

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

Mauricio Neira¹, Catalina Gómez², Diego A. Gómez¹, Marcela Hernández Hoyos¹, Jaime E. Forero-Romero³, Pablo Arbeláez²

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

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

²*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 ATRANCCATA (Annotated TRANSient Catalina CATALog) 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 variableS, 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 ATRANCCA will include images and benchmarks with deep learning models. All our code and data is available to the community at <https://github.com/MachineLearningUniandes/ATRANCCATA>.

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 near to 20 petabytes every night.

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 com-

puted 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 Perceptron 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 dataset. In Section ?? we describe its main features together with the repository structure gathering the files and some Python code to explore it. In Section ?? we show how this dataset can be used to perform some tests using some classical ML tests 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 magnitudes and its corresponding uncertainty (Stetson 1996).

The CRTS dataset already provides a classification. 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 pos-

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.

sible classes are found for a single event e.g. SN/CV. Table 1 summarizes the number of objects in each class.

We obtain the lightcurves from the transient dataset directly from the CRTS project. To compile the non-transient dataset we retrieve 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. Our dataset includes all lightcurves with at least one observation.

3 REPOSITORY DESCRIPTION

The repository that contains the light-curves and some code to process the data can be found in <https://github.com/MachineLearningUniandes/ATRANCCATA>.

The repository has three main folders

- **data:** Contains all the lightcurves and labels.
 - **lightcurves:** raw light-curves as csv files.
- **nb-explore:** jupyter notebooks with examples of simple light-curve exploration.
- **ML:** Contains the data sets and notebooks to perform the ML tests. It is divided into three different directories:
 - **notebooks.** Contains all the code for the ML tests. There are three numbered jupyter notebooks summarizing our experimental framework. The documentation is included into each notebook. There are various .py files with the required python routines.
 - **features:** resulting features of pre-processing the light curves for both transients and non-transients.
 - **inputs** data-set split into training/test for the different machine learning classification tasks. Data is split into different meta-parameters: 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*)
 - **results.** Products of the ML test. 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.

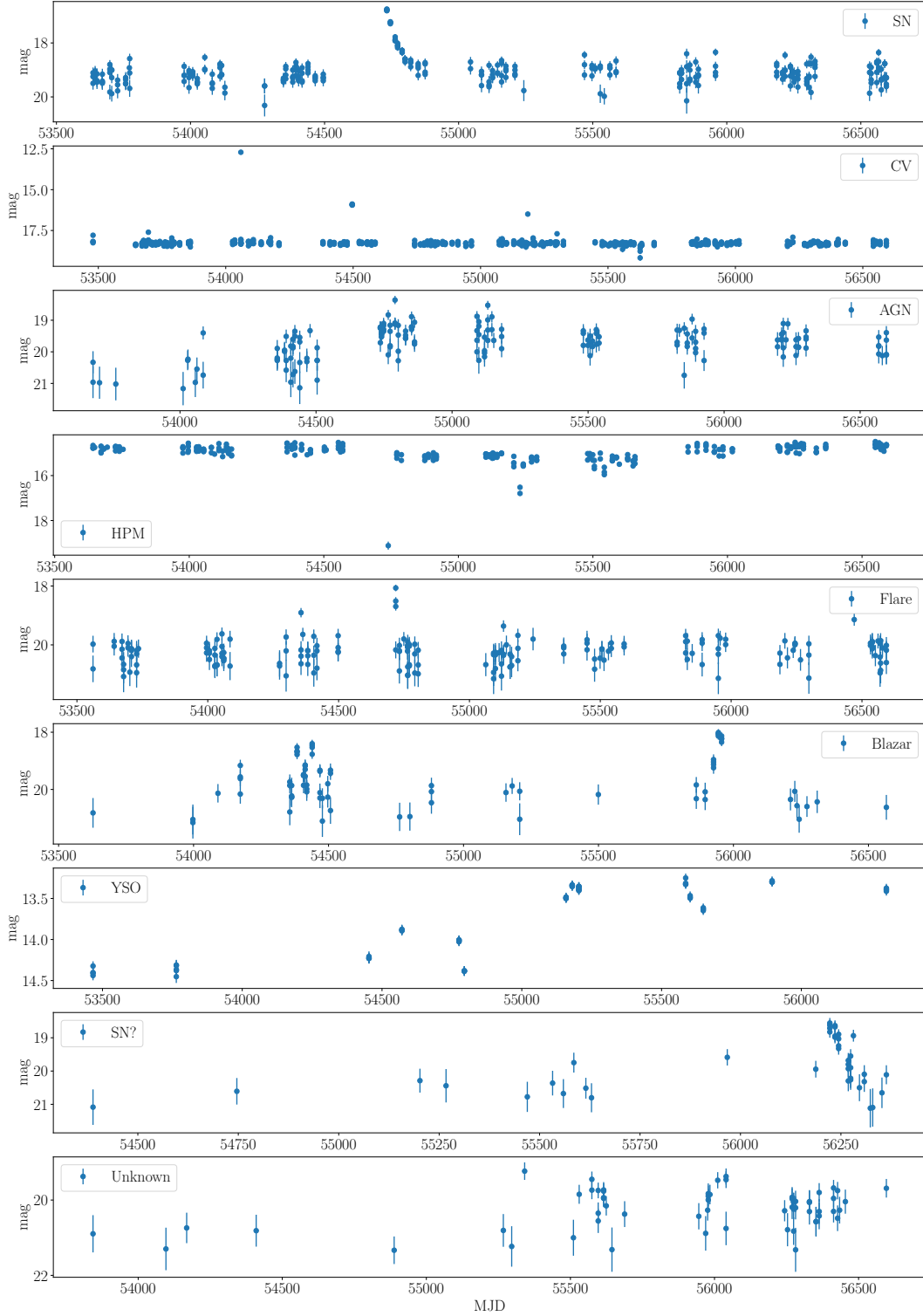


Figure 1. Randomly selected lightcurves for the most represented classes as compiled in the database we present in this paper.

4 CLASSICAL MACHINE LEARNING TESTS

4.1 Data Preprocessing

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

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

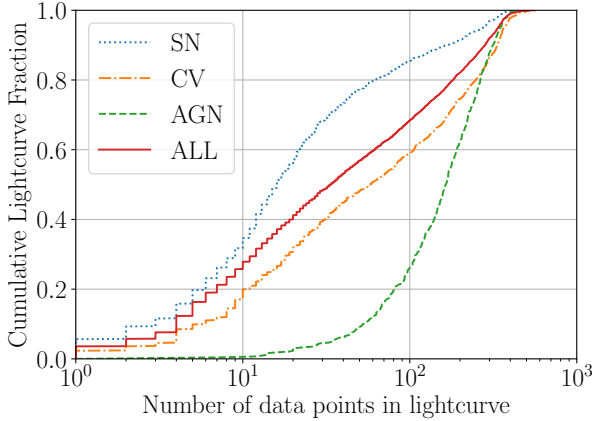


Figure 2. Cumulative number of data curves (expressed as a fraction) sampled with at least N data points. This includes information for the three most representative classes (SN, CV, AGN) and the whole database (ALL).

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 of 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 derived 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 a 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 (α): 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:

	Precision	Recall	f1-Score	Support
Non-Transient	94.13	94.13	94.13	3798
Transient	79.10	79.10	79.10	1067

Table 2. Precision, Recall and f1-score for the Binary Classification Task with Regular inputs.

	non-transient	transient
non-transient	3575	223
transient	223	844

Table 3. Confusion Matrix for the best performing model in the Binary task.

- Kernel: Radial Basis Function (RBF).
- Kernel Coefficient (γ): 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 3 shows the confusion matrix of the best performing algorithm and Table 2 summarizes the scores. These results imply that non-transients are better classified overall.

Figure 3 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 5 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

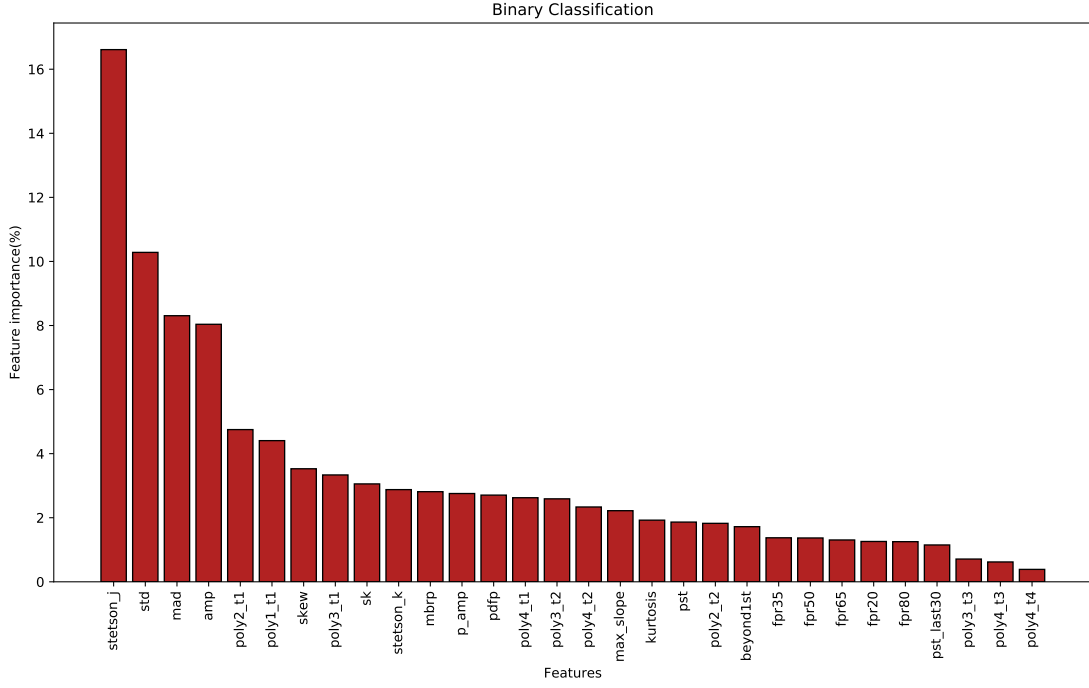


Figure 3. Feature importance rank for the best Random Forest classifier for the Binary classification task. Feature importance is represented with percentages.

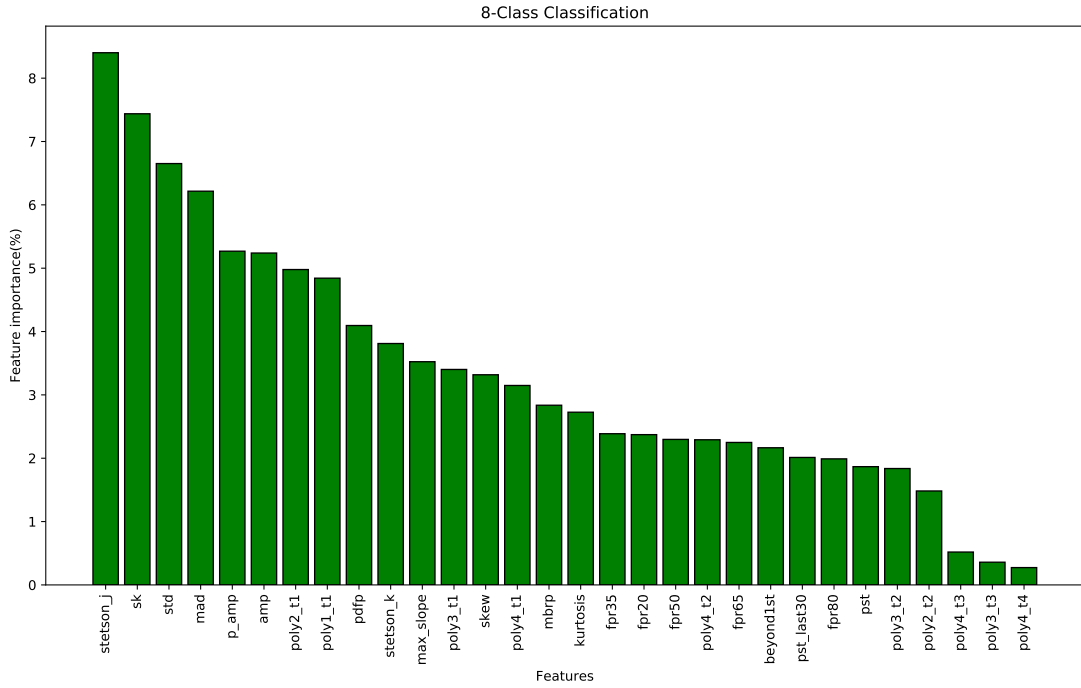


Figure 4. Feature importance rank for the best Random Forest classifier for the best 8-Class classification task. Feature importance is represented with percentages.

	Precision	Recall	f1-Score	Support
SN	48.82	51.39	50.07	323
CV	66.96	70.69	68.77	215
AGN	48.14	85.84	61.69	106
HPM	25.19	86.84	39.05	76
Blazar	46.77	49.15	47.93	59
Flare	7.00	41.17	11.96	51
Other	31.11	44.01	36.46	234
Non-transient	96.06	79.69	87.12	3798
avg/total	46.25	63.59	50.38	4862

Table 4. Precision, Recall and f1-score for the 8-Class Classification Task with Regular inputs.

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 4 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 motivates the development and exploration of automatized ways to detect transient sources. Such an effort requires the compilation of publicly available databases to train and test new algorithms. In this paper we present the results of such compilation based on data from the Catalina Real-Time Transient Survey (CRTS). The dataset we presented in this paper summarizes more than 10000 light curves for six different classes of transients. The dataset is publicly available at <https://github.com/MachineLearningUniandes/AstroTransientReference>.

We illustrated how to use this database by extracting characteristic features to use them as inputs to train three different machine learning algorithms (Random Forests, Neural Networks and Support Vector Machines) for classification tasks. The features extracted from light curves were either statistical descriptors of the observations, or polynomial curve fitting coefficients applied to the light curves. Overall the best classifier for all tasks was the Random Forest. In this model the most important feature was always `stetson_j`.

In a second paper we will present another reference dataset for astronomical transient event recognition based

on images of the CATALINA survey. The corresponding tests will use state-of-the-art deep learning techniques for transient classification.

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.

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](https://arxiv.org/abs/1111.2566)), [doi:10.1017/S1743921312000889](https://doi.org/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](https://arxiv.org/abs/0805.2366))
- Jurić M., et al., 2015, preprint, ([arXiv:1512.07914](https://arxiv.org/abs/1512.07914))
- Klencki J., Wyrzykowski L., 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

	SN	CV	AGN	HPM	Blazar	Flare	Other	Non-Transient
SN	166	25	0	0	7	5	40	97
CV	17	152	0	1	5	3	12	37
AGN	1	2	91	0	10	1	35	49
HPM	5	0	0	66	0	0	5	186
Blazar	8	13	4	0	29	0	6	2
Flare	16	5	0	0	3	21	4	251
Other	53	12	7	1	3	3	103	149
Non-Transient	57	6	4	8	2	18	29	3027

Table 5. Confusion Matrix for the best performing model in the 8-Class task.