

Original citation:

Armstrong, David J., Kirk, J., Lam, K. W. F., McCormac, J. J., Osborn, H. P., Spake, J., Walker, S., Brown, D. J. A., Kristiansen, M. H., Pollacco, Don, West, Richard G. and Wheatley, P. J.. (2016) K2 variable catalogue – II. Machine learning classification of variable stars and eclipsing binaries in K2 fields 0–4. *Monthly Notices of the Royal Astronomical Society*, 456 (2). pp. 2260-2272.

Permanent WRAP URL:

<http://wrap.warwick.ac.uk/80249>

Copyright and reuse:

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions. Copyright © and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable the material made available in WRAP has been checked for eligibility before being made available.

Copies of full items can be used for personal research or study, educational, or not-for profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

Publisher's statement:

This article has been accepted for publication in *Monthly Notices of the Royal Astronomical Society* © 2016 The Authors Published by Oxford University Press on behalf of the Royal Astronomical Society. All rights reserved.

Link to final published version: <http://dx.doi.org/10.1093/mnras/stv2836>

A note on versions:

The version presented in WRAP is the published version or, version of record, and may be cited as it appears here.

For more information, please contact the WRAP Team at: wrap@warwick.ac.uk

K2 variable catalogue – II. Machine learning classification of variable stars and eclipsing binaries in K2 fields 0–4

D. J. Armstrong,^{1,2★} J. Kirk,¹ K. W. F. Lam,¹ J. McCormac,¹ H. P. Osborn,¹ J. Spake,^{1,3} S. Walker,¹ D. J. A. Brown,¹ M. H. Kristiansen,⁴ D. Pollacco,¹ R. West¹ and P. J. Wheatley¹

¹University of Warwick, Department of Physics, Gibbet Hill Road, Coventry CV4 7AL, UK

²ARC, School of Mathematics & Physics, Queen's University Belfast, University Road, Belfast BT7 1NN, UK

³Astrophysics Group, School of Physics, University of Exeter, Stocker Road, Exeter EX4 4QL, UK

⁴DTU Space, National Space Institute, Technical University of Denmark, Elektrovej 327, DK-2800 Lyngby, Denmark

Accepted 2015 December 1. Received 2015 November 20; in original form 2015 October 9

ABSTRACT

We are entering an era of unprecedented quantities of data from current and planned survey telescopes. To maximize the potential of such surveys, automated data analysis techniques are required. Here we implement a new methodology for variable star classification, through the combination of Kohonen Self-Organizing Maps (SOMs, an unsupervised machine learning algorithm) and the more common Random Forest (RF) supervised machine learning technique. We apply this method to data from the K2 mission fields 0–4, finding 154 ab-type RR Lyraes (10 newly discovered), 377 δ Scuti pulsators, 133 γ Doradus pulsators, 183 detached eclipsing binaries, 290 semidetached or contact eclipsing binaries and 9399 other periodic (mostly spot-modulated) sources, once class significance cuts are taken into account. We present light-curve features for all K2 stellar targets, including their three strongest detected frequencies, which can be used to study stellar rotation periods where the observed variability arises from spot modulation. The resulting catalogue of variable stars, classes, and associated data features are made available online. We publish our SOM code in PYTHON as part of the open source PYMVPA package, which in combination with already available RF modules can be easily used to recreate the method.

Key words: methods: data analysis – techniques: photometric – catalogues – binaries: eclipsing – stars: variables: general.

1 INTRODUCTION

Data flows from new and planned astronomical survey telescopes are steadily increasing. This shows no sign of stopping, with Large Synoptic Survey Telescope (LSST) starting operations in ~ 2020 . There is clearly a need for accurate, fast, automated classification of photometric light curves to maximize the scientific returns from these surveys. Even when later spectroscopic followup is required, finding which targets to prioritize is a necessary first step.

The literature contains multiple examples of such classification, using a wide variety of techniques. These include a variety of supervised machine learning applications (e.g. Eyer & Blake 2005; Mahabal et al. 2008; Blomme et al. 2010; Debosscher et al. 2011; Brink et al. 2013; Nun et al. 2014). Recently Random Forests (RFs) have begun to gain popularity, due to their robustness and appli-

cability to different sets of data, extracted light-curve properties, and classification schemes (e.g. Richards et al. 2011a, 2012; Masci et al. 2014). Several improvements have been proposed, in areas such as parametrizing light curves with maximal information retention (Kügler, Gianniotis & Polsterer 2015), and adjusting for training set deficiencies (Richards et al. 2011b). One method of *unsupervised* machine learning is a Kohonen Self-Organizing Map (SOM; Kohonen 1990) demonstrated by Brett, West & Wheatley (2004) in an astronomical context. Here we adopt a novel technique based on a combination of SOM and RF machine learning. SOMs can efficiently parametrize light curve shapes without resorting to specific light-curve features, and RFs are capable of placing objects into classes.

In this work, we apply these techniques to data from the K2 mission, the repurposed *Kepler* satellite (Borucki et al. 2010). K2 and its predecessor *Kepler* have left a lasting mark in studies of variable stars, showing that most δ Scuti and γ Dor stars show pulsations in both the p-mode and g-mode frequency regimes (Grigahcène

* E-mail: d.j.armstrong@warwick.ac.uk

et al. 2010). Many studies have been performed on *Kepler* variable stars (e.g. Blomme et al. 2010; Balona & Dziembowski 2011; Balona et al. 2011; Debosscher et al. 2011; Uytterhoeven et al. 2011; Tkachenko et al. 2013; Bradley et al. 2015), but few so far on K2. Balona et al. (2015) studied B star variability in *Kepler* and K2, and found that K2 data presented some new challenges from the original mission. Despite these, it has for example discovered the several RR Lyrae stars known outside our own Galaxy (Molnár et al. 2015). LaCourse et al. (2015) have also produced a catalogue of eclipsing binary stars in K2 field 0.

The initial version of this catalogue (Armstrong et al. 2015) classified several thousand K2 variable stars in K2 fields 0 and 1. This classification was based on an interpretation of light-curve periodicity, and split objects into periodic, quasi-periodic, and aperiodic variables. Here we improve on this initial work, by applying an automated technique to classify variables into more usual classes. We extend the classification to K2 fields 0–4, and will release updates as more K2 fields become available.

2 DATA

2.1 Source

Data are taken from the K2 satellite (Howell et al. 2014). K2 is the repurposed *Kepler* mission, and provides light-curve flux measurements at a 30 min ‘long’ cadence continuously for 80 d per target. Targets are organized into campaigns, with each campaign spanning an ~ 80 d period and covering several thousand objects. A much smaller number of targets (a few tens per campaign) are available at the ‘short’ cadence of ~ 1 min. For the purposes of this work, we restrict ourselves to long cadence data only, to preserve uniformity in the data. At the time of writing, five campaigns had been released to the public (covering fields 0–4), with more due as the mission continues. Four of these campaigns cover ~ 80 d, with the first campaign 0 covering ~ 40 d. We take data for these campaigns from the Michulski Archive for Space Telescopes (MAST) website,¹ limiting ourselves to objects classified as stars in the MAST catalogue. This cut primarily removes a small number of Solar system bodies and extended sources from the analysis. At this point, we have 68 910 object light curves.

For the purposes of training the classifier, we also use data from the original *Kepler* mission. In these cases a single quarter of long cadence *Kepler* data is randomly selected. This covers ~ 90 d, and hence is similar to a single K2 campaign in duration and cadence. *Kepler* does however have different noise properties than K2, particularly in regards to the ~ 6 h thruster firing, which is present in K2 but not in *Kepler*. *Kepler* data were also downloaded from MAST, and the Presearch Data Conditioning (PDC) detrended light curves (Smith et al. 2012; Stumpe et al. 2012) used.

2.2 Extraction and detrending

K2 data show instrumental artefacts not previously seen in the original *Kepler* mission. The strongest of these is a signal at ~ 6 h, which is the time-scale on which the satellite thrusters are fired to adjust the spacecraft pointing. This pointing adjustment is necessary due to drift associated with the new mode of operations, and is explained fully in the K2 mission papers. It has the unfortunate effect of causing systematic noise, due to aperture losses and

Table 1. Times of pointing characteristic change, used to split the K2 data before detrending.

Campaign	Split time BJD 2454833
0	N/A
1	2016.0
2	2101.41
3	N/A
4	2273.0

inter-pixel sensitivity changes. A number of techniques have been put forward for removing this noise (Vanderburg & Johnson 2014; Aigrain et al. 2015; Lund et al. 2015), including one in the previous version of this catalogue (Armstrong et al. 2015). Each has advantages and disadvantages; our experience has been that while overall most techniques perform comparably, for individual objects the differences can be large. We use an updated version of our own extraction and detrending method here, which is fully described in Armstrong et al. (2015). The only change from that publication is the performing of a polynomial fit to the light curve, prior to detrending. This fit is performed by considering successive 0.3 d long regions of the light curve, and fitting third degree polynomials to 4 d regions centred on these. Outlier points more than 10σ from the initial fit are masked, and the fit redone without these points. The 10σ masking and refitting is repeated for 10 iterations. Masked points are not cut from the final light curves. The final fit is removed, detrending is performed, and the fit then added back in. This step was added to improve preservation of variability signals, a notable improvement on the first method. Light curves detrended using this method are publicly available at the MAST website.

It is important to note that, as described in Armstrong et al. (2015), our detrending method works best when performed separately on each half of the light curve (the exact split can be a few days from the precise halfway time). This is due to a change in the pointing characteristics of the spacecraft near the middle of each campaign, possibly the result of a change in orientation to the Sun. The precise times used to split the data are given in Table 1. Before conducting the analysis presented later in this work, we normalize each light curve half by performing a linear fit.

With the release of campaign 3, the K2 mission team began to release its own detrended light curves (these are not available for earlier campaigns at the time of writing). Similarly to the other methods, we find that these perform well overall but are by no means the best choice for every object. We will apply the classifier to both our light curves (hereafter the ‘Warwick’ set) and the K2 team light curves (hereafter the ‘PDC’ set) for campaigns 3 and 4. The comparison is complicated by the fact that the above-mentioned change in pointing characteristics does not occur in the usual way for these campaigns. Rather than change once in the middle of the campaign, in campaign 3 the change occurs twice, at roughly one third intervals. We do not adjust our detrending method for this, as introducing the option for another split adds an additional layer of complexity, and reduces the number of points available in each section (a risky option, as these points form the base surface used to decorrelate flux from pointing). Instead, we perform the detrending with no split at all. For campaign 4, we split at time 2273 (BJD 2454833, as given in the K2 data files), and cut points up to the first change in pointing at 2240.5. This shortens each campaign 4 light curve by 11 d, but results in improved detrending. We do

¹ <https://archive.stsci.edu/k2/>

not perform such an adjustment for campaign 3 as even more data would need to be cut.

3 CLASSIFICATION

3.1 Methodology

We employ a classification scheme using two distinct components. These are SOMs, otherwise known as Kohonen maps, and an RF classifier. Each is described below.

3.1.1 SOM

SOMs have been tested in an astrophysical context before (Brett et al. 2004; Tornaiainen et al. 2008; Carrasco Kind & Brunner 2014), but are rarely to date applied in astronomy in practice. As such we outline their methodology here.

A SOM is a form of dimensionality reduction; data consisting of multiple pieces of information can be condensed into a pre-defined number of dimensions, and is grouped together according to similarity. In our case, the SOM takes phase-folded light curve shapes and groups similar shapes into clusters, in one or two dimensions. The great strength of an SOM is in the unsupervised nature of its clustering algorithm. The user need not specify what groups or labels to look for; any set of similar input data, including for example previously unseen variability classes, will form a cluster in the resulting map. Similar clusters will lie near each other, those that are the same according to the input data will overlap. Furthermore, the input parameters for the algorithm are quite insensitive to small variations, making the clustering process robust (Brett et al. 2004).

The key component of an SOM is the Kohonen layer. This can be N -dimensional, but we will consider 2D layers here for clarity. The layer consists of pixels, each of which represents a template against which the input data is compared. The size of the layer is unimportant as long as it is sufficiently large to express the variation present in the input data. Once trained on a set of data, the Kohonen layer becomes a set of templates, representing the observed data features that it was trained on. These templates can be examined to spot interesting features in the data set, such as variation within an already known class. Further data (or the original data itself) can be compared to the trained layer and the closest matching template found. In this way, an object is placed on to the map.

The specific implementation of SOMs used here is described in Section 3.3, with an example of their use shown. The result is a map against which any input K2 phase-folded light curve can be compared. The location of the light curve on the map gives us its similarity to certain shapes, such as the distinctive light curve of an eclipsing binary star.

3.1.2 Random forest

The SOM allows us to classify and study the *shape* of a given phase curve, and the sets of similar shapes found within a data set. It does not place an object into a specific variability class. For that we utilize an RF classifier (Breiman 2001). These have been used in a number of previous variable star studies cited above. To use an RF classifier the light curve must be broken down into specific features, which represent the data (see Section 3.4 for those used here). These features are then paired with known classes in a training set of known variables, and the classifier fit to this set. For a given object, the RF classifier can then map sets of features

to probabilities for class membership, giving the likelihood for an unclassified object to be in each class.

RF classifiers are ensemble methods, in that they give results based on a large sample of simple estimators, in this case decision trees. In this way they can reduce bias in estimation. The core components of an RF are these decision trees. See Richards et al. (2011a) for a concise discussion of the underlying trees and how they are constructed. The specific parameters and implementation used here are discussed in Section 3.7.

3.2 Automated period finding

Our classification methodology relies heavily on the phase-folded light curves of our targets. This requires knowledge of the target's dominant period. Such knowledge is available for some known variables, but not for the general K2 sample at the time of writing. As such we use the K2 photometry to determine frequencies for each target.

There are a number of methods popularly used for determining light-curve frequencies. The most common is the Lomb–Scargle (LS) periodogram (Lomb 1976; Scargle 1982), which performs a fit of sinusoids at a series of test frequencies. Other available methods include the autocorrelation function (ACF; see e.g. McQuillan, Mazeh & Aigrain 2014) and wavelet analyses (Torrence & Compo 1998). We use LS here, due to its provenance and simplicity of implementation. The same arguments can be made for the ACF, which for stellar rotation periods has been shown to be more resilient than LS at detecting dominant frequencies (McQuillan, Aigrain & Mazeh 2013). However we find removing unwanted power from frequencies and harmonics, and detecting multiple frequencies from the same light curve, to be simpler for the LS method, at least in the implementations that we had available. In future, utilizing the ACF alone or in combination with the LS may be possible.

We use the fast LS method of Press & Rybicki (1989), with an oversampling factor of 20 run up to our Nyquist frequency of 24.5 d^{-1} . To avoid excessive human interference (and maintain the ‘automated’ status of this classification), the dominant frequencies for a target must be found without supervision. To avoid frequencies commonly associated with thruster firing noise in K2 (see Section 2.2), we remove frequencies within 5 per cent of 4.0850 d^{-1} and their 1/2, first, second, third, and fourth harmonics from the periodogram, by removing the best-fitting sinusoid of form

$$z = a \sin(2\pi ft) + b \cos(2\pi ft) + c \quad (1)$$

at each of these frequencies. In this model f represents the frequency being removed, t and z the time and flux data, and a , b and c free parameters of the model. We then cut these frequencies altogether before extracting the dominant period. We also remove frequencies associated with the data cadence which commonly show power in our periodogram ($48.943\,558\,19$ and $20.394\,709 \text{ d}^{-1}$) and their 1/2 frequency harmonics, by similarly fitting and removing a sinusoid at these frequencies and then cutting the frequencies from the periodogram. We did not find it necessary to remove other harmonics of the data cadence frequencies, as doing so provided little improvement. Finally, periods above 20 d (10 d in campaign 0) are cut, as the data baseline is not long enough to reliably determine them without the introduction of spurious noise related frequencies. At this point the most significant peak in the LS periodogram is taken.

To extract other significant frequencies, we remove the dominant frequency using a fit of the model of equation (1), then recalculate the LS periodogram, again ignoring thruster firing and cadence related frequencies as above. The remaining most significant peak

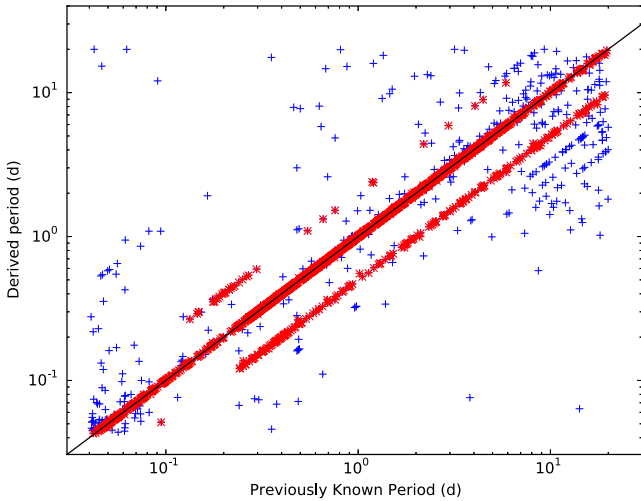


Figure 1. Periods determined using our method compared to previously known period, for a set of known variables in *Kepler*. An acceptance rate of 70.3 per cent is obtained, 82.2 per cent if half and double periods are included, and 90.8 per cent if second and third detected frequencies are included. Variables lying at the correct, half or double frequency are plotted as red stars.

is taken. To compare the power of different peaks, we calculate their amplitude A using $A = (a^2 + b^2)^{1/2}$. This is used to produce the frequency amplitude ratios used later in this work. We repeat this process to extract a total of three frequencies from each light curve.

A common weakness in period-finding algorithms occurs for eclipsing binary stars, a significant variability class. The LS periodogram often gives its highest power for half the true binary orbital period (i.e. when the primary and secondary eclipses occur at the same phase). This error is simple to spot by inspection, but harder to correct automatically. We account for this potential error source by introducing a check into the automated period finder. This phase folds each light curve at double its LS-determined dominant period. The phase-folded light curve is then binned into 64 bins, and the bin values at the minimum bin and the lowest bin value between phases 0.45 and 0.55 from this minimum found. We perform two checks on these two bin values. If the initial period is correct, they should be the same. We first check for an absolute difference between the two, finding that 0.0025 in relative flux works well as a threshold. We further test that the difference between them is greater than 3 per cent of the range of the un-phase-folded light curve. We calculate this range by taking the difference between the median of the largest 30 and median of the lowest 30 flux points in the light curve, to avoid unwanted outlier effects. If the difference between the two tested bin values is greater than both of these thresholds, the object period is doubled. If the doubled period would be over the 20 d upper period limit already applied (10 d for campaign 0), the doubling is not allowed. Similar adjustments have been made in previous variability studies (e.g. Richards et al. 2012). Only the dominant extracted period may be adjusted in this way.

To test the efficacy of our automated period finding software, we trial it against a known sample of variable stars from the *Kepler* data. See Section 3.6 for a full description of this set, which is also used as a training set for our classifier. We use one randomly selected quarter of *Kepler* data, to give data with a similar baseline and cadence to a single K2 campaign. There are 2128 training objects with previously determined periods (after removing objects with periods below our Nyquist period of 0.0408 d). Fig. 1 shows the comparison between

our dominant determined periods and the previously known ones. The acceptance rate is 70.3 per cent, rising to 82.2 per cent if half and double periods are included. In 90.8 per cent of the sample, one of our three determined periods finds either the previously known period or its half or double harmonic. In the remaining light curves, we find that either the noise obscures the known period (due possibly to different quarters with differing noise properties being used by us and previous studies) or that the dominant period has changed.

3.2.1 Phase curve template preparation

The SOM element of our classifier requires phased light-curve shapes to function. We create these using the periods determined in Section 3.2. Each light curve is phase-folded on this period. For known training set objects (see Section 3.6), the literature period is used. Once phase-folded, the light curve is binned into 64 equal width bins, and the mean of each bin used to form the phase curve that will be passed to the classifier. The exact number of bins is unimportant, as long as it gives enough resolution to see any variability in the phase curve. Brett et al. (2004) used 32 bins and found satisfactory results, we use 64 as the performance decrease is small and it reduces the chances of missing rapidly changing variability such as eclipses.

It is essential that the phase curves be on the same scale and aligned, so that the classifier can spot similarities between them (see next section). As such we normalize each phase curve to span between 0 and 1, and shift it so that the minimum bin is at phase 0. Each phase curve then consists of 64 elements, with the first being at (0,0).

3.3 Training the SOM

There are variations in the literature on how precisely to train the SOM. Here we run through the procedure followed for this work. The input parameters are the initial learning rate, α_0 , which influences the rate at which pixels in the Kohonen layer are adjusted, and the initial learning radius, σ_0 , which affects the size of groups. Initially, each pixel is randomized so that each of its 64 elements lies between 0 and 1, as our phase curves have been scaled to this range. For each of a series of iterations, each input phase curve is compared to the Kohonen layer. The best matching pixel in the layer is found, via minimizing the difference between the pixel elements and the phase curve. Each element in each pixel in the layer is then updated according to the expression

$$m_{xy,k,\text{new}} = \alpha e^{\frac{-d_{xy}^2}{2\sigma^2}} (s_k - m_{xy,k,\text{old}}), \quad (2)$$

where $m_{xy,k}$ is the value m of the pixel at coordinates x,y and element k in the phase curve, d_{xy} is the Euclidean distance of that pixel from the best matching pixel in the layer, and s_k is the k th element of the considered input phase curve. This expression is specific to two-dimensional SOMs, but can be easily adapted for one-dimension by setting the size of the second dimension to be 1. Note that distances are continued across the Kohonen layer boundaries, i.e. they are periodic. Once this has been performed for each phase curve, α and σ are updated according to

$$\sigma = \sigma_0 e^{\left(\frac{-i \log(r)}{n_{\text{iter}}}\right)} \quad (3)$$

$$\alpha = \alpha_0 \left(1 - \frac{i}{n_{\text{iter}}}\right), \quad (4)$$

where i is the current iteration, and r is the size of the largest dimension of the Kohonen layer. This is then repeated for n_{iter} iterations.

It is possible to use different functional forms for the evolution of α and σ ; typically a linear or exponential decay is used. Brett et al. (2004) found that the performance of the SOM was largely unimpeded by the choice of form or initial value, as long as the learning rate does not drop too quickly. We find satisfactory results for the expressions above and values of $\alpha_0 = 0.1$ and $\sigma_0 = r$, as can be seen in the below example. The code used in this study was initially adapted from the SOM module of the open source PYMVPA package² (Hanke et al. 2009), and has now been contributed as an update to that package by the authors. As such any readers wishing to use this code should look to the given reference. Note that the functional form of equations (3) and (4) are slightly different in the online version of the code, to preserve compatibility with older versions of the module. The formulae described here are the ones used in this work.

As an example we train an SOM on the K2 data from campaigns 0–2, as well as *Kepler* data used for training the classifier (see Section 3.6 for a full description of the data set). We use a 40×40 Kohonen Layer. K2 data were only used if the range of variation in the phase curve before normalization was greater than 1.5 times the overall mean of the standard deviations of points falling in each phase bin (see previous section). This cut was imposed to avoid essentially flat light curves from impacting the SOM, removing ~ 40 per cent of the K2 light curves. The majority of these were classified as ‘Noise’ or ‘AP’ in Armstrong et al. (2015), showing that we are not removing many periodically varying sources. We note that the SOM is robust enough to work without this cut, and it is imposed only to increase the purity of the training set.

We take the known *Kepler* variables, along with ‘OTHPER’ other periodic and quasi-periodic objects from K2, and plot them on the resulting SOM in Fig. 2. Clear groups can be seen, with eclipsing binary types well differentiated but bordering each other, as would be expected. RR Lyraes are very well grouped, and δ Scuti variables cluster but more weakly. Example templates from the Kohonen layer are shown in Fig. 3, representing the major clusters seen. Note that the size of a group is determined by a number of factors, including the number of input objects matching it, and the extent of small variations within the group. As there are many more sinusoidal variables than eclipsing binaries or RR Lyraes, the δ Scuti, γ Doradus and ‘OTHPER’ groups fill most of the map. Different regions within these groups show for example slight skews from a pure sinusoid, and may represent interesting intraclass differences. δ Scutis lying near the eclipsing binary groups have likely been mapped using double their true period, and so look similar to a contact binary star. They may also have been previously misclassified. It is also interesting to see that δ Scutis and ‘OTHPER’ objects overlap, as would be expected given that their phase curve shapes are not particularly distinctive to their respective classes. ‘OTHPER’ objects also overlap with the RR Lyrae cluster, and likely mark out newly discovered RR Lyrae stars.

The SOM used for final classification is the same as that described above, but using only one dimension of 1600 pixels. This produces the same clustering results, but is less useful for visualization. We use only one dimension so that the other part of our classifier (the RF) can more easily make use of the information contained within the SOM.

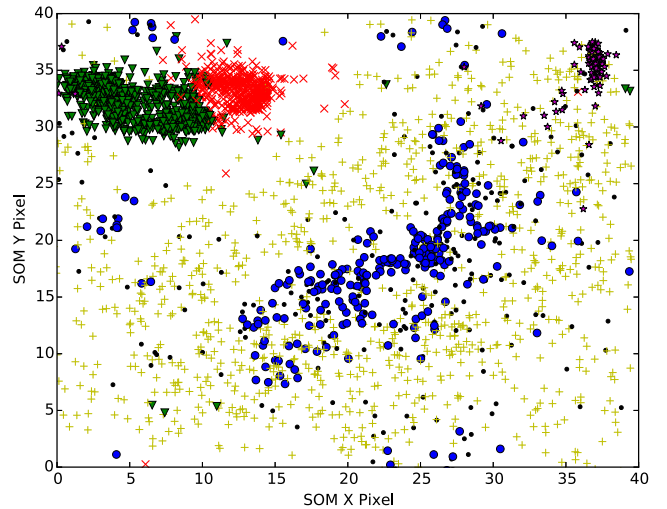


Figure 2. Known variables placed on to an SOM. Random jitter within each pixel has been added for clarity. Green triangles = ‘EA’ (detached eclipsing binaries), red crosses = ‘EB’ (semidetached and contact eclipsing binaries), pink stars = ‘RRab’ (ab-type fundamental mode RR Lyraes), blue circles = ‘DSCUT’ (δ Scuti variables), black dots = ‘GDOR’ (γ Dor variables) and yellow pluses = ‘OTHPER’ (other periodic and quasi-periodic objects). See Section 3.5 for more detail on these variability classes.

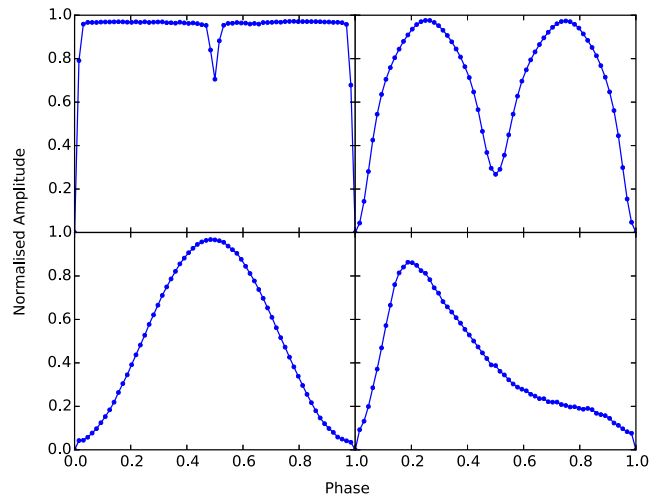


Figure 3. Template phase curves from the Kohonen layer of the SOM in Fig. 2. Clockwise from top left, templates are for pixel [13,34] (EA), [6,32] (EB), [37,35] (RRab) and [25,19] (DSCUT). See Section 3.5 for a description of the classes. Note that templates do not have to span the range 0–1, even if the input phase curves do. Note also that all these templates were found from initially random pixels without any human guidance or input.

3.4 Data features

For the classification of variables into classes, we use a number of specific features of each light curve. This is common practice in general classification problems (e.g. Richards et al. 2011a). However, there is a subjective element to selecting features, and it can be desirable to minimize this if possible (see e.g. Kügler et al. 2015). We do so through the use of the SOM. This encodes the shape of the phase curve into one parameter (the location of the closest pixel in the SOM to the light curve in question), rather than a series of features, none of which may capture the desired shape properties.

² <http://www.pymvpa.org>

Table 2. Data features.

Feature name	Description
Period	Most significant period (Section 3.2).
Amplitude	Max–min of phase curve.
period_2	Second detected period (Section 3.2).
period_3	Third detected period (Section 3.2).
ampratio_21	period_2 to period amplitude ratio.
ampratio_31	period_3 to period amplitude ratio.
SOM_index	Index of closest pixel in 1D SOM.
SOM_distance	Euclidean distance to closest pixel in 1D SOM.
p2p_98perc ^a	98th percentile of point to point scatter in light curve.
p2p_mean ^a	Mean of point to point scatter in light curve.
phase_p2p_max	Maximum point to point scatter in binned phase curve.
phase_p2p_mean	Mean of point to point scatter in binned phase curve.
std_ov_err ^a	Whole light curve standard deviation over mean point error.

Note. ^aadjusted between data sets, see text.

There are however other features which are useful and which are uninformed by the SOM. A key example is the dominant (most significant) period of the light curve. Other significant frequencies can also be used, and in some cases many more have been studied. We only use the three most significant periods here.

The full range of features used is described in Table 2. These features are incorporated largely to separate out light curves which show purely noise, something which is generally uninformed by the SOM, as well as those without one particularly dominant frequency. We take the potentially controversial step of adjusting some of the noise related features between the *Kepler* and K2 data sets, due to the differing noise properties between each set. This is unavoidable here, as the scatter and increased noise in K2 causes catastrophic errors in the classifier if *Kepler* light curves are used as they come. In this case the general result is that the vast majority of K2 objects are classified as Noise. This problem is solved by multiplying the marked features in Table 2 by a factor to align their median values with those of K2. These features are those driven primarily by data set noise, rather than those associated with periodicity (noise-related periodicity is assumed to have been removed by the procedure in Section 3.2). As the *Kepler* data used all comes from known variable stars, the median of the features is not strictly comparable to K2, where the data comes from the whole target list. As such we set the multiplication factor so that the median of the non-eclipsing binary *Kepler* data features is increased to equal the median of the ‘OTHPER’ K2 data features. Eclipsing binaries are left alone, as their features are in our case dominated by the binary eclipses.

A similar problem arises when studying the PDC light curves. These have different characteristics to the Warwick light curves. Assuming that the intrinsic distribution of stellar variability should be the same across fields, this difference is due to the differing detrending methods. We adjust for it in the same way and to the same features as above, marked in Table 2. As we do not have prior classifications for fields 3 and 4, the factor is applied to the whole data set, and set so as to match the medians of these features between the PDC campaigns 3 and 4 and the Warwick campaigns 0–2. Each PDC campaign is adjusted separately.

It would be desirable to use colour information as a feature to aid classification of variability types connected to specific stellar spectral types. However, colours are not uniformly available for

the K2 sample, although some can be found through a cross-match with the TESS input catalogue (Stassun et al. 2014). As such we do not use them, as doing so would mean large fractions of the K2 targets would need to be disregarded. This has consequences for the variability classes we use, see Section 3.5.

3.5 Classification scheme

An important decision is in which variability classes to use. We experimented with classifying RR Lyrae (subtype ab), δ Scuti, eclipsing binary (split into detached, subtype EA, and semidetached or contact, subtype EB), γ Dor, and so-called ROT variables, a class applying to likely rotationally modulated light curves seen in Bradley et al. (2015). We also attempted to split the γ Dor class into symmetric, asymmetric, and ‘MULT’ classes, as defined in Balona et al. (2011). This approach had varied success; RR Lyrae ab, δ Scuti, γ Dor and eclipsing binary classes performed well, but we found that the γ Dor subtypes were not well constrained by our available features. This may be because we lack sufficient training objects to reliably map the range of features offered by these subtypes. This problem could be navigable when an increased sample of objects is available through K2, and we plan to address this in later work.

Similarly, we found that the ‘ROT’ class was not very coherent – the classifier struggled to identify regions in parameter space corresponding to these variables. This likely arises due to the tendency of this class to have an indistinct cluster of low-frequency peaks rather than one clear signal (Bradley et al. 2015). Rather than use the ROT class by itself, we make use of the previous version of this catalogue, which contained a ‘QP’ quasi-periodic variable class. This class contains a number of variable types, but is characterized by periodic variability that is not strictly sinusoidal, and changes in amplitude and/or period. We use this as a variable classification, to catch interesting variables of astrophysical origin which are not one of the five other classes (RR Lyrae ab, EA, EB, δ Scuti, γ Dor). It is likely dominated by spot-modulated stars, but also contains other variables such as Cepheids. We rename this class to ‘OTHPER’ for ‘other periodic’ to avoid confusion, as variables which are strictly periodic but not in another class can be classified by this group.

We considered including other variable classes, such as Cepheids, the other RR Lyrae subtypes (first-overtone or multimode RR Lyraes), and Mira variables. We could not find sufficient training set objects in any of these classes (less than 20 in each case). While it is possible to attempt classification with small training sets, rather than present a weak or unreliable classification for these classes we prefer to wait for more K2 data. As more fields are observed, more training set objects will become available. We intend to include more classes in future versions of this catalogue.

Finally, we include ‘Noise’, non-variable light curves, as a class label. This leaves seven classes, DSCUT (δ Scuti), GDOR (γ Doradus), EA (detached eclipsing binaries), EB (semidetached and contact eclipsing binaries), OTHPER (other periodic and quasi-periodic variables), RRab (RR Lyrae ab type) and Noise. It is important to note that as we do not have colour information, there will be degeneracy in the DSCUT class between true δ Scutis and β Cep variables, as in Debosscher et al. (2011). This is also true for slowly pulsating B stars, which are degenerate with γ Dor variables.

3.6 Training set

Although the SOM described is unsupervised and so requires no training set, the RF classifier we use for final classification does. An ideal training set would consist of a set of known variable stars from

the K2 mission, to which we can fit the classifier. Some previous classification work on K2 has been done (for B stars Balona et al. 2015, for eclipsing binaries LaCourse et al. 2015, and in the previous version of this catalogue). These sources however suffer from either small numbers, only being applicable to a few variable types, or in the Armstrong et al. (2015) case using variability classes derived from the light curves rather than externally recognized types. We cross-matched the observed K2 targets in fields 0–3 (4 was not available at that time) with catalogues of known variable stars, including those from AAVSO,³ GCVS (Samus et al. 2009) and ASAS (Richards et al. 2012). This led to a small number of targets (a few tens of each class at best), not enough for a full training set. As such, we turned to the original *Kepler* mission. Much classification work has been done on the *Kepler* light curves. The data have differing noise properties to K2 data, but the same cadence, instrument, and if only one 90 d quarter of data is used a similar baseline to a K2 campaign.

Although multiple works are available offering classified variable stars in *Kepler*, we limit ourselves to a small number of relatively large-scale catalogues, in order to maintain homogeneity among classification methods and simplify the process. We began by taking the EA, EB, DSCUT classes from Bradley et al. (2015). We also took ROT, SPOTM and SPOTV, low-frequency variables likely due to rotational modulation, reclassifying these objects as OTHPER. We supplemented the DSCUT set with those from Uytterhoeven et al. (2011). The bulk of our eclipsing binary training set come from the Kepler Eclipsing Binary Catalogue (Prsa et al. 2011; Slawson et al. 2011). We removed all heartbeat binaries (Thompson et al. 2012) and those where the primary eclipse depth was less than 1 per cent. A threshold of 1 per cent was implemented in order to avoid shallow, likely blended binary eclipses from being included in the training set and hence increase training set purity. This also avoids the problem of noisy light curves with instrumental systematics of the order of a per cent being misclassified as eclipsing binaries. Binaries were then classified as EA or EB based on a morphology threshold of 0.5 (see Matijević et al. 2012 for a discussion of morphology in this context). For RR Lyrae stars we use the list in Nemec et al. (2013). Fundamental mode subtype ab stars were labelled RRab, and the first-overtone subtype c stars classified as OTHPER. To increase this relatively small RR Lyrae sample we used the results from the K2 AAVSO cross-match, taking fundamental mode RR Lyraes and adding them to the RRab training set. The B-star catalogue of Balona et al. (2015) was also used, with the SPB class reclassified as GDOR (given the degeneracy between GDOR and SPB present without temperature information) and the ROT class being reclassified as OTHPER.

For the OTHPER and Noise classes, we also use our previous catalogue. This contained 5445 OTHPER (QP in the original catalogue) and 29 228 Noise objects in fields 0–2, with labels assigned by human eyeballing. To avoid having an excessive disparity between training set classes, we downsample this set to 1000 of each class, selected randomly, which are then added on to the *Kepler* OTHPER set above. This also makes the results on fields 0–2 more independent, as we can compare previously classified OTHPERs (the majority of which are now not in the training sample) with newly found ones. To reduce the impact of potential mistakes in the previous catalogue, we removed the small number of objects in the OTHPER training set which were in an initial run of this classifier reclassified as another class. Objects with a probability of being

Table 3. Training Set.

Class	N objects
RRab	91
DSCUT	278
GDOR	233
EA	694
EB	759
OTHPER	1992
Noise	976

in the RRab class of greater than 0.2 were also removed, as the probabilities for the RRab class are not well calibrated (see Section 3.8). These cuts caught ~50 objects misclassified as OTHPER and ~30 objects misclassified as Noise out of the 1000 each initially selected.

The final classes and number of objects in each training set are shown in Table 3.

3.7 RF implementation

We use the implementation of RFs in the `scikit-learn` PYTHON module.⁴ There are several input parameters for an RF classifier. The key ones are the number of estimators, the maximum features considered at each branch in the component decision trees, and the minimum number of samples required to split a node on the tree, which controls how far each tree is extended. In a typical case, increasing the number of estimators always leads to improvement in performance but with decreasing returns and increasing computation time. The theoretical optimum maximum features for a classification problem is the square root of the total number of features, in our case 3. We optimize the parameters using the ‘out-of-bag’ score of the RF. When training, the classifier uses a random subset of the total data sample given to it for each tree, to reduce the chance of bias. The left out data are then used to test the performance of the tree – its known class is compared to the predicted class, giving a performance metric between 0 (for absolute failure) and 1 (for perfect classification). Maximizing this metric allows us to optimize the parameters. We find the best results for 300 estimators, a maximum of three features, and five samples to split a node. These parameters are used for classification. Additionally we apply weights to the training set, so that each class is inversely weighted according to its frequency in the training set (input option `class_weight=‘auto’`). This makes sure that classes with more members (such as OTHPER and Noise) do not drown out other classes, and in effect imposes a uniform prior on the class probabilities.

There are several random elements in our method. These are the selection of the OTHPER and Noise training sets, as well as certain elements of the RF. Random subsets of training objects and features are selected for each decision tree as part of the RF method, to avoid bias. To minimize any effects of this randomness (especially the OTHPER and Noise selection), we train 50 classifiers with the above parameters and repeat the selection for each, applying each classifier to the K2 data set. The average class probability across the classifiers gives the final result.

To explore the power of the SOM method, we trial the RF on only the SOM map location (SOM_index). The classifier is cross-validated by taking one training set member and training the classifier on the remaining members (so-called leave-one-out cross-validation). The left out object is then tested on the classifier, and

³ www.aavso.org⁴ <http://scikit-learn.org/stable/>

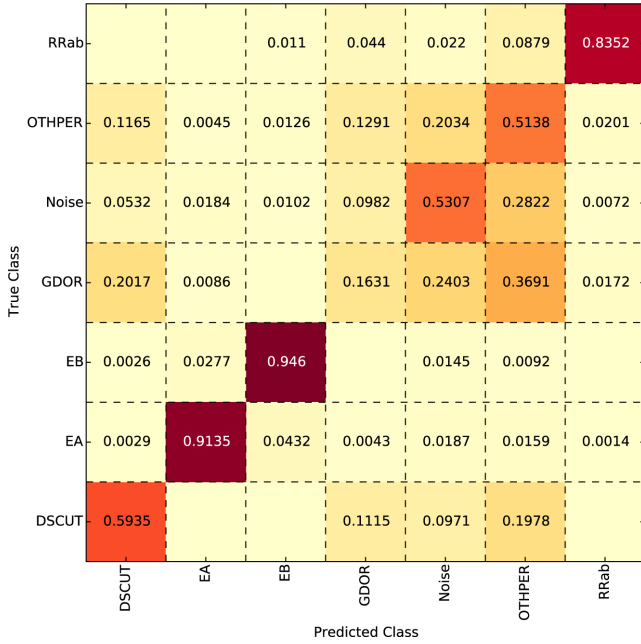


Figure 4. Confusion matrix for an RF considering only SOM map location, generated using leave-one-out cross-validation. Text shows the percentage of each sample which was classified into the relevant box. Correct classification lies on the diagonal.

the process repeated for each member. The performance of the classifier is best described by a ‘confusion matrix’, shown in Fig. 4. This shows what proportion of training members in each class were assigned to which other classes. In the ideal case each object is predicted correctly. Here we can see clearly which classes are well informed by the SOM. RRab, EA and EB classes are strongly recovered, as expected from their strong localization in Fig. 2. The DSCUT class is also recovered although less so. On the other side, OTHPER and Noise classes are found more weakly, and GDOR barely at all, due to the often multiple pulsation frequencies in this class combining to produce no distinctive phase curve shape. This demonstrates the power of the SOM alone to classify certain classes of variable stars.

Moving on to the full classification scheme, we test the RF in a similar manner. All seven classes are used, and the classifier cross-validated as before. The resulting confusion matrix is shown in Fig. 5. It highlights some interesting cases. First, the classifier works well, with an overall success rate of 92.0 per cent. There is some porosity between the two eclipsing binary classes, with objects of one class being placed into the other. As there is no rigid boundary in light-curve shape between them, this is to be expected. Similarly there is some spread between OTHPER and Noise. This is not desirable, but the numbers involved are low, and represent objects with either variability only just emerging above the noise or objects with unusual noise properties. The biggest misclassification occurs between the GDOR and OTHPER classes. This arises due to the less distinct nature of the OTHPER class – it acts as a ‘catch-all’ class to find any periodic or quasi-periodic variables which do not fit the other classes. GDOR objects can in some circumstances present similar light-curve features to for example fast rotating stars, leading to some confusion between the classes.

One advantage of RF classifiers is the ability to estimate feature importance. The classifier naturally measures which features have more descriptive power, through for example how often those

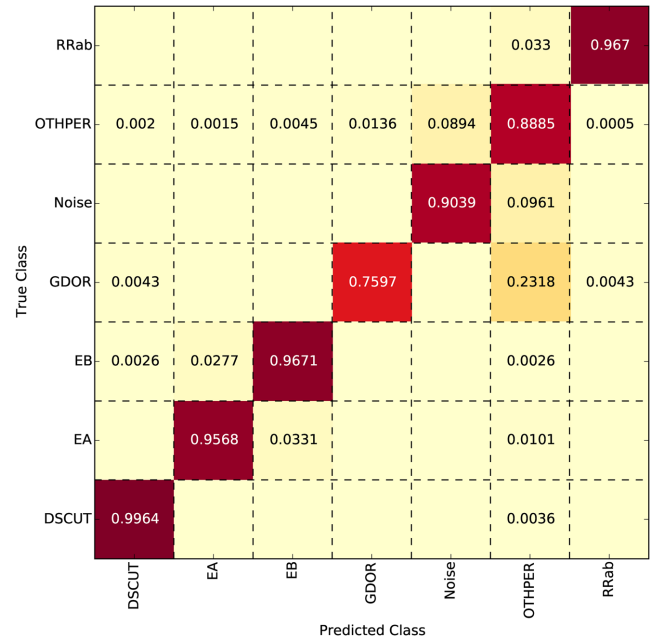


Figure 5. Confusion matrix for an RF considering all features and classes, generated using leave-one-out cross-validation. Text shows the percentage of each sample which was classified into the relevant box. Correct classification lies on the diagonal.

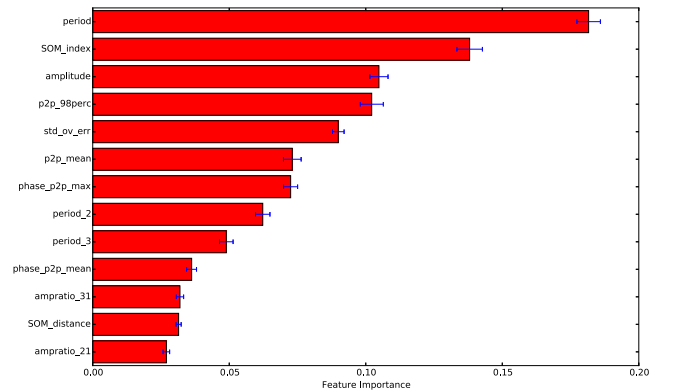


Figure 6. Relative importance of features to the RF. Values and errors arise from the mean and standard deviation of the feature importances extracted from 100 trained classifiers.

features are used in the decision trees, or through the reduction in performance that would be observed if a feature was replaced by a randomly sampled distribution. This allows for model refinement, and is of great use in developing a classifier. We plot the importance of our features in Fig. 6. These are found through training the classifier 100 times, and extracting the mean and standard deviation of the feature importances for each classifier.

3.8 Class posterior probability calibration

The RF classifier automatically generates class probabilities (through the proportion of estimators classifying an object into each class). These probabilities are not necessarily accurate. Although it is true that higher class probability means more likelihood of an object being in that class, the probabilities can need calibrating to ensure that they are true posterior probabilities. This is where, if a

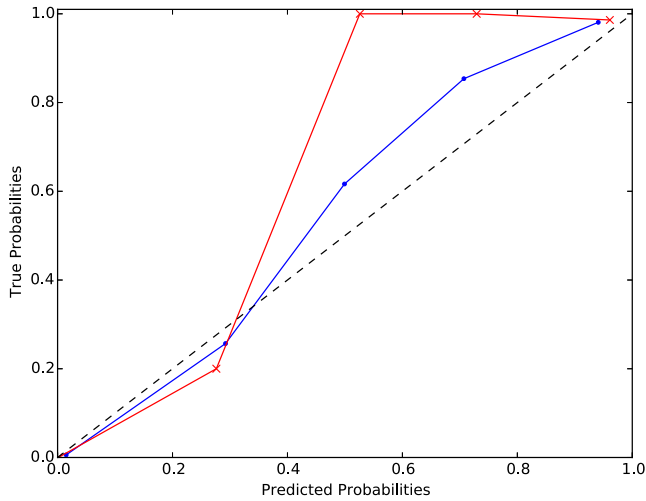


Figure 7. Overall classifier predicted probability against true probability for the RRAb class (crosses) and the average of all other classes (dots). The straight black dashed line represents the ideal case.

set of objects have probability p that they are in a certain class, the same proportion of them actually are of that class.

Initially we test the calibration of our ‘raw’ class probabilities. Fig. 7 shows the class probabilities found from the cross-validated training set data created as described in Section 3.7. This allows the predicted class probabilities for each training set object to be compared to their known classes. They are clearly not true posterior probabilities, especially for the RRAb class, where essentially every object with class probability >0.5 is a true class member. For the other classes the given probabilities are closer, but still show some departure from the ideal case.

One common way of testing classifier performance in this way is the Brier score (Brier 1950). Our raw probabilities have a Brier score of 0.1336. We attempted a number of methods of calibrating them (and so reducing this score). The most usual methods are sigmoid and isotonic regression, which fit certain functions to the calibration curve to transform the probabilities. Similarly to Richards et al. (2012), we find that these methods are not effective in our case. We attempted the method of Bostrom (2008) to transform the initial class probabilities, but also found the results to be unsatisfactory. Rather than present an incomplete calibration, we give the class

probabilities as they are. Users should be aware of this, and avoid interpreting class probabilities as true posterior probabilities.

As the training set will not be representative of the true K2 distribution, biases may exist. As the priors are not well known, and the distribution of training sources by no means matches the underlying distribution of variables in K2, true posterior probabilities are impossible to create. Hence the given class probabilities, even if calibrated, would only be posterior probabilities under the assumption that each class has a uniform probability of arising.

4 CATALOGUE

4.1 Overview

The full catalogue for K2 fields 0–4 inclusive is given in Table 4. This table contains classifications using the Warwick light curves, as described in Section 2.2. The features used to classify these objects are given in Table 5. We also run the classifier on the PDC light curves produced by the K2 mission team. These were only available for campaigns 3–4. The resulting classifications are given in Table 6, and their associated features in Table 7.

The total number of objects found in each class is given in Table 8, at various probability cuts. Note that for RRAb class objects in particular, most objects with class probability >0.5 are real classifications. In the other cases the probability calibration is better, but these probabilities should still not be interpreted as posterior probabilities.

We find that the classifier works well on all fields. The RRAb class performs well throughout, due to the distinctive shape of their phase curves. These are well characterized by the SOM. There are however some distinct features unique to fields 3 and 4. The EA class has a tendency to pick up noise-dominated light curves in these fields, primarily because their point to point scatter is much higher than in fields 0–2. In these cases the class probability, although highest for EA, is still relatively low however. Similarly for DSCUT objects, there are a higher proportion of objects in these fields with many anomalous points, possibly due to flaring or instrumental noise. These points can cause biases in the phase curve, resulting in an artificial sinusoid, which when combined with a short period results in a DSCUT classification. Again these noise objects have a lower probability than real DSCUT light curves. One final interesting property is the split between OTHPER and Noise light curves. This

Table 4. Catalogue table for our Warwick detrended light curves. Fields 0–4 are included. Only an extract is shown here for guidance in form. The full table is available online.

K2 ID	Campaign	Class	Class probabilities							Anomaly
			DSCUT	EA	EB	GDOR	Noise	OTHPER	RRab	
202059070	0	Noise	0.004 195	0.120 507	0.016 615	0.005 925	0.604 636	0.246 088	0.002 034	0.023 891
...

Table 5. Data features for our Warwick detrended light curves. Fields 0–4 are included. Only an extract is shown here for guidance in form. The full table is available online.

K2 ID	Campaign	SOM_index	period (d)	period_2 (d)	period_3 (d)	SOM_distance	phase_p2p_mean rel. flux	phase_p2p_max rel. flux	amplitude rel. flux	ampratio_21	ampratio_31
202059070	0	1544	4.764 370	1.241 680	0.174 448	1.180 831	0.003 801	0.487 419	0.042 283	0.629 987	0.548 721
...
p2p_mean rel. flux	p2p_98perc rel. flux	std_ov_err									
0.016 326	0.047 548	1.310 764									
...

Table 6. Catalogue table for PDC detrended light curves. Fields 3–4 only. Only an extract is shown here. The full table is available online.

K2 ID	Campaign	Class	Class probabilities							Anomaly
			DSCUT	EA	EB	GDOR	Noise	OTHPER	RRab	
205889250	3	Noise	0.000 067	0.000 000	0.000 000	0.000 030	0.966 544	0.033 359	0.000 000	0.000 000
...

Table 7. Data features for PDC detrended light curves. Fields 3 and 4 only. Only an extract is shown here. The full table is available online.

K2 ID	Campaign	SOM_index	period (d)	period_2 (d)	period_3 (d)	SOM_distance	phase_p2p_mean rel. flux	phase_p2p_max rel. flux	amplitude rel. flux	ampratio_21	ampratio_31
205889250	3	0630	19.754 572	12.803 889	2.281 881	1.179 035	0.003 795	0.421 976	0.008 715	0.741 302	0.592 596
...
p2p_mean rel. flux	p2p_98perc rel. flux	std_ov_err									
0.005 249	0.017 133	1.371 857									
...

Table 8. Total objects in each class.

Class	Total	Prob >0.5	Prob >0.7	Prob >0.9
RRab	248	154	72	25
DSCUT	750	562	377	166
GDOR	451	264	133	37
EA	607	308	183	99
EB	463	392	290	186
OTHPER	22 428	18 698	9399	3547
Noise	43 963	38 609	21 210	6018

is good for fields 0–2. In fields 3 and 4, while OTHPER light curves are recognized, several Noise light curves can be classified as OTHPER. Probability cuts remove the worst of these, but there is no way to distinguish between quasi-periodic instrumental noise and astrophysical variability in this scheme. These issues all lead to the conclusion that the classifier has more trouble with fields 3 and 4, due to a pattern of increased noise. We expect this issue to improve as K2 detrending methods become more robust.

4.2 Detrending method comparison

Table 9 shows the numbers of variable stars found using each data set. At first glance the numbers in Table 9 seem to imply significant differences between detrending methods. The discrepancy in RRab numbers is largely a result of differing probability calibration – the same stars are found in both data sets, but those in the Warwick set given lower probabilities (although still higher than all other classes). Other major discrepancies are in the GDOR and EA classes. For GDOR, we find that the PDC set gives better results. Several GDOR light curves are misclassified in the Warwick set due to poor detrending masking the true variability. In some cases the PDC GDOR classification is inaccurate, but this is rare for the class probability >0.7 objects. For the EA objects, the reverse is true. Several PDC light curves are misclassified as EA due to a higher number of light curves in the PDC set with very significant remnant outliers. These lead to a high point-to-point scatter, which is interpreted by the classifier as an eclipse. Here the Warwick set is more reliable. The largest absolute difference in the variable classes is in the OTHPER objects, where ~1000 light curves extra pass the high probability cut for the PDC set. This is partly a result of a similar effect as for the RRab objects, where similarly classified objects are given lower probabilities in the Warwick set. However, there are also several objects found in the PDC set which are missed in

the Warwick set, due to increased noise levels. The converse is also true, with some light curves found in the Warwick set but missed by the PDC. Overall, the two detrending methods perform comparably well, and can be used to reinforce each other when studying variable classes.

4.3 Anomaly detection

Due to the limited classification scheme used, it is inevitable that some objects will not fit any of the given classes (Protopapas et al. 2006). Due to the inclusion of Noise and OTHPER as classes, this is not a large problem as each class is quite broad. However it is worth noting any particular anomalies. One way of doing this is already intrinsic to the SOM – the Euclidean distance of a phase curve to its nearest matching pixel template. However this metric only works for periodic sources, and can flag high for noisy sources. We perform a check for anomalies following the method of Richards et al. (2012). This works by extracting the proximity measure, ρ_{ij} between each tested object i and each object j in the training set. The proximity measure is the proportion of trees in the classifier for which each object ends at the same final classification. It is close to unity for similar objects, and close to zero for dissimilar ones. From the proximity the discrepancy d is calculated, via

$$d_{ij} = \frac{1 - \rho_{ij}}{\rho_{ij}}. \quad (5)$$

The anomaly score is then given by the second smallest discrepancy of an object to the training set. High anomaly scores represent objects which are not well explained by any object in the training set, and are hence outliers.

We find that in this case, the highest few percentiles of anomalous objects are a mixture of noise-dominated light curves, unusual eclipsing binaries and variability which does not fit into the used classification scheme. We leave a full analysis of these unusual light curves to future work.

4.4 Eclipsing binaries

Encouragingly, we identify 139 (96 at class probability >0.7) of the 165 EPIC, non-M35 eclipsing binaries identified by LaCourse et al. (2015) in field 0 as either ‘EA’ or ‘EB’ type, despite automating the process and not focusing on exclusively eclipsing binaries. The majority of the remainder are identified as ‘OTHPER’ or ‘DSCUT’, and are discussed below. We further identify an additional 61 EPIC,

Table 9. Total objects in each class in fields 3 and 4, split by detrending method (W=Warwick, PDC=K2 Team released light curves).

Class	Total W	Total PDC	Prob >0.5 W	Prob >0.5 PDC	Prob >0.7 W	Prob >0.7 PDC
RRab	141	152	95	115	48	83
DSCUT	280	266	180	201	116	148
GDOR	198	382	122	238	61	101
EA	255	413	97	223	54	102
EB	168	150	140	131	106	105
OTHPER	11 402	9102	8709	8034	3522	4565
Noise	17 143	19 126	13 012	17 919	3625	11 566

Table 10. EPIC IDs for 29 visually identified eclipsing binaries classified as ‘OTHPER’ by our classifier, from fields 0–4.

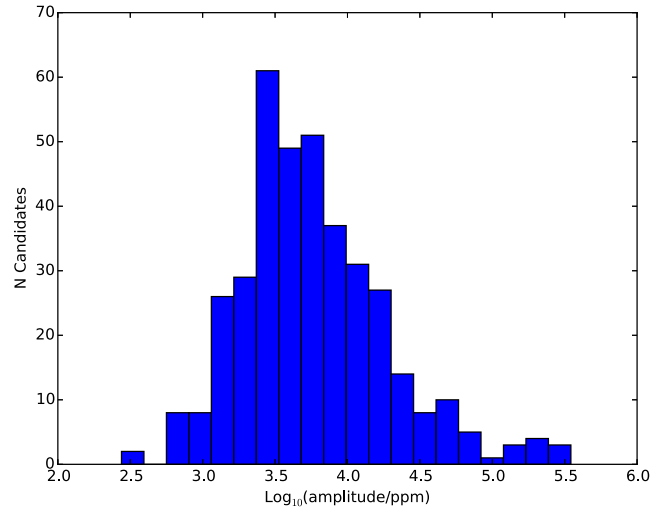
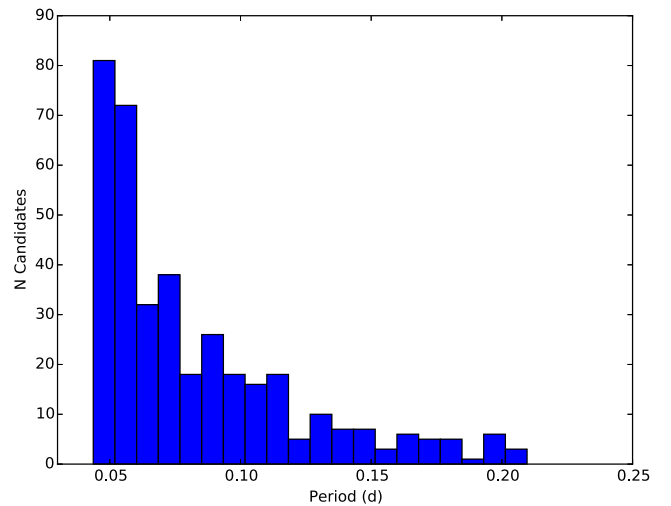
201158453	201173390	201569483	201584594
201638314	202072962	202137580	203371239
203476597	203637922	204043888	204193529
204328391	204411840	205510143	205919993
205985357	205990339	206047297	206060972
206066862	206226010	206311743	206500801
210350446	210501149	210766835	210945342
211093684	211135350		

non-M35 objects in field 0 as ‘EA’ or ‘EB’ at class probability >0.7 , although as our identification is automated rather than visual some of these may be misidentified by the classifier. Many more eclipsing binaries are found in the other fields.

The previously labelled, but not identified by our classifier, eclipsing binaries fall into three main groups. The first show near-sinusoidal short period light-curves, and are generally identified as ‘DSCUT’. In these cases, it is difficult to reliably assign a class with the information available. These objects may be actual δ Scuti stars, or contact eclipsing binaries. The other and largest group, with 14 members, are identified as ‘OTHPER’, and show pulsations or spot-modulation in addition to the known eclipses. We note that the classifier will assign a class based largely on the dominant period and phase curve at this period, hence performs as expected in these cases. Pulsating stars in eclipsing binaries are useful objects, and so while a detailed study of these objects is beyond the scope of this paper, we provide a list of such objects in Table 10. These are eclipsing binaries identified by a visual check of the light curves performed ourselves (as the LaCourse et al. 2015 catalogue only covered field 0), which are classified as ‘OTHPER’ by our classifier. Some may be blended signals, and hence the pulsator or spot-modulated star may not be a member of the eclipsing binary system.

4.5 δ Scuti stars

We have a sample of 377 δ Scuti candidates, using a class probability cut of 0.7. The majority of these candidates were previously unknown. It is interesting to study their frequency and amplitude distribution. Note that here we use amplitude defined as in the max-min of the binned phase curve, and semi-amplitude as half this value. The distribution of amplitudes for the 377 δ Scuti candidates is shown in Fig. 8. We see a number of HADS (high amplitude δ Scutis). Using an amplitude threshold of 10^4 ppm as used by Bradley et al. (2015), 104 of our candidates are HADS. Included in this sample are 11 candidates with an amplitude greater than 10^5 ppm. The period distribution of the whole sample is shown in Fig. 9, and covers the expected range for δ Scuti variables, limited by our Nyquist sampling frequency.

**Figure 8.** The distribution of phase curve amplitude for DSCUT classified objects. Several high-amplitude candidates are visible.**Figure 9.** The distribution of pulsation periods for DSCUT classified objects. The cutoff at the low-period end is imposed by our Nyquist sampling frequency.

As has been mentioned, the DSCUT classified objects are degenerate with β Ceph variables due to the lack of colour information available. There is a catalogue of estimated K2 temperatures available for some objects (Stassun et al. 2014) which could be used to make probable distinctions if necessary.

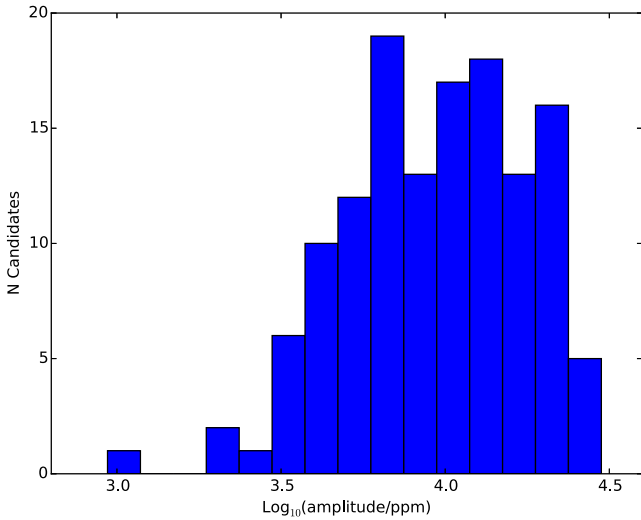


Figure 10. The distribution of phase curve amplitude for GDOR classified objects.

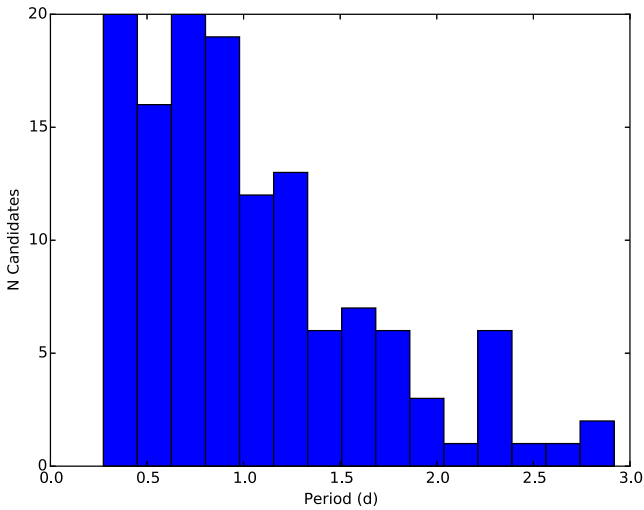


Figure 11. The distribution of pulsation periods for GDOR classified objects.

4.6 γ Doradus stars

We have a sample of 133 γ Doradus candidates, using a class probability cut of 0.7. We plot the amplitude and period distributions in Figs 10 and 11, following the same definition of amplitude as for the δ Scuti sample. Note that this amplitude is only for the dominant period phase curve, and so does not include the other significant frequencies often present in γ Doradus light curves. The period distribution covers the expected range for γ Doradus variables. Due to the lack of colour information available, γ Doradus objects are degenerate with slowly pulsating B stars.

4.7 RR Lyrae ab-type stars

As the RRab class has less well-calibrated probability (almost all candidates with Prob(RRab) greater than 0.5 seem to be real), we use an adjusted class probability threshold of 0.5 to study this class. This leaves 154 candidates. Their amplitude distribution is shown in Fig. 12, and peaks at significantly higher amplitude than that of the DSCUT and GDOR candidates as would be expected. Most of

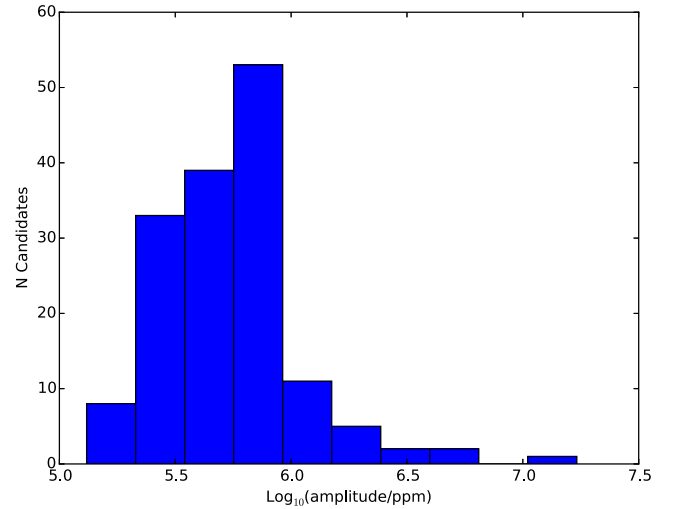


Figure 12. The distribution of phase curve amplitude for RRAb classified objects.

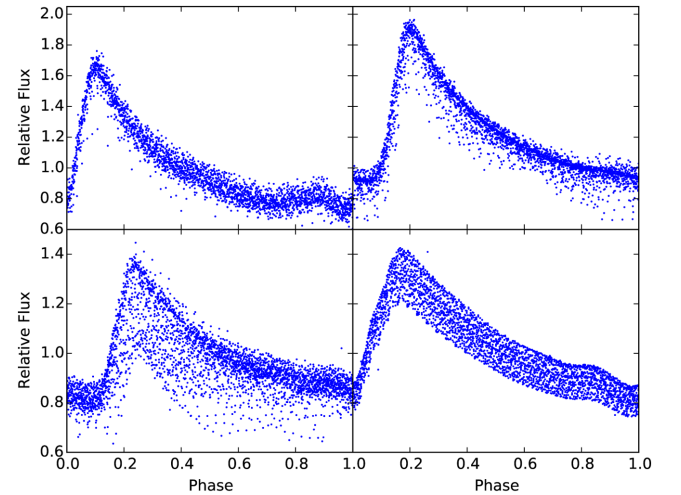


Figure 13. Four phase-folded RRAB classified light curves. Clockwise from top-left, the EPIC IDs are 210830646, 206409426, 211069540 and 203692906.

these candidates are previously known; we find that 129 of them are in K2 proposals focused on RR Lyrae stars. These proposals contain both known and candidate RR Lyraes; in the candidate cases our classification provides some support for them truly being RR Lyrae variables. Assuming these proposals were comprehensive (reasonable, given the multiple teams involved), this leaves 25 candidates as potential new discoveries by this catalogue. However, as these objects are those not in the proposals, there is a selection effect in favour of misclassified non-RR Lyrae objects. We performed a visual examination of each of these 25 light curves, which resulted in 8 of the 25 being confirmed as real RR Lyrae candidates (the others being either misclassified outbursting stars or particularly high-amplitude noise). An additional two candidates were found by using the PDC light curve set and checking objects in both sets with class probability between 0.4 and 0.5, resulting in 10 total new candidates. These objects may still be blends of true RR Lyraes, hence the candidate designation. We plot the phase-folded light curves for two new discoveries and two known RR Lyrae stars in Fig. 13. Some amplitude modulation can be seen, due to some of these

targets exhibiting the Blazhko effect (Blažko 1907). RR Lyraes are immensely useful objects, allowing studies of the evolution of stellar populations throughout the Galaxy and in other nearby galaxies. Due to an absolute magnitude–metallicity relation (Sandage 1981), it is possible to use them for distance estimation.

5 CONCLUSION

We have implemented a novel combined machine learning algorithm, using both SOMs and RFs to classify variable stars in the K2 data. We consider fields 0–4, and intend to update the catalogue as more fields are released. As more data builds up, it may become possible to implement new variability classes, and study the effect of different detrending methods on the catalogue performance. We obtain a success rate of 92 per cent using out of bag estimates on the training set.

We train the classifier on a set of Kepler and some K2 data from fields 0–2. As such it is applied completely independently to the majority of the K2 data, and the whole of fields 3–4. That we obtain good results for fields 3–4 bodes well for application of the classifier to future data.

Algorithms like this will become an increasingly important step in processing the data volumes expected from future astronomical surveys. To maximize scientific return, it is critical to select interesting candidates, and do so rapidly and with minimal input. We hope that this method will contribute to the growing body of work attempting to address this issue.

ACKNOWLEDGEMENTS

The authors thank the anonymous referee for a helpful review of the manuscript. This paper includes data collected by the Kepler mission. Funding for the Kepler mission is provided by the NASA Science Mission directorate. The data presented in this paper were obtained from the Mikulski Archive for Space Telescopes (MAST). STScI is operated by the Association of Universities for Research in Astronomy, Inc, under NASA contract NAS5-26555. Support for MAST for non-*HST* data is provided by the NASA Office of Space Science via grant NNX13AC07G and by other grants and contracts. We acknowledge with thanks the variable star observations from the AAVSO International Database contributed by observers worldwide and used in this research.

REFERENCES

- Aigrain S., Hodgkin S. T., Irwin M. J., Lewis J. R., Roberts S. J., 2015, *MNRAS*, 447, 2880
- Armstrong D. J. et al., 2015, *A&A*, 579, A19
- Balona L. A., Dziembowski W. A., 2011, *MNRAS*, 417, 591
- Balona L. A., Guzik J. A., Uytterhoeven K., Smith J. C., Tenenbaum P., Twicken J. D., 2011, *MNRAS*, 415, 3531
- Balona L. A., Baran A. S., Daszyńska-Daszkiewicz J., De Cat P., 2015, *MNRAS*, 451, 1445
- Blažko S., 1907, *Astron. Nachr.*, 175, 325
- Blomme J. et al., 2010, *ApJ*, 713, L204
- Borucki W. J. et al., 2010, *Science*, 327, 977
- Bostrom H., 2008, in *Proc. 7th Int. Conf. Machine Learning and Applications, Calibrating Random Forests*. IEEE, p. 121
- Bradley P. A., Guzik J. A., Miles L. F., Uytterhoeven K., Jackiewicz J., Kinemuchi K., 2015, *AJ*, 149, 68
- Breiman L., 2001, *Mach. Learn.*, 45, 5
- Brett D. R., West R. G., Wheatley P. J., 2004, *MNRAS*, 353, 369
- Brier G. W., 1950, *Mon. Weather Rev.*, 78, 1
- Brink H., Richards J. W., Poznanski D., Bloom J. S., Rice J., Negahban S., Wainwright M., 2013, *MNRAS*, 435, 1047

- Carrasco Kind M., Brunner R. J., 2014, *MNRAS*, 438, 3409
- Debosscher J., Blomme J., Aerts C., De Ridder J., 2011, *A&A*, 529, A89
- Eyer L., Blake C., 2005, *MNRAS*, 358, 30
- Grigahcène A. et al., 2010, *ApJ*, 713, L192
- Hanke M., Halchenko Y. O., Sederberg P. B., Hanson S. J., Haxby J. V., Pollmann S., 2009, *Neuroinformatics*, 7, 37
- Howell S. B. et al., 2014, *PASP*, 126, 398
- Kohonen T., 1990, *Proc. IEEE*, 78, 1464
- Kügler S. D., Gianniotis N., Polsterer K. L., 2015, *MNRAS*, 451, 3385
- LaCourse D. M. et al., 2015, *MNRAS*, 452, 3561
- Lomb N. R., 1976, *Ap&SS*, 39, 447
- Lund M. N., Handberg R., Davies G. R., Chaplin W. J., Jones C. D., 2015, *ApJ*, 806, 30
- McQuillan A., Aigrain S., Mazeh T., 2013, *MNRAS*, 432, 1203
- McQuillan A., Mazeh T., Aigrain S., 2014, *ApJS*, 211, 24
- Mahabal A. et al., 2008, *Astron. Nachr.*, 329, 288
- Masci F. J., Hoffman D. I., Grillmair C. J., Cutri R. M., 2014, *AJ*, 148, 21
- Matijević G., Prsa A., Orosz J. A., Welsh W. F., Bloemen S., Barclay T., 2012, *AJ*, 143, 123
- Molnár L., Pál A., Plachy E., Ripepi V., Moretti M. I., Szabo R., Kiss L. L., 2015, *ApJ*, 812, 2
- Nemec J. M., Cohen J. G., Ripepi V., Derekas A., Moskalik P., Sesar B., Chadid M., Bruntt H., 2013, *ApJ*, 773, 181
- Nun I., Pichara K., Protopapas P., Kim D.-W., 2014, *ApJ*, 793, 23
- Press W. H., Rybicki G. B., 1989, *ApJ*, 338, 277
- Protopapas P., Giammarco J. M., Faccioli L., Struble M. F., Dave R., Alcock C., 2006, *MNRAS*, 369, 677
- Prsa A. et al., 2011, *AJ*, 141, 83
- Richards J. W. et al., 2011a, *ApJ*, 733, 10
- Richards J. W. et al., 2011b, *ApJ*, 744, 192
- Richards J. W., Starr D. L., Miller A. A., Bloom J. S., Butler N. R., Brink H., Crellin-Quick A., 2012, *ApJS*, 203, 32
- Samus N. N., Durlevich O. V., Kazarovets E. V., Kireeva N. N., Pastukhova E. N., Zharova A. V., 2009, *VizieR On-line Data Catalog: B/gcvs*, 102025
- Sandage A., 1981, *ApJ*, 248, 161
- Scargle J. D., 1982, *ApJ*, 263, 835
- Slawson R. W. et al., 2011, *AJ*, 142, 160
- Smith J. C. et al., 2012, *PASP*, 124, 1000
- Stassun K. G., Pepper J. A., Paegert M., De Lee N., Sanchis-Ojeda R., 2014, preprint ([arXiv:1410.6379](https://arxiv.org/abs/1410.6379))
- Stumpe M. C. et al., 2012, *PASP*, 124, 985
- Thompson S. E. et al., 2012, *ApJ*, 753, 86
- Tkachenko A. et al., 2013, *A&A*, 556, A52
- Tornaiainen I. et al., 2008, *A&A*, 482, 483
- Torrence C., Compo G. P., 1998, *Bull. Am. Meteorol. Soc.*, 79, 61
- Uytterhoeven K. et al., 2011, *A&A*, 534, A125
- Vanderburg A., Johnson J. A., 2014, *PASP*, 126, 948

SUPPORTING INFORMATION

Additional Supporting Information may be found in the online version of this article:

Table 4. Catalogue table for our Warwick detrended light curves.

Table 5. Data features for our Warwick detrended light curves.

Table 6. Catalogue table for PDC detrended light curves.

Table 7. Data features for PDC detrended light curves.

(<http://www.mnras.oxfordjournals.org/lookup/suppl/doi:10.1093/mnras/stv2836/-/DC1>).

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 \LaTeX file prepared by the author.