





## Article

# The miniJPAS and J-NEP surveys: Machine learning for star-galaxy separation

Ana Paula Jeakel<sup>1,‡</sup> , Gabriel Vieira dos Santos<sup>2,‡</sup> , Valerio Marra<sup>2,3,4,\*</sup> , Rodrigo von Marttens<sup>5</sup> , Siddhartha Gurung-López<sup>6,7</sup>, Raul Abramo<sup>8</sup>, Jailson Alcaniz<sup>9</sup>, Narciso Benitez<sup>10</sup>, Silvia Bonoli<sup>11,12</sup>, Javier Cenarro<sup>12,13</sup>, David Cristóbal-Hornillos<sup>12</sup>, Simone Daflon<sup>9</sup>, Renato Dupke<sup>9</sup>, Alessandro Ederoclite<sup>12,13</sup>, Rosa M. González Delgado<sup>10</sup>, Antonio Hernán-Caballero<sup>12,13</sup>, Carlos Hernández-Monteagudo<sup>15</sup>, Jifeng Liu<sup>15</sup>, Carlos López-Sanjuan<sup>12,13</sup>, Antonio Marín-Franch<sup>12,13</sup>, Claudia Mendes de Oliveira<sup>16</sup>, Mariano Moles<sup>12</sup>, Fernando Roig<sup>9</sup>, Laerte Sodr   Jr.<sup>16</sup>, Keith Taylor<sup>17</sup>, Jes  s Varela<sup>12</sup>, H  ctor V  zquez Ram  o<sup>12,13</sup>, Jos   M. Vilchez<sup>10</sup>, Christopher Willmer<sup>18</sup>, Javier Zaragoza-Cardiel<sup>12</sup>

- <sup>1</sup> PPGFis, Universidade Federal do Esp  rito Santo, 29075-910, Vit  ria, ES, Brazil
- <sup>2</sup> Departamento de F  sica, Universidade Federal do Esp  rito Santo, 29075-910, Vit  ria, ES, Brazil
- <sup>3</sup> INAF – Osservatorio Astronomico di Trieste, via Tiepolo 11, 34131 Trieste, Italy
- <sup>4</sup> IFPU – Institute for Fundamental Physics of the Universe, via Beirut 2, 34151, Trieste, Italy
- <sup>5</sup> Instituto de F  sica, Universidade Federal da Bahia, 40210-340, Salvador, BA, Brazil
- <sup>6</sup> Observatori Astron  mic de la Universitat de Val  ncia, Ed. Instituts d'Investigaci  , Parc Cient  fic. C/ Catedr  tico Jos   Beltr  n, n2, 46980 Paterna, Valencia, Spain
- <sup>7</sup> Departament d'Astronomia i Astrof  sica, Universitat de Val  ncia, 46100 Burjassot, Spain
- <sup>8</sup> Departamento de F  sica Matem  tica, Instituto de F  sica, Universidade de S  o Paulo, 05508-090, S  o Paulo, SP, Brazil
- <sup>9</sup> Observat  rio Nacional, 20921-400, Rio de Janeiro, RJ, Brazil
- <sup>10</sup> Instituto de Astrof  sica de Andal  c  a - CSIC, 18080, Granada, Spain
- <sup>11</sup> Donostia International Physics Center (DIPC), 20018 San Sebasti  n, Spain
- <sup>12</sup> Centro de Estudios de F  sica del Cosmos de Arag  n (CEFCA), 44001, Teruel, Spain
- <sup>13</sup> Unidad Asociada CEFCA-IAA, CEFCA, Unidad Asociada al CSIC por el IAA y el IFCA, Plaza San Juan 1, 44001 Teruel, Spain
- <sup>14</sup> Instituto de Astrof  sica de Canarias, 38205, La Laguna, Tenerife, Spain
- <sup>15</sup> National Astronomical Observatory of China, Chinese Academy of Sciences, Beijing 100012, China
- <sup>16</sup> Departamento de Astronomia, Instituto de Astronomia, Geof  sica e Ci  ncias Atmosf  ricas, Universidade de S  o Paulo, 05508-090, S  o Paulo, SP, Brazil
- <sup>17</sup> Instruments4, La Canada Flintridge CA, 91011, USA
- <sup>18</sup> Steward Observatory, University of Arizona, Tucson AZ, 85721, USA
- \* Correspondence: valerio.marra@me.com
- ‡ These authors contributed equally to this work.

**Abstract:** We present a supervised machine learning classification of sources from the Javalambre Physics of the Accelerating Universe Astrophysical Survey (J-PAS) Pathfinder datasets: miniJPAS and J-NEP. Leveraging crossmatches with spectroscopic and photometric catalogs, we construct a robust labeled dataset comprising 14 594 sources classified into extended (galaxies) and point-like (stars and quasars) objects. We assess dataset representativeness using UMAP analysis, confirming broad and consistent coverage of feature space. An XGBoost classifier, with hyperparameters tuned using automated optimization, is trained using purely photometric data (60-band J-PAS magnitudes) and combined photometric and morphological features, with performance thoroughly evaluated via ROC and purity–completeness metrics. Incorporating morphology significantly improves classification, outperforming the baseline classifications available in the catalogs. Permutation importance analysis reveals morphological parameters, particularly concentration, normalized peak surface brightness, and PSF, alongside photometric features around 4000 and 6900   , as crucial for accurate classifications. We release a value-added catalog with our models for



Received:

Revised:

Accepted:

Published:

**Citation:** Jeakel, Vieira dos Santos et al. The miniJPAS and J-NEP surveys: Machine learning for star-galaxy separation.

*Galaxies* **2025**, *1*, 0.

<https://doi.org/>

**Copyright:**    2025 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

star-galaxy classification, enhancing the utility of miniJPAS and J-NEP for subsequent cosmological and astrophysical analyses.

**Keywords:** Astronomy; Galaxies; Stars; Data Analysis; Machine Learning

## 1. Introduction

The Javalambre Physics of the Accelerating Universe Astrophysical Survey<sup>1</sup> (J-PAS) is a ground-based imaging survey that will map thousands of square degrees using a unique set of 56 filters, providing a low-resolution spectrum for every pixel of the sky ( $R \sim 60$ ),<sup>2</sup> without the target selection biases or fiber collision issues that typically affect spectroscopic surveys. Designed and optimized to achieve highly accurate photometric redshifts, its scientific objectives span from cosmology to galaxy evolution [1].

Here we focus on data from the Pathfinder camera, which was used to test the telescope performance and execute the first scientific operations. The camera features a  $9k \times 9k$  CCD, with a  $0.3\text{deg}^2$  field of view and  $0.225$  arcsec pixel size. This effort led to the miniJPAS [1] and J-NEP [2] surveys, which covered  $1\text{deg}^2$  of the AEGIS field (four pointings around  $(\text{RA}, \text{Dec}) = (215^\circ, 53^\circ)$ ) and  $0.29\text{deg}^2$  of the North Ecliptic Pole (one pointing centred at  $(\text{RA}, \text{Dec}) = (260.6^\circ, 65.8^\circ)$ ), respectively. Both surveys use the 56 J-PAS filters together with the  $u$ ,  $g$ ,  $r$ , and  $i$  broad bands<sup>3</sup>, for a total of 60 bands. However, J-NEP adopted longer exposure times, reaching between  $0.5$ – $1.0$  magnitudes deeper than miniJPAS, which itself reaches the target depth of J-PAS—namely,  $m_{\text{AB}} \approx 24$  in the broad-band filters ( $5\sigma$  in a  $3''$  aperture).<sup>4</sup>

While miniJPAS targeted the well-studied Extended Groth Strip field, an area observed by many experiments, J-NEP fully covers the time-domain field of the James Webb Space Telescope’s North Ecliptic Pole [3], a region chosen for its low galactic extinction, absence of bright stars, and continuous visibility from Webb’s northern hemisphere.

Our main goal here is to provide a machine learning-based classification of J-NEP sources as an alternative to the Bayesian stellar and galactic loci classifier (SGLC [4], in `jnep.StarGalClass`<sup>5</sup>). This aims both to validate the quality of the SGLC, possibly outperforming it, and to offer an independent classification that can be used to assess the robustness of subsequent scientific analyses. In particular, our approach incorporates the full photometric information, which is crucial for faint galaxies whose morphology closely resembles that of stars.

The SGLC is based on morphological information, specifically on the bimodal distribution observed in the concentration–magnitude diagram, which distinguishes compact point-like objects from extended sources. López-Sanjuan et al. [4] modeled these two populations to compute the probability that a given source is compact or extended. This model, combined with appropriate priors, was then used to estimate the Bayesian probability that a source is a star or a galaxy [1].

In Baqui et al. [5], we obtained the machine learning classification of miniJPAS sources, available in `minijpas.StarGalClass`. Here, taking advantage of new spectroscopic data

<sup>1</sup> <https://www.j-pas.org>

<sup>2</sup> The wavelength resolution  $R_\lambda$  is defined as  $R_\lambda = \lambda/\Delta\lambda$ , so that  $R \sim 60$  is the approximate value in the intermediate wavelength range in the J-PAS filter system.

<sup>3</sup> The J-PAS survey itself employs only the  $i$  broad band in addition to the 56 narrow bands.

<sup>4</sup> Here  $m_{\text{AB}}$  denotes AB magnitudes, defined such that a source with constant flux density of  $3631$  Jy has  $m_{\text{AB}} = 0$ .

<sup>5</sup> Throughout this work, for reproducibility, we explicitly reference the table names available at [j-pas.org/datareleases](https://j-pas.org/datareleases).

and an improved methodology, we revise the classification of miniJPAS sources and provide a machine learning classification of J-NEP sources. This paper is organized as follows. Section 2 describes the training data, Section 3 presents the feature configuration, Section 4 details the classification methodology and examines the representativeness of the labeled set, and Section 5 reports the results. Finally, conclusions are drawn in Section 6.

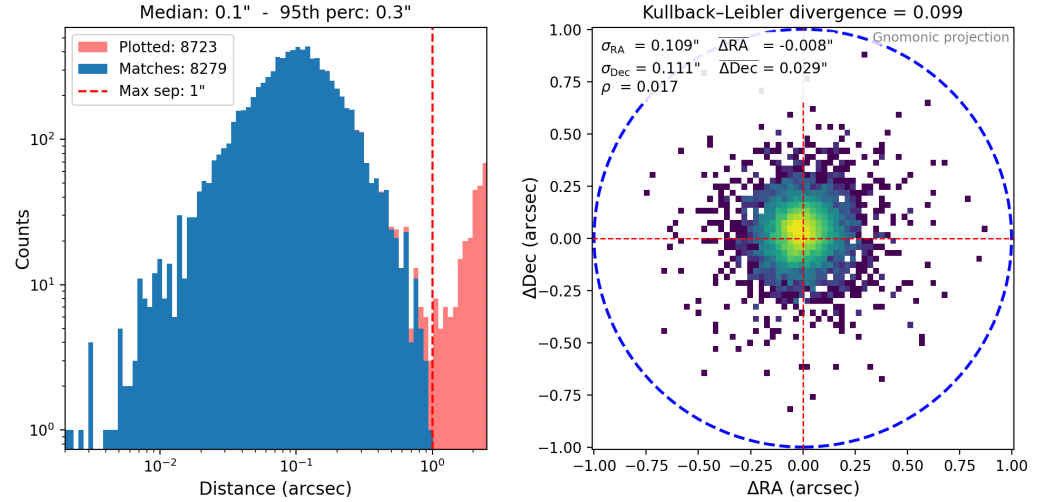
## 2. Data

In this work, we primarily used the `jnep.MagABDual0bj` and `minijpas.MagABDual0bj` tables from the data releases J-NEP-PDR202107 and MINIJ-PAS-PDR201912, respectively, available at [j-pas.org/datareleases](http://j-pas.org/datareleases). These catalogs provide photometric and morphological measurements for all sources detected in the *r* band, with forced photometry in the remaining filters. The *r* band was adopted as the detection image owing to its depth, image quality, and completeness, as detailed in Bonoli et al. [1]. All queries employed in this work are listed in Appendix A.

To ensure data quality, we required objects to satisfy `mask_flags=0`, which excludes detections affected by image boundaries, bright-star halos, or artifacts. For the `flags` field, which encodes issues encountered during photometric extraction, we retained objects with values 0 (clean, no problems detected), 1 (neighbor, the object is close to bright neighbors or contains bad pixels), 2 (blended, originally detected together with another source but successfully deblended), 4 (saturated, containing one or more pixels at or near saturation), and 6 (blended + saturated, a combination of the previous two cases). These situations can still yield reliable photometry and remain relevant for galaxy and cluster science, whereas more severe conditions (e.g., truncation at image edges, incomplete aperture measurements, or failed deblending) were excluded. This yields a total of 25,198 objects for miniJPAS and 7,290 objects for J-NEP within the range  $15 < r < 23.5$ .

Since training supervised machine learning models requires ground-truth labels, we cross-match miniJPAS and J-NEP data with surveys that provide reliable photometric or spectroscopic classifications into extended (galaxies) and point-like (stars and quasars) sources. This is a crucial step, as the performance of a supervised ML model depends critically on the quality and representativeness of the adopted training set. The surveys used for crossmatching in this work are:

- **SDSS DR12** (Section 2.1): the final Data Release (DR12) of the third-generation Sloan Digital Sky Survey [SDSS-III, 6], which collected imaging and spectroscopy from 2008 to 2014, providing wide-area multiband photometry and spectra of millions of sources.
- **Gaia EDR3** (Section 2.2): the Early Data Release 3 of the Gaia mission [7], delivering high-precision astrometry, broad-band photometry, and radial velocities for over 1.5 billion stars in the Milky Way.
- **DEEP3 DR4** (Section 2.3): spectroscopic redshift survey conducted with DEIMOS on the Keck II telescope in the Extended Groth Strip (EGS) [8]. It provides secure galaxy redshifts and spectra, forming part of the AEGIS survey program.
- **Binospec** (Section 2.4): spectroscopic data from Binospec [9], a high-throughput optical spectrograph on the 6.5-m MMT, covering 370–1000 nm with wide field of view, used to obtain redshifts for faint galaxies.
- **DESI DR1** (Section 2.5): the first data release of the Dark Energy Spectroscopic Instrument [10], based on the first 13 months of the main survey. It provides high-confidence redshifts for 18.7 million objects (13.1M galaxies, 1.6M quasars, 4.0M stars) over more than 9000 deg<sup>2</sup>, making it the largest extragalactic spectroscopic sample to date and a key resource for cosmology and large-scale structure studies.



**Figure 1.** Crossmatch between DEEP3 DR4 and miniJPAS sources with  $r < 23.5$ .

- **HSC-SSP PDR2** (Section 2.6): the second public data release of the Hyper Suprime-Cam Subaru Strategic Program [HSC-SSP, 11], providing deep, wide-field imaging in five broad bands (*grizy*), with exquisite seeing from the Subaru telescope.

### 2.1. Crossmatch with SDSS DR12

We used the table `minijpas.xmatch_sdss_dr12`, which contains the crossmatch with SDSS DR12 and includes both the photometric classification `class`, expected to be reliable for  $15 < r < 20$  [5]<sup>6</sup>, and the spectroscopic classification `spC1`. The photometric classification yields 802 galaxies and 1228 stars within  $15 < r < 20$ , for a total of 2030 labeled sources. The spectroscopic classification provides 355 galaxies and 220 point-like sources (94 stars and 124 quasars) for  $r < 23.5$ , for a total of 573 labeled sources.

We also used the analogous table `jnep.xmatch_sdss_dr12`. The photometric classification yields 156 galaxies and 724 stars within  $15 < r < 20$ , for a total of 880 labeled sources. No spectroscopic labels were available, since J-NEP lies outside the SDSS footprint.

### 2.2. Crossmatch with Gaia EDR3

We used the tables `minijpas.xmatch_gaia_edr3` and `jnep.xmatch_gaia_edr3`, which contain the crossmatches with Gaia EDR3<sup>7</sup>. Stars were identified as objects with a parallax measured at a significance greater than  $3\sigma$ , i.e., `parallax_over_error` > 3. This selection yields 989 stars for miniJPAS and 570 stars for J-NEP.

### 2.3. Crossmatch between miniJPAS and DEEP3 DR4

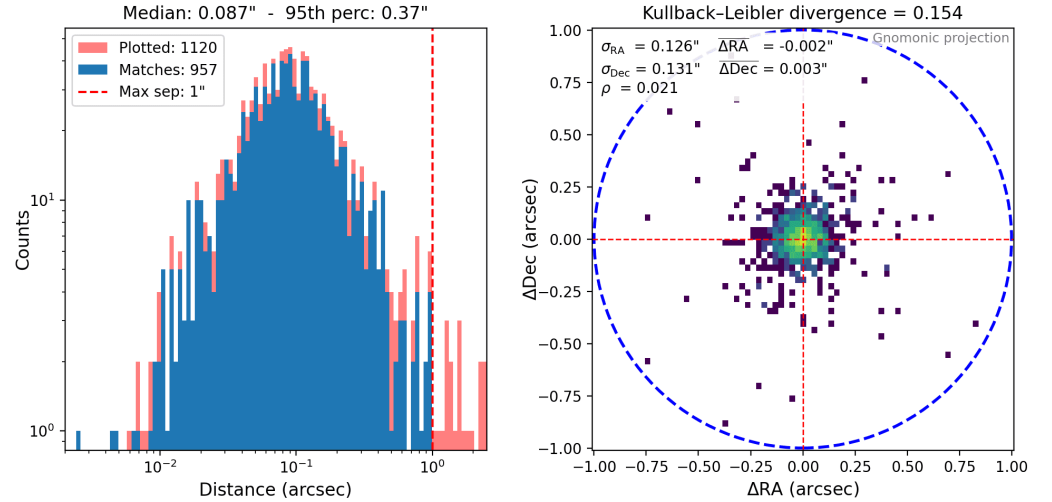
We crossmatched miniJPAS with  $r < 23.5$  with the DEEP3 DR4 catalog `alldeep.egs.2012jun13`<sup>8</sup>. In the case of duplicates, we retained the entry with the highest `ZQUALITY`, and, if necessary, the one with the smallest redshift error `Z_ERR`. The crossmatch was performed using the Python module `onexmatch`<sup>9</sup> [12], adopting a maximum separation of 1 arcsec and an ambiguity threshold of 0.5 arcsec, yielding 8279 matched objects as shown in Figure 1.

<sup>6</sup> We verified this by comparing DESI, DEEP3, SDSS (spectroscopic), and Gaia sources (taken as ground truth) against the SDSS photometric classification.

<sup>7</sup> As we are interested in parallax, Gaia EDR3 is essentially equivalent to the newer Gaia DR3.

<sup>8</sup> <https://sites.uci.edu/deep3/deep3zcat/>

<sup>9</sup> <https://github.com/valerio-marra/onexmatch>



**Figure 2.** Crossmatch between Binospec and J-NEP sources with  $r < 23.5$ .

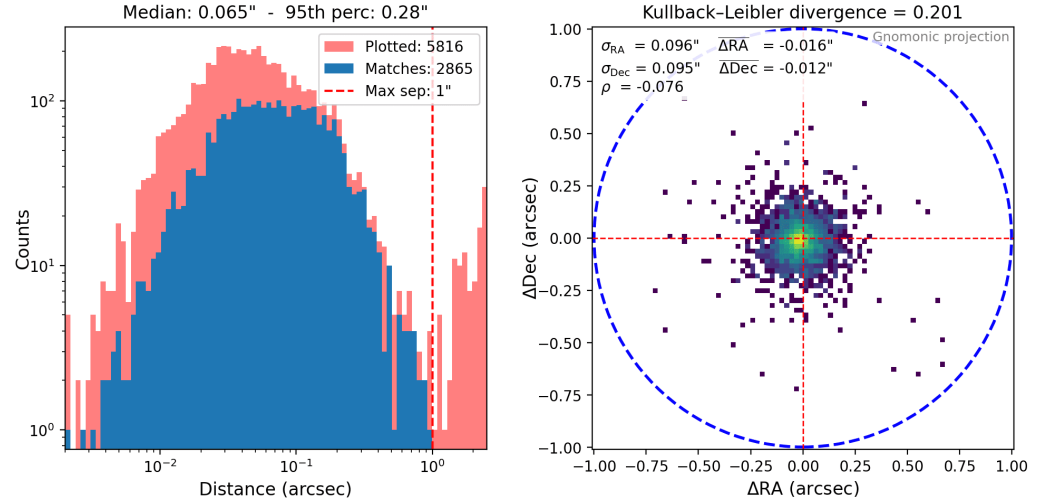
As shown by the diagnostic metrics, the positional offsets in Right Ascension and Declination are well described by a zero-mean isotropic 2D Gaussian whose standard deviation is equal to the average of the measured  $\sigma_{\Delta RA}$  and  $\sigma_{\Delta DEC}$ ; the resulting Kullback–Leibler divergence is  $\lesssim 0.1$ . The adoption of an ambiguity threshold of 0.5 arcsec removed 10 ambiguous DEEP3 matches associated with duplicated miniJPAS sources, whose positional differences were below this threshold—about twice the pixel scale and comparable to the PSF size. See Appendix C for more details.

After crossmatching, we retained only sources with  $ZQUALITY \geq 3$ , yielding 6583 galaxies, 433 active galactic nuclei (AGNs), and 5 stars. Unlike quasars (QSOs), many AGNs exhibit extended morphologies, which introduces ambiguity in the classification process. To mitigate this, we restricted point-like sources to objects with  $morph\_prob\_star > 0.9$  and redshift  $Z > 0.1$ , and stars to those with  $morph\_prob\_star > 0.9$  and  $Z < 0.01$ . The  $morph\_prob\_star$  values are provided by the `jnep.StarGalClass` table, which estimates stellar probability from morphological information [4]. This selection resulted in 36 point-like sources (35 AGNs and 1 star). This filtering step was necessary to construct a clean labeled dataset. Our study focuses on a binary star-galaxy classification, rather than a three-class star-AGN/QSO-galaxy scheme, which would require a significantly larger and more balanced AGN sample.

#### 2.4. Crossmatch between J-NEP and Binospec data

We crossmatched J-NEP sources with  $r < 23.5$  with data that was obtained by Christopher Willmer using Binospec at the MMT observatory. Details about the dataset are provided in Hernán-Caballero et al. [2]. The crossmatch was performed using the `onexmatch` code, adopting a maximum separation of 1 arcsec and an ambiguity threshold of 0.5 arcsec, and yielded 957 matched sources, as shown in Figure 2. The adoption of an ambiguity threshold of 0.5 arcsec removed 67 ambiguous Binospec matches associated with duplicated miniJPAS sources, mostly with positional differences below 0.2 arcsec (Figure C1).

After crossmatching, we retained only the 726 objects with  $ZQUALITY \geq 3$  and  $Z\_ERR < 0.01$ . Since the Binospec catalog does not provide spectroscopic classifications, we assigned class labels based on morphology and redshift: objects with  $morph\_prob\_star > 0.9$  and  $Z < 0.01$  were classified as stars, and those with  $morph\_prob\_star < 0.1$  and  $Z > 0.01$  as galaxies. This yielded 53 stars and 593 galaxies, for a total of 646 objects. There are 69 objects with  $morph\_prob\_star > 0.1$  and  $Z > 0.01$ . These could correspond either to QSOs or to compact galaxies; we exclude this small set of sources from the catalog.



**Figure 3.** Crossmatch between DESI DR1 and miniJPAS sources with  $r < 23.5$ .

**Table 1.** Confusion matrices in three  $r$ -band magnitude bins, comparing DESI, DEEP3, SDSS (spectroscopic) and Gaia sources (used as ground truth) against HSC-SSP classifications.

	18.5–20.5		20.5–22.5		22.5–23.5	
	pred. stars	pred. gal.	pred. stars	pred. gal.	pred. stars	pred. gal.
true stars	79	10	92	21	6	4
true gal.	0	385	5	1523	9	2350

### 2.5. Crossmatch between miniJPAS and DESI DR1

We crossmatched miniJPAS sources with  $r < 23.5$  with the DESI DR1 catalog.<sup>10</sup> The crossmatch was performed using the `onexmatch` code, adopting a maximum separation of 1 arcsec and an ambiguity threshold of 0.5 arcsec, and yielded 2865 matches, as shown in Figure 3. The adoption of an ambiguity threshold of 0.5 arcsec removed 1232 ambiguous DESI matches associated with duplicated miniJPAS sources, mostly with positional differences below 0.05 arcsec (Figure C1).

After crossmatching, as suggested in [10], only sources with `ZCAT_PRIMARY=True`, `ZWARN=0` and `DELTACHI2>25` were retained, for a total of 2640 objects. Of these, 2307 are galaxies and 333 are point-like sources (124 QSOs and 209 stars).

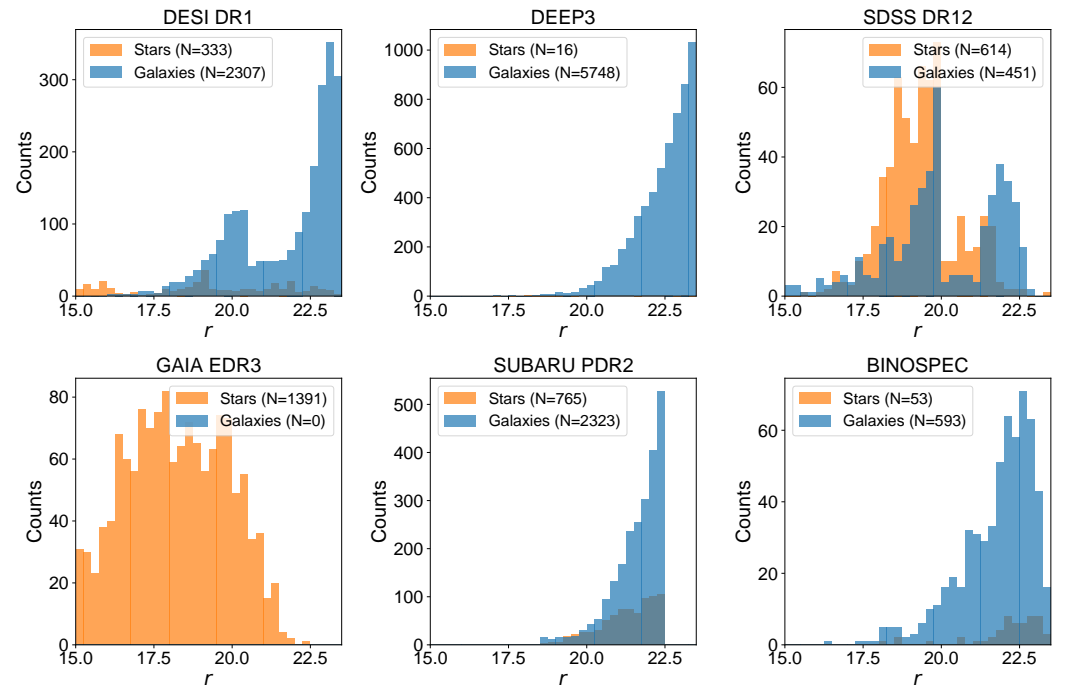
### 2.6. Crossmatch between miniJPAS and HSC-SSP PDR2

We used the photometric classification provided by `x_extendedness_value` (with `x=g, r, i, z, y`) from the catalog `minijpas.xmatch_subaru_pdr2`, which contains the crossmatch with the HSC-SSP PDR2 catalog. To ensure consistency, we retained only the 11 960 objects for which all five `x_extendedness_value` measurements agree.

We tested the HSC-SSP classification in the range  $18.5 < r < 23.5$  against DESI, DEEP3, SDSS (spectroscopic) and Gaia sources, used as ground truth; see Table 1. We find that for  $r > 22.5$  a non-negligible fraction of stars, though small in absolute numbers, are misclassified as galaxies. To minimize this effect, we restrict the HSC-SSP labels to the range  $18.5 < r < 22.5$ , where the bright limit avoids saturation. This selection yields 4406 galaxies and 946 stars, for a total of 5352 objects.

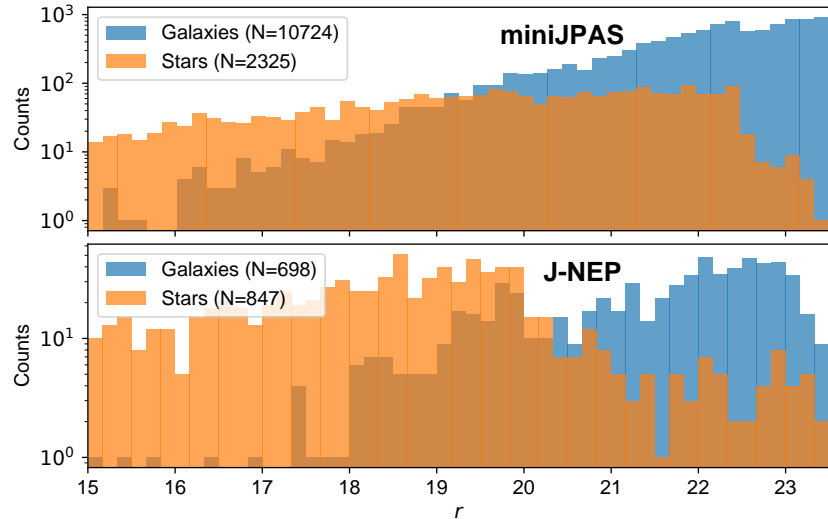
Crossmatch Source	miniJPAS (galaxies/stars)	J-NEP (galaxies/stars)
DESI DR1	2307/333 (2307/333)	–
Binospec	–	593/53 (593/53)
DEEP3 DR4	5748/16 (6583/36)	–
SDSS DR12 (spectro)	189/124 (355/218)	–
Gaia EDR3	–/826 (–/989)	–/565 (–/570)
HSC-SSP PDR2	2323/765 (4406/946)	–
SDSS DR12 (photo)	157/261 (802/1228)	105/229 (156/724)
Total	10724/2325	698/847

**Table 2.** Summary of external catalogs used to crossmatch and label sources in miniJPAS and J-NEP. Values in parentheses indicate the total number of matched sources before duplicate removal. When a source appears in multiple catalogs, labels are assigned based on a priority order corresponding to the table’s top-down sequence, with DESI given the highest priority. As a result, the final number of adopted labels is smaller than the per-catalog match count.



**Figure 4.**  $r$ -band magnitude distributions of labeled stars and galaxies across the catalogs used for crossmatching. SDSS spectroscopic and photometric classifications are separated at  $r = 20$ .





**Figure 5.**  $r$ -band magnitude distributions of labeled stars and galaxies in miniJPAS and J-NEP.

### 2.7. Final labeled set

The crossmatched catalogs from the previous sections partially overlap. When dealing with duplicate J-PAS sources, we assigned class labels according to the following priority order, reflecting the reliability of each classification: DESI, Binospec, DEEP3, SDSS (spectroscopic classification), Gaia, HSC-SSP, and SDSS (photometric classification). Table 2 summarizes the labeled dataset, comprising 14 594 unique sources. The corresponding  $r$ -band magnitude distributions of the retained sources from each survey are shown in Figure 4.

Figure 5 shows the  $r$ -band magnitude distributions of stars and galaxies for the miniJPAS and J-NEP surveys. In both surveys, galaxies dominate at  $r \gtrsim 21$ , while stars are more prevalent at  $r \lesssim 18.5$ . This results in a class imbalance in the training set, which can bias both the performance of the classifier and the interpretation of its metrics. To avoid misleading conclusions from this imbalance, we evaluate performance using purity-completeness curves, which provide a balanced view of classification quality independent of the relative proportions of stars and galaxies.

Finally, Figure 6 shows the histogram of the magnitude differences between the SDSS and J-PAS  $r$ -band total flux magnitudes for point sources. This provides an additional crosscheck of the quality of the crossmatch. We restrict the analysis to point sources because J-PAS does not provide a total-flux model magnitude for extended objects, as `MAG_AUTO` is measured within an aperture of two Kron radii and therefore misses a non-negligible fraction of the galaxy flux in the outskirts.

As stated earlier, we perform a binary classification into extended and point-like sources. Point-like sources include both stars and QSOs, which are morphologically similar at the J-PAS resolution. However, the number of spectroscopically confirmed QSOs in our sample is insufficient to support a reliable three-class classification [13–17] that would exploit the low-resolution J-PAS spectra. As is customary, we refer to extended sources as galaxies and to point-like sources as stars throughout this paper.

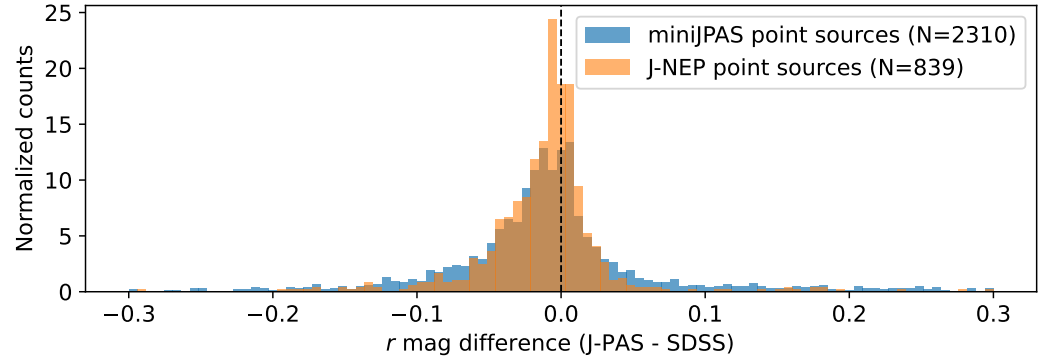
## 3. Feature configuration

### 3.1. Observational features

To characterize the observing conditions and data quality, we included four observational features:

<sup>10</sup> `zall-pix-iron.fits` at <https://data.desi.lbl.gov/public/dr1/spectro/redux/iron/zcatalog/v1>





**Figure 6.** Histogram of the differences between the SDSS  $r$ -band model magnitudes and the J-PAS  $r$ -band `MAG_APER_COR_3` magnitudes for the stars of the final labeled set. The latter corresponds to the total magnitude obtained from a 3 arcsec diameter aperture corrected for aperture effects.

1. Point Spread Function (PSF) of the tile,
2. Galactic extinction  $E(B - V)$ ,
3.  $E(B - V)$  error,
4. Object `flags` indicating observational issues.

The PSF is characterized as the mean full width at half maximum (FWHM) of point-like sources in each tile, estimated with `PSFEx` [18]. The color excess  $E(B - V)$  is the integrated reddening due to Milky Way dust along the line of sight, provided at the position of each source by the Bayestar17 three-dimensional dust map [19,20].

### 3.2. Photometric features

We used a total of 120 photometric features derived from the 60 J-PAS bands. Specifically, we considered `MAG_AUTO` magnitudes and their associated uncertainties in each of the 60 bands from `MagABDualObj`. Bands without detection were assigned a magnitude of 99, while the magnitude errors encapsulate valuable information about the strength and reliability of each detection.

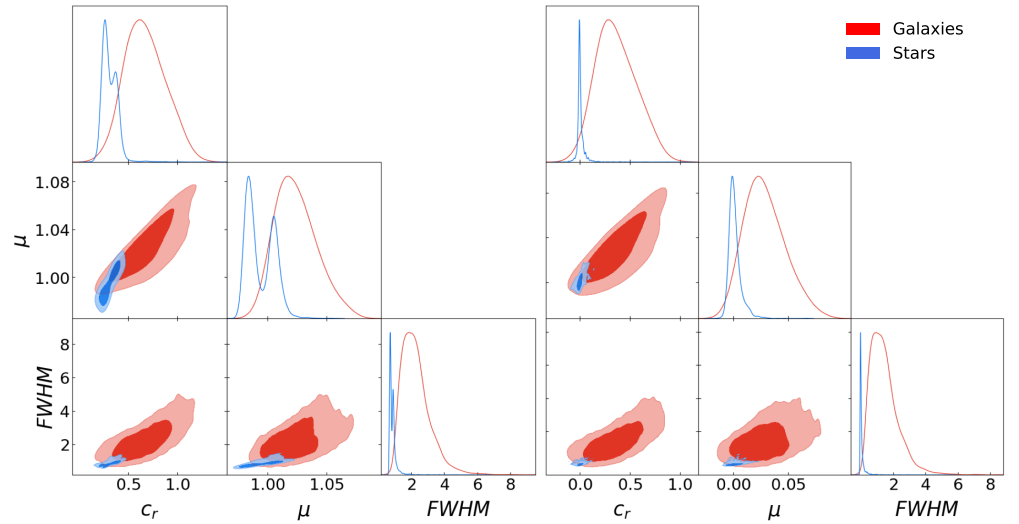
Although `MagABDualPointSources` provides PSF-corrected photometry (`MAG_APER_COR_3_0`) optimized for point-like objects—shown to outperform `MAG_AUTO` for spectral energy distribution (SED) fitting and photometric redshift estimation—our classification tests indicate a different trend. For the photometry-only model, `MAG_AUTO` yields better results, with ROC AUC = 0.992 and star purity-completeness AUC = 0.984, compared to 0.984 and 0.965 for `MAG_APER_COR_3_0`.<sup>11</sup> For galaxies, the improvement is smaller but still in favor of `MAG_AUTO` (AUC = 0.997 vs. 0.994). When morphology is included, the two definitions perform nearly identically—0.1% better performance for `MAG_APER_COR_3_0`. Given the superior performance of `MAG_AUTO` in the photometry-only case and its comparable performance when morphology is included, we adopt `MAG_AUTO` as our default choice in this work.

### 3.3. Morphological features

We adopted four morphological parameters derived from columns relative to the detection  $i$  band, available in the `MagABDualObj` catalog:

1. *Ellipticity*  $A/B = A\_WORLD/B\_WORLD$ , where `A_WORLD` and `B_WORLD` are the root mean square values of the light distribution along the major and minor axes of the source, respectively.

<sup>11</sup> ROC stands for receiver operating characteristic, and AUC denotes the area under the curve.



**Figure 7.** Distribution of morphological features before (left) and after (right) subtracting the tile-by-tile median stellar values. For each panel, the dark and light shades show the 68% and 95% contours of the distributions, respectively. The original distributions (left) are affected by PSF-induced multimodality across tiles.

2. *Concentration*  $c_r = \text{MAG\_APER\_1\_5} - \text{MAG\_APER\_3\_0}$  (in the  $r$ -band), calculated from magnitudes measured within circular apertures of 1.5 and 3.0 arcsec.
3. *Full Width at Half Maximum*  $\text{FWHM\_WORLD}$ , defined assuming a Gaussian intensity distribution.
4. *Normalized peak surface brightness*  $\mu = \text{MU\_MAX}/\text{MAG\_APER\_3\_0}$  (in the  $r$ -band), where  $\text{MU\_MAX}$  represents the peak surface brightness above the local background.

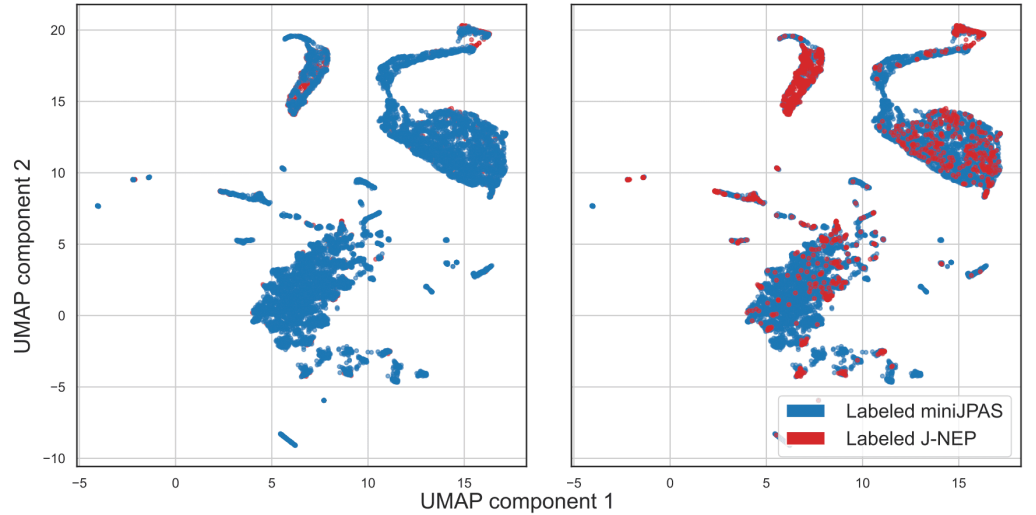
The spatial variation of the PSF across the J-PAS tiles introduces multimodality into the distributions of concentration, FWHM, and peak surface brightness. To correct for this effect and achieve uniform morphological features, we normalize each parameter by subtracting the corresponding median stellar value computed within the same tile:

$$X_i \rightarrow X_i - \langle X \rangle_{\text{stars, tile}}, \quad (1)$$

where  $X$  denotes concentration, FWHM, or peak surface brightness. This correction effectively removes the multimodal behavior, as illustrated in Figure 7, which compares the distributions before and after correction. In addition, the separation between stellar and galactic populations becomes more pronounced following this normalization.

We do not apply this correction to the ellipticity, as it does not display multimodal features and can take spurious values for stars. Due to the finite pixel scale, a narrow PSF may lead SExtractor to assign  $A/B \neq 1$  to unresolved sources, introducing artificial shape distortions caused by limitations in the PSF fitting procedure. See Bonoli et al. [1, App. C] for the SExtractor configuration adopted.

We use the normalized morphological parameters for the representativeness analysis presented in the next section. However, for the classification task, we retain the uncorrected catalog values. Our tests show that including the PSF as an additional feature is sufficient to account for its spatial variations and prevent confusion, while the normalization procedure introduces additional complexity to the pipeline without yielding a significant improvement in classification performance.



**Figure 8.** UMAP projection of labeled miniJPAS and J-NEP samples, showing consistent coverage across feature space. Left (right) panel shows miniJPAS (J-NEP) sources plotted on top to enhance visibility.

### 3.4. Feature sets

We perform three separate classifications. The first uses the full set of 128 features, which includes 60 magnitudes with their associated errors, four morphological parameters, the PSF information, Galactic extinction with its error, and detection flags. The second is restricted to 123 features, consisting of the 60 magnitudes with their errors, Galactic extinction with its error, and detection flags. We refer to these two configurations as ‘morphology+photometry’ and ‘photometry only’, respectively. The rationale for testing the latter is that morphological parameters may be biased in tiles produced by combining exposures taken under different observing conditions; in such cases, the tile-averaged PSF may not adequately represent individual sources, making a photometry-only classification potentially more robust.

In addition, we evaluate a third model based solely on morphological information, comprising the four morphological parameters, PSF, and detection flags. This setup serves as a benchmark to compare with the other models and, in particular, with the SGLC classifier, thereby quantifying the performance gain provided by machine learning.

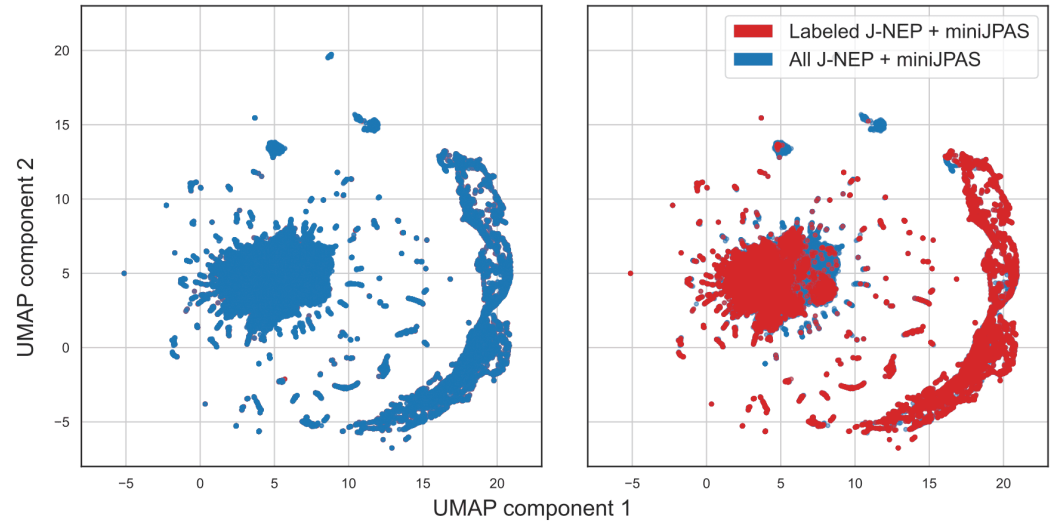
## 4. Supervised classification

### 4.1. Representativeness

The miniJPAS and J-NEP surveys were conducted using the same instrument and filter system, though under slightly different observing conditions. To evaluate the reliability and generalization capacity of the labeled datasets assembled in Section 2, we assess their representativeness along two complementary directions:

1. the internal consistency between the labeled miniJPAS and labeled J-NEP subsets;
2. the external coverage of the full miniJPAS and J-NEP datasets by the combined labeled set.

For this purpose, we adopt the Uniform Manifold Approximation and Projection (UMAP) algorithm [21], a non-linear dimensionality reduction technique designed to preserve both local and global data structure in low-dimensional projections. Importantly, our use of UMAP is diagnostic: we do not attempt classification in the embedded space, but rather evaluate whether the labeled samples span the same feature manifold as the full catalogs.



**Figure 9.** UMAP projection of the combined labeled set and full miniJPAS and J-NEP catalogs, showing representative coverage. Left (right) panel shows miniJPAS (J-NEP) sources plotted on top to enhance visibility.

We apply UMAP to the 124 photometric and morphological features described in Section 3, excluding the observational features listed in Section 3.1, as these are non-object-specific or discrete and could introduce spurious clustering. Figure 8 shows the two-dimensional UMAP projection of the labeled miniJPAS and J-NEP sources. The J-NEP subset is largely embedded within the distribution of miniJPAS, indicating that the two labeled samples are consistent and can be reliably combined for model training. Figure 9 compares the combined labeled dataset against the full miniJPAS and J-NEP populations.<sup>12</sup> The labeled set broadly spans the full feature space, confirming that it is sufficiently representative to serve as a training base for supervised classification of the entire catalog. It is important to note, however, that while the labeled set spans the overall feature space, this does not guarantee that labeled stars (or galaxies) fully cover the star (or galaxy) subspaces of the full catalogs.

#### 4.2. TPOT pipeline for XGBoost optimization

Based on the results of von Marttens et al. [13], who compared several machine learning algorithms, we adopt the eXtreme Gradient Boosting algorithm [XGBoost, 22] for classification. Given the relatively modest size of our labeled dataset (fewer than 15 000 objects), XGBoost is a suitable choice. While neural networks may outperform tree-based models for substantially larger datasets, they are unlikely to offer significant advantages in our case. XGBoost is an ensemble learning method based on decision trees and gradient boosting. It iteratively adds new trees to correct the errors of the existing ensemble while minimizing a regularized objective function, which balances prediction accuracy and model complexity. This process results in a robust and efficient classifier that is well-suited for structured tabular data. For technical details, see Chen and Guestrin [22].

The final hyperparameter configuration was selected using the Tree-based Pipeline Optimization Tool [TPOT, 23],<sup>13</sup> an automated machine learning framework that optimizes machine learning pipelines through genetic programming. The number of pipelines evaluated by TPOT depends on three key hyperparameters: *generations*, *population size*, and *offspring*

<sup>12</sup> Note that the UMAP projection differs from that in Figure 8, so the two plots are not directly comparable in terms of their coordinate positions.

<sup>13</sup> <https://epistasislab.github.io/tpot/latest/>

size. These respectively control the number of optimization iterations, the number of top-performing pipelines retained in each generation, and the number of new pipelines generated per generation.

The total number of pipelines analyzed is given by *Population Size + Generations*  $\times$  *Offspring Size*. In this work, we set all three hyperparameters to 200, resulting in the evaluation of 40,200 distinct pipelines. The final hyperparameter configuration is listed in Appendix B. Model selection was performed exclusively within the training set, comprising 80% of the combined miniJPAS and J-NEP labeled samples. TPOT used 5-fold cross-validation within this training set to optimize pipeline architectures and hyperparameters. After optimization, the best-performing pipeline was retrained on the same 80% split, while the remaining 20% was held out as an independent test set. This test set was used only once for final evaluation, ensuring that the reported performance provides an unbiased estimate of the model’s generalization ability.

#### 4.3. Performance metrics

The XGBoost classifier outputs a probabilistic score between 0 (star) and 1 (galaxy), requiring the choice of a classification threshold  $p_{\text{cut}}$  to separate the two classes. Once  $p_{\text{cut}}$  is specified, the model’s performance can be evaluated using a  $2 \times 2$  confusion matrix. From its entries, we derive standard classification metrics such as the receiver operating characteristic (ROC) curve and the purity–completeness curve, which provide graphical assessments of classifier performance.

The ROC curve is a parametric plot of the true positive rate (TPR) versus the false positive rate (FPR) as a function of  $p_{\text{cut}}$ :

$$\text{TPR}(p_{\text{cut}}) = \frac{\text{TP}}{\text{TP} + \text{FN}}, \quad \text{FPR}(p_{\text{cut}}) = \frac{\text{FP}}{\text{FP} + \text{TN}}, \quad (2)$$

where TP, FN, FP, and TN are the entries of the confusion matrix. TPR corresponds to recall, commonly referred to as completeness in astronomy. A commonly used scalar summary of the ROC curve is the area under the curve (AUC), which ranges from 0.5 (random classifier) to 1 (perfect classifier).

To account for class imbalance in the training set, we also consider the purity–completeness curve. It is a parametric plot of the purity (or precision) versus the completeness (or recall) as a function of  $p_{\text{cut}}$ :

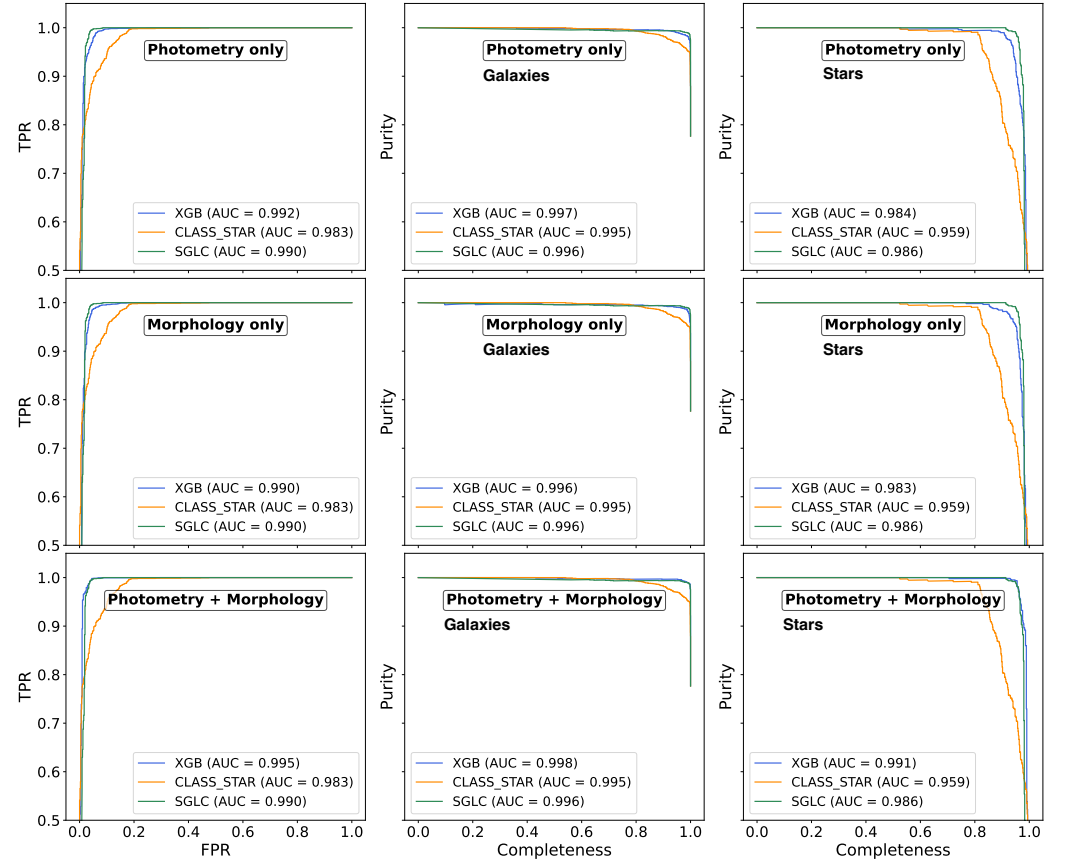
$$\text{Purity}(p_{\text{cut}}) = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{Completeness}(p_{\text{cut}}) = \frac{\text{TP}}{\text{TP} + \text{FN}}. \quad (3)$$

A compact summary of this curve is given by the average precision (AP), defined as the area under the purity–completeness curve. AP values range from 0 to 1, with higher values indicating better classification performance.

## 5. Results

### 5.1. Metrics

Figure 10 shows the ROC curves for the three XGBoost models: photometry only (123 features, top), morphology only (6 features, middle), and morphology+photometry (128 features, bottom). For comparison, we also include the performance of SGLC and CLASS\_STAR, the neural-network classifier implemented in SExtractor [24]. The photometry-only model outperforms CLASS\_STAR and achieves performance comparable to, and slightly better than, SGLC. It also surpasses the morphology-only model, which performs as SGLC. The morphology+photometry model is the best overall, yielding the highest ROC AUC and thus the most accurate classifications. The same figure also shows the purity and



**Figure 10.** ROC and purity-completeness curves for XGB, CLASS\_STAR, and SGLC.

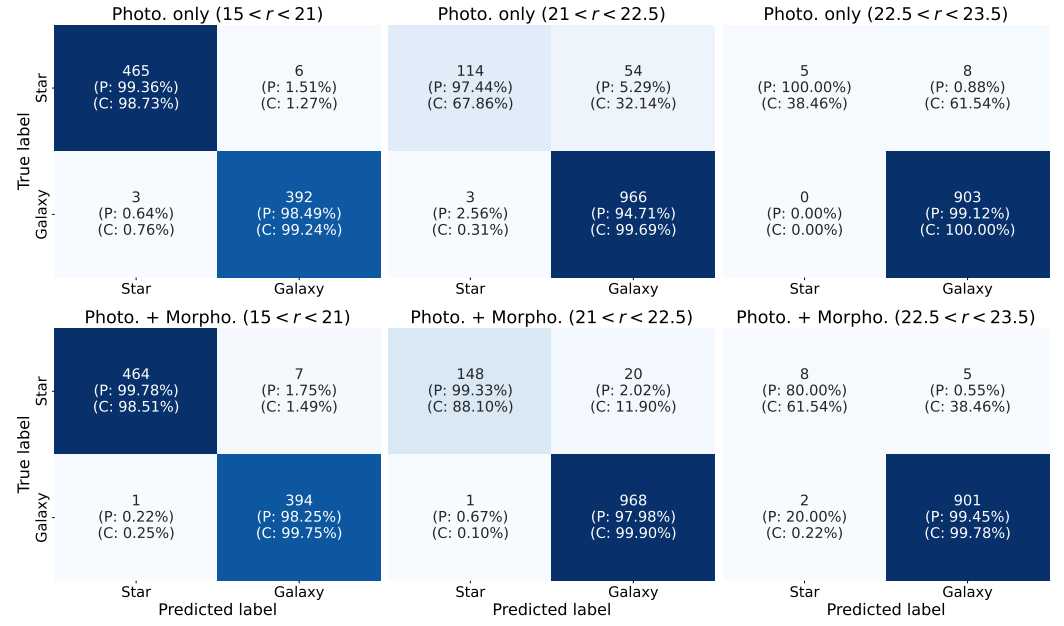
completeness curves for galaxies (second column) and stars (third column). Despite the class imbalance in the training set, stars are well classified, with only a modest reduction in average precision (area under the purity-completeness curve).

Figure 11 shows the confusion matrices for the models based on photometry only (top row) and on photometry combined with morphology (bottom row), evaluated in three  $r$ -band magnitude bins. A classification threshold of  $p_{\text{cut}} = 0.5$  is used throughout. As expected, classification performance deteriorates at fainter magnitudes, primarily due to stars being increasingly misclassified as galaxies. The inclusion of morphological features significantly mitigates this effect, confining most misclassifications to the faintest bin ( $r > 22.5$ ). In contrast, the classification of galaxies remains consistently accurate across all bins.

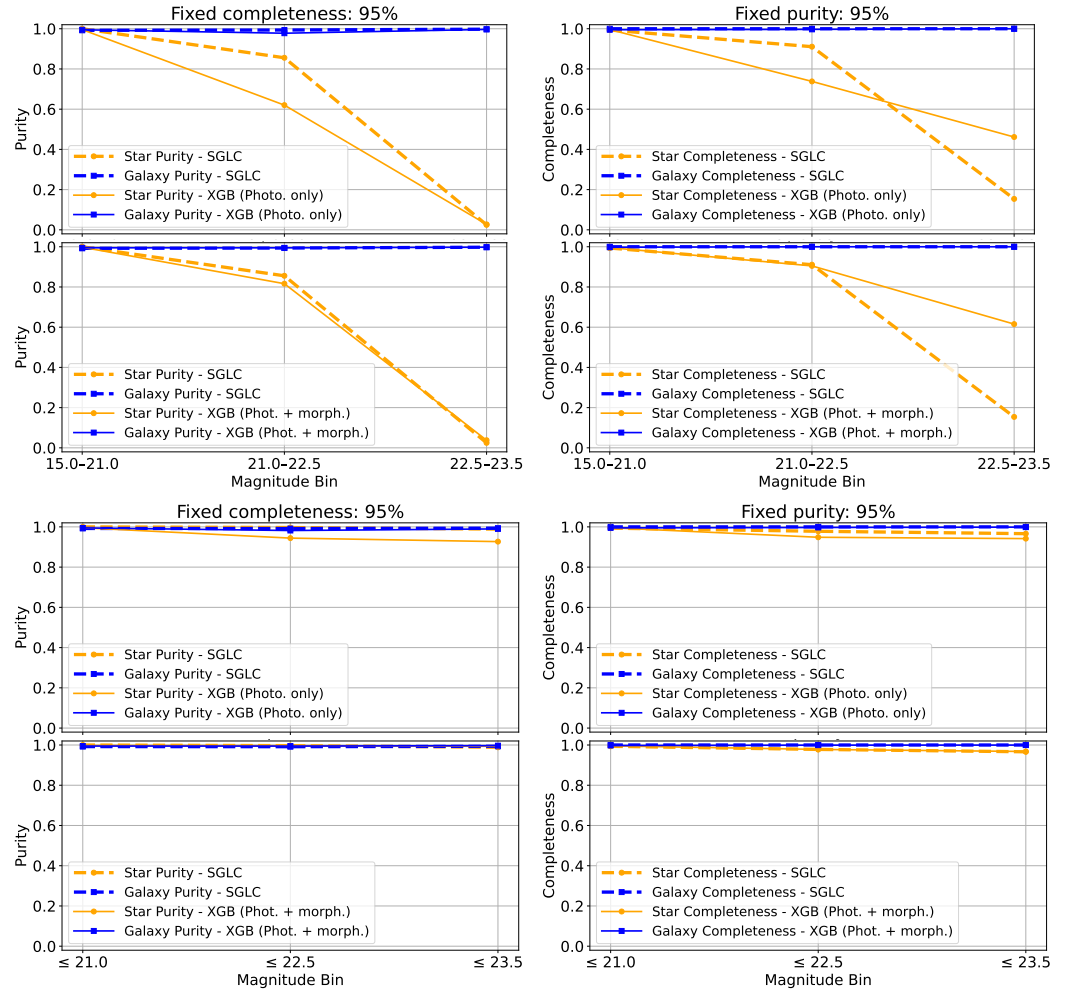
Figure 12 shows the purity and completeness for galaxies and stars as functions of the  $r$ -band magnitude, evaluated at a fixed target value of 95%. In each magnitude bin, this is achieved by selecting the probability threshold  $p_{\text{cut}}$  that yields the desired classification level. Results are presented in discrete magnitude bins (top panels) and cumulatively (bottom panels), for XGBoost models trained with photometry only and with photometry plus morphology, compared against SGLC. Galaxy classifications remain highly accurate across all magnitudes. For stars, performance degrades at faint magnitudes, but improves with the inclusion of morphology, surpassing SGLC. The apparent drop at  $r > 22.5$  in the binned results is explained by the very small number of faint stars in the training set (see Figure 5) and does not significantly affect the cumulative statistics.

## 5.2. Misclassifications

Misclassifications are dominated by faint stars ( $r > 21$ ) erroneously labeled as galaxies. Two effects drive this: (i) a brightness imbalance in the training set (relatively brighter stars and fainter galaxies) and (ii) color/SED degeneracy in photometry-only models at low

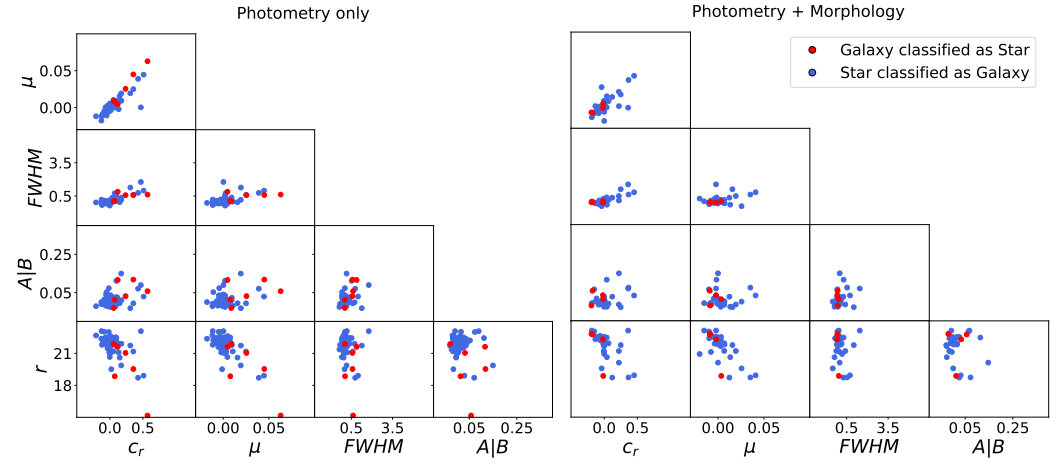


**Figure 11.** Confusion matrices for the XGBoost models using photometry only (top) and photometry plus morphology (bottom), shown for three bins of magnitudes. Each matrix displays the number of objects and associated purity (P) and completeness (C) for each class.



**Figure 12.** Purity (left) and completeness (right) for galaxies and stars as a function of  $r$ -band magnitude, evaluated at a fixed target value of 95%. Top panels show results in discrete magnitude bins, while bottom panels display the corresponding cumulative statistics.





**Figure 13.** Misclassified sources in the normalized morphological parameter space and  $r$ -band magnitude, for models using photometry only (left) and photometry with morphology (right).

SNR. Adding morphology largely mitigates the problem—star→galaxy errors drop sharply and galaxy→star errors are essentially removed. At the faint end, however, stellar and galactic morphologies converge, increasing ambiguity. This is evident in Figure 13, where the normalized peak surface brightness  $\mu$  drifts toward the stellar reference; the same trend appears in the photometry-only model. Thus, morphology reduces the rate of mistakes but does not fully break the underlying degeneracies, indicating that part of the confusion arises from photometric/SED similarity in faint sources. These results align with von Marttens et al. [13] and suggest that more sophisticated treatment of photometry (e.g., improved uncertainty handling) could yield further gains.

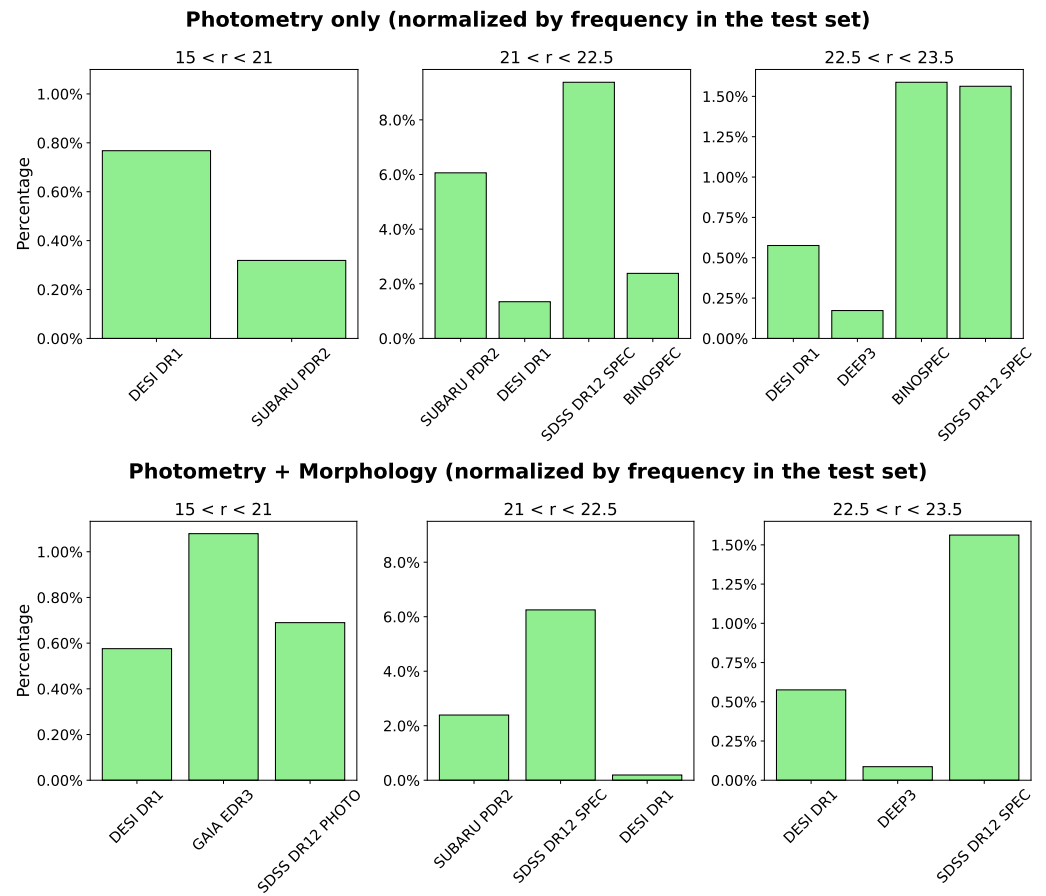
To investigate whether additional factors contribute to the misclassifications, Figure 14 shows the fraction of stars classified as galaxies by survey of origin within the test set. No clear trends are observed, excluding the possibility that stellar misclassifications are mainly due to incorrect truth labels. We also examined misclassifications as a function of **flags**, but the numbers were negligible. Finally, Appendix D presents visual inspections of representative misclassifications across both models and magnitude ranges.

### 5.3. Stellar locus

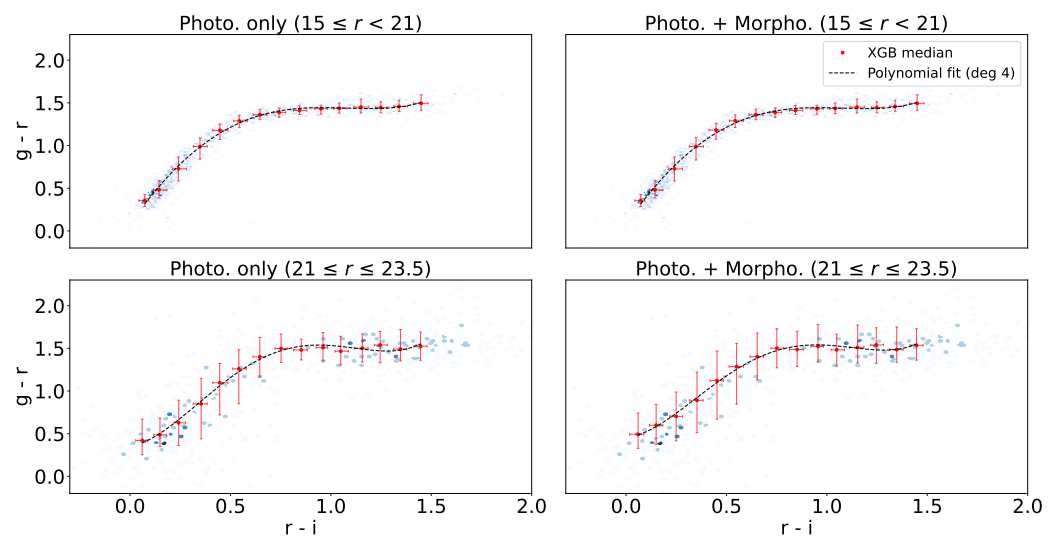
Figure 15 shows the stellar locus. The dashed line is a fourth-degree polynomial fit to all labeled stars, overplotted on the 2D histogram to guide the reader’s eye. Red markers show the XGBoost predictions from the full value-added catalog (VAC) using the model with photometry (left) and photometry and morphology (right): each point corresponds to the median value in a bin, with error bars indicating the 16th and 84th percentiles. The predicted locus closely follows the reference across the magnitude range, consistent with the high stellar purity reported in Figure 11.

### 5.4. Feature importance

To evaluate the relative contribution of each feature group, we computed permutation importance on the test set. This metric measures the change in ROC AUC when the values of a given feature group are randomly permuted, thereby breaking its relationship with the target. A larger drop corresponds to higher importance. Small negative values (i.e., a slight increase in AUC after permutation) do not imply that a feature ‘confuses’ the classifier; they typically arise from sampling noise on a finite test set and/or correlations with other features. We therefore interpret such negative importances as statistically consistent with zero (non-informative) importance.



**Figure 14.** Fraction of stars misclassified as galaxies in the test set, normalized by the number of objects per survey, shown in three  $r$ -band magnitude bins.



**Figure 15.** Stellar locus: dashed line shows the fit to labeled stars; red markers indicate XGBoost predictions using the model with photometry (left) and photometry and morphology (right) from the full VAC.

Figure 16 shows the permutation importance results for the XGBoost model trained with photometry only. The most informative features are the three broad bands  $i$ ,  $g$ , and  $r$ , together with three narrow bands: J0460 (4603 Å), J0680 (6812 Å), and J0390 (3904 Å). These results are consistent with those of Baqui et al. [5, Figure 13], who found peaks in feature importance around 4000 and 6900 Å. This suggests that the model exploits these wavelengths to estimate the slope of the spectral energy distribution.

Figure 17 shows the results for the XGBoost model trained with photometry and morphology. The most informative features are the concentration index  $c_r$  and the normalized peak surface brightness  $\mu$ , followed by the PSF. FWHM and photometry also contribute, while ellipticity ( $A/B$ ), extinction  $E(B - V)$  with its error, and detection flags have negligible impact. This reflects the fact that miniJPAS and J-NEP each cover only about  $1 \text{ deg}^2$ , where Milky Way extinction varies very little. The results therefore confirm that morphological features carry the dominant information for effective star–galaxy separation in these fields.

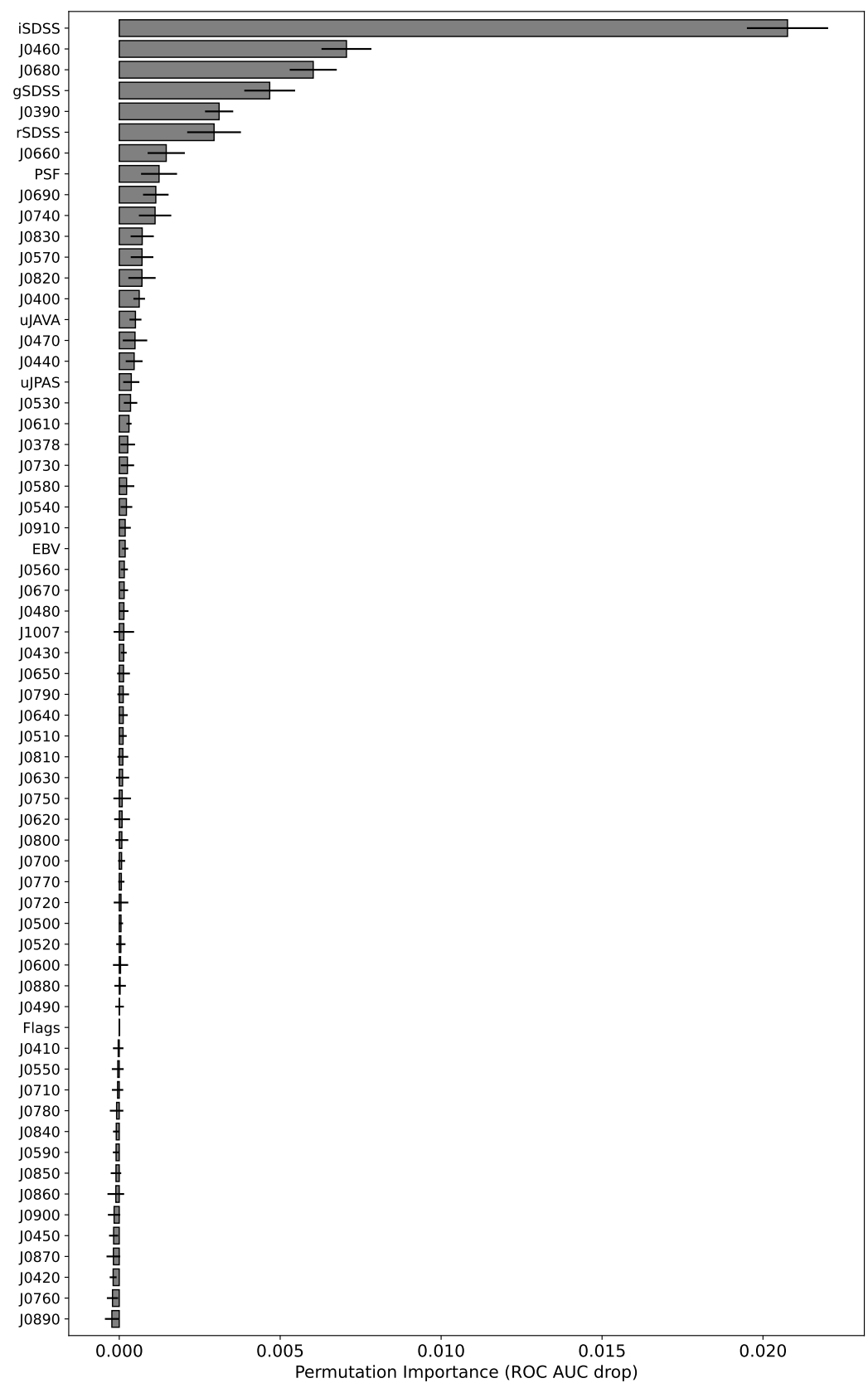
## 6. Conclusions

The machine learning-based classification of J-NEP and miniJPAS sources presented in this paper improves star–galaxy separation compared to previous methods. Using an XGBoost classifier with hyperparameters tuned through automated optimization, trained on a robust labeled dataset derived from crossmatching with multiple spectroscopic and photometric catalogs, we provide reliable classifications both from purely photometric data and from combined photometric and morphological features.

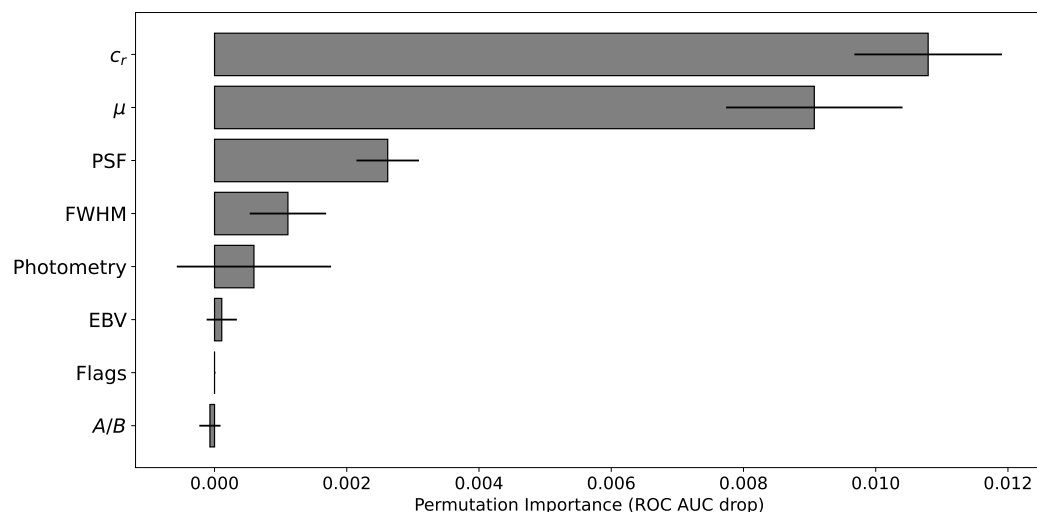
Our key findings can be summarized as follows:

1. The XGBoost model combining photometric and morphological features achieves a higher area under the ROC curve and outperforms both the SGLC and CLASS\_STAR classifiers across all magnitude ranges.
2. Using photometry alone yields competitive performance, slightly surpassing SGLC, but leads to misclassifications of faint stars as galaxies at  $r > 22$ . This ambiguity is reduced when morphological features are included.
3. The permutation importance analysis identifies the concentration, normalized peak surface brightness, and PSF as critical morphological parameters for classification. Photometric bands at approximately 3900, 4600, and 6800 Å are consistently among the most informative, highlighting their relevance for characterizing stellar and galactic spectral slopes.
4. The UMAP representativeness analysis demonstrates the robustness and consistency of the labeled training set, validating its applicability for reliable classification across the full miniJPAS and J-NEP datasets.
5. Our ML pipeline, publicly released as a value-added catalog (VAC), offers reliable star–galaxy classifications that can directly support downstream scientific analyses, including studies of galaxy evolution and cosmological investigations reliant on precise object classifications.

**Author Contributions:** Conceptualization, VM; methodology, APJ, GVS, VM and RvM; software, APJ and GVS; validation, APJ, GVS, VM and RvM; formal analysis, APJ, GVS, VM and RvM; investigation, APJ, GVS, VM and RvM; resources, R.A., J.A., N.B., S.B., J.C., D.C.H., S.D., R.D., A.E., R.M.G.D., A.H.C., C.H.M., J.L., C.L.S., A.M.F., C.M.O., M.M., F.R., L.S., K.T., J.V., H.V.R., J.M.V., C.W., J.Z.C.; data curation, APJ, GVS and VM; writing—original draft preparation, VM; writing—review and editing, APJ, GVS, VM, RvM, SGL, CW; visualization, APJ and GVS; supervision, VM; project administration, VM; funding acquisition, V.M., R.A., J.A., N.B., S.B., J.C., D.C.H., S.D., R.D., A.E., R.M.G.D., A.H.C., C.H.M., J.L., C.L.S., A.M.F.,



**Figure 16.** Permutation importance, measured as the decrease in ROC AUC when each group is shuffled, using the XGBoost model trained with photometry only. Each feature group includes both its value and its associated error.



**Figure 17.** Permutation importance of feature groups, measured as the decrease in ROC AUC when each group is shuffled, using the XGBoost model trained with photometry and morphology. Photometry refers to the combined importance of all 60 bands and their associated errors; EBV includes both the extinction value and its error.

C.M.O., M.M., F.R., L.S., K.T., J.V., H.V.R., J.M.V., J.Z.C. All authors have read and agreed to the published version of the manuscript.

**Funding:** APJ acknowledges financial support from CAPES (Brazil). GVS acknowledges financial support from UFES (Brazil). VM acknowledges partial support from CNPq (Brazil) and FAPES (Brazil). RvM is supported by Fundação de Amparo à Pesquisa do Estado da Bahia (FAPESB) grant TO APP0039/2023. SGL acknowledge the financial support from the MICIU with funding from the European Union NextGenerationEU and Generalitat Valenciana in the call Programa de Planes Complementarios de I+D+i (PRTR 2022) Project (VAL-JPAS), reference ASFAE/2022/025. SGL also acknowledge the research Project PID2023-149420NB-I00 funded by MICIU/AEI/10.13039/501100011033 and by ERDF/EU and the project of excellence PROMETEO CIPROM/2023/21 of the Conselleria de Educación, Universidades y Empleo (Generalitat Valenciana). RGD acknowledges financial support from the Severo Ochoa grant CEX2021-001131-S, funded by MICIU/AEI (10.13039/501100011033), and is also grateful for support from project PID2022-141755NB-I00.

**Data Availability Statement:** The data and code underlying this article are available from the corresponding author upon reasonable request. The value-added catalog will be publicly released in accordance with J-PAS collaboration policies, through the tables `jnep.StarGalClass` and `minijpas.StarGalClass` at [j-pas.org/datareleases](https://j-pas.org/datareleases).

**Acknowledgments:** This paper has undergone internal review by the J-PAS Collaboration. This work made use of computational resources provided by the [Sci-Com Lab](#) of the Department of Physics at UFES, supported by FAPES, CAPES and CNPq.

Based on observations made with the JST/T250 telescope and JPCam at the Observatorio Astrofísico de Javalambre (OAJ), in Teruel, owned, managed, and operated by the Centro de Estudios de Física del Cosmos de Aragón (CEFCA). We acknowledge the OAJ Data Processing and Archiving Unit (UPAD) for reducing and calibrating the OAJ data used in this work.

Funding for the J-PAS Project has been provided by the Governments of Spain and Aragón through the Fondo de Inversiones de Teruel; the Aragonese Government through the Research Groups E96, E103, E16\_17R, E16\_20R, and E16\_23R; the Spanish Ministry of Science and Innovation (MCIN/AEI/10.13039/501100011033 y FEDER, Una manera de hacer Europa) with grants PID2021-124918NB-C41, PID2021-124918NB-C42, PID2021-124918NA-C43, and PID2021-124918NB-C44; the Spanish Ministry of Science, Innovation and Universities (MCIU/AEI/FEDER, UE) with grants PGC2018-097585-B-C21 and PGC2018-097585-B-C22; the Spanish Ministry of

Economy and Competitiveness (MINECO) under AYA2015-66211-C2-1-P, AYA2015-66211-C2-2, and AYA2012-30789; and European FEDER funding (FCDD10-4E-867, FCDD13-4E-2685). The Brazilian agencies FINEP, FAPESP, FAPERJ and the National Observatory of Brazil have also contributed to this project. Additional funding was provided by the Tartu Observatory and by the J-PAS Chinese Astronomical Consortium.

**Conflicts of Interest:** The authors declare no conflicts of interest.

## Appendix A. Queries

### Appendix A.1. miniJPAS and J-NEP sources

Query used to select miniJPAS sources (an analogous query is used for J-NEP):

```
SELECT
    t1.NUMBER, t1.TILE_ID,
FROM
    minijpas.MagABDualObj t1
WHERE
    t1.flags[minijpas::rSDSS] IN (0,1,2,4,6)
    AND t1.mask_flags[minijpas::rSDSS] = 0
    AND t1.MAG_AUTO[minijpas::rSDSS] < 23.5
    AND t1.MAG_AUTO[minijpas::rSDSS] > 15
```

### Appendix A.2. SDSS DR12 photometric crossmatch

Query used to extract miniJPAS sources with SDSS DR12 crossmatch and reliable photometric classification (an analogous query is used for J-NEP):

```
SELECT
    t1.NUMBER, t1.TILE_ID,
    t2.class AS label
FROM
    minijpas.MagABDualObj t1
JOIN
    minijpas.xmatch_sdss_dr12 t2
    ON t1.TILE_ID = t2.TILE_ID AND t1.NUMBER = t2.NUMBER
WHERE
    t1.flags[minijpas::rSDSS] IN (0,1,2,4,6)
    AND t1.mask_flags[minijpas::rSDSS] = 0
    AND t1.MAG_AUTO[minijpas::rSDSS] < 20
    AND t1.MAG_AUTO[minijpas::rSDSS] > 15
```

### Appendix A.3. SDSS DR12 spectroscopic crossmatch

Query used for miniJPAS sources with SDSS spectroscopic crossmatch for  $r > 20$ :

```
SELECT
    t1.NUMBER, t1.TILE_ID,
    t2.spCl AS label
FROM
    minijpas.MagABDualObj t1
JOIN
    minijpas.xmatch_sdss_dr12 t2
    ON t1.TILE_ID = t2.TILE_ID AND t1.NUMBER = t2.NUMBER
WHERE
```

```

t1.flags[miniJPAS::rSDSS] IN (0,1,2,4,6)
AND t1.mask_flags[miniJPAS::rSDSS] = 0
AND t1.MAG_AUTO[miniJPAS::rSDSS] < 23.5
AND t1.MAG_AUTO[miniJPAS::rSDSS] > 20
AND t2.f_zsp = 0

```

#### Appendix A.4. Gaia EDR3 crossmatch

Query used to extract miniJPAS stars via Gaia parallax (also for J-NEP):

```

SELECT
    t1.NUMBER, t1.TILE_ID,
    t2.parallax_over_error AS label
FROM
    miniJPAS.MagABDualObj t1
JOIN
    miniJPAS.xmatch_gaia_edr3 t2
    ON t1.TILE_ID = t2.TILE_ID AND t1.NUMBER = t2.NUMBER
WHERE
    t1.flags[miniJPAS::rSDSS] IN (0,1,2,4,6)
    AND t1.mask_flags[miniJPAS::rSDSS] = 0
    AND t1.MAG_AUTO[miniJPAS::rSDSS] < 23.5
    AND t1.MAG_AUTO[miniJPAS::rSDSS] > 15
    AND t2.parallax_over_error > 3

```

#### Appendix A.5. HSC PDR2 crossmatch

Query used to extract miniJPAS sources crossmatched with HSC PDR2:

```

SELECT
    t1.NUMBER, t1.TILE_ID,
    t2.g_extendedness_value AS label_g,
    t2.r_extendedness_value AS label_r,
    t2.i_extendedness_value AS label_i,
    t2.z_extendedness_value AS label_z,
    t2.y_extendedness_value AS label_y
FROM
    miniJPAS.MagABDualObj t1
JOIN
    miniJPAS.xmatch_subaru_pdr2 t2
    ON t1.TILE_ID = t2.TILE_ID AND t1.NUMBER = t2.NUMBER
WHERE
    t1.flags[miniJPAS::rSDSS] IN (0,1,2,4,6)
    AND t1.mask_flags[miniJPAS::rSDSS] = 0
    AND t1.MAG_AUTO[miniJPAS::rSDSS] < 23.5
    AND t1.MAG_AUTO[miniJPAS::rSDSS] > 18.5
    AND t2.g_inputcount_value >= 4
    AND t2.g_cmodel_flag = false
    AND t2.g_pixelflags = false
    AND t2.r_inputcount_value >= 4
    AND t2.r_cmodel_flag = false
    AND t2.r_pixelflags = false
    AND t2.i_inputcount_value >= 4

```



```

AND t2.i_cmodel_flag = false
AND t2.i_pixelflags = false
AND t2.z_inputcount_value >= 4
AND t2.z_cmodel_flag = false
AND t2.z_pixelflags = false
AND t2.y_inputcount_value >= 4
AND t2.y_cmodel_flag = false
AND t2.y_pixelflags = false

```

## Appendix B. Hyperparameter configuration

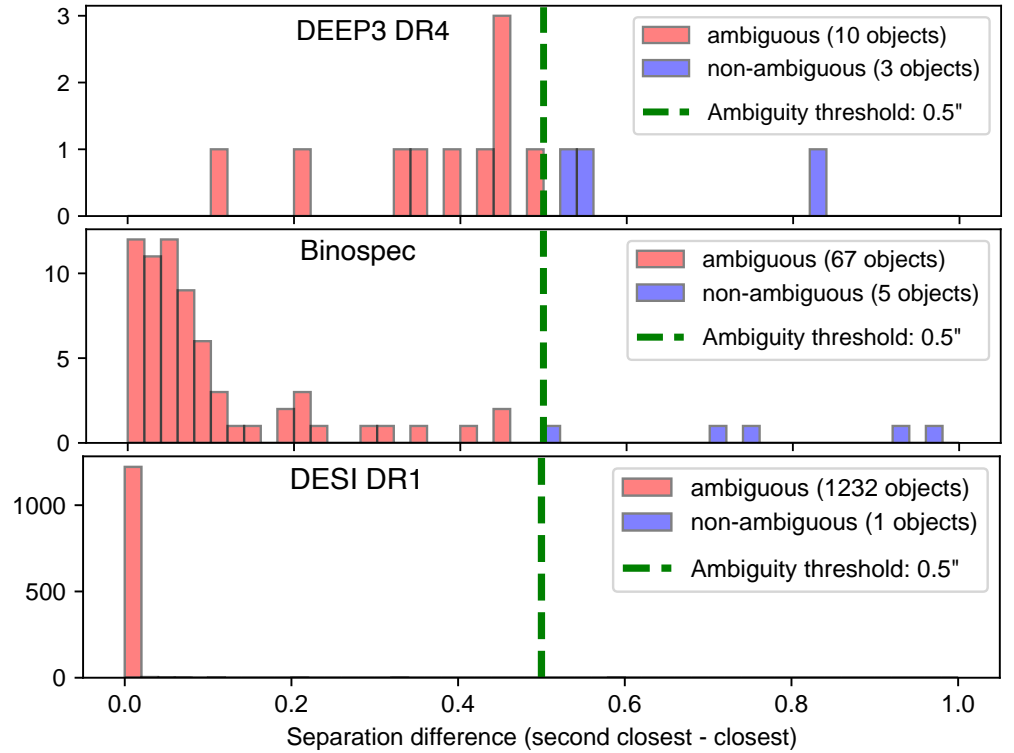
The best-performing pipeline identified by TPOT consists of a **MaxAbsScaler** for feature scaling followed by an **XGBClassifier**. The choice of XGBoost is further supported by von Marttens et al. [13], which—although based on 12 photometric bands in J-PLUS rather than the 60 in J-PAS—uses the same morphological features and comparable photometric coverage. In that study, TPOT systematically evaluated multiple machine learning algorithms and identified XGBoost as the top performer. The consistency between the two surveys thus provides a strong justification for adopting XGBoost in our analysis. The final hyperparameter configuration is as follows:

- **booster=gbtree**: specifies the tree-based boosting model, suitable for structured tabular data.
- **colsample\_bylevel=0.6**: fraction of features sampled for each tree level, reducing overfitting by introducing randomness.
- **colsample\_bynode=1.0**: fraction of features sampled for each split, set to 1 (all features used).
- **colsample\_bytree=0.9**: fraction of features sampled for each tree, promoting diversity among trees.
- **gamma=0.5**: minimum loss reduction required for further partitioning; larger values make the model more conservative.
- **learning\_rate=0.05**: step size shrinkage, balancing learning speed with generalization.
- **max\_depth=7**: maximum depth of each tree, controlling model complexity.
- **min\_child\_weight=5**: minimum sum of instance weights in a child node; higher values prevent overfitting small fluctuations.
- **n\_estimators=600**: number of boosting rounds (trees).
- **objective=binary:logistic**: specifies binary classification with probabilistic outputs.
- **reg\_alpha=0**: L1 regularization term, promoting sparsity in weights; set to zero here.
- **reg\_lambda=1**: L2 regularization term, penalizing large weights to reduce overfitting.
- **scale\_pos\_weight=3.46**: weighting factor to address class imbalance between stars and galaxies.
- **subsample=0.95**: fraction of training instances sampled for each tree, preventing overfitting by adding stochasticity.

This configuration balances model flexibility and regularization, effectively handling the imbalance between stellar and galactic classes while maintaining robustness against overfitting.

## Appendix C. Treatment of Ambiguous Crossmatches

In crossmatching miniJPAS and J-NEP with external spectroscopic catalogs (DEEP3, Binospec, and DESI), ambiguous matches may arise when multiple J-PAS sources fall within the matching radius of the same external catalog object. By default, **onexmatch** resolves duplicates by retaining only the closest J-PAS source for each catalog entry. However, when



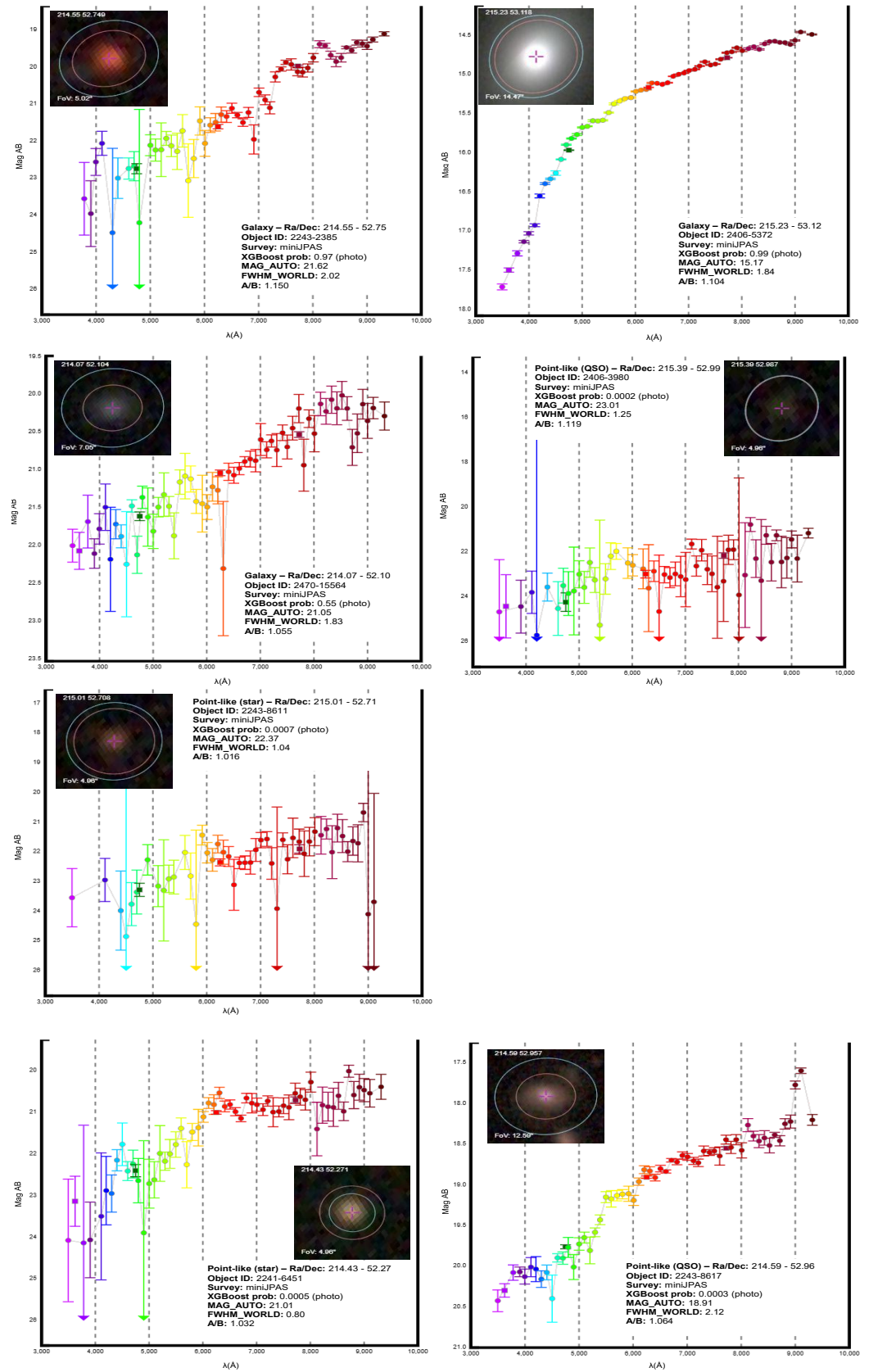
**Figure C1.** Distribution of ambiguous and non-ambiguous crossmatches. The horizontal axis shows the separation difference between the closest and second-closest sources matched to the same miniJPAS or J-NEP object. Red bars indicate ambiguous matches (within 0.5 arcsec), and blue bars show non-ambiguous cases.

two candidate sources are nearly equidistant from the same object, this procedure may not be reliable. To address such cases, we adopted an ambiguity threshold of 0.5 arcsec, which is approximately twice the pixel scale and comparable to the PSF size. If the positional difference between two candidate matches is below this threshold, the association is deemed ambiguous and the sources are discarded from the final crossmatch, see Figure C1. This conservative approach minimizes spurious associations, ensuring that the resulting training set retains only high-confidence matches.

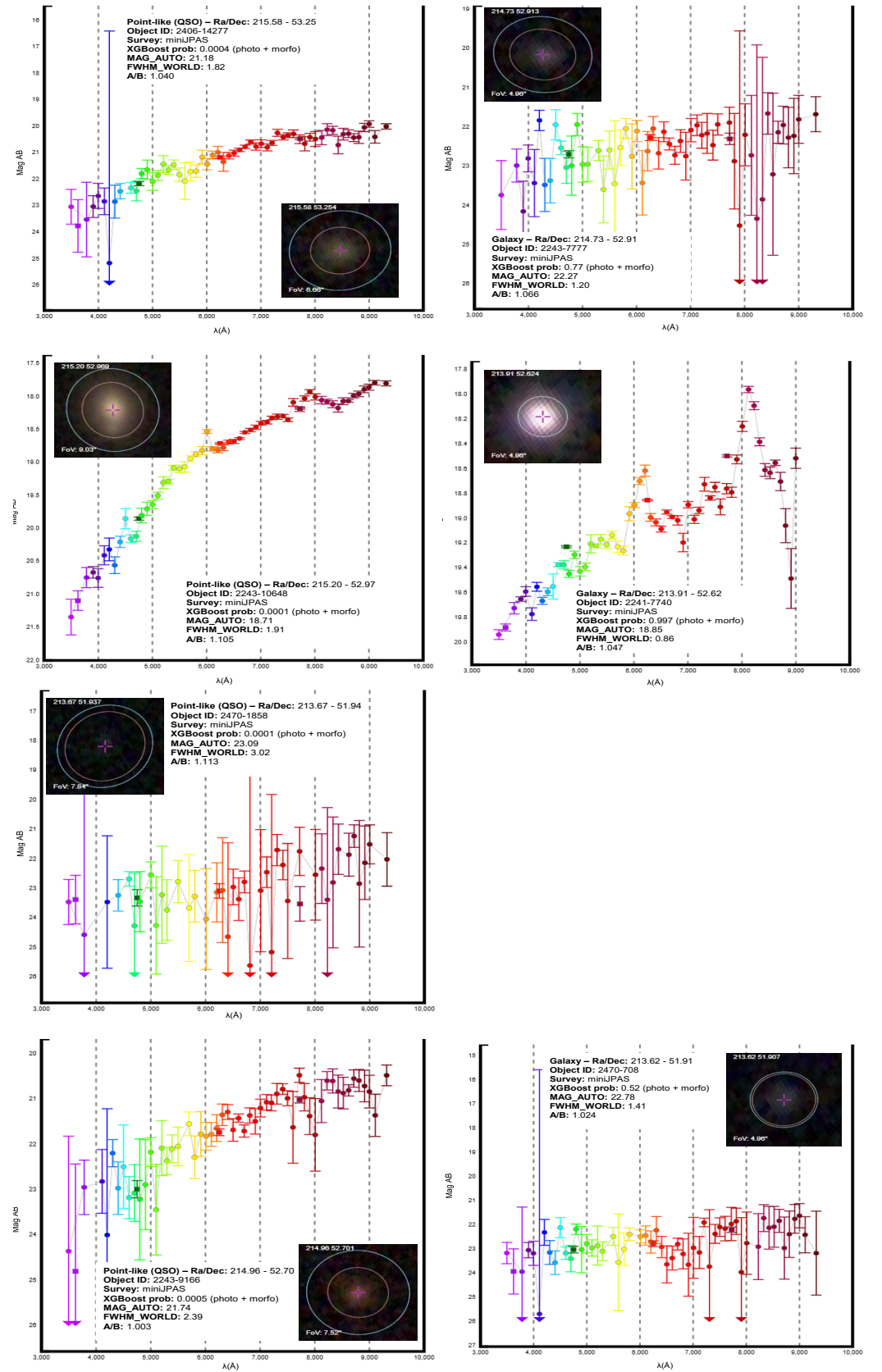
## Appendix D. Visual inspection of misclassifications

Here we present a visual inspection of misclassifications produced by the models trained with photometric information only and with combined photometric and morphological information. Representative examples were selected to cover the full magnitude range considered in this work. This qualitative analysis complements the statistical results by illustrating typical failure modes and offering insights into possible improvements for future modeling.

Each panel of Figures D1 and D2 displays the spectral energy distribution across the 60 J-PAS bands, with broad-band points shown as squares, together with an RGB cutout from the  $r$ ,  $i$ , and  $g$  bands (inset). The examples include galaxies misclassified as point-like sources and point-like objects (stars or QSOs) misclassified as galaxies. For each case, we report the object ID, sky coordinates, survey origin, classifier probability (1 = star, 0 = galaxy),  $r$ -band magnitude, PSF, FWHM, and ellipticity  $A/B$ .



**Figure D1.** Representative examples of misclassified sources by the model trained with photometric features only. These cases illustrate typical limitations of a purely photometric approach.



**Figure D2.** A Figure D1 but for the model trained with combined photometric and morphological features. These cases illustrate typical failure modes when morphology is included in the model.

## References

1. Bonoli, S.; et al. The miniJPAS survey: A preview of the Universe in 56 colors. *Astron. Astrophys.* **2021**, *653*, A31, [arXiv:astro-ph.CO/2007.01910]. <https://doi.org/10.1051/0004-6361/202038841>.
2. Hernán-Caballero, A.; et al. J-NEP: 60-band photometry and photometric redshifts for the James Webb Space Telescope North Ecliptic Pole Time-Domain Field. *Astron. Astrophys.* **2023**, *671*, A71, [arXiv:astro-ph.GA/2301.09623]. <https://doi.org/10.1051/0004-6361/202244759>.
3. Jansen, R.A.; Windhorst, R.A. The James Webb Space Telescope North Ecliptic Pole Time-domain Field. I. Field Selection of a JWST Community Field for Time-domain Studies. *PASP* **2018**, *130*, 124001, [arXiv:astro-ph.GA/1807.05278]. <https://doi.org/10.1088/1538-3873/aae476>.
4. López-Sanjuan, C.; Vázquez Ramió, H.; Varela, J.; Spinoso, D.; Angulo, R.E.; Muniesa, D.; Viironen, K.; Cristóbal-Hornillos, D.; Cenarro, A.J.; Ederoclite, A.; et al. J-PLUS: Morphological star/galaxy classification by PDF analysis. *A&A* **2019**, *622*, A177, [arXiv:astro-ph.GA/1804.02673]. <https://doi.org/10.1051/0004-6361/201732480>.
5. Baqui, P.; Marra, V.; et al. The miniJPAS survey: star-galaxy classification using machine learning. *A&A* **2021**, *645*, A87, [arXiv:astro-ph.IM/2007.07622]. <https://doi.org/10.1051/0004-6361/202038986>.
6. Alam, S.; Albareti, F.D.; Allende Prieto, C.; Anders, F.; Anderson, S.F.; Anderton, T.; Andrews, B.H.; Armengaud, E.; Aubourg, É.; Bailey, S.; et al. The Eleventh and Twelfth Data Releases of the Sloan Digital Sky Survey: Final Data from SDSS-III. *ApJS* **2015**, *219*, 12, [arXiv:astro-ph.IM/1501.00963]. <https://doi.org/10.1088/0067-0049/219/1/12>.
7. Gaia Collaboration.; Brown, A.G.A.; Vallenari, A.; Prusti, T.; de Bruijne, J.H.J.; Babusiaux, C.; Biermann, M.; Creevey, O.L.; Evans, D.W.; Eyer, L.; et al. Gaia Early Data Release 3. Summary of the contents and survey properties. *A&A* **2021**, *649*, A1, [arXiv:astro-ph.GA/2012.01533]. <https://doi.org/10.1051/0004-6361/202039657>.
8. Cooper, M.C.; Aird, J.A.; Coil, A.L.; Davis, M.; Faber, S.M.; Juneau, S.; Lotz, J.M.; Nandra, K.; Newman, J.A.; Willmer, C.N.A.; et al. The DEEP3 Galaxy Redshift Survey: Keck/DEIMOS Spectroscopy in the GOODS-N Field. *ApJS* **2011**, *193*, 14, [arXiv:astro-ph.CO/1101.4018]. <https://doi.org/10.1088/0067-0049/193/1/14>.
9. Fabricant, D.; Fata, R.; Epps, H.; Gauron, T.; Mueller, M.; Zajac, J.; Amato, S.; Barberis, J.; Bergner, H.; Brennan, P.; et al. Binospec: A Wide-field Imaging Spectrograph for the MMT. *PASP* **2019**, *131*, 075004, [arXiv:astro-ph.IM/1905.03320]. <https://doi.org/10.1088/1538-3873/ab1d78>.
10. Abdul Karim, M.; et al. Data Release 1 of the Dark Energy Spectroscopic Instrument **2025**. [arXiv:astro-ph.CO/2503.14745].
11. Aihara, H.; AlSayyad, Y.; Ando, M.; Armstrong, R.; Bosch, J.; Egami, E.; Furusawa, H.; Furusawa, J.; Goulding, A.; Harikane, Y.; et al. Second data release of the Hyper Suprime-Cam Subaru Strategic Program. *PASJ* **2019**, *71*, 114, [arXiv:astro-ph.IM/1905.12221]. <https://doi.org/10.1093/pasj/psz103>.
12. Marra, V. onexmatch, 2025. <https://doi.org/10.5281/zenodo.17148385>.
13. von Marttens, R.; Marra, V.; et al. J-PLUS DR3: Galaxy-Star-Quasar classification. *Mon. Not. Roy. Astron. Soc.* **2023**, *527*, 3347–3365, [arXiv:astro-ph.GA/2212.05868]. <https://doi.org/10.1093/mnras/stad3373>.
14. Chaini, S.; Bagul, A.; Deshpande, A.; Gondkar, R.; Sharma, K.; Vivek, M.; Kembhavi, A. Photometric identification of compact galaxies, stars, and quasars using multiple neural networks. *MNRAS* **2023**, *518*, 3123–3136, [arXiv:astro-ph.GA/2211.08388]. <https://doi.org/10.1093/mnras/stac3336>.
15. del Pino, A.; López-Sanjuan, C.; Hernán-Caballero, A.; Domínguez-Sánchez, H.; von Marttens, R.; Fernández-Ontiveros, J.A.; Coelho, P.R.T.; Lumbreras-Calle, A.; Vega-Ferrero, J.; Jimenez-Esteban, F.; et al. J-PLUS: Bayesian object classification with a stream of BANNJOS. *A&A* **2024**, *691*, A221, [arXiv:astro-ph.GA/2404.16567]. <https://doi.org/10.1051/0004-6361/202450503>.
16. Zeraatgari, F.Z.; Hafezianzadeh, F.; Zhang, Y.; Mei, L.; Ayubinia, A.; Mosallanezhad, A.; Zhang, J. Machine learning-based photometric classification of galaxies, quasars, emission-line galaxies, and stars. *MNRAS* **2024**, *527*, 4677–4689, [arXiv:astro-ph.GA/2311.02951]. <https://doi.org/10.1093/mnras/stad3436>.
17. Asadi, V.; Haghi, H.; Zonoozi, A.H. Semi-supervised classification of stars, galaxies and quasars using K-means and random-forest approaches. *A&A* **2025**, *700*, A259, [arXiv:astro-ph.GA/2507.14072]. <https://doi.org/10.1051/0004-6361/202555620>.
18. Bertin, E. Automated Morphometry with SExtractor and PSFEx. In Proceedings of the Astronomical Data Analysis Software and Systems XX; Evans, I.N.; Accomazzi, A.; Mink, D.J.; Rots, A.H., Eds., 2011, Vol. 442, *Astronomical Society of the Pacific Conference Series*, p. 435.
19. Green, G.M.; Schlafly, E.F.; Finkbeiner, D.; Rix, H.W.; Martin, N.; Burgett, W.; Draper, P.W.; Flewelling, H.; Hodapp, K.; Kaiser, N.; et al. Galactic reddening in 3D from stellar photometry - an improved map. *MNRAS* **2018**, *478*, 651–666, [arXiv:astro-ph.GA/1801.03555]. <https://doi.org/10.1093/mnras/sty1008>.
20. López-Sanjuan, C.; Varela, J.; Cristóbal-Hornillos, D.; Vázquez Ramió, H.; Carrasco, J.M.; Tremblay, P.E.; Whitten, D.D.; Placco, V.M.; Marín-Franch, A.; Cenarro, A.J.; et al. J-PLUS: photometric calibration of large-area multi-filter surveys with stellar and white dwarf loci. *A&A* **2019**, *631*, A119, [arXiv:astro-ph.IM/1907.12939]. <https://doi.org/10.1051/0004-6361/201936405>.

21. McInnes, L.; Healy, J.; Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *arXiv e-prints* **2018**, p. arXiv:1802.03426, [arXiv:stat.ML/1802.03426]. <https://doi.org/10.48550/arXiv.1802.03426>.
22. Chen, T.; Guestrin, C. XGBoost: A Scalable Tree Boosting System. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* **2016**, p. 785–794, [arXiv:cs.LG/1603.02754]. <https://doi.org/10.1145/2939672.2939785>.
23. Le, T.T.; Fu, W.; Moore, J.H. Scaling tree-based automated machine learning to biomedical big data with a feature set selector. *Bioinformatics* **2020**, *36*, 250–256.
24. Bertin, E.; Arnouts, S. SExtractor: Software for source extraction. *A&AS* **1996**, *117*, 393–404. <https://doi.org/10.1051/aas:1996164>.

**Disclaimer/Publisher’s Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.