

# Hierarchical Machine Learning Classification of Eclipsing Binaries in the OGLE Data

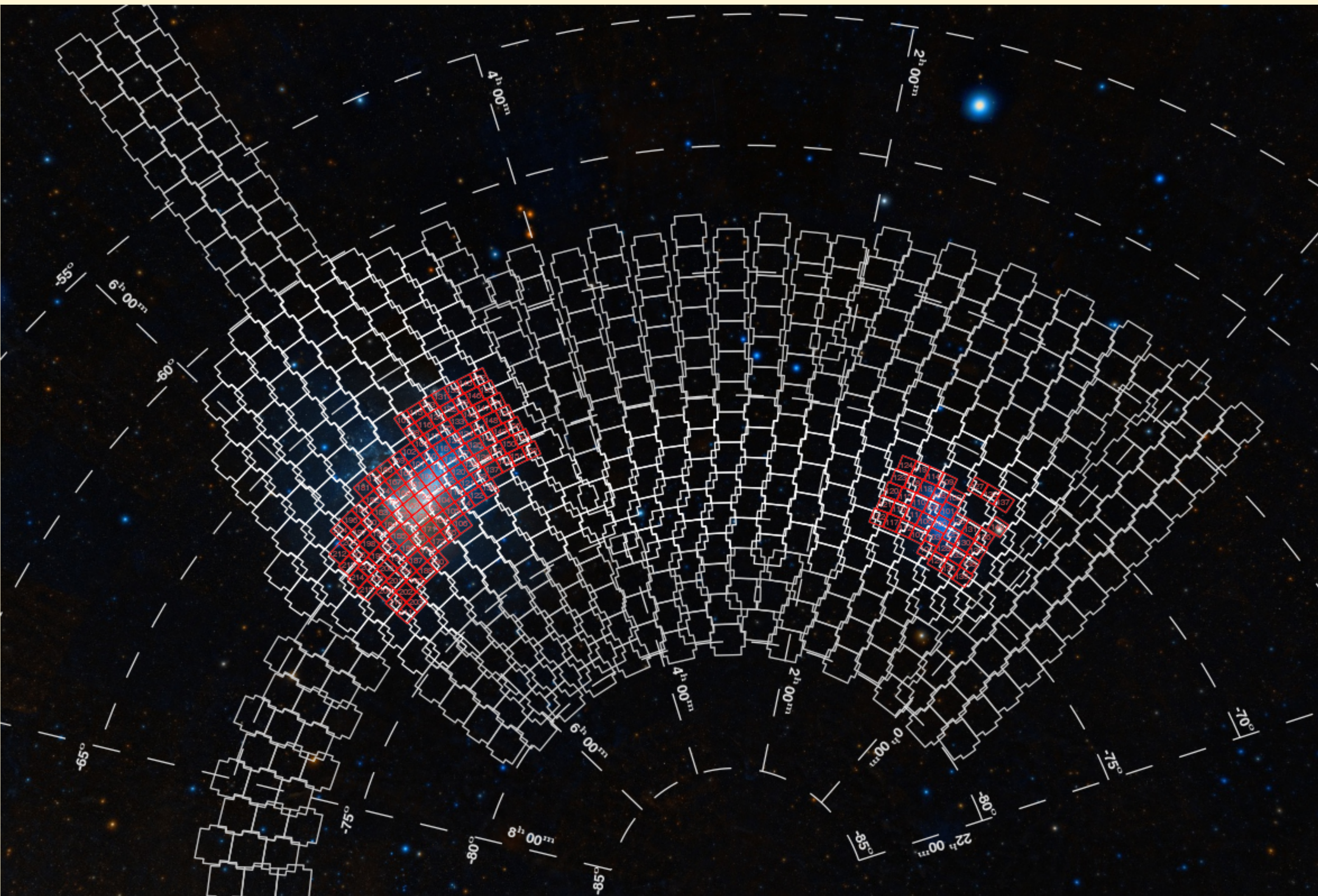


Michał Pawlak  
Warsaw University Observatory  
mpawlak@astrouw.edu.pl



## The OGLE Project

The OGLE project has been operating since 1992. Three phases of the project were successfully completed: OGLE-I (1992-1995), OGLE-II (1997-2000) and OGLE-III (2001-2009). Since 2010 OGLE-IV is in operation. During the first phase of the project, 1.0-m Swope Telescope was used. From the OGLE-II, 1.3-m Warsaw Telescope at the Las Campanas Observatory in Chile is dedicated for the project. The number of monitored object increased significantly from about one million in OGLE-I to about 40 million in OGLE-II and more than 400 million in OGLE-III. In OGLE-IV the number of observed sources reached 1 billion. Most of the observations (about 90%) are taken in *I* filter, while the remaining in *V* filter. The OGLE project original aim was the detection and analysis of microlensing events, however its scientific results cover much wider field of study. One of the most significant of them is discovery and analysis of hundreds of thousands of variables stars.



## OGLE-III and OGLE-IV

The figure above presents the comparison of the OGLE-III (red contours) and OGLE-IV (white contours) sky coverage in the Magellanic Clouds. The OGLE-III catalogue of variable stars contains more than 26 000 eclipsing binaries in the LMC and 6000 in the SMC, making a perfect training set for machine learning selection and classification of stars observed in OGLE-IV which covers much larger area.

## Hierarchical classification

### First Step

The OGLE-III catalog of eclipsing binaries in the LMC (Graczyk et al. 2011) was used to create the training set. The entire sample of known eclipsing binaries from one of the OGLE-IV fields containing about 1000 objects was complemented with a matching number of other, randomly selected stars from the same field. BLS algorithm (Kovacs et al. 2002) is used to determine the period for each of the objects, as well as series of statistics including signal to noise ratio of the periodogram, duration and depth of eclipse,  $\chi^2$  of the fit, white and red noise. The BLS statistics are used as features for Random Forest (Breiman 2001) classification algorithm. To reduce the number of false positives, the prior on probability is set to 80%. Recall of such classification is 93%.

### Second Step

For the next step of the classification, the training set is composed of the same eclipsing binaries as previously and other, non-eclipsing stars that passed the first step. The ratio of false to true detections after the first step is about 6:1. A part from BLS parameters a series of statistical features (standard deviation, skewness, kurtosis, third and first quartile difference) is used. Classification is again performed with Random Forest. Recall of this step, evaluated on the training set is 88% and precision 89%.

### Overall Result

The overall performance has been tested independently on one another OGLE-IV field not used in the training set. Recall of the entire method in respect to the OGLE-III catalog is estimated to be about **81.5%**. This is consistent with the recall derived from the cross-validation method which is **83%**.

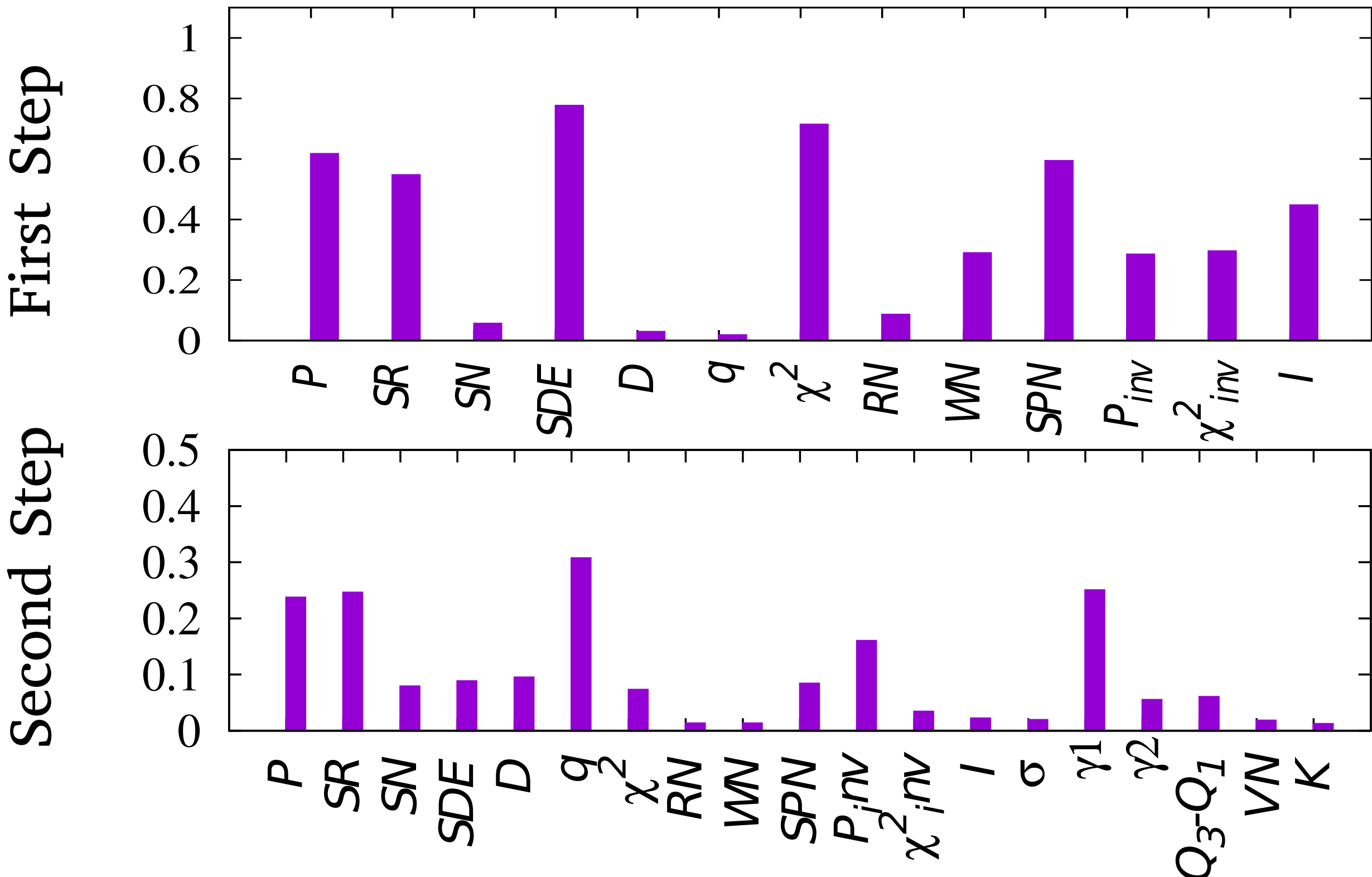
Parameter	Description
<i>P</i>	Period
<i>SR</i>	Signal Residue
<i>SN</i>	Signal to Noise Ratio
<i>SDE</i>	Signal Detection Efficiency
<i>D</i>	Depth of transit
<i>q</i>	Fraction of the phase in transit
$\chi^2$	$\chi^2$ of the transit model fit
<i>RN</i>	Red Noise
<i>WN</i>	White Noise
<i>SPN</i>	Signal to Pink Noise Ratio
<i>P<sub>inv</sub></i>	Period of inverse transit
$\chi^2_{inv}$	$\chi^2$ of the inverted transit model fit
<i>I</i>	mean <i>I</i> -band magnitude

Parameter	Description
$\sigma$	standard deviation
$\gamma_1$	skewness
$\gamma_2$	kurtosis
$Q_3 - Q_1$	difference between third and first quartile
<i>VN</i>	von Neumann index
<i>K</i>	Stetson <i>K</i> index

### Parametrization

Set of parameters used in the First Step (left) and Second step (right). Evaluation of the information gain given by the parameters shown in the figure below.

## Information gain



### References:

Breiman 2001, Machine Learning, 45, 5  
Graczyk et al. 2011, Acta Astronomica, 61, 103  
Kovacs et al. 2002, A&A, 391, 369  
Pawlak et al. 2013, Acta Astronomica, 63, 323

### Acknowledgement:

This work has been supported by Polish National Science Center grant no. 2014/13/N/ST9/00075