

Refining Exoplanet Detection Using Supervised Learning and Feature Engineering

SLIOIA (CLEI) and CLEI-EJ

M. Bugueño, F. Mena and M. Araya

Departamento de Informática
Universidad Técnica Federico Santa María

April 2019



Summary

- 1** Introduction
- 2** Background
 - State of the art
- 3** Data
 - What and Where
 - Inside the data
 - Data generation
- 4** Models and Methods
 - Feature extraction
 - Learning models
 - Evaluation
- 5** Experiments
 - Results
 - Detail of results
 - Tagging
- 6** Working on
- 7** Conclusion



Motivation



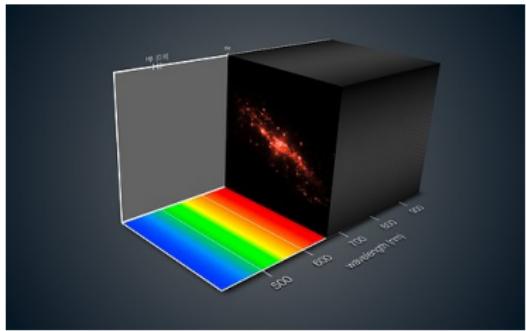
A lot of manual process and analysis



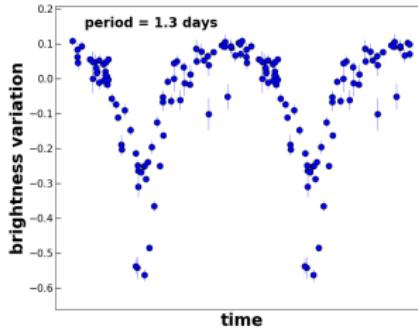
Different fields of application



Data types in astronomy



(a) Data Cube



(b) Time series

source_id	ra	ra_error	dec	dec_error	parallax	parallax_error	phot_g_mean_mag	tp_rp
	mag	mag		mag		mag		mag
428399982269110424	277.025621210121012	0.4410000090268175	5.4231830000360316	0.4823280000027175	0.9312620000001195	0.01916200000000121	19.29894	2.4483987
4283999876640975298	278.659500000018113	0.4435908000211309	5.4223057000245471	0.4806700000000406	0.9311240000000113	0.01915800000000121	19.24185	2.3798861
4284023040469360298	278.779900000003737	0.44210080002781394	5.4223300000018977	0.48051000000004514	0.93095000000003349	0.01915790000000121	19.09298	19.09298
4284023708030000216	278.717900000004738	0.4408400000190043	5.4223400000018944	0.48052000000003308	0.93095000000003268	0.01915790000000121	20.33848	2.14991103
4284023708030000344	277.242470000005048	0.40630000001949918	5.4233500000100055	0.48160000000001638	0.93084000000002068	0.01904700000000121	16.14209	1.864604
4284023708030000464	278.740800000005048	0.41667000001927618	5.4230040000100164	0.48130000000001841	0.93084000000002068	0.01904700000000121	16.03427	1.848207
4284023708030000560	278.760400000005048	0.41637000001927619	5.4230040000100164	0.48130000000001841	0.93084000000002068	0.01904700000000121	17.70862	1.770862
4294005072374546268	277.450246000004550	0.5223700000004433	1.76212700000005096	0.48077000000004134	0.93087000000002064	0.01904800000000121	20.11939	2.057945
429400509721013460	277.495500000004554	0.39420000001932299	5.42411791101507	0.48039000000003264	0.930740000000014798	0.0190887	1.693364	
429402502005331938	278.070500000005057	0.39604000001941668	5.4220128000173198	0.48031000000004981	0.930316000000017033	0.0190885	2.309092	
4294026012051093624	278.544274000003093	0.219840000015801454	5.4202180000098056	0.48212000000003123	0.944850000000017395	0.0220300000161176	19.36369	1.704170
4294026012051093732	278.579100000003093	0.219840000015801454	5.4202180000098056	0.48212000000003123	0.944850000000017395	0.0220300000161176	20.05995	1.7591270
4294026012051093732	278.7270050000027007	0.0011729820000302	5.42079431730492	1.0390300000021304	0.930918000000012107	0.0190879	2.5901270	
42940272294591772758	278.7093217371001	0.4696700000209287	5.4206040000093056	0.48243000000005064	0.930917000000014537	0.0190879	2.043394	1.909893
42940272294591772758	278.7771877371001	0.4696700000209287	5.4206040000093056	0.48243000000005064	0.930917000000014537	0.0190879	21.89919	1.883079
4304079050011354408	298.170490000004170	0.16430000000645508	1.8775971000000808	0.48194000000003879	0.929630000000014404	0.01904840000000121	27.85208	2.499528
4304079050011354408	298.170490000004170	0.16430000000645508	1.8775971000000808	0.48194000000003879	0.929630000000014404	0.01904840000000121	26.01023	2.196671
4304079050011354408	298.170490000004170	0.16430000000645508	1.8775971000000808	0.48194000000003879	0.929630000000014404	0.01904840000000121	28.79446	2.4893088
4313993480011360008	294.187973310000889	-0.45176000000003104	5.3908174000000506	0.48087110000004084	0.930500000000018308	0.0190807	18.793465	18.793465
4313993480011360008	295.0408560000042002	1.52165017441220	5.3908174000000506	0.48087110000004084	0.930500000000018308	0.0190807	20.417179	2.055513
43480570304857255472	205.2406801367943	0.45600600000502223	1.42374000000749285	0.50404000000280516	0.9098400000033764423	0.00864000000000121	19.267062	1.9879844

(c) Catalog



(d) Image

What is an exoplanet?

Planets orbiting stars outside our solar systems are called **extra-solar planets** or **exoplanets**.

Birth of a planet

The theory of extrasolar planets is under development until these days (since the first confirmed detection of the giant planet **51 Pegasi b**)

- » Current theories suggest that the dust particles of the protoplanetary disk begin to collapse by gravity forming larger grains.
- » If these discs survive to stellar radiation and comets, the matter continues compacting giving way to a planetoid.



Images by¹ and²

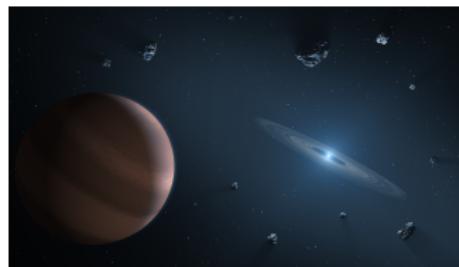


1. <http://www.spitzer.caltech.edu/news/1876-feature16-07-Light-Echoes-Give-Clues-to-Protoplanetary-Disk>
2. <https://www.space.com/19100-alien-planet-birth-alma-telescope.html>

Why is difficult to detect exoplanets ?

Detecting these planets is a challenging problem !

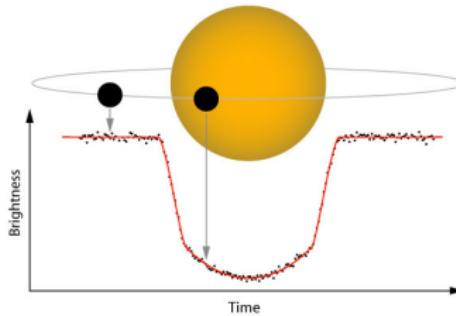
- They emit or reflect very dim magnitudes compared to their host stars
- They are very near to their host stars compared to the observation distance
- Fine-grained analysis is needed



Detection methods

The most successful detection mechanisms are the **indirect**, some of them :

- ♂ **Radial velocity**, which studies the speed variations of a star product of its orbiting planets, analyzing the spectral lines of this one through the Doppler effect.
Successful, but only effective on giant planets near its star
- ♺ **Transit photometry**, photometric observation of the star and detection of variations in the light intensity when an orbiting planet passes in front of it.
Efficient, detect high-volume planets independently of the proximity to its star

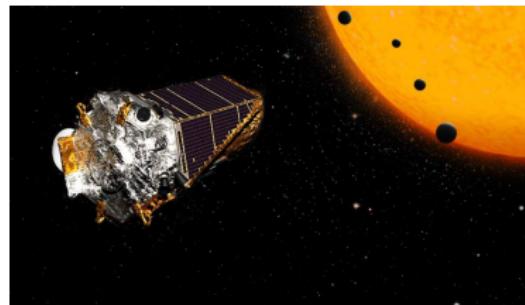


What is needed?

Problem : A large volume of data is being generated today !

- The use of automated methods could reproduce the astronomer analysis to decide if the data supports the existence of an exoplanet or not.
- 400 billions of planets waiting to be discovered³

Fortunately, technological advances in photometry have allowed experiments like the space observatory Kepler to have sufficient sensitivity for detecting a greater range of exoplanets.



Previous work I

A typical research are detections of **RR Lyrae Stars**, because their intensity varies through time independent of a planet.

- Richards et al. (2011) presents a catalog of variable stars and manually extracted light curve specialized features from simple statistics and other features based on the period and frequency analysis of a LombScargle fitted model
- Donalek et al. (2013) classify variable stars from the Catalina Real-Time Transient Survey (CRTS) and the Kepler mission extracting similar features from the light curve to Richards et al.
- Mahabal et al. (2017) also with the propose of classify variable stars, transformed light curves into an image (grid) that represent the variations of magnitude through the variations of time intervals.



Previous work II

- Hinners et al. (2018) presents different machine learning techniques and models with the objective of classify and predict features over the same data as we use. They extract some statistical features from the light curve but they were not interested on exoplanet detection.
- Thompson et al. (2015) used unsupervised learning to reduce dimensionality of Kepler light curve and cluster together similar shapes with the focus of find candidates object of interest. They use k-NN model to predict.
- The inspiration come from the **Autovetter** project (McCauliff et al. 2015) : used a Random Forest model based on features derived from Kepler pipeline statistics to identify candidates.

Our Objective : Exoplanet detection through transit light curve with ad-hoc feature extraction and machine learning.



Where the data come from ?

Kepler space observatory

- Data collected by Kepler Mission (launched in 2009 by NASA) with the goal of searching similar planet to Earth.
 - Around 65% of exoplanet discoveries have been detected thanks to Kepler Mission⁴
- The Kepler Objects of Interest (**KOI**⁵) dataset is provided by MAST (Mikulski Archive for Space Telescopes) archive.
 - It contains 8054 KOI's.

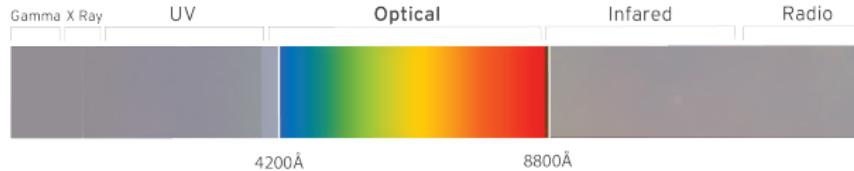


FIGURE – Kepler field of view of the spectrum.



4. https://exoplanetarchive.ipac.caltech.edu/docs/counts_details.html
5. http://archive.stsci.edu/search_fields.php?mission=kepler+koi

What are the data ?

The Labels

- Every record is associated to a Kepler Object of Interest labeled as :
 - *Confirmed* : those that have been confirmed as exoplanet, through extensive analysis.
 - *False Positive* : those that were initially selected as candidate exoplanets but there is additional evidence that shows they are not.
 - *Candidate* : those that are still under study.
- according to *Nasa Exoplanet Science Institute*⁶
- The reasons to catalog as a False Positive are observation that did not match with the star position on study. Also can be that the deep of the even transit was statistically different to the deep of the odd transits, showing a binary system.

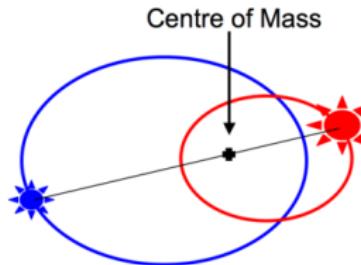


FIGURE – binary star system



What are the data ? - Class distribution

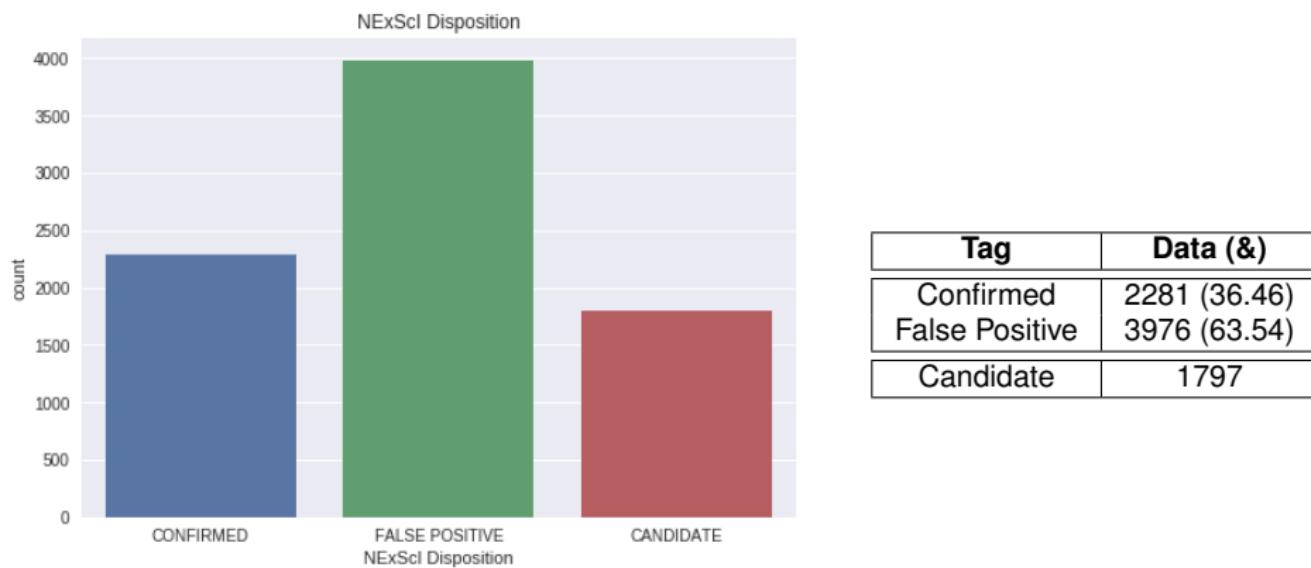


FIGURE – Categories distributions on collected data



What contain each KOI

Error

- The error associated to each measure

Time

- The time (in Julian Date - January 1, 4713 BC) when the measure was made

Raw light curve

- The raw measurements (intensity) of star light, with about 70000 measurements
- Frequency sampling rate of 30 minutes approximately (4 years of observations)
 - Measurements are not recorded uniform
- On average, the missing data is about 23% , this mean approximately about 55000 *effective* measurements.



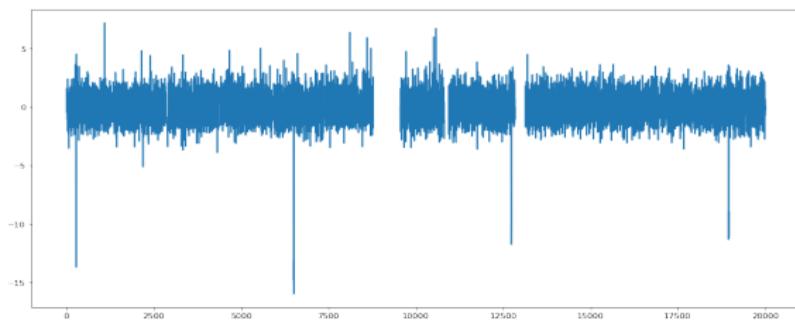
Fill missing values

- Fill with zeros

- Stationary state of the system (not eclipsing)
- Expected value, $E[I] \approx 0$

- Fill the gaps with linear interpolation

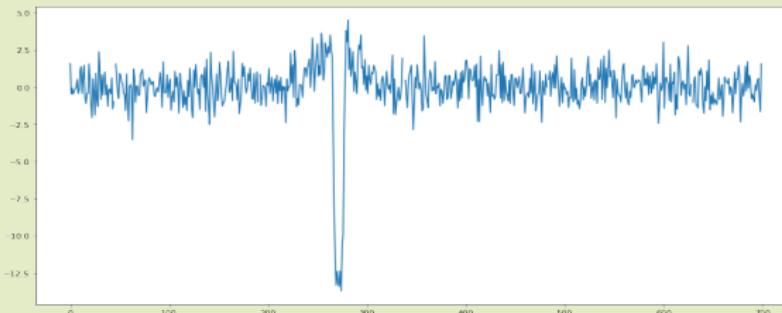
- To get a smooth version in extense continuous missing data



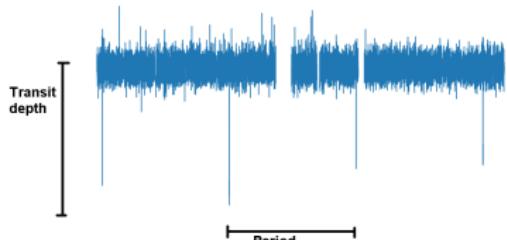
What contain each KOI

Filtered light curve

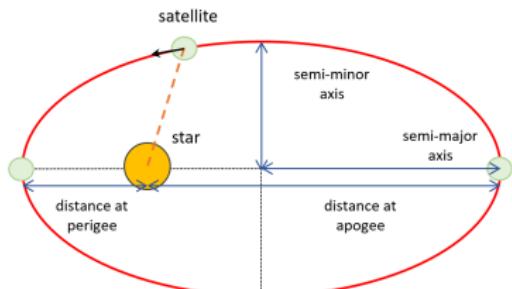
- Light curve with a **whitened filter**. This transform led a constant white noise on the light curve, i.e., the higher signal get amplified and give a more uniform signal.
- The white noise is a random signal that have the same intensity on different frequencies, which gives a spectral density of constant power.



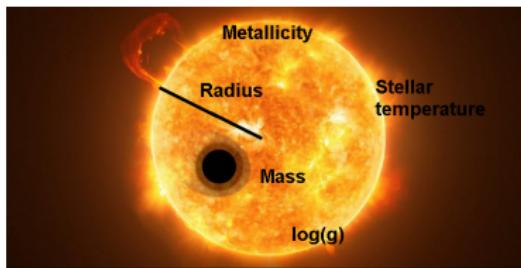
What does each KOI contain ? - Metadata



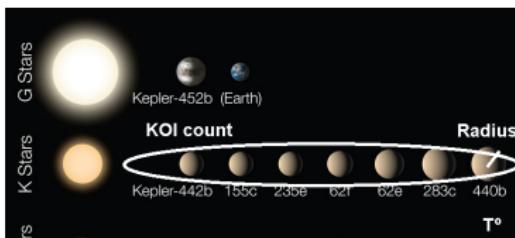
(a) Planet metadata



(b) Transit metadata



(c) Stellar metadata



(d) Planet metadata



Data process

- A system can contain various exoplanet orbiting the host star.
- There are systems that have from one to seven objects under study.

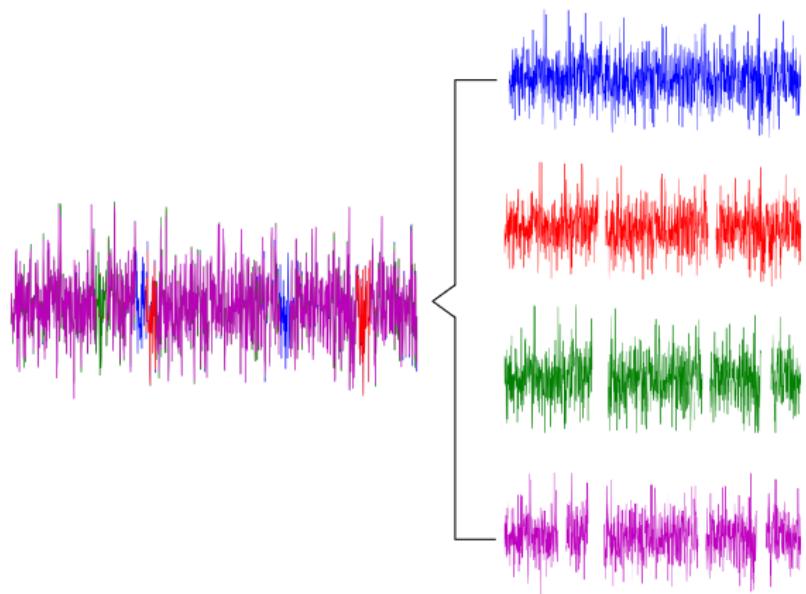


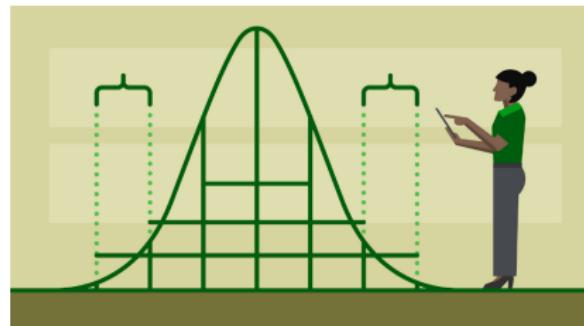
FIGURE – Light curve separation



Manual - baseline

We used extraction techniques specialized on time series, which in this case corresponds to measurements of intensity of the light along time, inspired on the library Feature Analysis for Time Series (**FATS**⁷) for Python.

- Amplitude
- Slope
- Max
- Mean
- Median
- Median absolute deviation
- Min
- Q1
- Q2
- Q31
- Residual bright faint ratio
- Skew
- Kurtosis
- Std



To this we decide to add the metadata.



7. <https://github.com/isadoranun/FATS>

Selected Metadata

Based on previous small knowledge	
Planet+Transit	Stellar
Radius	Radius
Temperature	Temperature
KOI count	Metallicity
Period	Mass
Transit Depth	$\log(g)$



Automatic

- We use : **Unsupervised Learning Methods**
the objective is to find intrinsic patterns among all the data independently from task.

Pre-process

- Firstly we applied a Discrete Fourier Transform¹ to the light curve.
 - It transforms the data from the time domain in which the measurements were obtained, to the frequency domain where the signal was generated.

1. *this method is designed to analyze periodic signals, which is exactly the case of transit light curves.*

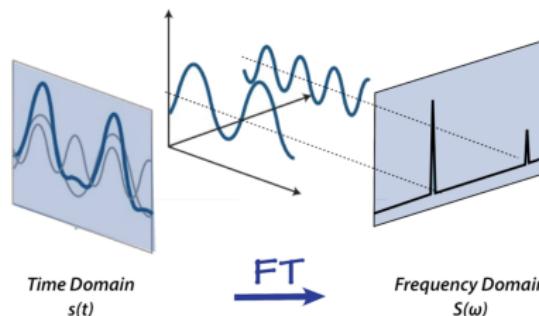


FIGURE – image by <http://mriquestions.com/fourier-transform-ft.html>



Automatic

PCA

- The **Principal Component Analysis**, as a linear method that projects data into a lower dimensional space (the higher variance vector).

- Why PCA ?

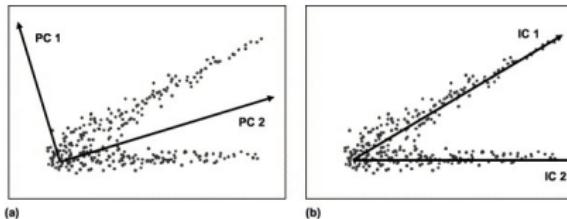
- Has been applied to several applications, obtaining particularly good results on time series
- Its great efficiency from specific optimizations over linear algebra methods with high dimensional data.

ICA

- The Fast iterative algorithm for **Independent Component Analysis** that finds statistically independent components of the data.

- Why ICA ?

- Is focused on the signal abstraction, since it tries to detect the independently sources that, mixed, produce the observed data.
- Differs to the uncorrelated componentes that finds PCA.



Machine Learning models

k-NN, regularization parameter : k

- Based on memory (i.e., non-parametric). Simple, but with good performance

Logistic Regression, regularization parameter : C

- Classify based on a probabilistic (logistic - sigmoid function) linear model.

SVM, regularization parameter : C

- A linear margin-based model... but with Kernel : RBF (Radial Basis Function)

Random Forest, regularization parameters : $depth$ and T

- A ensemble of decision trees.. not linear !



Metrics

Classification

The exoplanet problem is an instance of **unbalanced binary** classification problem.

- Precision : ability to label one class when the object effectively was from that class. (inverse to contamination)

$$P = \frac{T_p}{T_p + F_p}$$

- Recall : ability to include all objects that effectively are from one class. (similar to completeness)

$$R = \frac{T_p}{T_p + F_n}$$

- F_1 -score : as harmonic mean between P and R

$$F_1 = 2 \cdot \frac{P \cdot R}{P + R}$$

All these metrics reach their best values at 1 and worst at 0



Experiments

- Large amount of data to be processed! So it was necessary to use a cluster provided by **ChiVO**⁸ (Chilean Virtual Observatory)
- We have 6257 labeled data, corresponding to 121GB, and 1797 not labeled data (candidates), corresponding to another 33GB.
 - We used the (64/18/18)% of the labeled data for set-split. Approximately 4000, 1000 and 1000 registers were grouped as training, validation and testing sets. This last one represents the actual *target* unlabeled data that is unknown (Candidates).

Note : The selection of the hyper-parameters of the classifiers was not expensive in computational terms (dimensionality d much smaller than the original), unlike the feature extraction process.



8. <http://www.chivo.cl>

Experiments

- For automatic feature extraction techniques, fixed dimensions were experimented

	5	10	15	20	25	50
ICA	0.711	0.709	0.709	0.686	0.679	0.675
	5	10	25	55	100	255
PCA	0.713	0.701	0.701	0.699	0.702	0.689

TABLE – F1 score of the best classifier, Random Forest, in function of dimensionality.

- Surprisingly enough, completing missing data with zeros produces a consistent improvement of ~ 0.1 in the F_1 -score, while linear interpolation produced worse.



What about unbalanced data ?

- 1 **Undersampling technique**; the majority class was subsampled

- 2 **Weighting**; weights the different classes
 - This improved the results by ~ 0.1 on F_1 score metric.



Performance Results

Learners performance				
	<i>k-NN</i>	<i>Logistic Regression</i>	<i>SVM RBF</i>	<i>Random Forest</i>
Fourier + PCA	0.679	0.493	0.486	0.713
Fourier + ICA	0.679	0.493	0.486	0.711
OwnFATS	0.666	0.583	0.575	0.658
Planet metadata	0.825	0.848	0.848	0.870
Stellar metadata	0.766	0.718	0.751	0.766
OwnFATS + stellar & planet metadata	0.844	0.864	0.876	0.883

TABLE – F1 score on the classification of different models (learners) over the test set on the different representations generated.



Results analysis

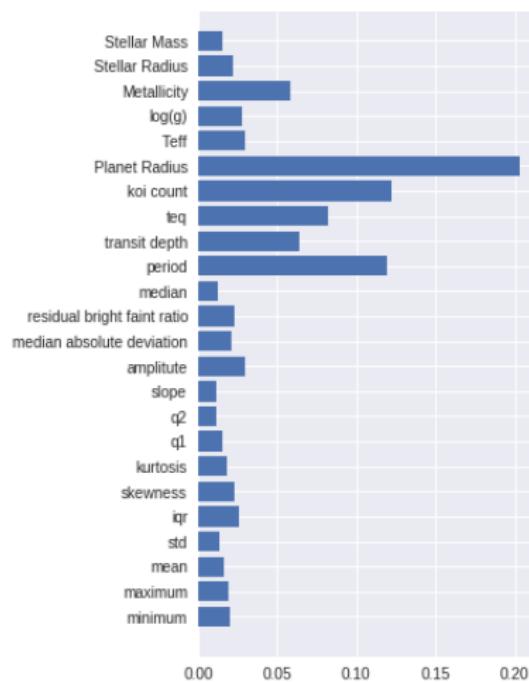


FIGURE – Random Forest feature importance

Performance analysis

- The best trained model, was Random Forest
- Performance of 88.3% for future classification by F_1 score metric
- In detail :
 - Radius, period and number of objects in the system as the most relevant features
 - The less important features are the extracted from the light curve (slope and second quartile)



Detail of results

Precision-Recall detail



FIGURE – Manual extraction features plus metadata

- The best model that can identify correctly the False Positive Class (FP) based on Precision and Recall is the SVM RBF.
- Classification on the Confirmed class shows lower scores than FP suggesting the difficulty on the exoplanet prediction.
 - Why ? Could be because all of them do not have very similar features on the light curve
- The best model in the task of confirmed exoplanets on the different representation of the data was Random Forest.



Final Results

Selected representation :

- OwnFATS + stellar & planet metadata

Selected model for *confirmed* class :

- Random Forest ($depth = 15, T = 15$)

Selected model for *false positive* class :

- SVM RBF ($C = 100$)

We show the classification over the Kepler Object of Interest that are still being studied by the staff of NexSci on September 2017 :

Total CANDIDATE	1791
Subtotal CONFIRMED	975
Subtotal FALSE POSITIVE	434
Unclassified	382

KOI name	Disposition	Confirmed on system	Star
K00601.02	<i>False Positive</i>	2/3	Kepler 619
K00750.02	<i>Unclassified</i>	1/3	Kepler 662
K01082.01	<i>Confirmed</i>		
K01082.02	<i>False Positive</i>	1/4	Kepler 763
K01082.04	<i>Confirmed</i>		
K01236.04	<i>Confirmed</i>	2/3	Kepler 279
K01358.01	<i>Confirmed</i>		
K01358.02	<i>Confirmed</i>	0/4	-
K01358.03	<i>Confirmed</i>		
K01358.04	<i>Confirmed</i>		
K01750.02	<i>Confirmed</i>	1/2	Kepler 948
K02064.01	<i>Unclassified</i>	0/1	-
K02420.02	<i>Confirmed</i>	1/2	Kepler 1231
K02578.01	<i>False Positive</i>	0/1	-
K02828.02	<i>False Positive</i>	1/2	Kepler 1259
K03444.03	<i>Unclassified</i>	0/4	-
K03451.01	<i>Unclassified</i>	0/1	-
K04591.01	<i>False Positive</i>	0/1	-
K05353.01	<i>False Positive</i>	0/1	-
K06267.01	<i>Confirmed</i>	0/1	-



github.com/FMena14/ExoplanetDetection

Publication

- SLIOIA 182772, CLEI publication :

Refining Exoplanet Detection Using Supervised Learning and Feature Engineering

Margarita Bugueño

Departamento de Informática
Univ. Técnica Federico Santa María
Santiago, Chile
margarita.bugueno.13@sansano.usm.cl

Francisco Mena

Departamento de Informática
Univ. Técnica Federico Santa María
Santiago, Chile
francisco.mena.13@sansano.usm.cl

Mauricio Araya

Departamento de Informática
Univ. Técnica Federico Santa María
Valparaíso, Chile
maray@inf.utfsm.cl

Abstract—The field of astronomical data analysis has experienced an important paradigm shift in the recent years. The automation of certain analysis procedures is no longer a desirable feature for reducing the human effort, but a must have asset for coping with the extremely large datasets that new instruments

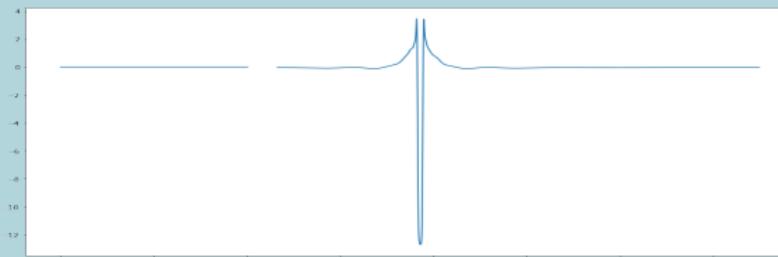
performance of supervised classification methods. To compare them, we used a few standard statistical learning criteria, yet the results motivated an hybrid selection of techniques that reduces the number of unconfirmed light curves.



Working on

Mandel-Agol fit

- The Mandel-Agol model the transit of a stratospheric planet around a stratospheric star, like an eclipse, assuming a uniform source.
 - It requires the distance from the planet to the parent star as well as the radius of each one of the bodies (ratio) and the period of the orbit. For a better fit could provide eccentricity and limb-darkening coefficients.
- The closeness of the planet to star (eclipse) is modeled as a quadratic polynomial.



Selected Metadata II

Based on previous not so small knowledge		
Transit	Stellar	Planet
Period	Radius	Radius
Transit Depth	Temperature	Temperature
Eccentricity	Metallicity	KOI count
Planet-Star Radius Ratio	Mass	
Orbit Semi-Major Axis	Log(g)	
Limb Darkening Coefs		
Transit Signal-to-Noise		
+Upper and lower uncertainty limit		



Dataset split

We selected a new test set which share certain similarity with the unlabeled objects

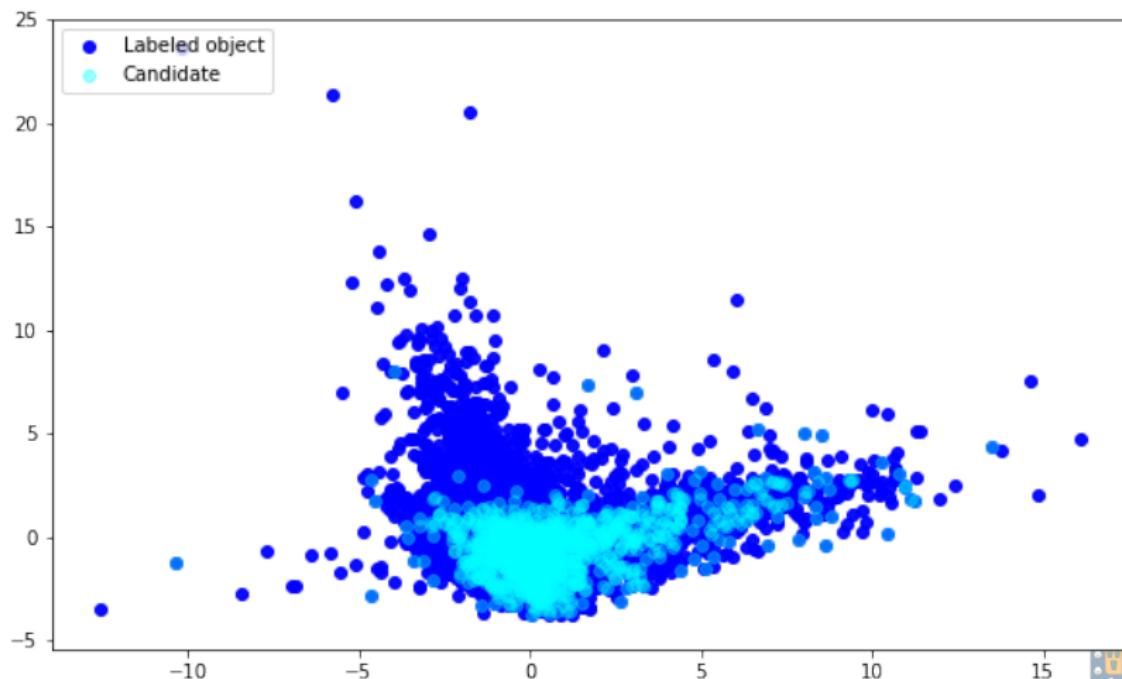


FIGURE – PCA over metadata

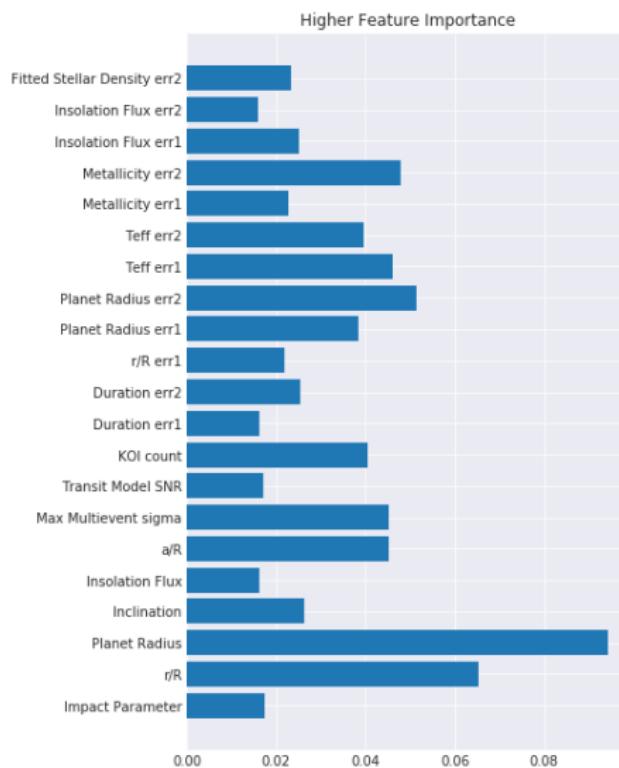
New and working on results

Learners performance new			
	<i>k-NN</i>	<i>SVM RBF</i>	<i>Random Forest</i>
Fourier-PCA (25) on Mandel Agol	0.764	0.440	0.793
Fourier-ICA (25) on Mandel Agol	0.773	0.440	0.804
All metadata + OwnFATS	0.889	0.916	0.932
(FSS) All metadata + OwnFATS			
All metadata + OwnFATS-MA	0.883	0.922	0.930
(FSS) All metadata + OwnFATS-MA			

TABLE – F1 score (weighted) on the classification of different models (learners) over the test set on the different representations generated.



Details



Performance analysis

- On the image 20 most important features for Random Forest.
- Different distribution on Planet Radius :

Statistic	Confirmed	False Positive
mean	2.470	9.415
std	2.027	53.021
min	0.440	0.220
25%	1.350	1.090
50%	1.940	1.580
75%	2.745	3.330
max	13.600	1543.740

Varying the dimensionality

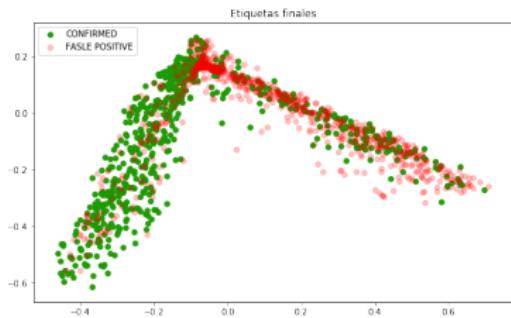
- Change on best dimensionality

	5	10	25	50
ICA	0.751	0.780	0.804	0.803
	5	10	25	50
PCA	0.748	0.764	0.793	0.793

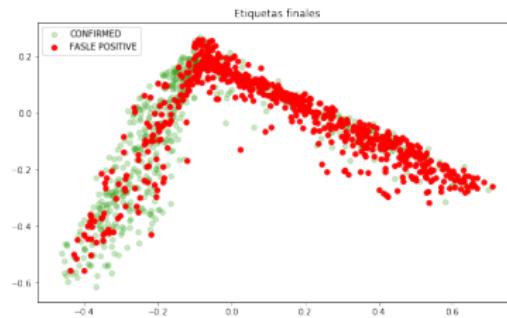
TABLE – F1 score of the best classifier, Random Forest, in function of dimensionality.



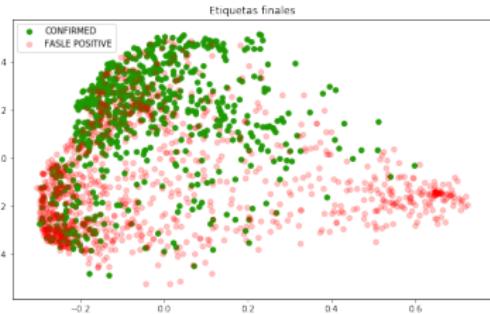
Tagging Visualization



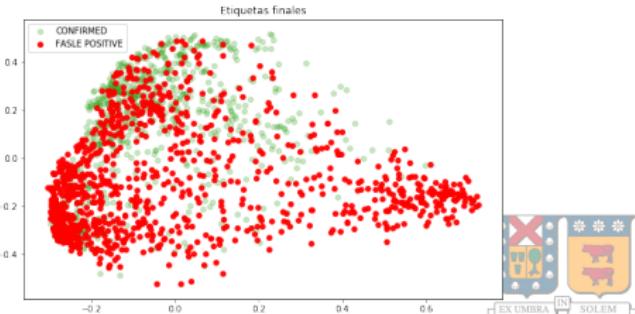
(a) One projection, Confirmed highlighted



(b) One projection, False Positive highlighted



(c) Another projection, Confirmed highlighted



(d) Another projection, False Positive highlighted

Conclusion

- We introduced a new refining method to decide if an object on study (KOI) is really an exoplanet using automatic learning and handling raw data
 - We reproduce the arduous and extensive work that experts perform on detection
- The results show that the automatic techniques used to extract information from the light curve was not good enough compared to the metadata
 - Maybe the not suitable methods for the feature extraction, or well, too simple for the complex problem that we faced
- Also the problem was complex regarding the execution time **on feature extraction**
- The assistance of an expert in this area could be of great value
- **Possible improvements** : Change input representation.

Acknowledgments : Thanks to Chilean Virtual Observatory, ChiVO. Also we thanks to the academic Ricardo Nanculef.



Refining Exoplanet Detection Using Supervised Learning and Feature Engineering

SLIOIA (CLEI) and CLEI-EJ

M. Bugueño, F. Mena and M. Araya

Departamento de Informática
Universidad Técnica Federico Santa María

April 2019

