Classification of Multiwavelength Transients with Machine Learning

K. Sooknunan¹, M. Lochner^{2,3,5}, Bruce A. Bassett^{1,2,3,4}, H. V. Peiris^{5,6}, R. Fender^{7,9}, A. J. Stewart^{7,8}, M. Pietka⁷, P. A. Woudt⁹, J. D. McEwen¹⁰, O. Lahav⁵

- ¹Department of Maths and Applied Maths, University of Cape Town, Cape Town, South Africa
- ²African Institute for Mathematical Sciences, 6 Melrose Road, Muizenberg, 7945, South Africa
- ³South African Radio Astronomy Observatory, The Park, Park Road, Pinelands, Cape Town 7405, South Africa
- ⁴ South African Astronomical Observatory, Observatory, Cape Town, 7925, South Africa
- ⁵ Department of Physics and Astronomy, University College London, Gower Street, London WC1E 6BT, UK
- ⁶ The Oskar Klein Centre for Cosmoparticle Physics, Department of Physics, Stockholm University, AlbaNova, 10691 Stockholm, Sweden
- ⁷ Astrophysics, Department of Physics, University of Oxford, Keble Road, Oxford OX1 3RH, UK
- ⁸Sydney Institute for Astronomy, School of Physics, The University of Sydney, NSW 2006, Australia
- ⁹Department of Astronomy, University of Cape Town, Cape Town, South Africa
- ¹⁰ Mullard Space Science Laboratory, University College London, Surrey RH5 6NT, UK

Accepted XXX. Received YYY; in original form ZZZ

ABSTRACT

With the advent of powerful telescopes such as the SKA and LSST, we are entering a golden era of multiwavelength transient astronomy. In order to cope with the dramatic increase in data volume as well as successfully prioritise spectroscopic follow-up resources, we propose a new machine learning approach for the classification of multiwavelength transients. The algorithm consists of three steps: (1) augmentation and interpolation of the data using Gaussian processes; (2) feature extraction using a wavelet decomposition; (3) classification with the robust machine learning algorithm known as random forests. We apply this algorithm to existing radio transient data, illustrating its ability to accurately classify most of the eleven classes of radio variables and transients after just eight hours of observations, achieving an overall accuracy of 73.5%. We show how performance is expected to increase as more training data are acquired, by training the classifier on a simulated representative training set, achieving an overall accuracy of 97.4%. Finally, we outline a general approach for including multiwavelength data for general transient classification, and demonstrate its effectiveness by incorporating a single optical data point into the analysis, which improves the overall accuracy by $\approx 22\%$.

Key words: Radio Transients – Machine learning

1 INTRODUCTION

In the coming years, radio astronomy will enter a new era of deep field surveys with the advent of the Square Kilometre Array¹ (SKA) and its precursors, MeerKAT² and the Australian Square Kilometre Array Pathfinder³ (ASKAP). These telescopes will achieve unprecedented sensitivity and resolution. Large MeerKAT science projects such as ThunderKAT (Armstrong et al. 2018) will dramatically increase the detected number of radio transients. In the past, radio

transient datasets have been small, allowing spectroscopic classification of all objects of interest. As the event rate increases, follow-up resources must be prioritised by making use of early classification of the radio data. Machine learning algorithms have proven themselves invaluable in this context (Ball & Brunner 2010).

There has been a substantial amount of work done with machine learning in astronomy over the last decade. This includes research done by Bailer-Jones (2001) in stellar classification, image-based classification of supernovae (Romano et al. 2006; Bailey et al. 2007), classifying variable stars (Richards et al. 2011), and photometric supernovae classification (Newling et al. 2011). In recent years, these algorithms have been used successfully in classifying optical

¹ www.skatelescope.org

² www.ska.ac.za/science-engineering/meerkat

³ www.atnf.csiro.au/projects/askap/index.html

transients, such as classification of transients in SDSS images (Buisson et al. 2015), supernovae (Lochner et al. 2016), variable sources (Farrell et al. 2016) or general optical transients (Mahabal et al. 2017). Some machine learning methods have been investigated for the upcoming ASKAP survey for Variables and Slow Transients (VAST) (Murphy et al. 2013), but these were only applied to optical data.

In the burgeoning era of multimessenger astronomy, incorporating data from different telescopes could dramatically improve classification of events. A prime example of this is the MeerLICHT ⁴ telescope (Bloemen et al. 2016), an optical telescope whose observing schedule is synchronised with that of the (night time) observations of the radio telescope MeerKAT, resulting in simultaneous optical and radio observations of transients. Alert streams from telescopes such as Fermi⁵ and LSST ⁶(LSST Science Collaboration et al. 2009) will also enable rapid coordination for multimessenger observations. Combining these data sources necessitates a new general framework for multimessenger machine learning.

In this paper we outline a method for the automatic classification of radio transients, based on Lochner et al. (2016), that makes use of multiwavelength data and machine learning. We define how to incorporate data in other wavelengths and alert streams from other telescopes. We test our method on existing radio transient light curves, exploring the effects of non-representative training sets (i.e. when the algorithm is tested on objects that are dissimilar to those in the training set). We also demonstrate, with an example, the effect of including optical data to improve classification accuracy.

1.1 Radio transients

A transient is an astronomical object observed to have a time-dependent brightness (flux). The time scales of variability range from milliseconds to a few years. Radio transients have emission frequencies in the radio regime (for a review on radio transients see Fender & Bell (2011)). Transient events are typically divided into two kinds: incoherent synchrotron events and coherent burst events.

Incoherent synchrotron events are thought to be caused by a high energy phenomenon, from which a large amount of energy is released over a longer time scale (on the order of minutes or larger). When in a steady state, this energy release is limited to a brightness temperature of $T \leq 10^{12} {\rm K}$ (Fender et al. 2015).

Coherent bursts occur on much shorter time scales (on the order of seconds) and can have brightness temperatures up to 10^{30} K. This type of emission can only be observed using a special mode on radio telescopes (Fender et al. 2015). For this study, therefore, we will only consider incoherent synchrotron transient events. We assume that a light curve is generated from measuring the fluxes from radio images. We include the following transients in our study: active galactic nuclei (AGN), algols, flare stars (FS), gamma ray bursts

(GRBs), kilonovae, magnetars, novae, RS canum venticorums (RSCVn), supernovae (SNe), tidal disruption events (TDEs) and X-ray binaries (XRBs).

Looking at the raw radio light curves of transients, while some objects exhibit obvious differences (such as binary star systems and AGN) others look more similar. Contextual information (such as the location of the object) can be invaluable in telling the difference between classes, as can multiwavelength data. We begin by studying how well we can distinguish between classes using radio data alone, and then show how including contextual and multiwavelength information can improve the classification accuracy.

1.2 Machine learning overview

Machine learning can broadly be split into two approaches: supervised and unsupervised learning. In this study we use supervised learning. A supervised machine learning algorithm automatically learns a model given a set of known inputs and outputs, called a training set. The now trained algorithm can be given new inputs, called a test set, that it will then map to an output. In terms of a classification problem, the inputs constitute objects to be classified, and the outputs are the class labels assigned to each object. For an in depth review of machine learning see Mitchell (1997) or MacKay (2003).

In this study we used the "random forests" algorithm (Ho 1995; Breiman 2001). Random forests have been shown to outperform other algorithms in a variety of cases (Caruana & Niculescu-Mizil 2006; Liu et al. 2013; Lochner et al. 2016). Deep learning algorithms⁷ such as LSTMs (Hochreiter & Schmidhuber 1997) would be an interesting approach to this problem. However we do not consider them here due to the requirement of large training sets and significant computational resources, especially in light of the excellent performance of faster traditional techniques.

1.2.1 Random Forests

Ensemble methods like random forests (Ho 1995; Breiman 2001) build robust classifiers out of a multitude of weak learners such as decision trees. A decision tree creates a mapping, by making a series of "yes/no" decisions, from an input vector (the feature vector) to an output label (the class) (Ball & Brunner 2010). The algorithm creates this mapping by making decisions based on whether or not a given component of the feature vector falls into some range. One of the main drawbacks of decision trees is the high variance on the labels it outputs.

This problem can be overcome by training many separate trees and taking the average of the output. Random forests perform an additional step, each tree in the "forest" is trained on a random subset of the total feature set. This leads to more robust overall predictions. We used the package Scikit

 $^{^4}$ www.meerlicht.uct.ac.za

 $^{^{5}\ \}mathrm{https://fermi.gsfc.nasa.gov}$

⁶ www.lsst.org

 $^{^{7}}$ For a review on deep learning algorithms see Vargas et al. (2017).

Learn⁸ (Pedregosa et al. 2011) to implement the random forest classifier.

1.2.2 Feature extraction in machine learning

Classical machine learning techniques can seldom use data in its raw format for classification. Feature extraction is a technique used to reduce the dimensionality of the data by summarizing the information contained in the original data. The features one uses should also be well-separated between classes. Taking a repeating light curve as an example, features one could extract are the frequency; the amplitude; the phase etc. For more on feature extraction see Li et al. (2016). An obvious choice of simple features for this problem would be changes of flux values over specific time periods. However we found these were inadequate to capture the variation between classes (see Sec. B1) and so we instead follow the feature extraction procedure used in Lochner et al. (2016). In Sec. 2.2, we outline the wavelet decomposition approach used, which resulted in much higher performance.

1.2.3 Visualising features

Visualising feature vectors is quite difficult because of their often high-dimensional nature. One tool commonly used to visualise higher-dimensional spaces is t-distributed Stochastic Neighbor Embedding or t-SNE (van der Maaten & Hinton 2008). It works by computing the probability that two points are similar in the higher dimensional space based on its Euclidean distance. It does this for every pair of points in the feature set, then attempts to find a lower dimensional representation of these points that preserves the probability distribution. Thus points clustered in this lower dimensional representation correspond to points clustered in the original higher dimensional space. t-SNE uses Student's tdistribution when determining the degree of similarity of two points. We stress that t-SNE plots are useful tools for visualisation purposes only, and cannot be used as classifiers themselves due to their stochastic nature.

1.2.4 Training, testing and cross validation

In machine learning, a dataset is generally split into two main subsets: the training set and the test set. The algorithm is given the training set from which to learn the parameters of the model. The test set is reserved in order to check if the algorithm has learned an accurate model, and is only presented to the algorithm after the training step is complete. In some cases an algorithm can perform very well on the training set but perform poorly on the test set. This occurs when the algorithm overfits the model parameters to the test set and hence the model will not generalise to the training set. One method for overcoming this is known as k-fold cross validation. Instead of splitting the dataset into two subsets, the dataset is split into k > 2 subsets. One of these subsets is then used as the validation set while the

others are used as the training set. This is then repeated k times until all k subsets has been used as a test set and average results are used for the model parameters.

Machine learning models have two types of parameters. The first are the parameters that are learned by the algorithm during training as mentioned above. The second type is known as hyperparameters. These are parameters that are not learned during training but are set by the user. These parameters can also be optimised by specifying a range for each hyperparameter, searching through this hyperparameter space and using cross-validation to choose the parameters with the best performance. The primary hyperparameter for the random forest algorithm is the number of decision trees in the forest. This was optimised using a 3-fold cross validation.

1.2.5 Evaluating machine learning results

The performance of a machine learning algorithm can be measured in different ways. We will evaluate our results using confusion matrices. A confusion matrix is a plot showing the true label of the object on one axis and the label predicted by the machine learning algorithm on the other. For the simplest classification problem, a binary classification problem, the confusion matrix would be a 2×2 matrix with the true positives and true negatives along the diagonal, and the false negatives and false positives on the off-diagonals. Therefore confusion matrices show how well the algorithm classifies each class. Classes classified correctly would appear on the diagonal, and incorrect classifications would appear on the off-diagonals.

2 GENERAL APPROACH TO MULTIWAVELENGTH TRANSIENT CLASSIFICATION

Drawing heavily from Lochner et al. (2016), we outline a general approach to classifying transients with multimessenger data. The approach is split into two main sections; the first deals with combining data from different sources and the second builds a machine learning classifier that uses light curve data of any wavelength. This creates a general approach useful for combining all sources of information that may be useful to classifying transients, in addition to classification using the light curves themselves.

2.1 Combining multiple data sources

The data used to classify a source need not be information extracted directly from the light curves. They can also be prior or external information about the sources, such as fluxes of the source in different wavelengths or, contextual information, such as position of the object in the sky. Information from alert streams from other observatories can also be added as external information (e.g. the presence of gamma ray emission or a gravitational wave detected by

⁸ www.scikit-learn.org

LIGO⁹ in the region of a new radio transient source can be highly discriminating).

There are two methods for incorporating information from other sources:

- Probabilistic Approach: most machine learning classification algorithms are capable of producing a score that can be interpreted as a probability of an object belonging to a particular class. To combine this with external information, such as the presence of a coincident alert at another wavelength, we can calculate the prior probability, P(C), of the object being in a certain class C, given all prior information. This probability, P(C), would then be multiplied by the probability given by the classifier to give a final probability of some object being in class C.
- Extra Features: the second method is to use the information as an extra feature in the machine learning process. For example, if one has a flux measurement at any other wavelength, one could add that flux as a feature. The advantage of this approach is that correlations between the different features are learned automatically by the machine learning algorithm, potentially resulting in improved classification accuracy.

The disadvantage of this latter approach is that machine learning algorithms do not intrinsically deal well with missing data. This could happen if, for instance, MeerKAT detects a transient during a daytime observation when MeerLICHT cannot observe. While feature imputation techniques exist (Quinlan 1993), a more interpretable approach may be to combine the probabilities where, if data are missing, a default probability based on prior observations (for example, known transient rates) can be used.

The specific setup of the problem will dictate which approach is more appropriate, but formalising this process is a step towards automated multimessenger machine learning pipelines, that can then be fed into downstream analysis including spectroscopic follow-up prioritisation.

2.2 General approach to transient classification with light curves

Our method for transient classification follows identically the technique outlined in Lochner et al. (2016), which was used for classification of supernovae light curves at optical wavelengths. The technique is applicable to any transient (or indeed, almost any time series data).

Step 1: Interpolation

Given some light curve data, \mathcal{D} , the first step is to interpolate the data so that it is on a uniform grid. This is done using Gaussian Processes (GPs; Rasmussen & Williams (2005)), since the mean function derived from a GP is an extremely robust interpolator in the presence of noisy data.

Step 2: Wavelet decomposition

Time series data can be decomposed into a linear combination of basis functions

$$f(x) = \sum_{k} a_k \phi_k(x) , \qquad (1)$$

where $\phi_k(x)$ are orthogonal basis functions and a_k are the respective coefficients.

This is a common approach in signal processing and can be a powerful tool for feature extraction, to obtain the set of coefficients used as the features with a machine learning algorithm. One widely used form of this is a Fourier decomposition, where a signal can be decomposed into the component frequencies. However the Fourier decomposition loses all localisation information and is thus mostly applicable to regular, repeating signals.

By contrast, in transient classification, the object may be observed at any point in its light curve and the algorithm must determine its class in this setting. Thus, we require a decomposition method that is translation-invariant but still sensitive to the intrinsic shape of the curve. A form of decomposition that is approximately scale and translation-invariant is known as the stationary wavelet transform (Mallat 2009; Holschneider et al. 1989). Following its successful use in Lochner et al. (2016) and Narayan et al. (2018), we make use of the stationary wavelet transform with the symlet family, as implemented in the package PyWavelets 10.

Step 3: Dimensionality reduction with PCA

The stationary wavelet transform produces a large number of redundant features, too many for standard machine learning techniques. Principal Component Analysis (PCA) (Pearson 1901) is a dimensionality reduction technique. It is a linear transformation that decorrelates a set of correlated variables by calculating the covariance matrix of the dataset. The eigenvalues and eigenvectors of this matrix are computed. The total variance of the data can be quantified by calculating the sum total of the eigenvalues. The eigenvectors with the largest eigenvalues, which describe the majority of the variability in the dataset are stored; smaller eigenvalues are disregarded. The number of eigenvectors kept are decided by the fraction of variability one would want to retain in the dataset; variability is defined as the sum of all eigenvalues. For example, if we want to keep 90% variability in our dataset, then we would keep the corresponding eigenvectors of the largest eigenvalues (in descending value) until their sum equals 90% of the sum of all eigenvalues. The large number of coefficients that we obtain from PyWavelets can be projected onto the stored eigenvectors, producing a new set of eigenvalues.

We use GPs as described in Sec. 3.2 for interpolation.

⁹ www.ligo.org

¹⁰ https://github.com/PyWavelets/pywt

Table 1. Breakdown of radio transient data into the relevant types and data sources. GBI refers to data collected from the Green Bank Interferometer. From Lit. refers to data collected from the literature.

Type	From Lit.	GBI	Total	
AGN	17	13	30	
Algol	1	2	3	
FS	5	0	5	
GRB	4	0	4	
Kilonova	1	0	1	
Magnetar	1	0	1	
Nova	8	0	8	
RSCVn	0	2	2	
SN	13	0	13	
TDE	2	0	2	
XRB	11	9	20	
Totals	63	26	89	

3 APPLICATION TO RADIO TRANSIENTS

Motivated by the expected increase in new transient detections with modern radio telescopes, we apply our general approach to existing radio transient data. Because this data is limited, we use a data augmentation technique described in Sec. 3.2 to artificially increase the number of light curves for training and testing. We follow the feature extraction method described in Sec. 2 and also illustrate the effect of incorporating additional data by including contextual information and optical data.

3.1 Radio data

The radio transient data used were collected by Pietka et al. (2017) except for the kilonova light curve, which was collected by Dobie et al. (2018).

Most of the data are obtained the literature and the rest is from the Green Bank Interferometer (GBI). Data collected from the GBI have a much higher cadence than the data obtained from the literature. The data consists of time series light curve data (radio flux as a function of time). These radio transient light curves can be separated into eleven different classes or types. The total number of light curves in each type and their sources are shown in Table 1. An example light curve for each class is shown in Fig. 1.

It is important to note that all the light curves have different lengths. The length of the light curve is correlated with the type of object, due to observational biases. Some objects are observed over years (e.g. AGN) and others are observed over only a few hours (e.g. FS). Figure 2 shows the number of light curves for each class as a function of the total length of observation for that object. Because of this bias, we restrict our study to a timescale of eight hours, which is the longest observation time for which we have measurements for all classes (see Fig. 2). Classification on this timescale will also allow relatively prompt follow-up triggers. The technique is

11 https://public.nrao.edu/telescopes/green-bank-interferometer

applicable on even shorter timescales, even if there are very few flux measurements, although classification accuracy will likely decrease.

3.2 Augmentation

Data augmentation is a way of creating new data from existing data. One widely used method of augmenting data is to use a model to simulate new data. Since models are not available for all classes of radio transient, we use Gaussian processes (GPs) (Rasmussen & Williams 2005) to augment our data, similar to Revsbech et al. (2018).

GPs are a set of indexed random variables, defined such that every finite subset of a GP follows a multivariate normal distribution. Thus every subset of the GP can be characterised fully by its mean and covariance functions. This allows a GP to be fit to data and used to generate new data that follow the same distribution as the original data, as the mean and variance are calculated from the underlying distribution. GP regression was performed on our dataset using George (Ambikasaran et al. 2014).

Each class of light curve has different general characteristics. This necessitated a combination of a few different kernel (or covariance) functions to be used for regression, which are defined below.

The first was a regular radial base function. The kernel, k, is given by

$$k_{\rm rad}(r) = \sigma^2 \exp(-r^2/2), \qquad (2)$$

where σ^2 is a hyperparameter and r^2 is the squared distance defined under some general metric C,

$$r^{2}(x_{i}, x_{i}) = (x_{i} - x_{i})^{T} C^{-1}(x_{i} - x_{i}),$$
(3)

where x_i and x_j are x data points. An isotropic metric was used, hence, $r^2(x_i, x_j) = |x_i - x_j|^2$.

The second kernel used was the exponential function given by

$$k_{\exp}(r) = \exp(-r), \qquad (4)$$

where r is defined in Eq. 3.

Lastly, we used the exponential sine squared function, given by

$$k_{\text{sine}}(r) = \exp\left(-\Gamma \sin^2\left[\frac{\pi}{P}|x_i - x_j|\right]\right),$$
 (5)

where Γ and P are hyperparameters, x_i and x_j are x data points.

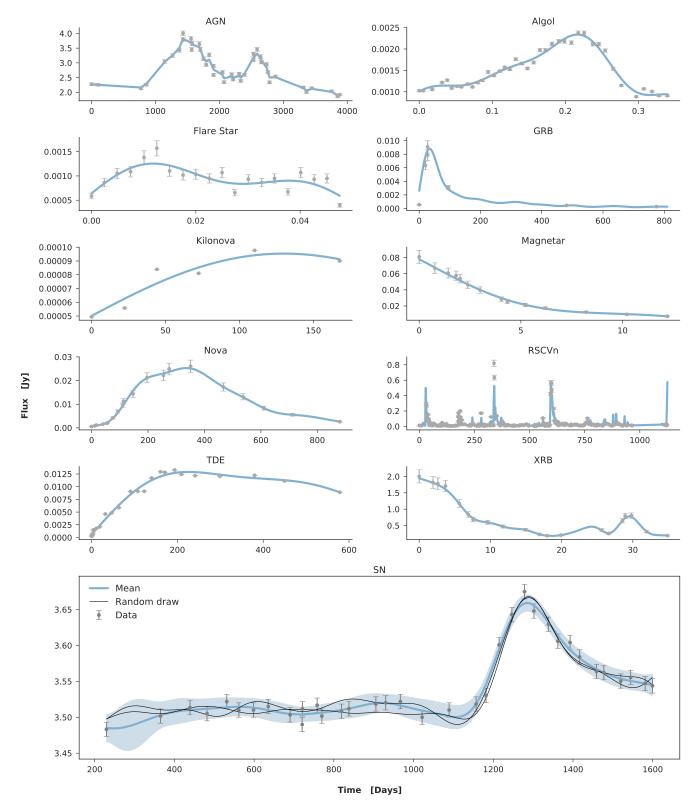


Figure 1. An example light curve for each of the the radio transient types plotted as flux (in Jy) as a function of time (in days). It is clear that the data are taken on extremely different timescales. The original data points are shown in grey and the mean of the Gaussian Process is shown in blue. For the SN light curve one standard deviation away from the mean is shown as a blue envelope. The black lines shown in the SN light curve are two random draws from the GP. It can be seen that these lines are different from, but still consistent with the original data.

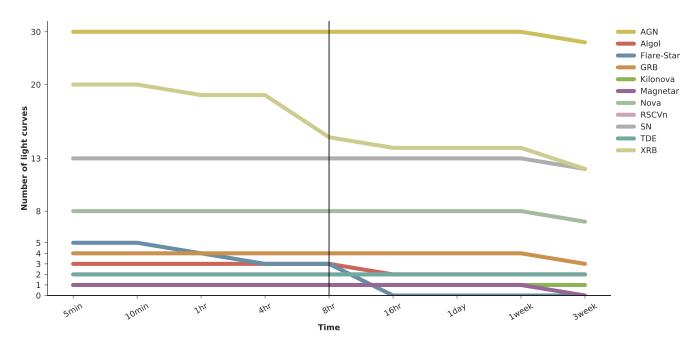


Figure 2. The number of light curves for each class as a function of the length of observation for that that source (in other words, there are y objects with light curves that were observed for at least time x). It can be seen that 8 hrs is the shortest observation period for which all the classes have a non-zero number of light curves.

Given the general characteristics of the different classes, it was found that one of three combinations of these kernels fit the data best. These combinations were defined as follows

$$K_1 = \omega_1 k_{\text{rad}}(r)$$

$$K_2 = \omega_1 k_{\text{rad}}(r) + \omega_2 k_{\text{rad}}(r) k_{\text{sine}}(r)$$

$$K_3 = \omega_1 k_{\text{rad}}(r) + \omega_2 k_{\text{exp}}(r)$$

where ω are weights on each kernel.

GP regression was performed using each of these combinations on each of the light curves. At the start of this regression, the hyperparameters were randomly initialised. The negative log likelihood was calculated using these hyperparameters and the data. The negative log likelihood was then minimised using the Limited memory Broyden-Fletcher-Goldfarb-Shanno (L-BFGS) optimization algorithm from the SciPy ¹² package, to obtain the best set of hyperparameters for each of the kernel combinations. The combination of kernels with the lowest negative log likelihood was used to construct the GP from which to sample. An example GP for a SN light curve is shown in Fig. 1.

3.3 Feature extraction

The feature extraction method used is described in Sec. 2.2. First GP regression was performed on the original data set. From this, many example light curves can be generated, each statistically consistent with the original data, which allows us to generate a realistic synthetic dataset of any

size. To simulate the fact that the transient may be detected at any point on the light curve, a reference time, t_0 , was then drawn at random to be somewhere within the curve. We then sampled 100 flux values between t_0 and t_0+8 hrs from the GP. This approximates an 8 hour radio observation where an image is produced every five minutes. We then used PyWavelets to perform a two-level wavelet decomposition on these 100 points, which returns 400 coefficients. PCA was performed on these, keeping 20 coefficients which corresponds to retaining 99% of the variability in the dataset.

3.4 Incorporating external information

Some classes can have similar light curves but are generally found in different parts of the sky. For example we are more likely to find novae in the Galactic plane, while SNe are more likely to be extragalactic.

In order to break the degeneracies between these classes we added a feature that characterises where the object is in the sky. Telescopes will always have access to the position in the sky in which it is pointing, hence the drawback of adding features that may sometimes be missing, as outlined in Sec. 2.1, will not be an issue.

The RA and Dec coordinates were obtained for all the objects in our dataset using a combination of Simbad ¹³ and NED ¹⁴. These coordinates were then converted to Galactic coordinates. We defined a feature that specified whether or

¹² www.scipy.org

¹³ www.simbad.u-strasbg.fr

¹⁴ https://ned.ipac.caltech.edu

not the object was in the Galactic plane. Any object with a Galactic declination of above 10° or below -10° was considered to be out of the Galactic plane.

4 RESULTS

4.1 The effect of training sets

4.1.1 Fully representative training data

The radio transient data currently available are unfortunately too small to use as a training set for machine learning. So, in order to test the performance of the classifier, we simulated a representative dataset by training the classifier on samples from all the light curves in our dataset.

GP regression was performed on the original 87 light curves. From these, 10,000 "simulated" light curves were generated for each of the eleven classes, to create a balanced training set. Wavelet feature extraction was then performed on each of the simulated light curves resulting in 400 wavelet coefficients. After performing PCA on these coefficients, the 20 most important components were selected and used as features.

In order to show the effect of adding contextual information two separate classifiers were trained using the method outlined in Sec. 3. One classifier was trained without any contextual feature and the other with contextual information as described in Sec. 3.4. The results are shown in Fig. 3. It can be seen that without any contextual information the classifier confuses the classes of XRB, SNe, Novae and GRBs. However after the contextual feature is added, the accuracy for the class of XRBs increases by 13%, SNe increases by 10%, Novae increases by 22% and GRBs by 13%. We expect the confusion matrix with contextual information to characterise the performance of the classifier in practice; thus this contextual feature was used in the rest of this work.

4.1.2 Non-representative training set

Because our original data set is limited, we expect the results outlined in Sec. 4.1.1 to represent an idealised case. In practice, we anticipate that any test set would consist of some objects not present in our training set. To construct a more realistic test of the classifier performance with the data currently available, instead of training the classifier on samples from all the light curves, we trained it on a subset of light curves. As can be seen from Fig. 2, only five classes have greater than four light curves. In order to ensure the training subset contained all classes, we only removed light curves from these five classes. We still included the classes with only one object in both the training and test sets, because these objects can cause confusion between classes which would be artificially removed if they were excluded.

We removed 25% of the light curves in these five classes at random, GP regression was performed on the remaining 75% of the original light curves. From these, 7500 simulated light curves were generated for each of the 11 classes. Wavelet

feature extraction was then performed on each of the simulated light curves followed by PCA as before. This was used as the training set. A further 2500 simulated light curves were drawn from the remaining 25% of the light curves. GP regression, wavelet extraction and PCA were performed as before. This was used as the test set. This was repeated 20 times, each time removing 25% of the light curves at random. The results of this is shown in Figs. 4 and 5. It can be seen that the performance of our classifier has decreased by $\approx 22\%$, which is to be expected given the non-representativeness of the training set. In particular, well-represented classes like AGN still perform well but classes such as SN, where the training set is highly diverse with many dissimilar objects, are poorly classified.

4.1.3 Single curve testing

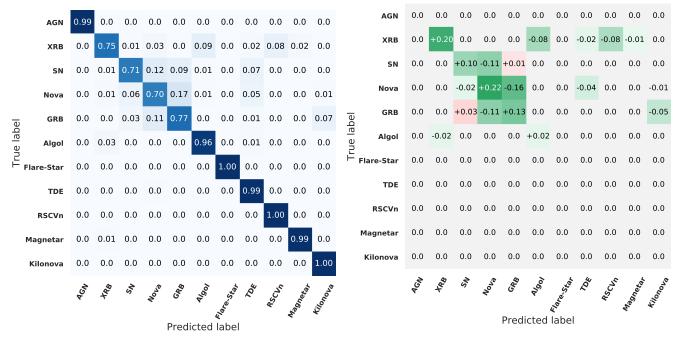
To further demonstrate the effect of a non-representative training set, we tested the classifier's performance on light curves not represented in a training set. We focused on the five classes (AGN, SNe, XRBs, Novae and GRBs) that contain four or more light curves.

We trained the classifier on samples from all of the original light curves, except for one. GP regression was thus performed on 86 of the original light curves and from these simulated light curves were drawn. Wavelet feature extraction was then performed on each of the simulated light curves followed by PCA as before. We then tested the classifier on simulated light curves from the one that was left out of the training set in the training process. This process was then repeated, each time omitting a different light curve in the training set, until every light curve in the dataset had been removed during the training of a classifier at least once. The accuracy for each of the dropped out curves is shown in Table A1. The results for these five main classes are summarised in the violin plots shown in Fig. 6. The colored outline shows the smoothed kernel count distribution of the accuracies. The thick central black line represents the interquartile range. The thin central black line shows the 95% interval. From this is can be seen that the classifier performs well for AGN as most of the accuracies are above 80%. It can also be seen that the classifier performs poorly for SNe as most of the accuracies are below 50%.

As described in Sec. 1.2.3, the t-SNE plots show the distribution of features in high order feature space in two dimensions. It can be seen from Fig. 7 that features from classes where the classifier performs well are relatively separated in feature space. However data for classes which the classifier confuses overlap in feature space hence the classifier is unable to tell the difference between these objects.

4.2 Adding optical data

Simultaneous optical observations do not exist for all the classes in our dataset. To show how the addition of multi-wavelength data could help in classification, an optical flux measurement had to be simulated. Optical data for four classes, namely: AGN, XRB, SNe and GRB, were collected



- (a) Confusion matrix without contextual information
- (b) Confusion matrix showing the difference when contextual information is added

Figure 3. The normalised confusion matrix for wavelet features extracted from eight hours of data. The y-axis shows the true label of the object (true class). The x-axis shows the label which the algorithm predicts for the object (predicted class). The right panel has two colour schemes. The first corresponds to the diagonals. If the values along the diagonal increase they will show in green; if they decrease they will show in red. The second corresponds to the off-diagonals. If the values along the off-diagonals increase they will show in red; if they decrease they will show in green. From the left panel we see that without any contextual information, the classifier confuses the classes of XRB, SNe, Novae and GRBs. The classifier is greatly improved with the added feature of the object's position on the sky. The four classes are no longer confused as most of the off-diagonals are green. The accuracy for the class of XRBs increased by 13%, SNe increased by 10%, Novae increased by 22% and GRBs by 13% (right panel).

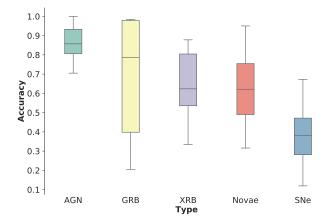


Figure 4. Summary of results from non-representative training set are shown in these box and whisker plots. The bar in each box shows the accuracy for each class averaged over the 20 runs (on each run a classifier was trained on 75% of the total dataset and tested on the remaining 25%, this process was repeated 20 times, each time randomising the training and testing sets). The coloured boxes show the interquartile range. The grey bars (whiskers) show the minimum and maximum accuracies in the 20 runs. It can be seen that the classifier performs poorly for SNe as it has the lowest average accuracy of $\approx 40\%$.

from Stewart et al. (2018). The data consisted of single flux measurements in both optical and radio wavelengths for each source. The dataset had total of 11,882 measurements of which 11,782 were AGN. Figure 8 shows the radio-optical flux distributions for each of the four classes. A two dimensional Gaussian was used to model the distribution all of the classes except for AGN which can be clearly seen to be highly non-Gaussian. The Gaussian fits are shown as contours. These Gaussian distributions were used to sample new optical fluxes for the three classes, new optical fluxes were sampled directly from the distribution for AGNs as there are $\approx 11\,000$ data points.

Using the class of GRB as an example, the process for simulating simultaneous optical and radio observations is is follows. First, a GRB radio light curve was simulated as described previously; the peak flux of this light curve was found. The radio-optical flux distribution were marginalised over, given the peak radio flux, as shown in black in Fig. 8. An optical flux was then drawn from this marginal distribution. Finally, this optical flux was added as an extra feature in the machine learning process. For the class of AGN, a peak radio flux was found as before, points in the radio-optical flux distribution was then binned, centered on the peak radio flux with a bin width of 0.1 mJy. These binned points, now marginalised over the peak radio flux follow a Gaussian distribution similar to that of the other classes. An

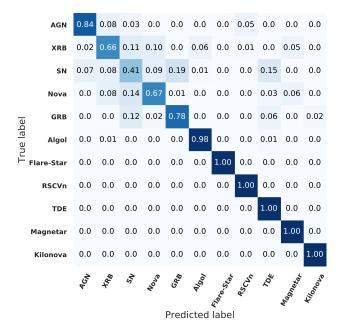


Figure 5. Confusion matrix showing the results of the classifier with a non-representative training set averaged over 20 runs. The y-axis shows the true label of the object (true class). The x-axis shows the label which the algorithm predicts for the object (predicted class). It can be seen that the classes that the classifier performs poorly on are SNe and Novae.

optical flux was then drawn from this marginal distribution and added an extra feature in the machine learning process.

The method outlined in Sec. 4.1.2 was repeated using the four classes mentioned above. The average confusion matrices with and without an optical feature are shown in Fig. 9. It can be seen that the addition of this optical feature significantly improves the performance of the classifier for classes such as SNe and GRBs, with the accuracies increasing by 44% and 25% respectively. This improvement in performance demonstrates the value of simultaneous optical observations, such as from MeerLICHT in classifying radio transients.

5 CONCLUSIONS

We have presented a general formalism for multiwavelength transient classification with machine learning. We outlined the different approaches that can be taken to include data from multiple telescopes, as well as additional information such as source location.

Extending Lochner et al. (2016), we developed a machine learning pipeline for classifying radio transients, and illustrated how to include contextual information (such as whether or not the source is in the Galactic plane) and simultaneous observations from an optical telescope.

We tested our pipeline using existing radio transient light curves gathered from literature. Because these light curves were few in number, we artificially augmented the dataset using Gaussian processes. Features were extracted from this augmented dataset using a wavelet decomposition, followed by dimensionality reduction. We found that the wavelet features have much higher performance than a simpler set of features based on the change in flux of the light curve over different time periods.

If the training data are representative of the test data, the performance of the algorithm is excellent, achieving an average accuracy of 97.4%. However, because our training set is so small, it is highly likely that sources observed by new telescopes such as MeerKAT would not be similar to anything in the training set.

If we remove 25% of the objects from the training set and then test on those removed light curves, we achieve an average accuracy of only 73.5%. From the investigation of the effect of dropping out single individual light curves we find that, while most of the light curves are still classified well, several light curves in the dataset are poorly classified. This is because they are dissimilar to anything in the training set. This effect is most pronounced for supernovae, which are generally easily confused with other classes such as GRBs or novae.

The classification accuracy increases by 22% when contextual information is included as an additional feature. The improvement is particularly noticeable for novae, objects typically found within the galactic plane, which are otherwise confused with GRBs.

Finally, we illustrate the effect of including optical data. Unfortunately simultaneous optical and radio transient data are scarce, so we simulate optical fluxes for the radio light curves, drawing on existing distributions. We find that the average accuracy increases from 72.9% to 94.7% using the more realistic "dropped out" training set.

These results indicate that by including multiwavelength information and making use of a sophisticated machine learning approach, we can expect highly accurate classification of radio transients with even a small training set created with early MeerKAT and MeerLICHT data.

ACKNOWLEDGEMENTS

We acknowledge support from SKA Africa and the National Research Foundation (NRF) towards this research. Opinions expressed and conclusions arrived at, are those of the authors and are not necessarily to be attributed to the NRF. We acknowledge support from STFC for UK participation in LSST through grant ST/N00258X/1 and travel support provided by STFC for UK participation in LSST through grants ST/L00660X/1 and ST/M00015X/1. HVP was supported by the European Research Council under the European Community's Seventh Framework Programme (FP7/2007-2013)/ERC grant agreement no 306478-CosmicDawn. PAW acknowledges support from the University of Cape Town.

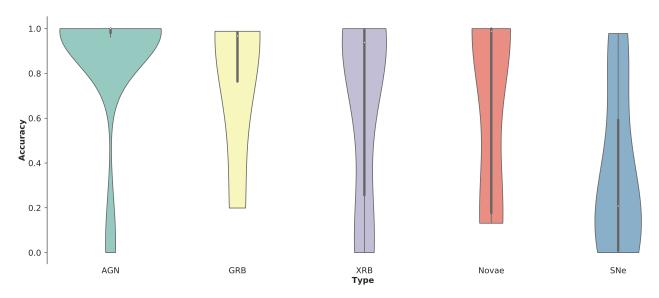


Figure 6. The single curve testing results for the five main classes are summarised in these violin plots. Each violin plot represents a different class as shown. The colored outline shows the smoothed kernel count distribution of the accuracies. The thick central black line represents the interquartile range. The thin central black line shows the 95% interval. From this is can be seen that the classifier performs well with AGN as most of the accuracies are above 80%. It can also be seen that the classifier performs poorly with SNe as most of the accuracies are below 50%.

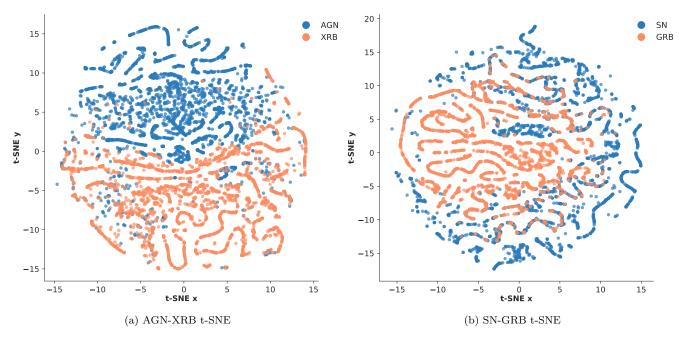


Figure 7. The t-SNE plots for the four classes of radio transients are shown. t-SNE plots show the high dimensional feature space embedding in two dimensions, hence the x and y axes are arbitrary. The left panel shows the t-SNE plot for two classes for which the classifier performs well. The blue points represent feature vectors of AGN objects, orange points represent the feature vectors of XRB objects. The right panel shows the t-SNE plot for two classes for which the classifier performs poorly. The blue points represent feature vectors of SN objects, orange points represent the feature vectors of GRB objects. It can be seen that features from classes that the classifier does well with are relatively separated in feature space (left panel) where as features for classes which the classifier confuses overlap in feature space (right panel).

APPENDIX A: TABLES

Table A1 shows the results of the individual dropout test. Samples from one light curve was removed during training. These samples were then used to test the classifier. The ac-

curacy for the samples of each of the light curves for the four classes are shown. From left to right the classes are: AGN, XRB, SN, Nova and GRB. It can be seen that the classifier is very accurate for AGNs, XRBs and GRBs but misclassifies

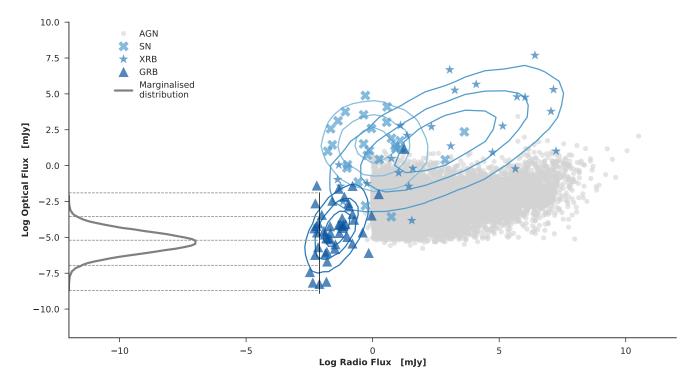


Figure 8. Optical - radio flux distributions for the classes of AGN, XRB, GRB and SNe. Two dimensional gaussian fits to the distributions of three of the classes are shown with contours. These fit the data for the three classes reasonably well, however it can be seen that the distribution for AGNs is highly non-Gaussian. An example peak GRB radio flux measurement is shown as a black vertical line. The GRB radio-optical flux distribution marginalised over this peak flux is shown in dark grey. Optical fluxes were sampled from these marginals.

SNe and Novae. This shows that the intrinsic variabilities of these two classes are not captured by our dataset.

APPENDIX B: FLUX FEATURES

B1 Flux feature extraction

For the first attempt at classifying the radio transients a very simple feature set was used. The feature vector, ΔF , was defined to be

$$\Delta F = F_t - F_{t_0} \,, \tag{B1}$$

where ΔF is the difference in flux between a reference flux, F_{t_0} , chosen at random from the light curve and the flux, F_t , at time $t > t_0$.

In anticipation of a fast imager on MeerKAT we chose the minimum difference between successive flux measurements to be two seconds. The maximum time difference was chosen to be three months to account for the objects that vary on very long time scales such as AGN. The complete set of t where chosen to be:

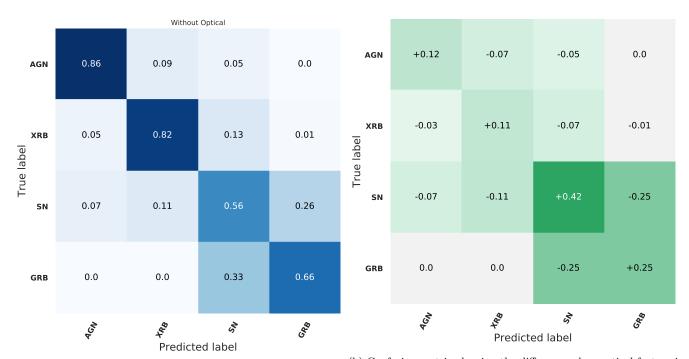
$$t = [2 \sec, 1 \min, 5 \min, 10 \min, 30 \min, 1 \ln, 2 \ln, 4 \ln, 6 \ln, 8 \ln, 12 \ln, 18 \ln, 1 day, 2 day, 4 day, 1 day, 2 week, 3 week, 1 month, 1.5 month, 2 month, 3 month] (B2)$$

The feature extraction method was then as follows. First, GP regression was performed on the original data set. A reference time, t_0 , was then drawn at random to be somewhere within the curve. The GP was then sampled at this t_0 to obtain an F_{t_0} . The GP was then sampled at points (t_0+t) . F_{t_0} was then subtracted from these fluxes to obtain the feature vector ΔF . This process was repeated multiple times for each light curve, each time generating a random t_0 . This was done to simulate the fact that the transient may be detected at any point on the light curve.

B2 Results

Once the flux features were extracted as described in Sec. B1, different subsets of the total feature set was used to train different classifiers. The subsets used were created by truncating the features at different timescales, i.e a classifier was trained on all features up to a maximum of 5 min, then another was trained on all features up to a maximum of 10 min and so on for all t in Eq. B2. The accuracy of each of these classifiers is shown in Fig. B1.

It can see from Fig. B1 that the classifier performs well on long timescales. We would like to investigate is how well the classifier performs on short time scales for each of this classes. It can be seen from Fig. 2 that the number of FS light curves in the dataset goes to zero, hence 8 hrs is the longest timescale at which we have a complete set of classes. Thus the time vector from Sec. B1 is changed to:



- (a) Confusion matrix without optical feature
- (b) Confusion matrix showing the difference when optical feature is added

Figure 9. The normalized confusion matrix for wavelet features extracted from 8hrs of data. The y-axis shows the true label of the object (true class). The x-axis shows the label which the algorithm predicts for the object (predicted class). The right panel has two colour schemes. The first corresponds to the diagonals. If the values along the diagonal increase they will show in green; if they decrease they will show in red, if they decrease they will show in green. From the left panel it can be seen that without the optical feature the classes of GRBs and SNe are confused. With the added optical feature the performance of the classifier is greatly improved as all the off-diagonals are green. From the diagonals it can be seen that all the classes sees an increase in accuracy with the most significant being the classes of SNe and GRBs which increase in accuracy by 42% and 25% respectively (right panel).

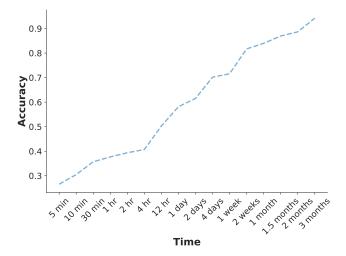


Figure B1. The accuracy of different random forest classifiers each trained on increasing timescales of the total feature set. The y-axis shows the accuracy and the x-axis shows the time observed. It can be seen that as we increase the time observed the accuracy of the classifier also increases.

t = [2 sec, 1 min, 5 min, 10 min, 30 min, 1 hr, 2 hr, 4 hr, 6 hr, 8 hr]

APPENDIX C: VARIABLE AND TRANSIENT CONFUSION MATRICES

From Fig. 1 it can be seen that radio transients can be split into two types: transient that have multiple outbursts on short time scales (e.g. FS) and transients that have one outburst that then decays over longer time scales (e.g. SNe). These two groups are called variables and transients respectively. The dataset was split into these two groups. A classifier was trained on the binary classification of transient vs variable. Classifiers were then trained on the individual groups of variables and transients. The results of this is shown in Fig. C1

It can be seen from Fig. C1 that the classifier can distinguish between variables and transients extremely well, with an overall accuracy of 98.5%. By making this split we increase the accuracy of the XRB and SN classifications slightly when compared to Fig. 3. The class of GRB, however, sees no improvement. As with Fig. 3, SNe are still being confused with TDEs. It can be seen that this split is unnecessary as it achieves little to no improvement in the performance of the classifier.

Table A1. The results of the individual dropout test. Samples from one light curve was removed during training. These samples were then used to test the classifier. The accuracy for the samples of each of the light curves for the five main classes are shown. From left to right the classes are: AGN, XRB, SN, Nova and GRB. From these tables it can be seen that the classifier performs well for AGNs, XRBs and GRBs but poorly for SNe and Novae. This shows that the intrinsic variabilities of these two classes are not captured by our dataset.

Name Acc		Name Acc		Name Acc		Name	Acc	Name	Acc
	Acc	Name	Acc	Name	Acc	Name	Acc	Name	Acc
NGC7213	0.0	B1259-63	0.0	SN1993J	0.0	V1974Cyg	0.131	GRB030329	0.199
0850-121	0.0	MAXIJ1836	0.007	SN2008iz	0.0	RSOph	0.152	GRB970508	0.952
0954 + 658	0.206	ScoX-1	0.014	SN1980K	0.0	V1500Cyg	0.184	GRB060418	0.983
0224 + 671	0.784	CygX-2	0.148	SN1988z	0.009	V407Cyg	0.976	GRB110709B	0.988
2005 + 403	0.824	XTEJ1550	0.581	SN2003L	0.147	Sco2012	1.0		
NRAO530	0.931	aqlX1	0.707	SN1998bw	0.195	TPyx	1.0		
2223 - 052	0.951	GROJ1655	0.927	SN2003bg	0.208	SSCyg	1.0		
1413 + 135	0.961	1909 + 048	0.948	SN2011dh	0.307	V1723Aql	1.0		
0528 + 134p	0.98	SS433	0.954	SN2004dk	0.593				
1622 - 297	1.0	0236 + 610	0.981	SN1994I	0.593				
CTA102	1.0	1915 + 105	0.986	SN2008ax	0.861				
0336-019	1.0	GX17+2	1.0	SN2004cc.	0.931				
3C345	1.0	CygX-1	1.0	SN2004gq	0.978				
B0605-085	1.0	CICam	1.0						
3C454.3	1.0	CirX-1	1.0						
1328 + 254	1.0								
3C120	1.0								
3C273	1.0								
0458-020	1.0								
3C279	1.0								
1237 + 049	1.0								
0851 + 202	1.0								
2200 + 420	1.0								
0528 + 134	1.0								
PKS2004-447	1.0								
1803 + 784	1.0								
0954 + 65	1.0								
1749 + 096	1.0								
NGC4278.	1.0								
AO0235+164	1.0								

REFERENCES

Ambikasaran S., Foreman-Mackey D., Greengard L., Hogg D. W., O'Neil M., 2014

Armstrong R., et al., 2018, PoS, MeerKAT2016, 013 Bailer-Jones C. A. L., 2001

Bailey S., Aragon C., Romano R., Thomas R. C., Weaver B. A., Wong D., 2007, ApJ, 665, 1246

Ball N. M., Brunner R. J., 2010, International Journal of Modern Physics D, 19, 1049

Bloemen S., et al., 2016, Proc. SPIE Int. Soc. Opt. Eng., 9906, 990664

Breiman L., 2001, Mach. Learn., 45, 5

Buisson L. d., Sivanandam N., Bassett B. A., Smith M., 2015, Mon. Not. Roy. Astron. Soc., 454, 2026

Caruana R., Niculescu-Mizil A., 2006, in Proceedings of the 23rd international conference on machine learning. ACM, pp

Dobie D., et al., 2018, Astrophys. J., 858, L15

Farrell S. A., Murphy T., Lo K. K., 2016, VizieR Online Data Catalog, 181

Fender R. P., Bell M. E., 2011, Bulletin of the Astronomical Society of India, 39, 315

Fender R., Stewart A., Macquart J.-P., Donnarumma I., Murphy T., Deller A., Paragi Z., Chatterjee S., 2015. (arXiv:1507.00729), https://inspirehep.net/record/ 1381138/files/arXiv:1507.00729.pdf

Ho T. K., 1995, in Proceedings of 3rd International Conference on Document Analysis and Recognition. pp 278-282 vol.1,

doi:10.1109/ICDAR.1995.598994

Hochreiter S., Schmidhuber J., 1997, Neural Comput., 9, 1735 Holschneider M., Kronland-Martinet R., Morlet J., Tchamitchian P., 1989, in Combes J.-M., Grossmann A., Tchamitchian P., eds, Wavelets. Time-Frequency Methods and Phase Space. p. 286

LSST $_{
m et}$ Science Collaboration al., 2009, preprint, (arXiv:0912.0201)

Li J., Cheng K., Wang S., Morstatter F., Trevino R. P., Tang J., Liu H., 2016, preprint, (arXiv:1601.07996)

Liu M., Wang M., Wang J., Li D., 2013, Sensors and Actuators B: Chemical, 177, 970

Lochner M., McEwen J. D., Peiris H. V., Lahav O., Winter M. K., 2016, Astrophys. J. Suppl., 225, 31

MacKay D. J. C., 2003, Information theory, inference, and learning algorithms. Cambridge University Press

Mahabal A., Sheth K., Gieseke F., Pai A., Djorgovski S. G., Drake A., Graham M., the CSS/CRTS/PTF Collaboration 2017, preprint, (arXiv:1709.06257)

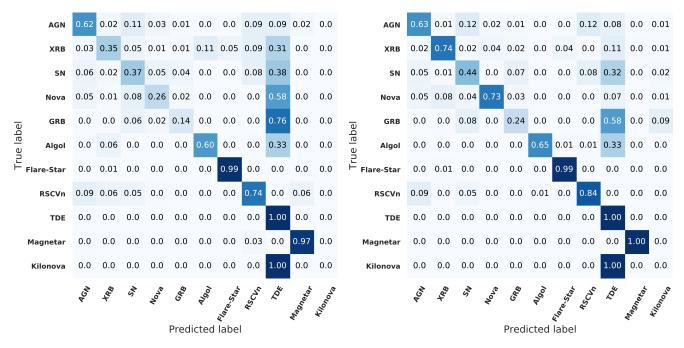
Mallat S. G. S. G., 2009, A Wavelet Tour of Signal Processing: The Sparse Way, third edn

Mitchell T. M., 1997, Machine Learning, 1 edn. McGraw-Hill, Inc., New York, NY, USA

Murphy T., et al., 2013, Publ. Astron. Soc. Australia, 30, e006 Narayan G., et al., 2018, Astrophys. J. Suppl., 236, 9

Newling J., et al., 2011, MNRAS, 414, 1987

Pearson K., 1901, The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science, 2, 559



(a) Confusion matrix without contextual information

(b) Confusion matrix with contextual information

Figure B2. The normalized confusion matrix for flux features extracted from 8hrs of data. The y-axis shows the true label of the object (true class). The x-axis shows the label which the algorithm predicts for the object (predicted class). It can be seen that while contextual information does improve the performance of the classifier this method of feature extraction does not do well to separate the classes.

Pedregosa F., et al., 2011, Journal of Machine Learning Research, 12, 2825

Pietka M., Staley T. D., Pretorius M. L., Fender R. P., 2017, Mon. Not. Roy. Astron. Soc., 471, 3788

Quinlan J. R., 1993, C4.5: Programs for Machine Learning. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA

Rasmussen C. E., Williams C. K. I., 2005, Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning). The MIT Press

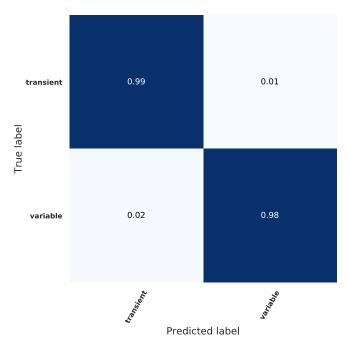
Revsbech E. A., Trotta R., van Dyk D. A., 2018, MNRAS, 473,

Richards J. W., et al., 2011, The Astrophysical Journal, 733, 10 Romano R., Aragon C., Ding C., 2006

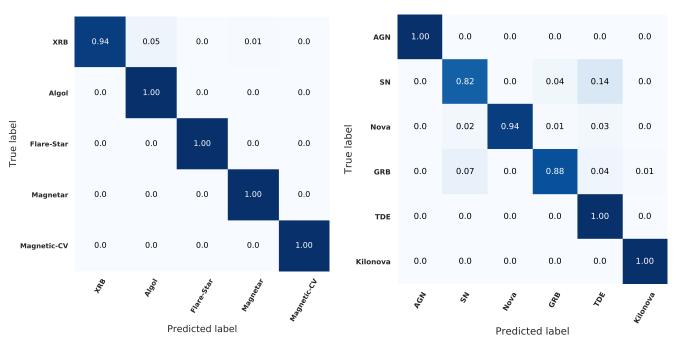
Stewart A. J., MuÃśoz-Darias T., Fender R. P., Pietka M., 2018, l 10.1093/mnras/sty1671

Vargas R., Mosavi A., Ruiz R., 2017, 5

van der Maaten L., Hinton G., 2008, Journal of Machine Learning Research, 9, 2579



(a) Confusion matrix for binary Transient/Variable classification



(b) Confusion matrix for variable classification

(c) Confusion matrix averaged over all runs

Figure C1. The normalized confusion matrix for wavelet features extracted from 8hrs of data. The y-axis shows the true label of the object(true class). The x-axis shows the label which the algorithm predicts for the object(predicted class). It can be seen from these figures that the classifier can distinguish between variables and transients extremely well. By making this split we increase the accuracy of the XRB and SN classifications. The class of GRB, however, sees no improvement. As with before SNe are still being confused with TDEs. It can be seen that this split is unnecessary as it achieves little to no improvement in the performance of the classifier.