

# Classifying *Kepler* light curves for 12 000 A and F stars using supervised feature-based machine learning

Nicholas H. Barbara,<sup>1,2★</sup> Timothy R. Bedding<sup>1,2★</sup>, Ben D. Fulcher,<sup>1</sup> Simon J. Murphy<sup>1,2</sup> and Timothy Van Reeth<sup>1,2,3</sup>

<sup>1</sup>*Sydney Institute for Astronomy, School of Physics, University of Sydney, Sydney, NSW 2006, Australia*

<sup>2</sup>*Stellar Astrophysics Centre, Department of Physics and Astronomy, Aarhus University, Aarhus, DK-8000, Denmark*

<sup>3</sup>*Institute of Astronomy, KU Leuven, Celestijnenlaan 200D, B-3001 Leuven, Belgium*

Accepted 2022 May 27. Received 2022 May 5; in original form 2022 February 24

## ABSTRACT

With the availability of large-scale surveys like *Kepler* and *TESS*, there is a pressing need for automated methods to classify light curves according to known classes of variable stars. We introduce a new algorithm for classifying light curves that compares 7000 time-series features to find those that most effectively classify a given set of light curves. We apply our method to *Kepler* light curves for stars with effective temperatures in the range 6500–10 000 K. We show that the sample can be meaningfully represented in an interpretable 5D feature space that separates seven major classes of light curves ( $\delta$  Scuti stars,  $\gamma$  Doradus stars, RR Lyrae stars, rotational variables, contact eclipsing binaries, detached eclipsing binaries, and non-variables). We achieve a balanced classification accuracy of 82 per cent on an independent test set of *Kepler* stars using a Gaussian mixture model classifier. We use our method to classify 12 000 *Kepler* light curves from Quarter 9 and provide a catalogue of the results. We further outline a confidence heuristic based on probability density to search our catalogue and extract candidate lists of correctly classified variable stars.

**Key words:** asteroseismology – methods: data analysis – binaries: eclipsing – stars: oscillations – stars: variables: general.

## 1 INTRODUCTION

The use of machine learning is becoming increasingly common in astronomy (Ball & Brunner 2010; Graff et al. 2014; Baron 2019; Ivezić et al. 2019). In particular, large-scale photometric surveys are producing light curves in numbers too large for humans to manually inspect and analyse. Considerable efforts have gone into using machine learning to classify light curves from large ground-based surveys (e.g. Carrasco-Davis et al. 2019; Tsang & Schultz 2019; Cabral et al. 2020; Hosenie et al. 2020; Jamal & Bloom 2020; Johnston et al. 2020; Szklenár et al. 2020; Bassi, Sharma & Gomekar 2021; Zhang & Bloom 2021). Such techniques have also been applied to light curves from NASA’s *Kepler* and K2 missions (e.g. Blomme et al. 2010, 2011; Debusscher et al. 2011; Armstrong et al. 2016; Bass & Borne 2016; Hon, Stello & Yu 2017, 2018a, b; Johnston et al. 2019; Kgoadi, Whittingham & Engelbrecht 2019; Le Saux et al. 2019; Giles & Walkowicz 2020; Kuszlewicz, Hekker & Bell 2020; Audenaert et al. 2021; Paul & Chattopadhyay 2022).

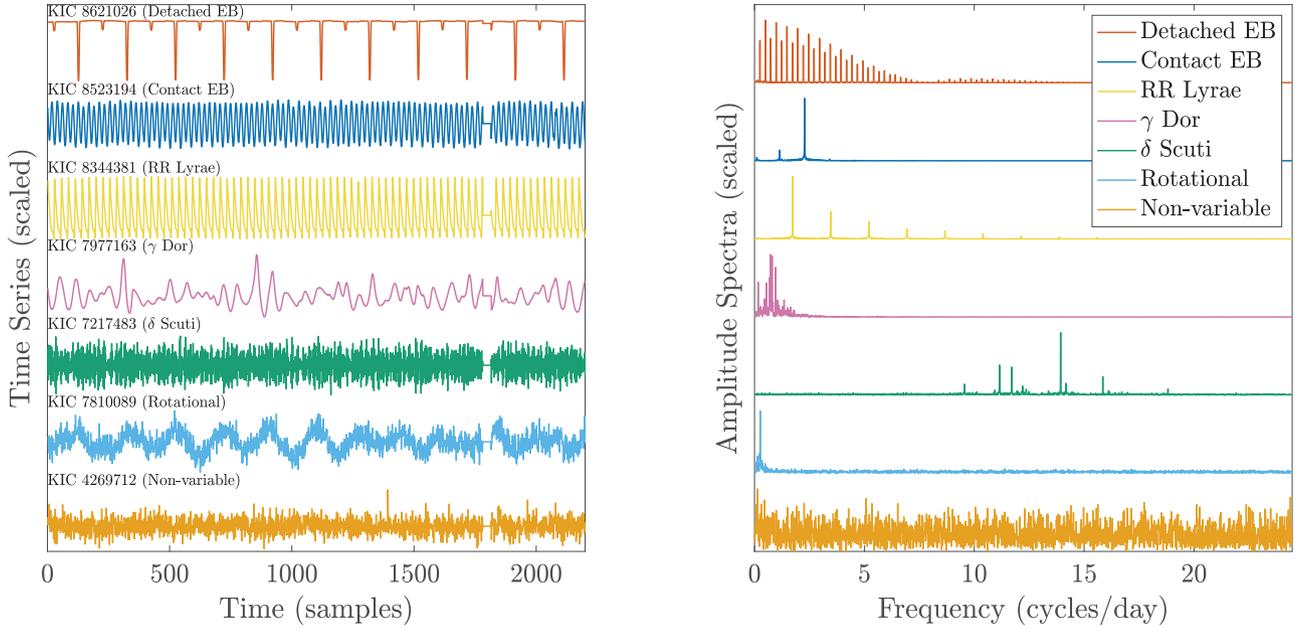
A range of algorithms have been proposed to classify light-curve data bases according to known classes of stars, but these algorithms often rely on black-box machine-learning methods, which limits their interpretability and hence ability to drive scientific understanding. Those algorithms that are more interpretable rely on manually selected temporal or spectral features of light curves (e.g. Pashchenko, Sokolovsky & Gavras 2017), with minimal comparison to the

performance of alternatives from across a highly interdisciplinary time-series analysis literature. Here, we introduce a new algorithm for classifying light-curve data bases that searches over 7000 time-series features to automatically find interpretable features relevant to classifying a given set of light curves. Our aim is to develop a simple, efficient classifier that uses these interpretable features to give us new insight into the classification of variable stars, while maintaining comparable performance with existing methods.

Over the course of its 4-yr mission, the *Kepler* spacecraft collected light curves for nearly 200 000 stars, most of which show variability. Subsets of *Kepler* stars have been classified systematically, resulting in catalogues of about 2900 eclipsing binaries (EBs; Kirk et al. 2016), 16 000 oscillating red giants (e.g. Yu et al. 2018, and references therein), 2000 pulsating  $\delta$  Scuti stars (Murphy et al. 2019), and over 600  $\gamma$  Doradus pulsators (Li et al. 2020). Independently, Balona (2018) used visual inspection of light curves and power spectra to classify over 20 000 A and F *Kepler* stars. Finally, Audenaert et al. (2021) have recently classified 167 000 light curves from one quarter of *Kepler* data (see Section 4.4).

In this paper, we present our methods and provide classifications based on a single 3-month quarter (Quarter 9, Q9) for approximately 12 000 *Kepler* stars with effective temperatures in the range 6500–10 000 K. Our primary interest is in pulsating stars and our chosen temperature range covers the classical instability strip, which has the richest variety of variability (e.g. Kurtz 2022). There are relatively few stars in the *Kepler* sample that are hotter than this range, while pulsations on the cooler side are dominated by a single class (solar-like oscillations) that have already been extensively classified and studied (see Jackiewicz 2021 for a recent review). We note that we

\* E-mail: [nicholas.barbara@sydney.edu.au](mailto:nicholas.barbara@sydney.edu.au) (NHB);  
[tim.bedding@sydney.edu.au](mailto:tim.bedding@sydney.edu.au) (TRB)



**Figure 1.** Examples of variable stars in the detached binary, contact binary, RR Lyr,  $\gamma$  Doradus,  $\delta$  Scuti, rotational, and non-variable stellar classes selected in the training set, respectively. The left-hand panel shows the first 2200 samples (at approximately 30 min per sample) of the light curves of each star, or half a *Kepler* quarter of data. The right-hand panel shows the corresponding Fourier transforms up to the Nyquist frequency. The vertical axes are scaled for ease of viewing.

have previously used our method to identify samples of  $\delta$  Scuti stars (Murphy et al. 2020) and  $\gamma$  Doradus stars (Li et al. 2020). Finally, we compare our classification of *Kepler* stars to labels assigned by a less interpretable, performance-focused classifier in Audenaert et al. (2021), where we find similar results.

## 2 TRAINING DATA

### 2.1 A selection of variable stars

Selecting an adequate training set to train a feature-based classifier for all 200 000 stars in the *Kepler* field is the most challenging and time-consuming aspect of developing a general classification algorithm. Many classes of variable stars are rare, while others remain poorly understood, and still others have not yet been identified. Indeed, new classes are occasionally proposed (Deboscher et al. 2011; Bass & Borne 2016; Pietrukowicz et al. 2017). For this reason, a good training set must be representative of a wide range of variable stars to construct a suitably general feature-space representation of *Kepler* light curves. Rather than attempting to compile a collection of all known classes of variable and non-variable stars, as attempted with limited success by Deboscher et al. (2011), we focused our research on a subset of seven well-studied classes, within the temperature range  $6500 \text{ K} \leq T_{\text{eff}} \leq 10\,000 \text{ K}$ , as an initial demonstration. The seven chosen classes are  $\delta$  Scuti stars,  $\gamma$  Doradus stars, RR Lyrae stars, rotational variables, contact EBs, detached EBs, and non-variable stars. Typical light curves and power spectra for each class are included in Fig. 1. These classes are commonly represented in the wider *Kepler* data (with the exception of the RR Lyrae class), are interesting stars that we wish to study in further detail, and are not intermediate or hybrid classes. We excluded hybrid stars to avoid having light curves in the training data that belong to multiple classes.

Two of our seven classes are subclasses of EB systems, which were catalogued in *Kepler* data by Kirk et al. (2016). EBs can be classed as one of the following: detached binaries, where the two stars are far from each other to give highly separated eclipses; contact binaries,

where there is no space between the two stellar envelopes, producing almost sinusoidal eclipse patterns; and semidetached binaries, the intermediate class of the two extremes. In accordance with our decision to exclude hybrid classes, only the detached and contact subclasses have been included in our training set.

Rotational variables are most commonly found among cooler stars, whose star-spots present darker patches on the surface that rotate in and out of view, but rotational variability is also seen by *Kepler* across the effective temperature range of our sample (Balona 2013; Sikora, Wade & Rowe 2020). Typically, the spots have lifetimes not much longer than the rotation period, and they may occur at different latitudes, so the variability is only quasi-periodic (Nielsen et al. 2013; McQuillan, Mazeh & Aigrain 2014). The  $\alpha^2$  CVn stars are hotter stars whose strong dipolar magnetic fields concentrate certain elements into spots. These also rotate with the star, but occur near the magnetic poles and are much longer lived, leading to light curves that do not change rapidly in period, amplitude, or shape (Wolff 1983). In our chosen temperature range, examples of both are found.

Three classes of pulsating variable star were included (for a recent review of pulsating stars, see Kurtz 2022). RR Lyr variables are bright, evolved stars burning helium in their cores. As they traverse the horizontal branch, they cross the instability strip and pulsate periodically with a characteristic phase curve. Their use as standard candles has allowed measurements of the distance to the Galactic centre and to globular clusters (Oort & Plaut 1975; Walker 1992). The two other pulsating star classes,  $\gamma$  Doradus and  $\delta$  Scuti variables, both comprise A- or F-type stars on or near the main sequence, and embody two distinct types of oscillation: g modes, or buoyancy-driven modes sensitive to the near-core region of a star; and p modes, pressure-driven modes most sensitive to the envelope.  $\gamma$  Doradus stars are multiperiodic g-mode pulsators with periods between approximately 0.3 and 3 d (Kaye et al. 1999). Despite having periods similar to the RR Lyr variables, the multiperiodic  $\gamma$  Doradus stars do not have simple phased curves. There are several hundred in the *Kepler* field (Li et al. 2020), and they have seen substantial recent

**Table 1.** Breakdown of stars in the training set.

Class	No. stars
Contact EB	171
Detached EB	83
$\delta$ Scuti	411
$\gamma$ Doradus	262
Non-variable	201
Rotational	166
RR Lyrae	25
Total	1319

attention because of their ability to probe internal rotation (Van Reeth et al. 2018; Ouazzani et al. 2019), diffusive mixing (Bouabid et al. 2013), and core overshooting (Mombarg et al. 2019).

Finally,  $\delta$  Scuti stars are the most common class of pulsating star at A and F spectral types, with approximately 2000 known in *Kepler* data alone (Murphy et al. 2019; Guzik 2021). These stars are p-mode oscillators, and with periods between 18 min and 8 h they are the highest frequency variables in our sample. For unknown reasons, even in the middle of the  $\delta$  Scuti instability strip only half of the stars pulsate as  $\delta$  Scuti stars (Murphy et al. 2019); hence, we include a non-variable class in this work. We note that some  $\delta$  Scuti stars are known to lie outside of the instability strip (e.g. Bowman & Kurtz 2018), but our classifications are based only on the *Kepler* light curves and not on parameters such as effective temperature.

## 2.2 Preparing the training set

We created a training set across all seven classes by hand-picking 1319 stars from candidate lists according to specific criteria. Stars were restricted to the temperature range  $6500 \text{ K} \leq T_{\text{eff}} \leq 10000 \text{ K}$  using effective temperatures from Mathur et al. (2017). We examined one quarter of long-cadence *Kepler* photometry for each star to prepare the training data. Quarter 9 (Q9) was chosen because it has no prolonged gaps in observation, such as those arising from telescope safe mode events, and no anomalies in data quality. We used light curves made with simple aperture photometry (SAP), downloaded from the Kepler Asteroseismic Science Operations Center (KASOC) website (Data Release 25).<sup>1</sup> The choice to examine a single quarter was made to reduce computation time, but this also precludes the analysis of variability on time-scales longer than a typical 90-d quarter. While we certainly recommend the investigation of 4-yr data in future research focused on a wider range of *Kepler* variables, this will not have a great effect on the stars chosen for our investigation. Of the seven classes, only a handful of detached binaries are known to have a period greater than 90 d (Kirk et al. 2016), and we did not include these in the training set.

The training stars were chosen from lists of possible candidates for each of the seven classes by visually inspecting their light curves and Fourier transforms. This laborious process embodies the motivation for automated variable star classification, and was a necessary task to ensure that the training data were accurate and would not mislead the automated feature-selection process. In the following paragraphs, we describe the selection of stars in training sets for each class. Table 1 summarizes the class-specific number of stars in the training set. The full list of stars is provided as supplementary material, with a sample shown in Table 2.

**Table 2.** The training set of 1319 *Kepler* stars. An extract of 14 stars is shown, with the full table provided in the supplementary material.

KIC ID	Class
10855535	Contact EB
9612468	Contact EB
3836439	Detached EB
9711751	Detached EB
9331207	$\delta$ Scuti
8376471	$\delta$ Scuti
4755510	$\gamma$ Doradus
1996456	$\gamma$ Doradus
1864603	Non-variable
2156425	Non-variable
1164109	Rotational
1435836	Rotational
3733346	RR Lyrae
3864443	RR Lyrae

EB systems were selected from the *Kepler* Eclipsing Binary Catalogue (Kirk et al. 2016), restricted to periods  $< 90 \text{ d}$  and a morphology index of  $0 \leq c \leq 0.5$  (detached binaries) or  $0.75 \leq c \leq 1.0$  (contact binaries), as recommended by Matijević et al. (2012).

Our selection of  $\delta$  Scuti stars began with 2405 stars manually identified as variable at frequencies above  $7 \text{ d}^{-1}$  from a preliminary version of the Murphy et al. (2019) catalogue. We randomly selected 1000 of these, and further refined this list to remove any stars that were also  $\gamma$  Doradus stars (i.e.  $\gamma$  Doradus/ $\delta$  Scuti hybrids) by manual inspection. From the same source, we also chose 500 stars that were not  $\delta$  Scuti pulsators, and removed stars with low-frequency variability to arrive at the 201-star non-variable class.

We selected the  $\gamma$  Doradus sample from the Deboscher et al. (2011) catalogue by choosing stars with a label confidence of  $> 95$  per cent and an effective temperature in the appropriate range. While the Deboscher et al. (2011) catalogue is known to have errors, this approach was taken due to a lack of an available list of  $\gamma$  Doradus stars exhibiting a broad range of oscillatory behaviours characteristic of the class – that is, a sample not restricted to neat and well-studied  $\gamma$  Doradus stars from which scientific inference has been made (references in Section 2.1). The addition of rigorous manual inspection ensured that the  $\gamma$  Doradus stars included in the final sample were significantly more likely to be correctly classified than in the Deboscher et al. (2011) catalogue, and that hybrid pulsators were removed.

Unlike the other classes, RR Lyr variables are not common in the *Kepler* data set. Of the 47 *Kepler* RR Lyr stars we found in the literature (Nemec et al. 2013; Molnár et al. 2018; Murphy et al. 2018), only 25 were observed in Q9. We admitted all 25 of these in the hope that we might discover additional RR Lyrae variables when classifying the remainder of the *Kepler* field (we did not).

The rotational variables were selected after trialling a preliminary version of our classifier, trained on the other six classes, on a test sample of *Kepler* stars. When visually inspecting the classification results across the six classes, we found that rotational variables constituted a considerable fraction of stars (approximately 15–20 per cent). From these, we added a list of 166 rotational variables to the training set after a second manual verification.

## 2.3 Processing *Kepler* data

Starting with SAP fluxes from Q9 light curves, we processed the data to remove instrumental variability by eliminating long-period trends

<sup>1</sup><http://kasoc.phys.au.dk/>

in the light curve of each star. Such variability can arise from physical drift of the telescope, causing changes in the flux levels falling in the aperture mask, as well as other instrumental effects distinct from stellar variability. Our processing involved division by a smoothed version of each light curve (smoothed using a Savitzky–Golay filter), removal of single-point outliers more than  $3\sigma$  from the mean of the smoothed light curve, and converting units to magnitudes.

Any gaps in the data of more than an hour (corresponding to two 29.45-min integrations) were padded with either the mean of the time-series for long gaps of four or more integrations, or the mean of the points on either side of smaller gaps. Even in high-quality quarters, long gaps arise from standard telescope operations such as the data downlinks that happen for approximately 24 h twice every quarter, while small gaps may be caused by e.g. cosmic ray events. Most machine-learning tools operate as functions of array index rather than explicit functions of time, hence it is imperative that these gaps are filled.

### 3 FEATURE-BASED LIGHT-CURVE CLASSIFICATION

Having constructed a training set, we next aimed to build a classifier to accurately predict the class of a star from features of its light-curve time-series. Our approach involved four steps: (i) mapping each light curve to a large feature vector, where each feature is a single, real-valued summary statistic that captures some interpretable property of the light curve; (ii) learning a classification rule that maps from a reduced subset of extracted features to the class label on a labelled training set; (iii) evaluating the performance of the learned classification rule on an independent test set; and (iv) applying this rule to classify the full *Kepler* data set.

#### 3.1 Feature extraction

The task of selecting relevant properties of a time-series for a given application, like light-curve classification, is commonly a manual one performed by a given researcher (e.g. Pashchenko et al. 2017). An alternative approach, termed ‘highly comparative time-series analysis’ (Fulcher, Little & Jones 2013; Fulcher & Jones 2014), is to include a large and comprehensive candidate set of possible time-series features, and take a data-driven approach to selecting those that are most relevant to the task at hand. To extract features from a light curve, we used a comprehensive candidate set of over 7000 time-series features from the HCTSA software package (v0.96) (Fulcher & Jones 2017).<sup>2</sup> The HCTSA feature set encompasses a wide range of time-series analysis methods, from properties of the distribution of time-series values, linear and nonlinear autocorrelation, entropy and complexity measures, stationarity, time-series model fits, wavelet and Fourier basis-function decompositions, and others (Fulcher et al. 2013). This approach allowed us to represent a set of  $L$  light curves as an  $L \times F$  matrix, where  $F$  is the number of features; applying HCTSA to our training data set yielded a  $1319 \times 7873$  light curve  $\times$  feature matrix, where each row is labelled according to one of the seven classes listed in Table 1. After performing feature extraction, we excluded features that contained special values (NaN, Inf), returned an error, or produced near-constant outputs (within  $10 \times$  machine precision) across all 1319 time-series, resulting in 6492 features after filtering. As a preprocessing prior to classification, feature values

were normalized to the unit interval using a scaled, outlier-robust sigmoidal transformation (Fulcher et al. 2013).

#### 3.2 Training and evaluating a classifier

In modern applications of machine learning, complexity is often introduced at the level of the classifier. In this work, we instead focused on selecting from a large candidate set of complex features, but using simple classifiers. This has the advantage of yielding features that can provide clear scientific interpretation, and follows the approach of Timmer et al. (1993): ‘The crucial problem is not the classifier function (linear or non-linear), but the selection of well-discriminating features. In addition, the features should contribute to an understanding’. For classification, we used a Gaussian mixture model (GMM; McLachlan & Peel 2000) on a labelled time-series  $\times$  feature matrix (described above). We fitted a single-Gaussian component to each of the seven training classes in feature space, combining them with equal prior probabilities to form a seven-component probability density function (PDF). While all classes are not equally common, equal priors are the simplest choice without knowing the true distribution of variable stars in the *Kepler* field. Classification was performed by evaluating the (posterior) probability of a star belonging to each class using the trained PDF, and selecting the class with highest probability. This GMM approach was substantially faster (by factors of approximately 10–100) than alternative algorithms such as nearest-neighbour clustering or support vector machines, but achieved similar classification performance on our training set.

We evaluated classification performance as the average balanced accuracy computed using 10-fold stratified cross-validation (Hastie, Tibshirani & Friedman 2009). Balanced accuracy,  $C_{\text{bal}}$ , accounts for class imbalance (the unequal number of observations in each class) in our data set and is defined as

$$C_{\text{bal}} = \frac{1}{m} \sum_{i=1}^m \frac{t_i}{c_i}, \quad (1)$$

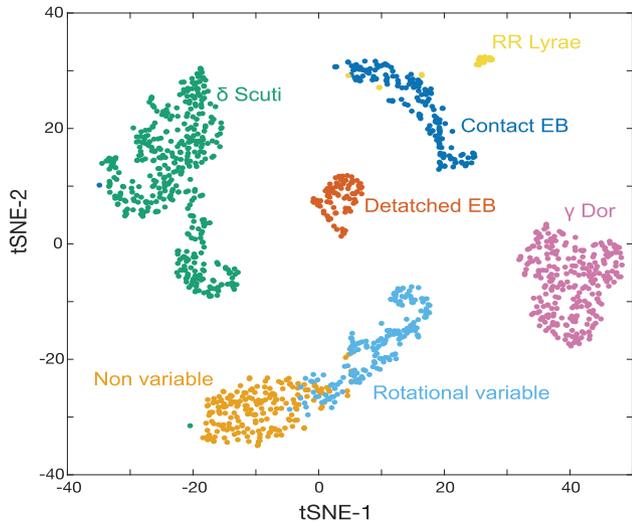
where  $m$  is the number of classes,  $t_i$  is the number of successfully identified time-series in the  $i$ th class, and  $c_i$  is the total number of time-series in this class.

#### 3.3 Feature subset selection

To extract a small number of HCTSA features that are most informative of the class labels, we used greedy forward feature selection (Hastie et al. 2009; Fulcher & Jones 2014). This simple algorithm iteratively builds a feature set, one feature at a time, with the objective of maximising the balanced classification accuracy,  $C_{\text{bal}}$ , at each iteration. That is, at iteration  $k$ , the algorithm searches across all individual features for the feature that maximizes  $C_{\text{bal}}$  when used in combination with the features selected in the  $k - 1$  previous iterations.

HCTSA was developed to encompass a comprehensive sample of the interdisciplinary time-series analysis literature, and thus contains groups of features with highly correlated behaviour (Fulcher et al. 2013; Henderson & Fulcher 2021). When multiple features exhibit similar classification performance, we implemented a simple heuristic constraint to favour the selection of faster-to-compute features: at each iteration, of the features with an accuracy within a margin of 1 per cent of the best-performing feature, the feature with the fastest computation time was selected. The iterative procedure was terminated when the improvement in training-set  $C_{\text{bal}}$  from adding

<sup>2</sup><https://github.com/benfulcher/hctsa>



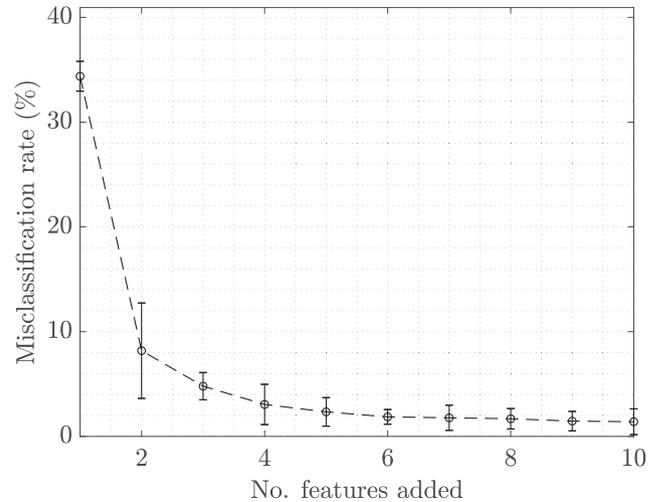
**Figure 2.** A 2D  $t$ -SNE projection of our *Kepler* training set of 1319 stars in the 6492D HCTSA feature space, where each light curve is coloured by its class label. Most stars form clear clusters that match their class identity, indicating that the HCTSA features provide a useful space in which to represent *Kepler* light curves.

another feature dropped below 1 percent. Note that applying this algorithm to the full training data set has the potential to overfit, since the selection step at each iteration (despite using cross-validation for each feature) uses the training set itself to select the best-performing feature. Accordingly, we evaluate the performance of our reduced feature set on an independent test set in Section 4.

## 4 RESULTS AND DISCUSSION

### 4.1 Representing light curves in a high-dimensional feature space

We first investigated the structure of the seven labelled classes of 1319 *Kepler* stars in the 6492D HCTSA feature space. We found that the HCTSA feature space is able to capture characteristic properties of the seven labelled classes of stars, obtaining a high mean 10-fold cross-validated balanced accuracy of 95.9 percent [using a linear support vector machine (Hastie et al. 2009), compared with a chance rate for seven classes of 14.3 percent]. This indicates that each type of star displays distinctive dynamics in ways that can be detected by the features included in HCTSA. To better understand the structure of light curves in the high-dimensional HCTSA feature space, we inspected a 2D  $t$ -SNE visualization ( $t$ -distributed stochastic neighbour embedding; Van Der Maaten & Hinton 2008). The result is shown in Fig. 2, where each point is a light curve, and light curves with similar features tend to be positioned closely in the space. While  $t$ -SNE is an unsupervised technique (Fig. 2 was constructed without knowledge of the class labels), stars are meaningfully organized according to their labelled class, with most stars clustering with other stars of the same type. Consistent with the high classification results reported above, this indicates that the HCTSA feature space captures distinctive dynamical properties of the light curves corresponding to the seven different types of stars. The plot also reveals scientifically meaningful structure between classes, such as the continuum from non-variable (light orange) stars to rotational-variable (light blue) stars. We also see a small overlap between RR Lyr stars and contact



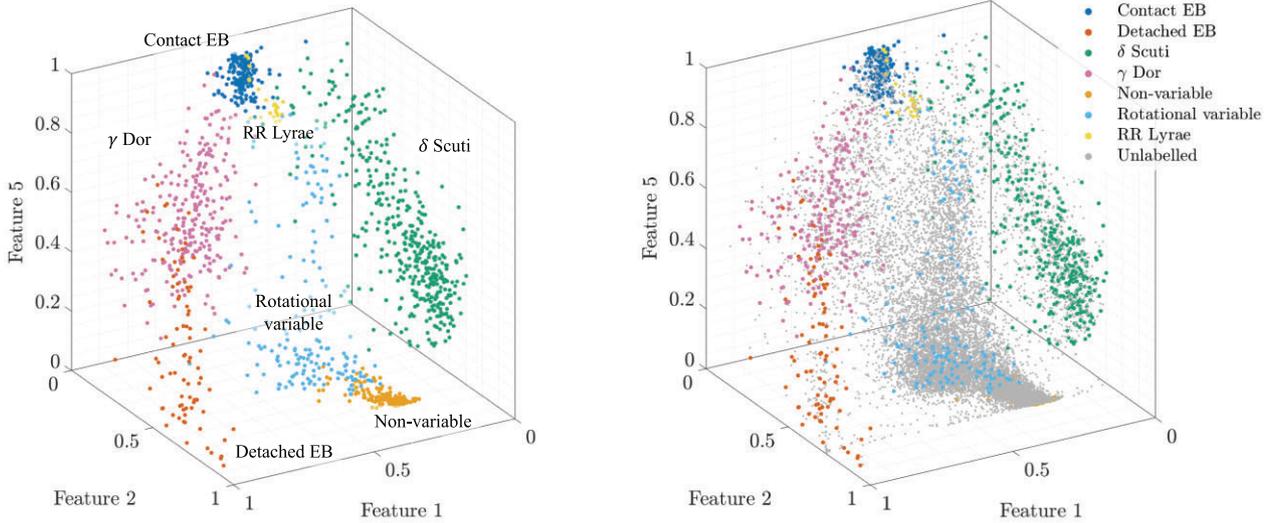
**Figure 3.** Classification performance as a function of the number of time-series features. Balanced misclassification rate on the training set (using a GMM classifier) is plotted as a function of selected features, shown as the mean and standard deviation across 10-fold cross-validation.

EBs, which reflects the similar morphologies of their light curves (see Fig. 1).

### 4.2 Representing light curves in a reduced feature space

The results above demonstrate that time-series properties in HCTSA can capture differences in light-curve dynamics between different types of stars. But which types of individual time-series features are most informative of these differences? To address this question, we aimed to construct a reduced set of HCTSA features that display strong classification performance using greedy forward selection (see Section 3.3 for details). The cross-validated balanced misclassification rate on the training set is shown as a function of the number of selected features in Fig. 3. This plot reveals that strong in-sample classification performance can be obtained with a relatively small set of well-chosen time-series features, e.g. a balanced accuracy of 95.2 percent with just three features. According to our termination criterion – when an additional feature provides  $< 1$  per cent marginal improvement in balanced accuracy – we obtained an informative 5D feature space in which to represent *Kepler* light curves.

To visualize how stars are organized in the reduced feature space, we plotted the training set in the space corresponding to three of the selected features in Fig. 4 (left). Despite a dramatic dimensionality reduction of each time-series – from the 4767 data points in a typical Q9 time-series to just three extracted summary statistics – the space meaningfully organizes all seven training classes in this low-dimensional feature representation, with each occupying a characteristic region of the space. Much like the  $t$ -SNE construction in Fig. 2, the relative positions of each class are consistent with what we would intuitively expect from their light curves and power spectra in Fig. 1. For example, detached binaries are highly separated from the other classes, as their light curves are the most distinct; non-variable stars blend with rotational variables when the rotations are weak and difficult to distinguish by eye, such that the light curves are almost non-varying; the  $\gamma$  Doradus and  $\delta$  Scuti stars lie on opposite sides of the space, reflecting their contrasting low- and high-frequency pulsations; and the contact binaries are close to the  $\gamma$  Doradus and RR Lyrae clusters, which all characteristically exhibit regular low-frequency variability.



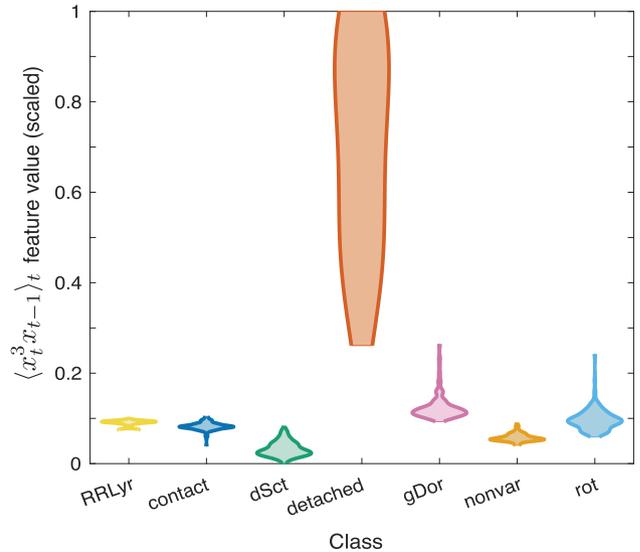
**Figure 4.** Separation of different classes of stars in the (normalized) space of three time-series features selected from HCTSA using greedy forward-feature selection. The two figures show (left) the training set, and (right) the training set alongside unlabelled *Kepler* data with  $6500 \text{ K} \leq T_{\text{eff}} \leq 10000 \text{ K}$ . The three features correspond to `AC_n1_001`, `MF_steps_ahead_ar_best_6_mabserr_5`, and `SP_Summaries_welch_rect_peakPower_5`, as described in detail in Section 4.2.1.

#### 4.2.1 The reduced feature set

We have demonstrated the usefulness of representing *Kepler* light curves in a low-dimensional feature space, but what types of properties are these features measuring, and what can that tell us about how light-curve dynamics differ between the seven classes of stars? In this section, we explain the five features in order of their selection by our greedy forward selection algorithm. Noting the small marginal improvements in accuracy after approximately three features (shown in Fig. 3), we focus in particular on these features. In the following discussion, note that the time-series were converted to magnitudes, so that positive excursions correspond to decreases in stellar flux, and vice versa.

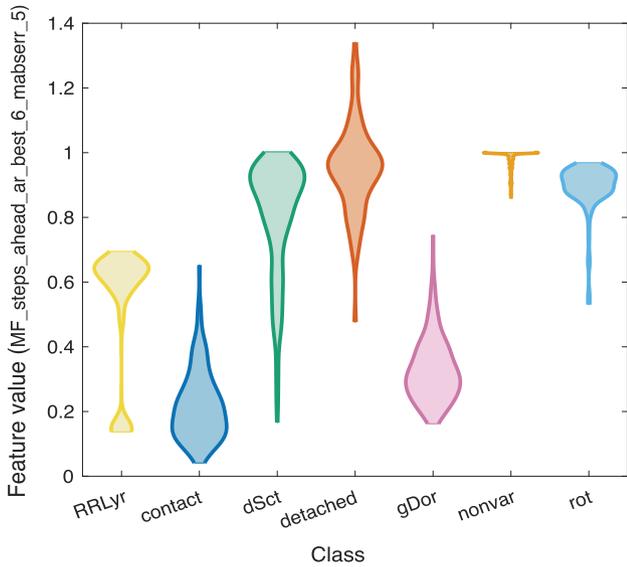
The first selected feature (labelled `AC_n1_001` in HCTSA), is a non-linear autocorrelation statistic that computes the time-average,  $\langle x_t^3 x_{t-1} \rangle_t$ , of the  $z$ -scored time-series  $x_t$ , with a time lag of 1 sample (approximately 30 min in the time domain). Similar to a lag-1 autocorrelation,  $\langle x_t x_{t-1} \rangle_t$ , it gives high values to highly autocorrelated light curves, but the modification ( $x_t^3$ ) accentuates large deviations from the mean. The distribution of this feature’s (sigmoid-normalized) values across the seven classes of stars is shown in Fig. 5. Detached binaries have the largest values of this statistic, driven by large positive excursions from the mean (since the time-series are in magnitudes). Autocorrelation arising from slower periodic patterns, as in  $\gamma$  Doradus, rotational stars, RR Lyr and contact binaries, lead to moderate positive values of `AC_n1_001`, while the non-variable stars have low values (raw values near-zero). The high-frequency oscillations seen in some  $\delta$  Scuti stars (e.g. Balona, Holdsworth & Cunha 2019; Bedding et al. 2020) resulted in negative values of `AC_n1_001` (the lowest normalized values).

Feature 2 (labelled `MF_steps_ahead_ar_best_6_mabserr_5` in HCTSA) uses a linear autoregressive (AR) model to measure how predictable a time-series is. This statistic captures how well an AR model (of optimal order, selected in the range 1–10 using Schwartz’s Bayesian Criterion) can predict 5 time-steps ahead in the time-series. This is measured relative to simple benchmark forecasting methods (including simple mean forecasts and a constant global-mean



**Figure 5.** Values for Feature 1, which computes  $\langle x_t^3 x_{t-1} \rangle_t$  (see Section 4.2.1). The violin plots show the normalized output of `AC_n1_001` across the seven classes of stars in the training set. Sigmoidal normalization scaled to the unit interval (see methods) was used to aid visualization of the large range of raw values of this feature.

forecast), calculated as the mean absolute error. The distribution of feature values across the seven classes of stars is shown in Fig. 6. Values near-zero indicate strong prediction performance of the AR model relative to simple benchmarks, while values greater than 1 indicate relatively inferior model performance. We see high values for the non-variable stars, detached binaries, rotational stars, and most of the  $\delta$  Scuti stars, with RR Lyr stars displaying intermediate values (a few RR Lyr stars with highly symmetric light curves have low values). The  $\gamma$  Doradus and contact binary light curves exhibit a strong linear correlation structure that allowed the AR models to



**Figure 6.** Distribution of Feature 2 values by class (see Section 4.2.1), which measures how predictable the time-series is using a linear autoregressive (AR) model; high values (near 1) are given to light curves for which the AR model performs worse than simple benchmarks, whereas values near 0 are given when the AR model strongly outperforms the benchmarks. Violin plots are shown for the distribution of this feature across the seven classes of stars.

make strong forecasts of these time-series, yielding low values for this statistic.

Feature 3, labelled `CO_trev_3_num` in HCTSA, evaluates the following time-average:  $\langle (x_t - x_{t-3})^3 \rangle_t$ . This statistic, using a time-lag  $\tau = 3$ , can be thought of as capturing asymmetry in the size of increases ( $x_t - x_{t-3} > 0$ ) versus decreases ( $x_t - x_{t-3} < 0$ ). For example, time-series with sudden increases but gradual decreases (at a time lag of three samples) will have large values of this feature. RR Lyr are distinguished by negative values of `CO_trev_3_num`, due to the characteristic asymmetry in the shapes of their light curves (e.g. Catelan & Smith 2015).

Feature 4, labelled `ST_LocalExtrema_n100_medianmax` in HCTSA, captures how positive outliers are distributed through the time-series. Operating on the  $z$ -scored time-series, this algorithm computes the maximum value in each of 100 overlapping windows (each containing 47 samples corresponding to approximately 23 h), and outputs the median of these local maxima. For time-series with relatively infrequent large positive excursions (like the light curves from many detached binaries, recalling that the calculations are done with magnitudes), most windows will have very low maxima, and thus the median of the maxima will be a low value. However, for time-series with maxima spaced more evenly throughout time, like most non-variable and  $\delta$  Scuti stars, high values are obtained for this statistic.

Feature 5, labelled `SP_Summaries_welch_rect_peakPower_5` in HCTSA, uses Welch’s method and a rectangular window to estimate the power spectrum and returns the proportion of power captured by the five most prominent identified peaks. Broadly, this feature gives high values to time-series that are well captured by a relatively small number of dominant frequencies. The lowest values for this feature were found for non-variable stars and rotational variables, while high values were obtained for contact binaries and RR Lyr stars.

**Table 3.** The test set of 515 *Kepler* stars. An extract of 12 stars is shown, with the full table provided in the supplementary material.

KIC ID	Class
8282730	Contact EB
6957185	Contact EB
8953296	Detached EB
5090690	Detached EB
8585472	$\delta$ Scuti
3648131	$\delta$ Scuti
6041803	$\gamma$ Doradus
8739181	$\gamma$ Doradus
5616145	Non-variable
8153411	Non-variable
3847563	Rotational
3967219	Rotational

GMM labels	contact	43 8.3%	3 0.6%	2 0.4%	1 0.2%	0 0.0%	6 1.2%	0 0.0%	78.2% 21.8%
	detached	0 0.0%	41 8.0%	0 0.0%	1 0.2%	18 3.5%	3 0.6%	0 0.0%	65.1% 34.9%
	dSct	3 0.6%	0 0.0%	72 14.0%	1 0.2%	0 0.0%	3 0.6%	0 0.0%	91.1% 8.9%
	gDor	2 0.4%	0 0.0%	0 0.0%	62 12.0%	0 0.0%	2 0.4%	0 0.0%	93.9% 6.1%
	nonvar	0 0.0%	0 0.0%	0 0.0%	0 0.0%	93 18.1%	3 0.6%	0 0.0%	96.9% 3.1%
	rot	3 0.6%	2 0.4%	3 0.6%	31 6.0%	7 1.4%	105 20.4%	0 0.0%	69.5% 30.5%
	RR Lyr	2 0.4%	0 0.0%	3 0.6%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0.0% 100%
		81.1% 18.9%	89.1% 10.9%	90.0% 10.0%	64.6% 35.4%	78.8% 21.2%	86.1% 13.9%	0.0% 0.0%	80.8% 19.2%
		Test labels							

**Figure 7.** Confusion matrix summarizing GMM classification performance. The GMM and test labels are the classifier predictions and manually assigned truth labels (respectively) for stars in our test set. Summaries in grey on the right of the matrix correspond to the (unbalanced) percentage of correct predictions, while summaries at the bottom are the (unbalanced) percentage of each class that was correctly classified. The raw classification accuracy is shown in blue (balanced accuracy 81.6 per cent). There were no RR Lyrae stars in our test set.

#### 4.2.2 Evaluation on a test data set

Having computed an informative low-dimensional space in which to represent *Kepler* light curves, we investigated its effectiveness in classifying variable stars outside our training set. We manually compiled a test set of 515 stars in the *Kepler* field belonging to classes of variable stars in our training set, and with  $6500 \text{ K} \leq T_{\text{eff}} \leq 10000 \text{ K}$ . The full list of test stars is provided as supplementary material, with a sample shown in Table 3.

To evaluate classification performance on the test set in the trained 5D feature space, we constructed a GMM consisting of seven Gaussian components, one fitted to each class in our training set (with uniform prior class probabilities), and used it to classify each of the test stars. Fig. 7 summarizes our results on the test set as a confusion matrix.

**Table 4.** Extract of GMM posterior class probabilities and probability density  $p(x)$  for 12 088 unlabelled stars in the *Kepler* field. The first 10 lines are shown, with the full table provided in the supplementary material.

KIC ID	Contact EB	Detached EB	$\delta$ Scuti	$\gamma$ Dor	Non-variable	Rotational variable	RR Lyrae	$p(x)$
757280	0.00	0.00	0.00	0.00	0.00	1.00	0.00	0.73
892667	0.00	0.00	0.00	0.00	0.00	1.00	0.00	9.12
892828	0.00	0.00	0.00	0.00	0.93	0.07	0.00	307.19
893234	0.00	0.00	0.00	0.00	0.01	0.99	0.00	4.12
893944	0.00	0.00	0.00	0.00	1.00	0.00	0.00	3802.55
1026133	0.00	0.00	0.00	0.00	0.00	1.00	0.00	1.02
1026255	0.00	0.00	0.00	0.00	0.51	0.49	0.00	0.71
1026475	0.00	0.00	0.00	0.00	0.00	1.00	0.00	9.54
1026861	0.00	0.00	0.00	1.00	0.00	0.00	0.00	0.44

The confusion matrix can be interpreted as follows. Labels on each row were assigned by the GMM classifier, while labels on each column correspond to manually assigned labels from our test set. Each cell  $(i, j)$  of the confusion matrix shows the number of stars (and percentage of all stars considered) that were classified as category  $i$  by the GMM, and as category  $j$  in our test set. For example, there were 18 stars classified as detached binaries by the GMM but labelled as non-variable in the test set. Diagonal elements of the matrix (in green) correspond to correctly classified stars. Summaries in grey on the right of the matrix correspond to the (unbalanced) percentage of correct predictions, while summaries at the bottom are the percentage of each class that was correctly classified. The raw classification accuracy is shown in blue in the bottom-right corner.

Our classifier achieved a balanced accuracy of 81.6 per cent on the test set and performed well on all classes, with two understandable exceptions highlighted in Fig. 7:

(i) Non-variable stars are commonly misclassified as detached binaries (18 misclassifications). Most have sharp transitions in their light curves at the beginning or end of the quarter, or just before or after the *Kepler* telescope paused observation for data transmission. These transitions appear as sharp peaks or troughs, and are represented similarly to eclipses in our feature space.

(ii)  $\gamma$  Doradus stars are commonly misclassified as rotational variables. Both classes have low-frequency variations (e.g. Li et al. 2019) and even for an expert eye, it can be difficult to resolve  $\gamma$  Doradus oscillations from a single quarter of *Kepler* data. The behaviour may therefore look similar to rotation in our feature space.

Apart from these exceptions, our approach yielded high overall classification accuracies despite relying on very simple methods (greedy forward feature selection and GMM classification), demonstrating the usefulness of the comprehensive HCTSA feature space in highlighting high-performing interpretable features for a given problem. We expect that repeating the feature selection and classification procedures with more sophisticated algorithms, while still working with a rich set of interpretable features, would further improve the accuracy reported here using simple methods. However, as discussed in Section 4.3, our method is already a useful tool for classifying and searching large data sets.

### 4.3 Classifying the *Kepler* field

In this section, we use the low-dimensional feature space learned from the training set and validated on the test set to classify variable stars across the entire *Kepler* field. Our full classification catalogue is provided as supplementary material (Table 4).

#### 4.3.1 Classifying unlabelled stars

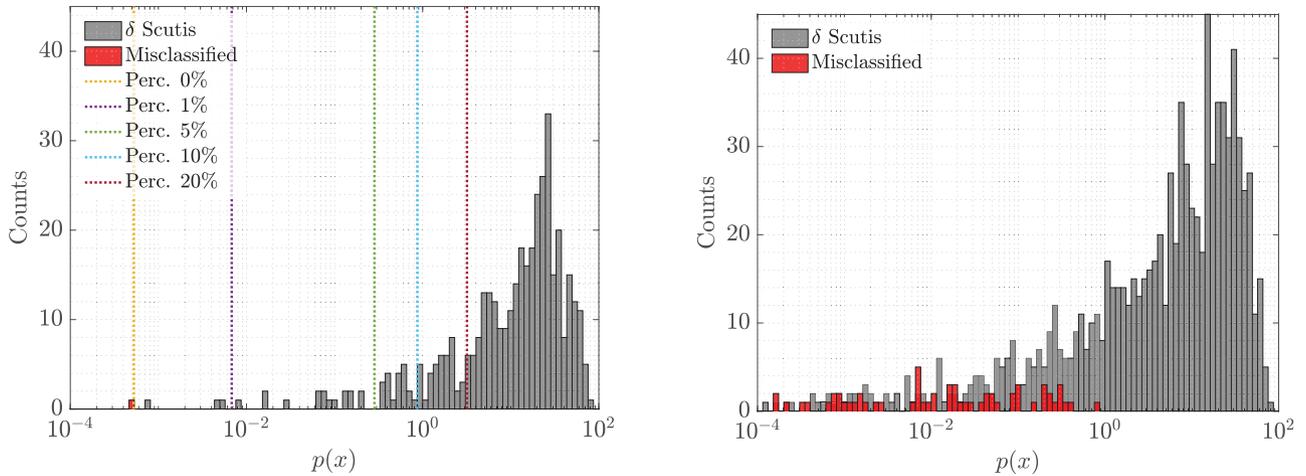
We computed the 5D feature-space representation of all 12 088 stars with Q9 data in the *Kepler* field and  $6500 \text{ K} \leq T_{\text{eff}} \leq 10\,000 \text{ K}$  (excluding our training set of 1319 stars). These are plotted in the right-hand panel of Fig. 4 as unlabelled stars (grey). Each feature vector was normalized using the same scaled robust sigmoid transformation (including its coefficients) as used on the training set, preserving the structure of our normalized feature space. The grey unlabelled stars are clearly clustered around the coloured training groups, with the majority of stars residing near the non-variable cluster. This clustering occurs naturally because of our choice of feature space. Intuitively, we might expect that (i) unlabelled stars near each training cluster belong to that respective class; (ii) stars midway between two groups are hybrids of both classes; and (iii) stars far from any group are new classes of variable stars unaccounted for in our training set. We have already verified the first of these hypotheses by applying our GMM classifier to the test set with reasonably high accuracy in Section 4.2.2. We leave investigation of the remaining two claims as future work.

We evaluated our trained GMM classifier on all 12 088 unlabelled stars to generate a catalogue of posterior probabilities, giving a predicted classification for each star. The first 10 lines of this catalogue are shown in Table 4, with the full catalogue provided as supplementary material. For each star in the catalogue, its classification is the class with maximum posterior probability. The catalogue is intended as a useful tool in searching for candidate variable stars of interest. We note that this catalogue and our broader methodology have already proven useful in identifying new  $\gamma$  Doradus stars (Li et al. 2019) and  $\delta$  Scuti stars (Murphy et al. 2020) in the *Kepler* field. We provide suggestions for searching our catalogue in Section 4.3.2.

#### 4.3.2 Using probability density as a confidence heuristic

Close examination of Fig. 4 reveals that many unlabelled stars lie in areas between the training set clusters, far from where the GMM classifier was trained. We may therefore ask: how does our classification accuracy improve if we restrict the test set to stars ‘near’ the training distributions? We define  $p(x)$  as the probability density of a star represented by feature vector  $x$ , where the PDF is the 7-class GMM used for classification. Stars close (in feature space) to the centre of the multivariate Gaussians will have large probability densities,  $p(x)$ , while those far from the class centroids will have low  $p(x)$ . In this sense, we can use  $p(x)$  as a heuristic measure of how likely a star is to belong to any of the training classes.  $p(x)$  is provided in the final column of Table 4.

As an example, the left-hand and right-hand panels of Fig. 8 show the distributions of  $p(x)$  for all  $\delta$  Scuti stars in the training set and



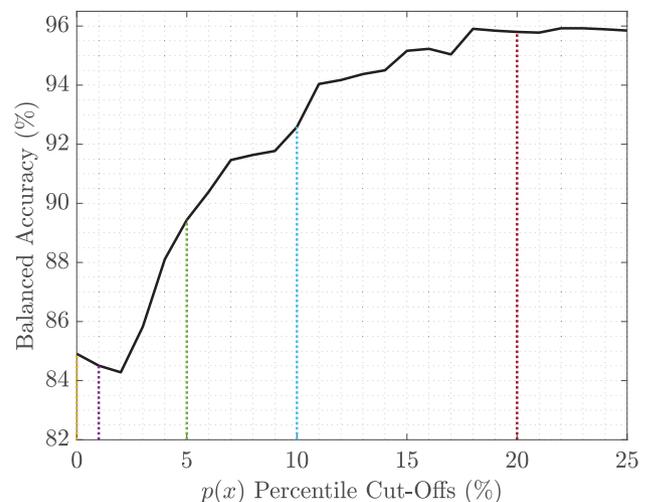
**Figure 8.** Histograms of GMM probability density  $p(x)$  for all stars classified as  $\delta$  Scuti in (left) the training set, and (right) the remainder of the *Kepler* field with  $6500 \text{ K} \leq T_{\text{eff}} \leq 10\,000 \text{ K}$ . Lines in the training distribution indicate  $p(x)$  percentile cut-offs. For example, 90 per cent of stars classified as  $\delta$  Scuti in the training set lie above the blue line. Stars in the right-hand panel are classified according to Murphy et al. (2019).

the rest of the *Kepler* field (in our temperature range of interest), respectively. Classifications in the right-hand panel of Fig. 8 are from the (manually compiled) Murphy et al. (2019) catalogue. As anticipated, misclassification is far more common at low densities. Interestingly, the distribution in Fig. 8 shows that above  $p(x) \approx 1$ , all predictions of  $\delta$  Scuti stars are correct. We would intuitively expect similar  $p(x)$  cut-offs for the other classes, above which we have high confidence in the GMM predictions. However, defining exact cut-offs is impossible without a full labelled catalogue of the *Kepler* field. Instead, we can define values for  $p(x)$  representing regions of increasing proximity to our trained distribution. The vertical lines in the left-hand panel of Fig. 8 show  $p(x)$  percentile cut-offs, above which a certain percentage of the training data fall. For example, only the 90 per cent ‘closest’  $\delta$  Scuti training stars to the  $\delta$  Scuti centroid (in terms of probability density) lie above the blue line in Fig. 8.

The results above, for  $\delta$  Scuti stars, suggest that our predictions are more accurate in higher confidence areas of the feature space, corresponding to areas with higher modelled density for the training set. To test whether this holds more generally, we computed the balanced classification accuracy (across all classes) on the test set for a range of  $p(x)$  percentile thresholds. As shown in Fig. 9, we find that accuracy improves with more stringent restrictions on  $p(x)$ , demonstrating the usefulness of  $p(x)$  as a proxy for prediction confidence. Even small restrictions in  $p(x)$ , such as the 95<sup>th</sup> percentile cut-off (green line), improve the classification performance on our test set to approximately 90 per cent accuracy. This is an example of a useful way to search our catalogue and obtain a list of confidently classified variable stars for further analysis – as  $p(x)$  increases for each class, so too does the confidence of our predictions. We once again stress that such intuitive search criteria are a direct consequence of our choice of feature space and simple classification algorithm. One could achieve even more accurate results with more sophisticated approaches, but this may come at the expense of interpretability of our low-dimensional feature space.

#### 4.4 Comparison with Audenaert et al. (2021)

When our paper was in the final stages of preparation, a new classification of *Kepler* light curves was published by Audenaert et al. (2021, hereafter Aud21). Their work was done as part of



**Figure 9.** Balanced classification accuracy as a function of  $p(x)$  percentile cut-offs. Applying the classifier to stars close to regions of feature space that we trained on significantly improves the overall accuracy. Dotted lines correspond to the same percentile cut-offs overlaid in Fig. 8.

efforts to design an automated classification algorithm for the *TESS* mission. Given the complementary nature of Aud21 and our own study, especially given that both were based on Q9 data, it is worthwhile to carry out a brief comparison. We should keep in mind that the emphasis in Aud21 was on providing a high-performance classification pipeline from existing methods, whereas ours involved designing an interpretable classifier from a rich library of time-series features.

The classification by Aud21 included about 167 000 *Kepler* Q9 light curves, regardless of effective temperature, whereas our work is restricted to about 12 000 stars with  $6500 \text{ K} \leq T_{\text{eff}} \leq 10\,000 \text{ K}$ . We have compared our classifications with Aud21 in Fig. 10. There is an obvious mapping between most of our classes and those used by Aud21, with the following differences:

- (i) Aud21 combined contact EBs and rotational (spotted) variables into a single class.

Audenaert et al. (2021) labels	contact EB/spots	5989 49.5%	76 0.6%	95 0.8%	153 1.3%	3094 25.6%	1 0.0%	63.7% 36.3%
	transit/eclipse	55 0.5%	77 0.6%	3 0.0%	8 0.1%	2 0.0%	0 0.0%	53.1% 46.9%
	dSct/bCep	106 0.9%	2 0.0%	825 6.8%	10 0.1%	3 0.0%	7 0.1%	86.6% 13.4%
	gDor/SPB	637 5.3%	1 0.0%	71 0.6%	555 4.6%	7 0.1%	0 0.0%	43.7% 56.3%
	constant	3 0.0%	2 0.0%	2 0.0%	0 0.0%	135 1.1%	0 0.0%	95.1% 4.9%
	RRLyrr/Ceph	11 0.1%	1 0.0%	1 0.0%	0 0.0%	0 0.0%	2 0.0%	13.3% 86.7%
	solar-like	104 0.9%	2 0.0%	0 0.0%	0 0.0%	33 0.3%	0 0.0%	0.0% 100%
	aperiodic	10 0.1%	4 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	0.0% 100%
		86.6% 13.4%	46.7% 53.3%	82.7% 17.3%	76.3% 23.7%	4.1% 95.9%	20.0% 80.0%	62.7% 37.3%
		contact detached /rot	dSct	gDor	nonvar	RRLyrr		
	This work (GMM labels)							

**Figure 10.** Confusion matrix comparing the results of the [Aud21](#) classifier to our own classifications for about 12 000 stars in the *Kepler* field with  $6500 \text{ K} \leq T_{\text{eff}} \leq 10\,000 \text{ K}$ . Much like in Fig. 7, the grey summary boxes on the right correspond to the percentage of [Aud21](#) labels that agree with our GMM predictions, while summaries at the bottom are the percentage of our predictions for each class that agree with [Aud21](#).

(ii) [Aud21](#) included  $\delta$  Scuti stars in a class with  $\beta$  Cephei stars. These have similar light curves but the  $\beta$  Cep pulsators have higher effective temperatures that lie outside the range of our sample. Similarly, [Aud21](#) combined  $\gamma$  Doradus stars with SPBs (slowly pulsating B stars), which are also hotter than our sample.

(iii) [Aud21](#) included a class for solar-like oscillators, which should not appear in our sample because they occur in stars whose effective temperatures fall below our range.

(iv) [Aud21](#) also included a class for aperiodic variables.

We see from the confusion matrix in Fig. 10 that there is generally excellent agreement between our results and those of [Aud21](#). We briefly discuss the areas with the greatest disagreement:

(i) 3094 stars that our classifier labelled as non-variable were classified by [Aud21](#) as rotational/contact EBs. We inspected 200 of these light curves (and their Fourier amplitude spectra) and found that most are non-variable, with some showing a weak rotation signal.

(ii) 637 stars that we labelled as contact EBs or rotational variables were classified by [Aud21](#) as  $\gamma$  Doradus pulsators. Inspection of 200 light curves shows that most are indeed  $\gamma$  Doradus stars. This may be a shortcoming of our specific feature space and classifier, particularly when considering Fig. 7, where the same disagreement occurs between our GMM classifications and our independent test set.

(iii) 153 stars were labelled by us as  $\gamma$  Doradus stars and by [Aud21](#) as contact EBs or rotational variables. Inspection of these shows that many are indeed  $\gamma$  Doradus stars, although it is sometimes difficult to be sure.

(iv) 139 stars in our sample were labelled by [Aud21](#) as having solar-like oscillations, which is not a class that we considered because these oscillations occur in stars below our temperature range.

Our classifications for these light curves were mainly as rotational variables, contact EBs, or non-variable. We inspected all 139 light curves and found that our classifications were mostly correct.

(v) 106 stars were labelled by us as contact EBs or rotational variables, and by [Aud21](#) as  $\delta$  Scuti stars. We inspected all light curves and found that most have  $\delta$  Scuti pulsations, but many also have low-frequency variability.

Finally, we note KIC 10024862, which is one of two stars listed by [Aud21](#) as non-variable and by our algorithm as a detached binary. In fact, Kawahara & Masuda (2019) identified this as a Jupiter-sized exoplanet in a long-period orbit that has only one transit during the 4-yr *Kepler* mission, which happened to be in Q9. This suggests that it might be worthwhile to look in more detail at groups for which classification methods are in disagreement for a small number of stars.

Much like our approach, [Aud21](#) assigned labels to each star according to the class with the highest posterior probability from their classifier. Fig. 10 therefore contains samples where either classifier may be confused – for example, a given light curve may have probabilities of 0.34, 0.35, and 0.01 split between three classes and the maximum probability (0.35) is relatively low. Not surprisingly, we found that by restricting to stars with a high maximum probability in both samples, the agreement increased between our classification labels and those from [Aud21](#). A detailed comparison of the two catalogues goes beyond the scope of this paper and would require a measure of label confidence from the [Aud21](#) classifier similar to the probability density heuristic from Section 4.3.2.

In general, we conclude that the two approaches produce results that generally agree well. The difference in point (i) reflects the subjectivity in drawing the line between variables and non-variables (and perhaps also different amounts of filtering applied to the light curves). Points (ii) and (iii) reflect the difficulty – especially with short data sets – in deciding whether low-frequency variability is due to pulsation or rotation (e.g. Briquet et al. 2007; Lee 2021; Kurtz 2022).

## 5 CONCLUSIONS

We have used a feature-based machine-learning algorithm to classify *Kepler* light curves for stars with effective temperatures in the range 6500–10 000 K. We first created a training set of 1319 light curves, which we classified into seven classes:  $\delta$  Scuti stars,  $\gamma$  Doradus stars, RR Lyrae stars, rotational variables, contact EBs, detached EBs, and non-variable stars. We built a classifier using features selected with the HCTSA package (highly comparative time-series analysis; Fulcher & Jones 2017), which includes over 7000 time-series features. We found that five features were sufficient to represent the training set with a balanced accuracy of 98 per cent, and a separate test set of 500 stars with a balanced accuracy of 82 per cent.

We used our method to classify *Kepler* light curves for all 12 000 stars with effective temperatures in the range 6500–10 000 K, and the results are tabulated in the online supplementary material (Table 4). We further outlined a confidence heuristic based on probability density to search our catalogue and extract candidate lists of correctly classified variable stars. We also compared our classifications to recent work on the same light curves by [Aud21](#) and generally found good agreement.

While many modern approaches to machine learning focus on performance over interpretability (resulting in the common description of being ‘black-box’ algorithms), here we favoured the selection of high-performing and interpretable features to meaningfully

represent *Kepler* light curves. Given the ease with which our five features can be computed for a large data base of light curves, comparing complex classification algorithms to our methods could provide an independent benchmark for general light-curve classification algorithms, much like we have shown with our comparison to Aud21.

Further extensions of this work might include using our catalogue to search for rare classes of variable stars, hybrid systems, and new stars entirely different from our training sample. In particular, we expect stars with roughly equal posterior probabilities between two classes to be hybrid systems, and very different stars to have much lower probability density scores than any other star in the *Kepler* field. Our methods could also be applied to individual classes of variable stars to try to identify interesting or unusual behaviour within a class, such as the recently discovered high-frequency  $\delta$  Scuti stars (Bedding et al. 2020). Another possibility is to extend our intuitive feature-based methods by adding more complex feature selection and classification algorithms. Such extensions are likely to improve our already strong classification performance, and strengthen results when applying our methods to even larger photometric surveys, such as that from *TESS*.

## ACKNOWLEDGEMENTS

We thank the *Kepler* team for providing such a wonderful data set. We gratefully acknowledge support from the Australian Research Council through DECRA grant DE180101104 and Discovery Project DP210103119, and from the Danish National Research Foundation (grant DNR106) through its funding for the Stellar Astrophysics Centre (SAC). TVR gratefully acknowledges support from the Research Foundation Flanders (FWO) under grant agreement number 12ZB620N.

## DATA AVAILABILITY

The data underlying this article are available at the Kepler Asteroseismic Science Operations Center (KASOC), at <http://kasoc.phys.uu.dk/>. The HCTSA software is available at <https://github.com/benfulcher/hctsa>.

## REFERENCES

Armstrong D. J. et al., 2016, *MNRAS*, 456, 2260  
 Audenaert J. et al., 2021, *AJ*, 162, 209 (Aud21)  
 Ball N. M., Brunner R. J., 2010, *Int. J. Mod. Phys. D*, 19, 1049  
 Balona L. A., 2013, *MNRAS*, 431, 2240  
 Balona L. A., 2018, *MNRAS*, 479, 183  
 Balona L. A., Holdsworth D. L., Cunha M. S., 2019, *MNRAS*, 487, 2117  
 Baron D., 2019, preprint ([arXiv:1904.07248](https://arxiv.org/abs/1904.07248))  
 Bass G., Borne K., 2016, *MNRAS*, 459, 3721  
 Bassi S., Sharma K., Gomekar A., 2021, *Front. Astron. Space Sci.*, 8, 168  
 Bedding T. R. et al., 2020, *Nature*, 581, 147  
 Blomme J. et al., 2010, *ApJ*, 713, L204  
 Blomme J. et al., 2011, *MNRAS*, 418, 96  
 Bouabid M. P., Dupret M. A., Salmon S., Montalbán J., Miglio A., Noels A., 2013, *MNRAS*, 429, 2500  
 Bowman D. M., Kurtz D. W., 2018, *MNRAS*, 476, 3169  
 Briquet M., Hubrig S., De Cat P., Aerts C., North P., Schöller M., 2007, *A&A*, 466, 269  
 Cabral J. B., Ramos F., Gurovich S., Granitto P. M., 2020, *A&A*, 642, A58  
 Carrasco-Davis R. et al., 2019, *PASP*, 131, 108006  
 Catelan M., Smith H. A., 2015, *Pulsating Stars*. Wiley-VCH, Weinheim

Debusscher J., Blomme J., Aerts C., De Ridder J., 2011, *A&A*, 529, A89  
 Fulcher B. D., Jones N. S., 2014, *IEEE Trans. Knowl. Data Eng.*, 26, 3026  
 Fulcher B. D., Jones N. S., 2017, *Cell Syst.*, 5, 527  
 Fulcher B. D., Little M. A., Jones N. S., 2013, *J. R. Soc. Interface*, 10, 20130048  
 Giles D. K., Walkowicz L., 2020, *MNRAS*, 499, 524  
 Graff P., Feroz F., Hobson M. P., Lasenby A., 2014, *MNRAS*, 441, 1741  
 Guzik J. A., 2021, *Front. Astron. Space Sci.*, 8, 55  
 Hastie T., Tibshirani R., Friedman J., 2009, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd edn. Springer-Verlag, Berlin  
 Henderson T., Fulcher B. D., 2021, *International Conference on Data Mining Workshops (ICDMW)*, p. 1032  
 Hon M., Stello D., Yu J., 2017, *MNRAS*, 469, 4578  
 Hon M., Stello D., Yu J., 2018a, *MNRAS*, 476, 3233  
 Hon M., Stello D., Zinn J. C., 2018b, *ApJ*, 859, 64  
 Hosenie Z., Lyon R., Stappers B., Mootooyaloo A., McBride V., 2020, *MNRAS*, 493, 6050  
 Ivezić Ž., Connelly A. J., Vanderplas J. T., Gray A., 2019, *Statistics, Data Mining, and Machine Learning in Astronomy*. Princeton Univ. Press, Princeton, NJ  
 Jackiewicz J., 2021, *Front. Astron. Space Sci.*, 7, 102  
 Jamal S., Bloom J. S., 2020, *ApJS*, 250, 30  
 Johnston K. B., Haber R., Caballero-Nieves S. M., Peter A. M., Petit V., Knotte M., 2019, *Comput. Astrophys. Cosmol.*, 6, 4  
 Johnston K. B., Caballero-Nieves S. M., Petit V., Peter A. M., Haber R., 2020, *MNRAS*, 491, 3805  
 Kawahara H., Masuda K., 2019, *AJ*, 157, 218  
 Kaye A. B., Handler G., Krisciunas K., Poretti E., Zerbi F. M., 1999, *PASP*, 111, 840  
 Kgoadi R., Whittingham I., Engelbrecht C., 2019, in Griffin R. E., ed., *Proc. IAU Symp. 339, Southern Horizons in Time-Domain Astronomy*. Cambridge Univ. Press, Cambridge, p. 310  
 Kirk B. et al., 2016, *AJ*, 151, 68  
 Kurtz D., 2022, preprint ([arXiv:2201.11629](https://arxiv.org/abs/2201.11629))  
 Kuzlewicz J. S., Hekker S., Bell K. J., 2020, *MNRAS*, 497, 4843  
 Le Saux A., Bugnet L., Mathur S., Breton S. N., García R. A., 2019, in Di Matteo P., Creevey O., Crida A., Kordopatis G., Malzac J., Marquette J. B., N'Diaye M., Venot O., eds, *SF2A-2019: Proceedings of the Annual meeting of the French Society of Astronomy and Astrophysics*. p. 437  
 Lee U., 2021, *MNRAS*, 505, 1495  
 Li G., Van Reeth T., Bedding T. R., Murphy S. J., Antoci V., 2019, *MNRAS*, 487, 782  
 Li G., Van Reeth T., Bedding T. R., Murphy S. J., Antoci V., Ouazzani R.-M., Barbara N. H., 2020, *MNRAS*, 491, 3586  
 McLachlan G. J., Peel D., 2000, *Finite mixture models. Probability and Statistics – Applied Probability and Statistics Section*, Vol. 299. Wiley, New York  
 McQuillan A., Mazeh T., Aigrain S., 2014, *ApJS*, 211, 24  
 Mathur S. et al., 2017, *ApJS*, 229, 30  
 Matijevič G., Prša A., Orosz J. A., Welsh W. F., Bloemen S., Barclay T., 2012, *AJ*, 143, 123  
 Molnár L., Plachy E., Juhász Á. L., Rimoldini L., 2018, *A&A*, 620, A127  
 Mombarg J. S. G., Van Reeth T., Pedersen M. G., Molenberghs G., Bowman D. M., Johnston C., Tkachenko A., Aerts C., 2019, *MNRAS*, 485, 3248  
 Murphy S. J., Moe M., Kurtz D. W., Bedding T. R., Shibahashi H., Boffin H. M. J., 2018, *MNRAS*, 474, 4322  
 Murphy S. J., Hey D., Van Reeth T., Bedding T. R., 2019, *MNRAS*, 485, 2380  
 Murphy S. J., Barbara N. H., Hey D., Bedding T. R., Fulcher B. D., 2020, *MNRAS*, 493, 5382  
 Nemec J. M., Cohen J. G., Ripepi V., Derekas A., Moskalik P., Sesar B., Chadid M., Bruntt H., 2013, *ApJ*, 773, 181

- Nielsen M. B., Gizon L., Schunker H., Karoff C., 2013, *A&A*, 557, L10  
Oort J. H., Plaut L., 1975, *A&A*, 41, 71  
Ouazzani R. M., Marques J. P., Goupil M. J., Christophe S., Antoci V., Salmon S. J. A. J., Ballot J., 2019, *A&A*, 626, A121  
Pashchenko I. N., Sokolovsky K. V., Gavras P., 2017, *MNRAS*, 475, 2326  
Paul S., Chattopadhyay T., 2022, preprint ([arXiv:2201.08755](https://arxiv.org/abs/2201.08755))  
Pietrukowicz P. et al., 2017, *Nat. Astron.*, 1, 0166  
Sikora J., Wade G. A., Rowe J., 2020, *MNRAS*, 498, 2456  
Szklenár T., Bódi A., Tarczay-Nehéz D., Vida K., Marton G., Mező G., Forró A., Szabó R., 2020, *ApJ*, 897, L12  
Timmer J., Gantert C., Deuschl G., Honerkamp J., 1993, *Biol. Cybern.*, 70, 75  
Tsang B. T. H., Schultz W. C., 2019, *ApJ*, 877, L14  
Van Der Maaten L., Hinton G., 2008, *J. Mach. Learn. Res.*, 9, 2579  
Van Reeth T. et al., 2018, *A&A*, 618, A24  
Walker A. R., 1992, *ApJ*, 390, L81  
Wolff S. C., 1983, *The A-type Stars: Problems and Perspectives*. NASA SP-463, Washington D.C.
- Yu J., Huber D., Bedding T. R., Stello D., Hon M., Murphy S. J., Khanna S., 2018, *ApJS*, 236, 42  
Zhang K., Bloom J. S., 2021, *MNRAS*, 505, 515

## SUPPORTING INFORMATION

Supplementary data are available at [MNRAS](https://www.mnras.org/onlineonly) online.

**table2-training-set.txt**

**table3-test-set.txt**

**table4-classification-posteriors.txt**

Please note: Oxford University Press is not responsible for the content or functionality of any supporting materials supplied by the authors. Any queries (other than missing material) should be directed to the corresponding author for the article.

This paper has been typeset from a  $\text{\TeX}/\text{\LaTeX}$  file prepared by the author.