

UNIVERSIDAD DE LOS ANDES

UNDERGRADUATE THESIS PROPOSAL

---

**Transient object classification using machine  
learning and deep learning techniques on real  
data**

---

*Author:*  
Mauricio Neira

*Supervisor:*  
Marcela Hernández Hoyos,  
Ph.D.

February 7, 2019

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Astronomy . . . . .	1
1.1.1	Light curves . . . . .	1
1.1.2	Supernovae . . . . .	2
1.2	LSST . . . . .	2
<b>2</b>	<b>Problem description</b>	<b>3</b>
2.1	Classification of transient objects from light curves . . . . .	3
2.2	Classification of transient objects from images . . . . .	4
<b>3</b>	<b>Project Background</b>	<b>4</b>
3.1	Diego's thesis - 2018-10 . . . . .	4
3.2	research internship?? No se como ponerle a esto 2018-2 . . . . .	6
3.2.1	Improvement on Diego's work . . . . .	6
3.2.2	PLAsTiCC - Kaggle competition . . . . .	6
<b>4</b>	<b>General objective</b>	<b>6</b>
<b>5</b>	<b>Specific objectives</b>	<b>6</b>
5.1	Improvement of the light curve feature space . . . . .	6
5.1.1	Feature pruning . . . . .	6
5.1.2	Addition of supernovae specific metrics . . . . .	7
5.2	Deep learning on light curves . . . . .	7
5.3	Deep learning on images . . . . .	8
5.3.1	CNNs . . . . .	8
5.3.2	CNNs and RNNs . . . . .	9
<b>6</b>	<b>Activities and schedule</b>	<b>9</b>
<b>7</b>	<b>Expected results</b>	<b>9</b>

# 1 Introduction

This project focuses on automatic transient object detection. This section will briefly outline the basic concepts needed to understand the motivation and reasoning behind it.

## 1.1 Astronomy

Astronomy is the field of science that studies celestial objects i.e. stars, planets, matter, etc. Outside of earth. These objects emit electromagnetic radiation all across the electromagnetic spectrum. Measurements of this radiation have been made with a wide variety of instruments and a lot of data has been recorded.

Usually, the data is recorded in *surveys* where various telescopes scan the sky through many days. One of these surveys is the subject of this paper's investigation known as the “Catalina Real Time Survey” [1]. In this database, the intensity of radiation of various objects was recorded at different time intervals. The group of points collected for a particular object is known as a light curve.

### 1.1.1 Light curves

After recording the intensity of emitted radiation of an object several times, it is possible to plot the result in a scatterplot. The following image is an example of a light curve of a type 1a supernova:

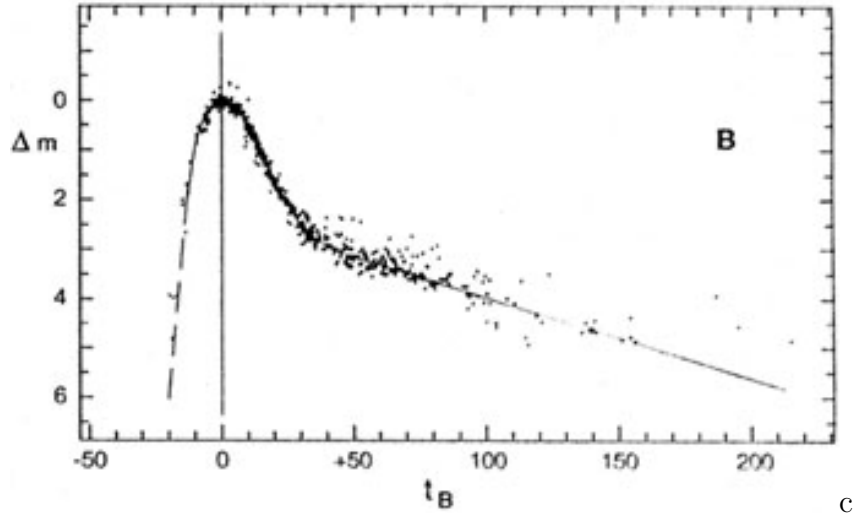


Figure 1: Light curve for a supernova type 1a [2].

Light curves hold a lot of information about the behavior of the object and can be used to identify the type of object [3, 4]. A big part of this project is dedicated to classifying objects based on their light curve.

### 1.1.2 Supernovae

Supernovae, specially type 1a supernovae, are of immense importance in present day astrophysics. They are used as “standard candles”<sup>1</sup> to estimate the rate of expansion of the universe [5]. So far, they have been used to pin down the cosmological constant (also known as Hubble’s constant) to demonstrate the accelerating expansion of the universe; “implying the existence of a nearly uniform component of dark energy” [6].

Finding and studying supernovae not only sheds light into their nature, it has implications that stretch out to the inner workings of the universe.

## 1.2 LSST

The Large Synoptic Survey Telescope is a project that will begin operating in 2020 and has four main goals[7]:

---

<sup>1</sup>A standard candle is a celestial object whose properties are optimal to measure large distances in the universe with high precision.

1. Probing dark energy and dark matter
2. Taking an inventory of the Solar System
3. Exploring the transient optical sky
4. Mapping the Milky Way

Unlike other telescopes of its nature, from the moment it begins operation, the LSST will release its data to the public. It plans on gathering 15 terabytes of data every night, generating alerts and transferring data with a 60 second delay [7].

With that amount of data, there simply aren't enough experts to classify all of the objects of interest. New automated techniques are needed to correctly classify all of the observed objects and notify the corresponding interested scientists.

Consequently, the motive behind this project is to provide an approximated solution to the problem outlined above.

## 2 Problem description

As stated in the previous section, the motivation behind this project is to create an algorithm that automates the classification of astronomical objects. The classification can be done directly on the images released by the CSS [1] or on the light curves extracted from the images.

### 2.1 Classification of transient objects from light curves

One of the main aims to develop an algorithm that correctly maps the light curve of an object to its corresponding category. It is unrealistic to claim that the function will correctly classify every single light curve it is provided. Nevertheless, it does aim to minimize classification errors.

There will be 2 main approaches to achieve this goal. The first will be to transfer the light curves into a feature space that attempts to extract many geometric properties of the light curves. After that, several classifiers will

be applied over the data points in the feature space. These include random forests and feed forward neural networks.

The second approach involves feeding the points to a recurrent neural network that will in turn be fed into a feed forward neural network. Both of these techniques will be explained in further detail in the upcoming sections.

## 2.2 Classification of transient objects from images

The second main aim of this proposed thesis is to develop an algorithm that correctly classifies astronomical objects from a series of *images*.

# 3 Project Background

## 3.1 Diego’s thesis - 2018-10

In the first semester of 2018, Diego Alejandro Gómez Mosquera worked on “Astronomical transient event recognition with machine learning” using random forests on a feature space calculated from light curves [8]. The author focused on creating a feature space robust enough to distinguish the different transient classes. This was achieved primarily through geometric parameters that were extracted from the light curves. The features used (as found on the author’s thesis) were:

- skew: Skewness.
- kurtosis: Kurtosis.
- small kurtosis: Small sample kurtosis.
- std: Standard deviation.
- beyond1std: Percentage of magnitudes beyond one standard deviation from the weighted mean. Each weights is calculated as the inverse of the corresponding photometric error.
- stetson j: The Welch-Stetson J variability index [39]. A robust standard deviation.

- stetson k: The Welch-Stetson K variability index [39]. A robust kurtosis measure.
- max slope: Maximum absolute slope (delta magnitude over deltatime) between two consecutive observations.
- amplitude: Difference between maximum and minimum magnitudes.
- median absolute deviation: from the median magnitude.
- median buffer range percentage: Percentage of points within 10% of the median magnitude.
- pair slope trend: Percentage of all pairs of consecutive magnitude measurements that have positive slope.
- percent amplitude: Largest percentage difference between the absolute maximum magnitude and the median.
- percent difference flux percentile: Ratio of F 5,95 and the median flux.
- flux percentile ratio mid20: Ratio F 40,60 / F 5,95
- flux percentile ratio mid35: Ratio F 32.5,67.5 / F 5,95
- flux percentile ratio mid50: Ratio F 25,75 / F 5,95
- flux percentile ratio mid65: Ratio F 17.5,82.5 / F 5,95
- flux percentile ratio mid80: Ratio F 10,90 / F 5,95
- poly1 a: Coefficient of the linear term in monomial curve fitting.
- poly2 a: Coefficient of the quadratic term in quadratic curve fitting.
- poly2 b: Coefficient of the linear term in quadratic curve fitting.
- poly3 a: Coefficient of the cubic term in cubic curve fitting.
- poly3 b: Coefficient of the quadratic term in cubic curve fitting.
- poly3 c: Coefficient of the linear term in cubic curve fitting.
- poly4 a: Coefficient of the quartic term in quartic curve fitting.

- poly4 b: Coefficient of the cubic term in quartic curve fitting.
- poly4 c: Coefficient of the quadratic term in quartic curve fitting.
- poly4 d: Coefficient of the linear term in quartic curve fitting.

The results from the authors thesis will not be presented here as a bug overrating the classification was found as will be discussed in the following section.

## 3.2 research internship?? No se como ponerle a esto 2018-2

### 3.2.1 Improvement on Diego's work

Diego's work was continued and improved during this semester with the mentorship of Marcela Hernández, Jaime Forero and Pablo Arbelaez.

### 3.2.2 PLAsTiCC - Kaggle competition

## 4 General objective

## 5 Specific objectives

### 5.1 Improvement of the light curve feature space

#### 5.1.1 Feature pruning

After seeing the results in 3.2.1, **REVISAR ESTO**, it was clear that the higher order coefficients in the polynomial fits did not contribute significantly to the correct classification in the feature space. In fact, these features could be harming the classification process. Thus, coefficients resulting from 3<sup>rd</sup> and 4<sup>th</sup> degree polynomial fitting will be removed and the random forest algorithm will be rerun. The classification metrics should improve but experimentation is needed to confirm the hypothesis.



### 5.1.2 Addition of supernovae specific metrics

To improve the classification of supernovae, metrics that target their specific light curve behavior are needed. In particular, there are functions that are known to approximate the light curve produced by a supernova. These functions are presented below:

**SALT2** is “an empirical model of Type Ia supernovae spectro-photometric evolution with time” [9]. This model should be better adjusted for type 1A supernovas than the rest of objects. The mathematical model of the function is the following[9]:

$$F(S, N, p, \lambda) = x_0 \times [M_0(p, \lambda) + x_1 M_1(p, \lambda) + \dots] \times \exp[cCL(\lambda)]$$

**Skewed Gaussian** fits of the form:

$$f^k(t) = A^k \frac{\exp - (t - t_0^k / \tau_{fall}^k)}{1 + \exp - (t - t_0^k) \tau_{rise}^k}$$

have been shown to resemble well a generalized supernova curve [10]. It should also approximate the shape of supernova light curves better than those that belong to other classes.

When these functions are fit to supernovae data, the resulting  $\chi^2$  should be considerably lower than the  $\chi^2$  calculated from non-supernovae fits. Consequently, the addition of the two  $\chi^2$  values from each of the functions should improve the binary classification of supernovae but as stated above, experimentation is needed for verification.

## 5.2 Deep learning on light curves

For the time being, the classification pipeline has had 3 steps:

1. Clean and filter light curves
2. Calculate features from clean light curves
3. Classify the objects on the feature space

When the feature space is calculated, a large quantity of information is lost. The features might be good descriptors of the objects but they will never be able to encapsulate the point-by-point data that the light curves have. Additionally, traditional machine learning methods, like the random forest classifier previously implemented, are unable to handle varying length input. To take advantage of all the information present in the light curves, a new approach is needed.

Recurrent neural networks (RNN) have been shown to be able to classify sequential data to a good extent. Speech recognition has been one of the fields with most progress[11]. A RNN like the long short term memory (LSTM) RNN or the gated recurrent unit (GRU) RNN are state of the art RNN's that seem promising for this purpose. Thus, several implementations of these networks will be carried out along with their hyperparameter tuning.

## **5.3 Deep learning on images**

### **5.3.1 CNNs**

So far, all the classification has been done on the objects' photometric light curves. The light curves, however, are not the raw data. These were calculated from the images that were taken by the telescopes. Specifically, the light curves were calculated from the catalina real time survey images [12]. Ideally, to minimize data loss, the classification process should be done on the raw data i.e. the images.

The state of the art algorithms used for classifying images are known as convolutional neural networks (CNNs) [13] and have had wide success on large variety of problems.

The first step to correctly classify the images will be to implement a baseline CNN algorithm. Once its classification metrics are established, a thorough search through the hyperparameter space will be carried out to maximize the classification metrics.

### 5.3.2 CNNs and RNNs

Implementing a succesful CNN is only half of the problem. Each object has multiple images taken at different instances in time. To fully exploit the available information, multiple images need to be used as input for each object.

Since the amount of images per object is not fixed, a RNN will need to be used. Thus, the CNN will extract descriptors from the images and then those descriptors will be fed in chronological order along with the date of when the image was taken into the RNN for classification. The architecture is depicted below:

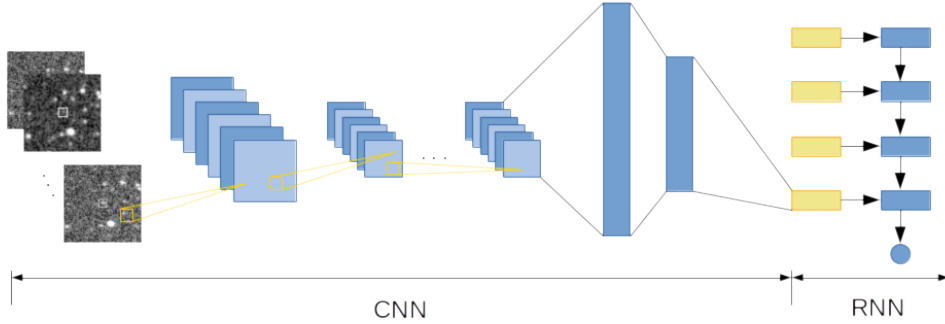


Figure 2: Architecture for classifying astronomical objects from multiple images.

## 6 Activities and schedule

## 7 Expected results

## References

1. Drake, A. J. *et al.*  
First Results from the Catalina Real-Time Transient Survey.  
*The Astrophysical Journal* **696**, 870 (2009).
2. *Type Ia Supernovae as Standard Candles*  
[https://ned.ipac.caltech.edu/level5/Branch2/Branch2\\_1.html](https://ned.ipac.caltech.edu/level5/Branch2/Branch2_1.html).
3. *Light Curves - Introduction*  
<https://imagine.gsfc.nasa.gov/science/toolbox/timing1.html>.
4. *About Light Curves — Aavso.Org*  
<https://www.aavso.org/about-light-curves>.
5. Goliath, M., Amanullah, R., Astier, P., Goobar, A. & Pain, R.  
Supernovae and the Nature of the Dark Energy.  
*Astronomy & Astrophysics* **380**, 6–18 (2001).
6. Perlmutter, S., Turner, M. S. & White, M. Constraining Dark Energy  
with Type Ia Supernovae and Large-Scale Structure.  
*Physical Review Letters* **83**, 670 (1999).
7. Ivezić, Z. *et al.* LSST: From Science Drivers to Reference Design and  
Anticipated Data Products. *arXiv preprint arXiv:0805.2366* (2008).
8. Gomez, D. *Astronomical Transient Recognition Using Machine  
Learning and Catalina Real Time Survey Dataset.:*  
*Diegoalejogm/Crts-Transient-Recognition* Sept. 2018.
9. Guy, J. *et al.* SALT2: Using Distant Supernovae to Improve the Use of  
Type Ia Supernovae as Distance Indicators.  
*Astronomy & Astrophysics* **466**, 11–21. ISSN: 0004-6361, 1432-0746  
(Apr. 2007).
10. Möller, A. *et al.* Photometric Classification of Type Ia Supernovae in  
the SuperNova Legacy Survey with Supervised Learning. *Journal of  
Cosmology and Astroparticle Physics* **2016**, 008–008. ISSN: 1475-7516  
(Dec. 2016).
11. Graves, A., Mohamed, A.-r. & Hinton, G.  
*Speech Recognition with Deep Recurrent Neural Networks*  
in *Acoustics, Speech and Signal Processing (Icassp), 2013 Ieee  
International Conference On* (IEEE, 2013), 6645–6649.

12. [http://nesssi.cacr.caltech.edu/catalina/CRTSII\\_SNCV.html#table26](http://nesssi.cacr.caltech.edu/catalina/CRTSII_SNCV.html#table26).
13. Wang, C. & Xi, Y.  
Convolutional Neural Network for Image Classification.  
*Johns Hopkins University Baltimore, MD* **21218**.