

# A reference dataset for astronomical transient event recognition I: lightcurves and tests on classical machine learning algorithms

Mauricio Neira<sup>1</sup>, Catalina Gómez<sup>2</sup>, Diego A. Gómez,<sup>1</sup> Marcela Hernández Hoyos<sup>1</sup>, Jaime E. Forero-Romero<sup>3</sup>, Pablo Arbeláez<sup>2</sup>

<sup>1</sup>*Systems and Computing Engineering Department, Universidad de los Andes, Cra. 1 No. 18A-10, Bogotá, Colombia*

<sup>3</sup>*Departamento de Ingeniería Biomédica, Universidad de los Andes, Cra. 1 No. 18A-10, Bogotá, Colombia*

<sup>2</sup>*Departamento de Física, Universidad de los Andes, Cra. 1 No. 18A-10, Bogotá, Colombia*

Accepted XXX. Received YYY; in original form ZZZ

## ABSTRACT

We introduce NNN an annotated dataset of more than 10000 transient and non-transient object light-curves built from the Catalina Real Time Transient Survey. This dataset provides a baseline to facilitate standardized quantitative comparison of astronomical transient event recognition algorithms. The classes included in the dataset are: supernovae, cataclysmic variable, active galactic nuclei, high proper motion stars, blazars and flares. As an example on how to use the dataset we experiment with multiple data pre-processing methods, feature selection techniques and classic machine learning algorithms (Support Vector Machines, Random Forests and Neural Networks). We assess quantitative performance in two classification tasks: binary (transient/non-transient) and eight-class classification. The best performing algorithm is a Random Forest Classifier for both classification experiments. The next release of this reference database will include images and benchmarks with deep learning models. All our code and data is available to the community at [www.github.com/XXX/YYY](http://www.github.com/XXX/YYY).

**Key words:** methods: data analysis, statistical

## 1 INTRODUCTION

The study and detection of astronomical variable sources is expected to occur on unprecedented scales with the new generation of forthcoming multi-epoch and multi-band (synoptic) astronomical surveys. For instance, projects like the Large Synoptic Survey Telescope (LSST) (Ivezić et al. 2008; Jurić et al. 2015) are expected to generate exuberant daily data-streams.

One of the main challenges that these datasets want to address is Real-Time Transient classification, i.e. flagging astrophysically relevant events whose luminosity varies in short duration relative in the timescale of the universe, from minutes to several years. Transients include phenomena such as supernovae, novae, neutron stars, blazars, pulsars, cataclysmic variable stars (CV), gamma ray bursts (GRB) and active galaxy nuclei (AGN). The time-domain dependency of these objects is one of the reasons why they are hard to classify: their data is usually heterogeneous, unbalanced, sparse, unevenly sampled and with missing information. This has motivated the use of Machine Learning (ML) algorithms to recognize and classify transient events.

There have been successful attempts to implement these algorithms using images as an input. For instance, data from the SkyMapper Supernova and Transient Survey and the High cadence Transient Survey (HiTS) have been used as inputs to automatic detection algorithms (Gieseke et al. 2017; Cabrera-Vives et al. 2017). Convolutional Neural Networks (CNN) have also achieved high accuracy in this binary classification task. Other studies have shown that artifacts can be detected using features extracted from raw images. Klencki et al. (2016) achieved reliable classification by transforming transient data from the OGLE-IV data-reduction pipeline and training it with machine learning algorithms such as Artificial Neural Networks, Support Vector Machines, Random Forests, Naive Bayes, K-Nearest Neighbors and Linear Discriminant Analysis. Similarly, Wright et al. (2015) used images from Pan-STARRS1 Medium Deep Surveys, and du Buisson et al. (2015) processed single-epoch multi-band images from the SDSS supernova survey for the same purpose.

A complementary approach uses the light curves computed from the images to perform the classification task. For instance D’Isanto et al. (2016) used classical machine learning algorithms such as Random Forest, MultiLayer Per-

ceptron and K-Nearest Neighbours light curves to classify transients from the Catalina Real Time Transient Survey; Lochner et al. (2016) used the same approach to find where supernovas from the Supernova Photometric Classification Challenge.

A crucial element in the developmet of any ML algorithm is the compilation of the training dataset. In astronomy this task has been facilitated by the publication of large databases from different observational projects. However, the publications that make use of these datasets still make extensive use of the internal know-how of the scientific collaboration. It is still difficult for another scientists to rebuild a training dataset and perform comparisons with published results.

To address this issue we compile and publish in easy-to-access files a dataset that can be used to train and test different ML algorithms for transient detection. We use public data from the Catalina Real-Time Transient Survey (CRTS) (Drake et al. 2012), an astronomical survey searching transient and highly variable objects as base for the dataset. In this paper we present light-curve data, in a future publication we will present an imaging dataset.

In Section ?? we present the CRTS and the steps we follow to build the curated dataset. In Section ?? we describe the main features of the dataset and the repository structure gathering the files and code useful to explore it. In Section ?? we show how this dataset can be used to perform some tests using some classical ML tests. We follow a similar approach as D’Isanto et al. (2016). We finalize in Section ?? with a summary of main features of our dataset and the results of our ML tests.

## 2 THE LIGHTCURVE DATASET

We use public data from the Catalina Real-Time Transient Survey (CRTS) (Drake et al. 2012), an astronomical survey searching transient and highly variable objects. It covered 33000 squared degrees of sky and took data since 2007. Three telescopes were used: Mt. Lemmon Survey (MLS), Catalina Sky Survey (CSS), and Siding Spring Survey (SSS). So far, CRTS has discovered more than 15000 transient events. We use light curves as measured with the CSS telescope of the CRTS, which is an f/1.8 Schmidt telescope located in the Santa Catalina Mountains, north of Tucson, Arizona and is equipped with a 111-megapixel detector, and covered 4000 square degrees per night, with a limiting magnitude of 19.5 in the V band. The public CRTS data base reports the source flux and its corresponding uncertainty (Stetson 1996).

All transient objects were classified in the CRTS dataset according to their type. The most relevant classes are: supernovae (SN), cataclysmic variable stars (CV), blazars, flares, asteroids, active galactic nuclei (AGN), and high-proper-motion stars (HPM). Though most objects in the transient object catalogue belong to a single class, there is some uncertainty in the categorization of some of them. In this case, an interrogation sign is used when a class is not clear e.g. SN? or sometimes multiple possible classes are found for a single event e.g. SN/CV. Table 1 summarized the number of objects in each class.

We use the light curves of 4269 unique transient events that contain at least 5 observations. We also use 15193 non-

Class	Object Count
SN	1293
CV	862
AGN	425
HPM	306
Blazar	237
SN?	236
Flare	207
AGN?	130
Unknown	114
CV?	55

**Table 1.** Top 10 transient classes, with their respective event count.

transient sources with at least 5 observations each. We obtained sources from the transient dataset directly from researchers of the CRTS project. Alternatively, we retrieved sources in the dataset from the CRTS online catalogue, by sampling light curves of objects within a 0.006 degree radius from CRTS detected transients, and removing known transient light curves from that set. Though this process should return only non-transient sources, it is possible that non-detected transients were captured and catalogued incorrectly as ‘non-transients’. Table 2 and Table 3 summarize some statistics on the number of observations available for each light curve for the transient and non-transient datasets, respectively.

## 3 REPOSITORY DESCRIPTION

The repository that contains the code, data and trained models we used in this project is found in the website <https://github.com/diegoalejogm/crts-transient-recognition>. For this project we used various python libraries: `jupyter` (1.0.0), `numpy` (1.14.3), `pandas` (0.22.0) and `scikit-learn` (0.19.1). Details other dependencies can be found in the included file `requirements.txt`.

We organized the repository in three main folders. The first one is `data`, and it contains the data we used in this project. We split data in three directories: `lightcurves`, `features` and `inputs`. The former contains pandas dataframes for raw and curated light-curves. The `features` subfolder contains the resulting features of pre-processing the light-curves, for both transients and non-transients. Finally, the `inputs` subfolder contains the already split training and testing inputs for the following hyper-parameter combinations: classification task, number of features, minimum number of observations, and balanced/unbalanced data. Each one of the input files contains a numpy tuple with the structure: *(train features, train labels, test features, test labels)*

The second main folder in the repository is `notebooks`. It contains all the code used for this project. The jupyter notebooks contain the pipeline of our proposed experimental framework, and they are numbered in sequential order. Moreover, each notebook is documented within itself. Two additional non-numbered notebooks are used for exploration purposes. Conversely, the `.py` files are used to contain the

Count	4269
Mean	102.810
Std	113.786
Min	5
25%	14
50%	48
75%	166
Max	564

**Table 2.** Transient’s data-set observation count per object.

Count	15193
Mean	118.37
Std	116.51
Min	5
25%	26
50%	72
75%	185
Max	537

**Table 3.** Non-Transient’s data-set observation count per object.

python code which is used in the notebooks. They’re separated by task type.

Finally, the **results** directory contains the products of running the code. One result type is named dataframes, which contains a pandas dataframe for each task. Each dataframe contains the testing and training results, including scores and confusion matrices. On the other hand, the results sub-folder contains the feature importance list figures presented in this document.

For more information on the repository, make sure to read the README file included, or email us directly.

The tables 2 and 3 in this annex contain statistical descriptors of the light curves’ observation count, for transient and non-transient curves that contain at least 5 observations. Table 2 makes reference to transient events, whereas table 3 refers to non-transient objects. Each of these tables present the mean, standard deviation, percentiles (25, 50 and 75) and maximum number of observations per light curve found in each dataset.

## 4 CLASSICAL MACHINE LEARNING TESTS

As an example on how the dataset can be used, we apply classical machine learning algorithms to perform different classification tasks.

### 4.1 Data Preprocessing

We do not feed directly the anotated lightcurves to the ML algorithms. There is preprocessing stage that follows six steps.

#### 4.1.1 Data Filtering

We discard light curves with few observations as they may not contain enough information to be classified correctly. Our nominal cut is 10 observations per light curve.

#### 4.1.2 Oversampling Transient Light Curves

The number of light curves per class is unbalanced. In order to have the same amount elements for each class we implement an oversampling step by artificially generating multiple mock light curves, each based on an observed one.

We generate a mock light curve from the observed light curve and then sample the observed magnitude from a Gaussian distribution centered on the observational apparent magnitude with the magnitude’s error as the standard deviation.

#### 4.1.3 Feature Extraction

Light curves are sampled at irregular time intervals and have different numbers of data points. Thus, it is challenging to directly use the time-series data for classification with traditional methods. We circumvent this problem by extracting a set of features for each light curve. These features are scalars derive from statistical and model-specific fitting techniques. The first features (moment-based, magnitude-based and percentile-based) were formally introduced in Richards et al. (2011), and have been used in other studies (Lochner et al. 2016; D’Isanto et al. 2016). We extend that list to include another set (polynomial fitting-based features. These groups are:

(i) Moment-based features, which use the magnitude for each light curve.

- Beyond1std (*beyond1std*): Percentage of observations which are over or under one standard deviation from the weighted average. Each weight is calculated as the inverse of the corresponding observation’s photometric error.
- Kurtosis (*kurtosis*): The fourth moment of the data distribution. Used to measure the heaviness or lightness in the tails of the statistical data.
- Skewness (*skew*): A measurement of the level of asymmetry from the normal distribution in a data distribution. Negative skewness is the property of a more pronounced left tail, while positive skewness is a characteristic that implies a more pronounced right tail.
- Small Kurtosis (*sk*): Small sample kurtosis.
- Standard deviation (*std*): The standard deviation of the magnitudes.
- Stetson J (*stetson\_j*): The Welch-Stetson J variability index Stetson (1996). A robust standard deviation.
- Stetson K (*stetson\_k*): The Welch-Stetson K variability index Stetson (1996). A robust kurtosis measure.

(ii) Magnitude-based features, which rely on the magnitude for each source.

- Amplitude (*amp*): The difference between the maximum and minimum magnitudes.
- Max Slope (*max\_slope*): Maximum absolute slope between two consecutive observations.
- Median Absolute Deviation (*mad*): The median of the difference between magnitudes and the median magnitude.
- Median Buffer Range Percentage (*mbrp*): The percentage of points within 10% of the median magnitude.
- Pair Slope Trend (*pst*): Percentage of all pairs of

consecutive magnitude measurements that have positive slope.

- Pair Slope Trend 30 (*pst\_last30*): Percentage of the last 30 pairs of consecutive magnitudes that have a positive slope, minus percentage of the last 30 pairs of consecutive magnitudes with a negative slope.

(iii) Percentile-based features, which use the sorted flux distribution for each source. The flux is computed as  $F = 10^{0.4\text{mag}}$ . We define  $F_{n,m}$  as the difference between the  $m$ -th and  $n$ -th flux percentiles.

- Percent Amplitude (*p\_amp*): Largest percentage difference between the absolute maximum magnitude and the median.
- Percent Difference Flux Percentile (*pdfp*): Ratio between  $F_{5,95}$  and the median flux.
- Flux Percentile Ratio Mid20 (*fpr20*): Ratio  $F_{40,60}/F_{5,95}$
- Flux Percentile Ratio Mid35 (*fpr35*): Ratio  $F_{32.5,67.5}/F_{5,95}$
- Flux Percentile Ratio Mid50 (*fpr50*): Ratio  $F_{25,75}/F_{5,95}$
- Flux Percentile Ratio Mid65 (*fpr65*): Ratio  $F_{17.5,82.5}/F_{5,95}$
- Flux Percentile Ratio Mid80 (*fpr80*): Ratio  $F_{10,90}/F_{5,95}$

(iv) Polynomial Fitting-based features, which are the coefficients of multi-level terms in a polynomial curve fitting. This is new set of features proposed in this paper.

- Poly1 T1: Linear term coeff. in monomial curve fitting.
- Poly2 T1: Linear term coeff. in quadratic curve fitting.
- Poly2 T2: Quadratic term coeff. in quadratic curve fitting.
- Poly3 T1: Linear term coeff. in cubic curve fitting.
- Poly3 T2: Quadratic term coeff. in cubic curve fitting.
- Poly3 T3: Cubic term coeff. in cubic curve fitting.
- Poly4 T1: Linear term coeff. in quartic curve fitting.
- Poly4 T2: Quadratic term coeff. in quartic curve fitting.
- Poly4 T3: Cubic term coeff. in quartic curve fitting.
- Poly4 T4: Quartic term coeff. in quartic curve fitting.

#### 4.1.4 Feature Scaling

We re-scale the magnitudes to have zero mean and unit variance.

## 4.2 Classification Tasks

We study two classification tasks.

### 4.2.1 Binary Classification

We use a balanced number of events from both classes in order to investigate the capability of distinguishing between Transients and Non-Transients.

### 4.2.2 8-Class Classification

We consider the unbalanced number of objects across classes to perform a classification into the following categories: AGN, Blazar, CV, Flare, HPM, Other, Non-Transient and Supernovae.

## 4.3 ML algorithms

We conduct experiments with three widely used families of supervised classification algorithms: Neural Networks (NNs), Random Forests (RFs) and Support Vector Machines (SVMs).

These algorithms are popular in published studies and are efficient for low dimensional feature datasets as is our case. We use sklearn (Pedregosa et al. 2011) Python's implementation of these algorithms. Details on the inner workings of these machine learning models can be found in Hastie et al. (2016).

The set of hyperparameter space explored for each algorithm is the following.

For Neural Networks:

- Learning Rate: Either constant vs adaptive.
- Hidden Layer Sizes: Single Layer with 100 nodes vs Two layers with 100 nodes each.
- L2 Penalty ( $\alpha$ ):  $10^{-1}$ ,  $10^{-2}$ ,  $10^{-3}$ ,  $10^{-4}$ .
- Activation Function: Logistic vs Relu.

For Random Forest:

- Number of Estimators: 200 or 700.
- Number of features Considered: Square Root or the Logarithm base 2 of the total number of features.

For Support Vector Machines:

- Kernel: Radial Basis Function (RBF).
- Kernel Coefficient ( $\gamma$ ):  $10^{-1}$ ,  $10^{-2}$ ,  $10^{-3}$ ,  $10^{-4}$ ,  $10^{-5}$
- Error Penalty ( $C$ ): 1 vs 10 vs 100 vs 1000.

## 4.4 Validation

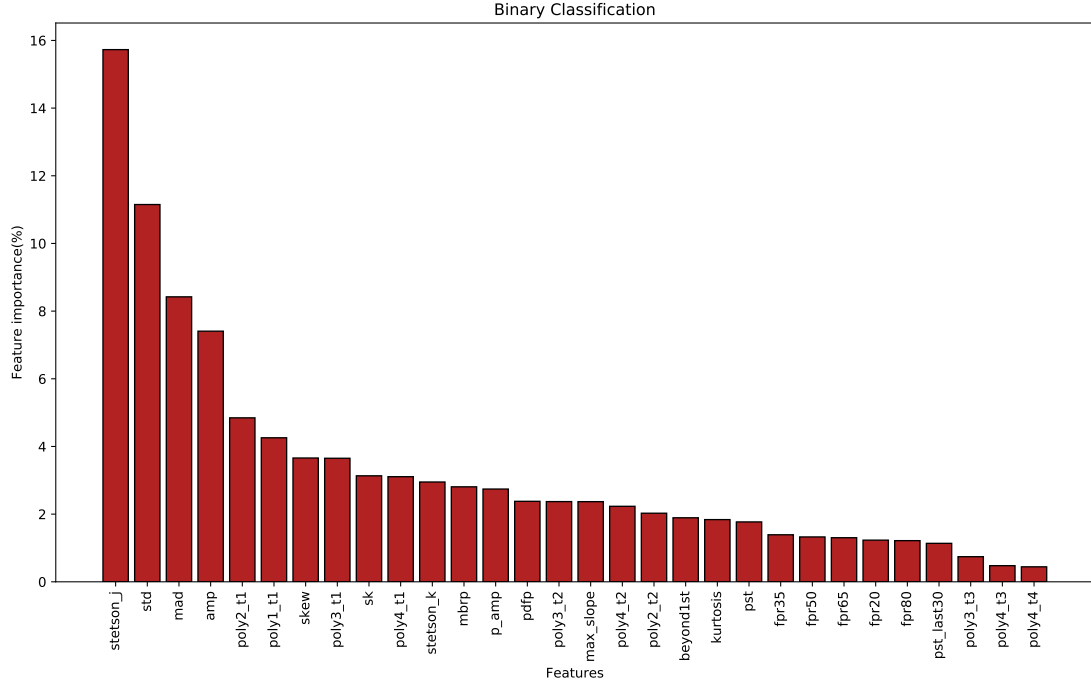
We split the input light curves in a training and test datasets. The test dataset contains only original light curves, without any oversampling. We use 2-fold cross-validation during training as evaluation protocol. Moreover, we use grid search during training to test multiple hyper-parameter configurations for each one of the possible algorithms. We use the F1-Score to assess the performance of a given model and we evaluate each task on the held-out test dataset.

## 4.5 Results

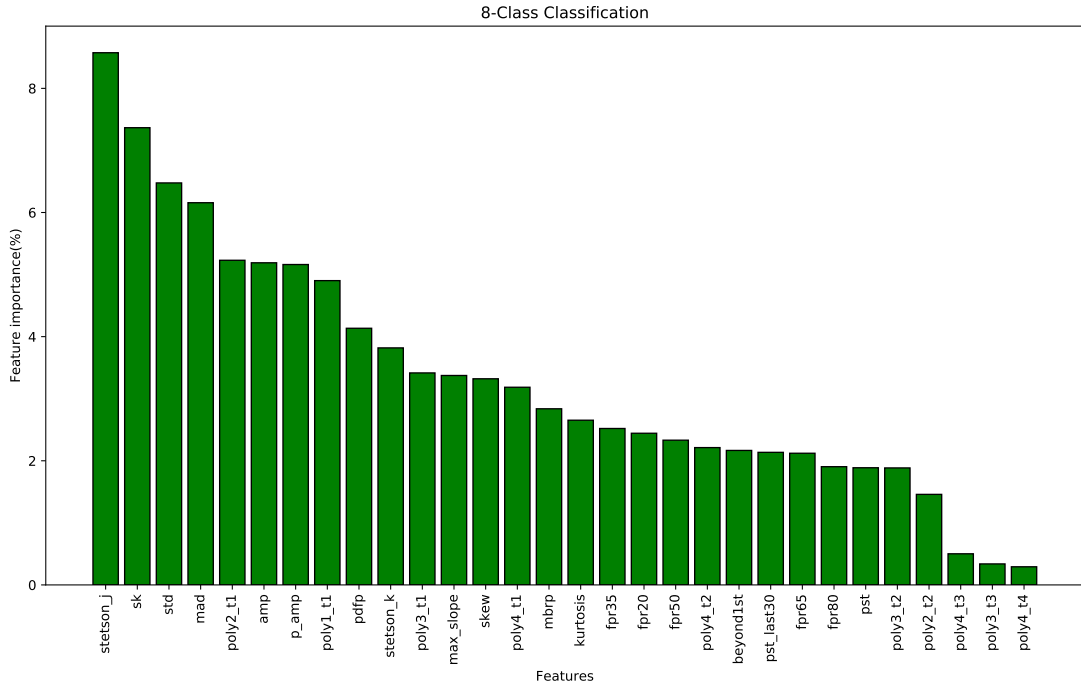
### 4.5.1 Binary Classification

The best algorithm in this task is RFs with a maximum F1-score of 87.69%. SVMs are the second best-performing model with the F1-score of 85.36%. Changing the number of features does not affect significantly the score. NNs are ranked third, although their scores are very similar to those of SVMs. The highest achieved score for NNs is 85.03%.

Table 5 shows the confusion matrix of the best performing algorithm and Table 4 summarizes the scores. These results imply that non-transients are better classified overall.



**Figure 1.** Feature importance rank for the best Random Forest classifier for the Binary classification task. Feature importance is represented with percentages.



**Figure 2.** Feature importance rank for the best Random Forest classifier for the best 8-Class classification task. Feature importance is represented with percentages.



	Precision	Recall	f1-Score	Support
Non-Transient	90.38	93.83	92.08	3798
Transient	93.59	90.02	91.77	3798

**Table 4.** Precision, Recall and f1-score for the Binary Classification Task with Regular inputs.

	non-transient	transient
non-transient	3564	379
transient	234	3419

**Table 5.** Confusion Matrix for the best performing model in the Binary task.

Figure 1 displays the most important features for the RFs classifier. The top five inputs for classification are *stetson\_j*, *std*, *mad*, *poly1\_t1* and *poly2\_t1*. The former achieved the highest importance with over 21%, compared to the following with values in the range 6% - 8%.

#### 4.5.2 Eight-Class Classification

RFs are the best classifier. The best f1-score is 66.05%. NNs are the second best. Its highest f1-score is 60.19%, while SVMs are the worst-performing model only achieving a maximum f1-score of 57.30%. Table ?? summarizes the results.

Table 7 presents the confusion matrix for the RF. The two classes with highest recall are HPM and Non-Transient, with a recall of 86.36% and 84.13%, respectively. The worst performing classes are Blazar, Flare and Other, with recall values in the range 36% - 40%. SN is the class with which most other class instances are incorrectly classified. Moreover, Flares have about 50% of the test samples classified as Non-Transients, AGNs have about 20% of their samples classified as Other, and Blazars and Other had most of its samples classified as AGN. Additionally, most incorrectly classified AGNs (~20.5%) are identified as Other and most Blazar instances are incorrectly categorized as either SN or AGN.

Figure 2 displays the feature importance ranking. This list ranks first *stetson\_j* with an 8% importance, followed by *amp*, *sk*, *std*, *mad*, with values around 6%. The lowest ranking features are the five high level polynomial: *poly4\_t1*, *poly4\_t2*, *poly3\_t3*, *poly4\_t3* and *poly4\_t4*.

## 5 CONCLUSIONS

The scope of forthcoming of large astronomical synoptic surveys such as the LSST (Ivezić et al. 2008) motivates the development and exploration of automatized ways to detect transient sources. In this paper we presented an approach for the automatic recognition of transient events with machine learning techniques.

The method we present is based on the study of light curves. We extracted its characteristic features to use them as inputs to train three different machine learning algorithms: Random Forests, Neural Networks and Support Vec-

	Precision	Recall	f1-Score	Support
SN	52.45	53.29	52.87	561
CV	71.66	68.98	70.29	561
AGN	71.81	75.40	73.56	561
HPM	90.21	88.77	89.48	561
Blazar	67.28	51.69	58.46	561
Flare	72.05	69.87	70.95	561
Other	45.91	45.09	45.50	561
Non-Transient	63.57	80.57	71.06	561
avg/total	66.87	66.71	66.52	4488

**Table 6.** Precision, Recall and f1-score for the 8-Class Classification Task with Regular inputs.

tor Machines. The features extracted from light curves were either statistical descriptors of the observations, or polynomial curve fitting coefficients applied to the light curves.

The machine learning algorithms performed two classification tasks, and multi-class classification of various transient classes including non-transients too. Overall, the best classifier for all tasks was the Random Forest. We studied the feature importance for this model. The most important feature was always *stetson\_j*. The proposed coefficients corresponding to the linear terms of the quadratic and monomial curve fitting are also useful in the classification task.

We provide the code and the datasets that were used in this project. The repository containing all this data may be found in the website <https://github.com/diegoalejogm/crts-transient-recognitionSection>.

In a continuation of this project we will present in a second paper a reference dataset for astronomical transient event recognition based on images of the CATALINA survey. We will preset tests using state-of-the art deep learning techniques for transient classification. lightcurves and tests on classical machine learning algorithms

## ACKNOWLEDGEMENTS

We thank Andrew Drake for sharing with us the CRTS Transient dataset used in this project. We thank Juan Pablo Reyes, Dominique Fouchez for helping with the research. We acknowledge funding from Universidad de los Andes in the call for project finalization. We also thank contributors and collaborators of the SciKit-Learn, Jupiter Notebooks and Pandas' Python libraries.

CRTS and CSDR2 are supported by the U.S. National Science Foundation under grant NSF grants AST-1313422, AST-1413600, and AST-1518308. The CSS survey is funded by the National Aeronautics and Space Administration under Grant No. NNG05GF22G issued through the Science Mission Directorate Near-Earth Objects Observations Program.

	SN	CV	AGN	HPM	Blazar	Flare	Other	Non-Transient
SN	299	64	7	0	89	21	75	15
CV	42	387	0	0	38	18	53	2
AGN	14	16	423	3	58	4	65	6
HPM	5	5	0	498	0	4	10	30
Blazar	24	27	44	0	290	9	37	0
Flare	43	8	18	1	13	392	35	34
Other	84	46	48	3	71	24	253	22
Non-Transient	50	8	21	56	2	89	33	452

**Table 7.** Confusion Matrix for the best performing model in the 8-Class task.

## REFERENCES

- Cabrera-Vives G., Reyes I., Förster F., Estévez P. A., Maureira J.-C., 2017, *ApJ*, **836**, 97
- D’Isanto A., Cavuoti S., Brescia M., Donalek C., Longo G., Riccio G., Djorgovski S. G., 2016, *MNRAS*, **457**, 3119
- Drake A. J., et al., 2012, in Griffin E., Hanisch R., Seaman R., eds, IAU Symposium Vol. 285, New Horizons in Time Domain Astronomy. pp 306–308 ([arXiv:1111.2566](#)), [doi:10.1017/S1743921312000889](#)
- Gieseke F., et al., 2017, *MNRAS*, **472**, 3101
- Hastie T., Tibshirani R., Friedman J., 2016, The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition (Springer Series in Statistics). Springer
- Ivezić Ž., et al., 2008, preprint, ([arXiv:0805.2366](#))
- Jurić M., et al., 2015, preprint, ([arXiv:1512.07914](#))
- Klencki J., Wyrzykowski Ł., Kostrzewa-Rutkowska Z., Udalski A., 2016, *Acta Astron.*, **66**, 15
- Lochner M., McEwen J. D., Peiris H. V., Lahav O., Winter M. K., 2016, *ApJS*, **225**, 31
- Pedregosa F., et al., 2011, *Journal of machine learning research*, **12**, 2825
- Richards J. W., et al., 2011, *ApJ*, **733**, 10
- Stetson P. B., 1996, *pasp*, **108**, 851
- Wright D. E., et al., 2015, *MNRAS*, **449**, 451
- du Buisson L., Sivanandam N., Bassett B. A., Smith M., 2015, *MNRAS*, **454**, 2026