in the first column. The numbers of objects of each SIMBAD category are exposed in the third column.

Main type	Associated SIMBAD types	Number of S-PLUS objects with SIMBAD match
Nebulae	PN	1
Supernovae	SN	1
Stellar binary system	CataclyV*, CV*_Candidate, EB*, EB*_Candidate	40
Star	star, WD*, WD*_Candidate, Blue, PM*, low-mass*	32
Variable star	RRLyr, Candidate_RRLyr, V*	7
Galaxy	EmG, Galaxy	42
QSO	QSO, QSO_Candidate	229
AGN	AGN, AGN_Candidate, Seyfert_1	31
Other type	UV, X, Unknown_Candidate, Radio(cm), Radio, EmObj, GICl	10
Total		393

### 3.1.3 Interacting binary systems

Following the classification available in SIMBAD, 27 known and 6 CV candidates were found by our analysis. CVs are interacting binary systems of very short orbital period, in which a low-mass, early-type star fills its Roche lobe and transfers mass to a white dwarf companion (Patterson 1984). For the sake of illustration, Fig. 12 (panel *b*) shows the S-PLUS photometry overlapped to the SDSS spectrum of 2SLAQ J204720.76+000007.7, a CV, which we correctly classified as a blue source. In fact, 31 of these CV are classified as blue sources by HAC. Eclipsing binaries are also listed on our catalog.

### 3.1.4 Stars

Table 1 shows that 32 objects in our sample of H $\alpha$  emitters are categorized, by SIMBAD, as: normal stars, white dwarf (WD\*), white dwarf candidates (Candidate\_WD\*), blue stars, high proper-motion stars (PM\*), variable stars of RR Lyr type and, low-mass star (low-mass\*, M<1M $_{\odot}$ ). Panel *c* of Fig. 12 shows the LAMOST and S-spectrum, as well as the coloured image of the source [ZJM2003] SA 95-2230 catalogued as star in SIMBAD. Its spectrum exhibits an emission line at the wavelength range covered by the *J*0660 filter. The spectroscopic redshift of this object is  $\sim$ 1.36, indicating that the excess on the narrow-band filter is caused by the Mg II 2798 Å emission line. This object could be a QSO misclassified in SIMBAD, which has been classified as blue sources for both HAC and HDBSCAN.

### 3.1.5 Galaxies

Galaxies are also included in our catalog: emission-line galaxies (EmG), Seyfert types-1 and other type of AGNs. Since, we are focusing on the H $\alpha$  emission line, the emission line galaxies in the local Universe ( $z \sim 0.02$ ) are of particular interest because their H $\alpha$  line still falls into the wavelength range covered by the *J*0660 S-PLUS filter.

Fig. 11 shows the redshift distribution of the galaxies in our sample that have SIMBAD correspondents. Of the 73 galaxies with SIMBAD counterparts, 65 objects have redshift measurements of which only 2 objects of them have small redshift values ( $z < 0.02$ ) showing that the emission detected in the *J*0660 filters is associated with the real H $\alpha$  λ6563 emission line. On the other hand, Fig. 11 also shows that there is an increment of galaxies with redshift between  $\sim$ 0.31 and  $\sim$ 0.38. This particular population of seemingly H $\alpha$  emitters is represented by AGN, Seyfert 1 galaxies and other emission line

galaxies. In fact, at the accumulative redshift range,  $0.306 < z < 0.376$ , represented by the filled area in the figure, with H $\beta$  and [O III] 4959, 5007 Å lines redshifted into the *J*0660 filter.

Panel *d* of Fig. 12 shows the S-PLUS photometry and the SDSS spectrum of a Seyfert type-I galaxy (2SLAQ J204340.03+002853.4). The spectrum clearly exhibits strong emission lines. The redshift of this object is around 0.317 indicating that the excess on the *J*0660 is due to the [O III] 4959, 5007 Å emission lines. Panel *e* exhibits the S-PLUS photometry and the SDSS spectrum a the galaxy 2SLAQ J031428.25+004506.6. This source has  $z \sim 0.334$  indicating, as the above Seyfert galaxy, that the excess in the *J*0660 filter is not due to the H $\alpha$  emission line. Note that 88% of these galaxies are classified as blue objects by HAC and HDBSCAN. However, 5% of the galaxies are found to fill up the population of red sources.

### 3.1.6 QSOs: redshifted lines mimicking the H $\alpha$ emission

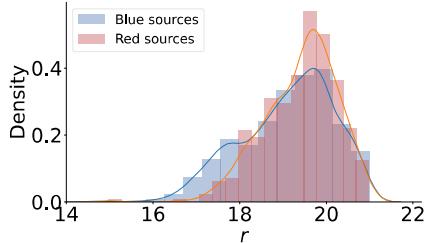
Following the classification available in the literature, about 13% of blue H $\alpha$  emitters sources in our sample are found to be QSOs. We have to point out here that the excess in the *J*0660 filter for QSOs is attributed to redshifted lines that fall in the wavelength range covered by that filter depending on the redshift of the QSOs – e.g., H $\beta$ , Mg II 2798 Å, C III] 1909 Å and C IV 1550 Å (see Gutierrez-Soto et al. 2020 and Nakazono et al. 2021). QSO 2SLAQ J220529.34-003110.6, shown in Panel *f* of Figure 12, is an example of a QSO at redshift  $\sim$ 2.45, for which the C III] line falls at the range covered by the *J*0660 filter.

### 3.1.7 Other classifications of objects

As it can be seen in Table 1, our sample also gathers a variety of objects without any previous classification. They may also be X and UV sources, globular cluster (GLCL), among others, indicating the richness of the sample in nature and in physical properties.

## 3.2 SDSS and LAMOST: a spectroscopic validation

Finally, we also cross-matched out a sample of H $\alpha$  emitters in the S-PLUS with the SDSS DR16 (Ahumada et al. 2020). For doing this, we adopted a 2 arcsec as the cross-matching radius. In the case of the cross-match with LAMOST (Wu et al. 2011), the same radius was considered, and we ended up with 218 sources belonging to both



**Figure 13.** Distribution of  $r$  magnitude for the blue and red sources of the sample of H $\alpha$  emitters. The histogram heights show density normalization scales. The smooth curves represent a Kernel density estimation for both samples.

catalogues and approximately 95% of them display emission lines spectra.

Most of the H $\alpha$  emitters with available spectroscopic information correspond to CVs, nova, emission-line galaxies, Seyfert I and AGN in general, and QSOs. However, we emphasize that more detailed analysis is necessary to check which other types of objects are included in these samples of spectra – what is not in the scope of this paper. Also, it is worth noticing that part of the objects does not have a conclusive classification.

The spectra from both SDSS and LAMOST projects provide a good validation to our approach, clearly showing that the methodology is actually effective for selecting sources with emission lines.

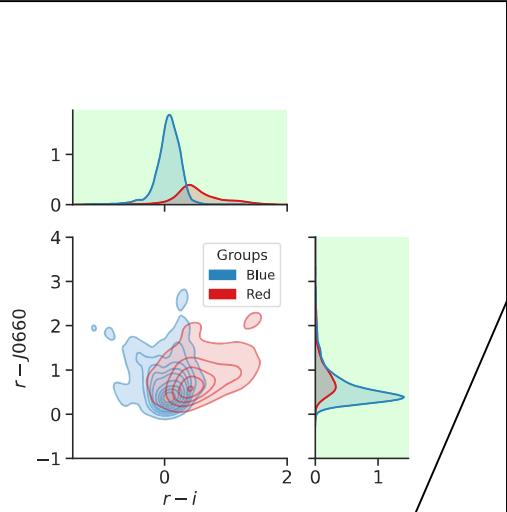
### 3.3 Magnitudes and colour distributions

In Fig. 13, we demonstrate the distribution of the blue and red population of S-PLUS H $\alpha$  emitters in term of their  $r$ -magnitude. Both, blue and red sources can be as bright as 16 mag in the  $r$  filter, while they show a peak at  $\sim 19.7$  mag. The fraction of blue sources in the  $16.0 \leq r \leq 19.0$  magnitude range is considerably higher in comparison with the red group. Therefore, the blue sample tends to be brighter than the red one in the  $r$ -band. In fact, the distribution of the blue sources is bimodal showing another peak at  $\sim 17.8$  mag.

Fig. 14 displays the  $(r - J0660)$  versus  $(r - i)$  diagram with the position of the two populations and the distributions of the blue and red H $\alpha$  sources in the  $(r - i)$  (upper) and  $(r - J0660)$  (side) colours. The  $(r - i)$  colour distribution of the blue and red sources have peaks at distinct values of 0.2 and 0.6, respectively. This result is consistent with that obtained from [Wevers et al. \(2017\)](#) who also used the  $(r - i)$  colour index to select blue outliers from the Galactic Bulge Survey (GBS; [Jonker et al. 2011](#)). Finally, the  $(r - J0660)$  colour index distribution of the blue and red objects peaks at 0.4 and 0.7, respectively. This result implies a strong H $\alpha$  emission in the red sources compared to the blue ones.

## 4 CONCLUSIONS

Here we exploited the capability of the S-PLUS project ([Mendes de Oliveira et al. 2019](#)) to survey H $\alpha$  emitters in the Southern Sky following a three-steps approach: identify H $\alpha$  emitters, distinguish



**Figure 14.**  $(r - J0660)$  versus  $(r - i)$  colour-colour diagram for the blue (blue contours) and red sources (red contours) of the sample of H $\alpha$  emitters. The side and upper panels show the respective colours distributions.

the blue and red populations as a first diagnostic about the nature of the sources, and validate the results through spectroscopic databases.

The H $\alpha$  emitters were identified by employing the  $J0660$  narrow-filter and  $r$  and  $i$  broad-filters available in the S-PLUS project. The  $(r - J0660)$  versus  $(r - i)$  colour-colour diagram was used to define the loci of the main-sequence and giant stars and disentangle objects with probability to be point sources in the local Universe with an H $\alpha$ -excess ( $r - J0660 > 0$ ) (see Fig. 4). 3,187 sources matched this criterion, with 273 of them claimed in the literature (SIMBAD) as QSOs and non-local galaxies, and therefore being false positive identifications of H $\alpha$  emission (see sections 3.1.5 and 3.1.6).

The  $(r - z)$  and  $(g - r)$  colour distributions of the H $\alpha$  emitters were found to be bimodal, indicating the presence of two distinct populations of bluer and redder sources with a narrow overlapping zone (Fig. 7). Two algorithms of unsupervised machine learning classification were used to distinguish the two populations: the HAC and HDBSCAN clustering algorithms. Both algorithms ended up to very similar clusters, on the  $(r - z)$  and  $(g - r)$  colour indices space.

Given that HDBSCAN is considered as a conservative algorithm, many objects were labeled as noise data points while they did not by the HAC algorithm (Section 2.2.3). To overcome this problem, a so-called “soft” clustering approach for HDBSCAN was employed and the probabilities of each data point to belong to the “blue” and “red” subgroup were computed. The results from the HAC and HDBSCAN algorithms are mutually consistent. We, therefore, reckon that the  $(r - z)$  and  $(g - r)$  colours are ideal for separating objects into the bluer and redder populations and correlate their colour to the nature of the sources. In particular, the bluer objects were found to be mainly CVs, PNe, dwarf compact galaxies, galaxies with redshifted spectra and QSOs, among others, while the redder sources are early type galaxies with emission lines, probably young/active late-type stars or even symbiotic stars (in fact, evolved binary systems hosting a red giant star).

Finally, we also cross-matched our catalog of H $\alpha$  emitters with

spectroscopic databases (SDSS and LAMOST; see Section 3.2). This exercise demonstrated that at least 95% of the objects with available spectroscopic information are genuine emission line sources, validating our approach to identify emission line objects in the S-PLUS project. The spectroscopic sample of H $\alpha$  emitters lists 21 sources of the local Universe (with  $z < 0.02$ ) indicating that the emission on the J0660 filter corresponds to the He I line, 197 sources with redshift larger than 0.02, indicating that they are very likely QSOs and AGN and/or non-local galaxies on which the excess of the J0660 filter is due to H $\beta$ , Mg II 2798 Å, C III] 1909 Å and C IV 1550 Å emission lines for the case of QSOs and H $\beta$  and [O III] 4959, 5007 Å emission lines for galaxies, those depending on their redshift.

As a practical result, here we make public a catalog from S-PLUS/DR3 that can be explored by the community in the identification and investigation of sources in twelve photometric bands in a systematic and homogeneous way.

## ACKNOWLEDGEMENTS

LAG-S acknowledges funding for this work from FAPESP grants 2019/26412-0. RLO acknowledges financial support from the Brazilian institutions CNPq (PQ-312705/2020-4) and FAPESP (#2020/00457-4). DGR acknowledges the CNPq (428330/2018-5; 313016/2020-8) and FAPERJ (269312) grants. F. A. -F. acknowledges funding for this work from FAPESP grants 2018/20977-2 and 2021/09468-1. C. C. is supported by the National Natural Science Foundation of China, No. 11803044, 11933003, 12173045. This work is sponsored (in part) by the Chinese Academy of Sciences (CAS), through a grant to the CAS South America Center for Astronomy (CASSACA). We acknowledge the science research grants from the China Manned Space Project with NO. CMS-CSST-2021-A05. AAC acknowledges support from the State Agency for Research of the Spanish MCIU through the “Center of Excellence Severo Ochoa” award to the Instituto de Astrofísica Andalucía (SEV-2017-0709). The authors would like to thank Amanda Reis Lopes for useful suggestions and comments.

The S-PLUS project, including the T80-South robotic telescope and the S-PLUS scientific survey, was founded as a partnership between the Fundação de Amparo à Pesquisa do Estado de São Paulo (FAPESP), the Observatório Nacional (ON), the Federal University of Sergipe (UFS), and the Federal University of Santa Catarina (UFSC), with important financial and practical contributions from other collaborating institutes in Brazil, Chile (Universidad de La Serena), and Spain (Centro de Estudios de Física del Cosmos de Aragón, CEFCFA). We further acknowledge financial support from the São Paulo Research Foundation (FAPESP), the Brazilian National Research Council (CNPq), the Coordination for the Improvement of Higher Education Personnel (CAPES), the Carlos Chagas Filho Rio de Janeiro State Research Foundation (FAPERJ), and the Brazilian Innovation Agency (FINEP).

Funding for the SDSS and SDSS-II has been provided by the Alfred P. Sloan Foundation, the Participating Institutions, the National Science Foundation, the U.S. Department of Energy, the National Aeronautics and Space Administration, the Japanese Monbukagakusho, the Max Planck Society, and the Higher Education Funding Council for England. The SDSS Web Site is <http://www.sdss.org/>.

The SDSS is managed by the Astrophysical Research Consortium for the Participating Institutions. The Participating Institutions are the American Museum of Natural History, Astrophysical Institute Potsdam, University of Basel, University of Cambridge, Case Western Reserve University, University of Chicago, Drexel Uni-

versity, Fermilab, the Institute for Advanced Study, the Japan Participation Group, Johns Hopkins University, the Joint Institute for Nuclear Astrophysics, the Kavli Institute for Particle Astrophysics and Cosmology, the Korean Scientist Group, the Chinese Academy of Sciences (LAMOST), Los Alamos National Laboratory, the Max-Planck-Institute for Astronomy (MPIA), the Max-Planck-Institute for Astrophysics (MPA), New Mexico State University, Ohio State University, University of Pittsburgh, University of Portsmouth, Princeton University, the United States Naval Observatory, and the University of Washington.

Guoshoujing Telescope (the Large Sky Area Multi-Object Fiber Spectroscopic Telescope LAMOST) is a National Major Scientific Project built by the Chinese Academy of Sciences. Funding for the project has been provided by the National Development and Reform Commission. LAMOST is operated and managed by the National Astronomical Observatories, Chinese Academy of Sciences.

Scientific software and databases used in this work include TOPCAT<sup>11</sup> (Taylor 2005), simbad and vixier from Strasbourg Astronomical Data Center (CDS)<sup>12</sup> and the following python packages: numpy, astropy, matplotlib, seaborn, scikit-learn.

## DATA AVAILABILITY

### REFERENCES

- Aggarwal C. C., 2015, Data Mining: The Textbook. Springer, Cham, doi:10.1007/978-3-319-14142-8
- Ahumada R., et al., 2020, *ApJS*, 249, 3
- Akras S., Guzman-Ramirez L., Leal-Ferreira M. L., Ramos-Larios G., 2019a, *ApJS*, 240, 21
- Akras S., Leal-Ferreira M. L., Guzman-Ramirez L., Ramos-Larios G., 2019b, *MNRAS*, 483, 5077
- Akras S., Guzman-Ramirez L., Gonçalves D. R., 2019c, *MNRAS*, 488, 3238
- Almeida-Fernandes F., et al., 2022, *MNRAS*, 511, 4590
- Barentsen G., et al., 2014, *MNRAS*, 444, 3230
- Benítez N., et al., 2014, arXiv e-prints, p. arXiv:1403.5237
- Bonoli S., et al., 2021, *A&A*, 653, A31
- Campeiro R. J. G. B., Moulaev D., Sander J., 2013, in Pei J., Tseng V. S., Cao L., Motoda H., Xu G., eds, Advances in Knowledge Discovery and Data Mining. Springer Berlin Heidelberg, Berlin, Heidelberg, pp 160–172
- Campello R., Moulaev D., Zimek A., Sander J., 2015, *A CM Transactions on Knowledge Discovery from Data*, 10, 1
- Cenarro A. J., et al., 2019, *A&A*, 622, A176
- Corradi R. L. M., Giannamico C., 2010, *A&A*, 520, A99
- Corradi R. L. M., et al., 2008, *A&A*, 480, 409
- Corradi R. L. M., Sabin L., Munari U., Cetrulo G., Englaro A., Angeloni R., Greimel R., Manpasó A., 2011, *A&A*, 529, A56
- Davies R. D., Elliott K. H., Meaburn J., 1976, *Mem. RAS*, 81, 89
- Drew J. E., et al., 2005, *MNRAS*, 362, 753
- Drew J. E., Greimel R., Irwin M. J., Sale S. E., 2008, *MNRAS*, 386, 1761
- Drew J. E., et al., 2014, *MNRAS*, 440, 2036
- Ester M., Kriegel H.-P., Sander J., Xu X., 1996, in Proc. of 2nd International Conference on Knowledge Discovery and Data Mining (KDD-96). pp 226–231
- Few D. J., 2008, PhD thesis, Department of Physics, Macquarie University, NSW 2109, Australia
- Fukugita M., Ichikawa T., Gunn J. E., Doi M., Shimasaku K., Schneider D. P., 1996, *AJ*, 111, 1748
- Gutiérrez-Soto L. A., et al., 2020, *A&A*, 633, A123
- Jacoby G. H., et al., 2010, *Publ. Astron. Soc. Australia*, 27, 156
- Jain A. K., Murty M. N., Flynn P. J., 1999, *ACM Comput. Surv.*, 31, 264

<sup>11</sup> <http://www.star.bristol.ac.uk/~mbt/topcat/>

<sup>12</sup> <https://cds.u-strasbg.fr/>

## No Comments.